

Introduction to Fourier Optics

McGraw-Hill Series in Electrical and Computer Engineering

SENIOR CONSULTING EDITOR

Stephen W. Director, **Carnegie Mellon University**

Circuits and Systems

Communications and Signal Processing

Computer Engineering

Control Theory

Electromagnetics

Electronics and VLSI Circuits

Introductory

Power and Energy

Radar and Antennas

PREVIOUS CONSULTING EDITORS

Ronald N. Bracewell, Colin Cherry, James F. Gibbons, **Willis W. Harman**, Hubert Heffner, Edward W. **Herold**, John G. Linvill, Simon Ramo, Ronald A. Rohrer, Anthony E. Siegman, Charles Susskind, Frederick E. **Terman**, John G. Truxal, Ernst Weber, and John R. Whinnery

Electromagnetics

SENIOR CONSULTING EDITOR

Stephen W. Director, Carnegie Mellon University

Dearhold and McSpadden: *Electromagnetic Wave Propagation*

Goodman: *Introduction to Fourier Optics*

Harrington: *Time-Harmonic Electromagnetic Fields*

Hayt: *Engineering Electromagnetics*

Kraus: *Electromagnetics*

Paul and Nasar: *Introduction to Electromagnetic Fields*

Plonus: *Applied Electromagnetics*

Introduction to Fourier Optics

SECOND EDITION

Joseph W. Goodman

Stanford University

THE McGRAW-HILL COMPANIES, INC.

New York St. Louis San Francisco Auckland **Bogotá** Caracas Lisbon
London Madrid Mexico City Milan Montreal New Delhi
San Juan Singapore Sydney Tokyo Toronto

McGraw-Hill

A Division of The McGraw-Hill Companies



INTRODUCTION TO FOURIER OPTICS

Copyright ©1996, 1968 by The McGraw-Hill Companies, Inc. Reissued 1988 by The McGraw-Hill Companies. All rights reserved. Printed in the United States of America. Except as permitted under the United States Copyright Act of 1976, no part of this publication may be reproduced or distributed in any form or by any means, or stored in a data base or retrieval system, without the prior written permission of the publisher.

This book is printed on acid-free paper.

3 4 5 6 7 8 9 0 FGR FGR 9 0 9 8 7

ISBN 0-07-024254-2

*This book was set in Times Roman by Publication Services, Inc.
The editors were Lynn Cox and John M. Morriss;
the production supervisor was Paula Keller.
The cover was designed by Anthony Paccione.
Quebecor Printing/Fairfield was printer and binder.*

Library of Congress Catalog Card Number: 95-82033

ABOUT THE AUTHOR

JOSEPH W. GOODMAN received the A.B. degree in Engineering and Applied Physics from Harvard University and the M.S and Ph.D. degrees in Electrical Engineering from Stanford University. He has been a member of the Stanford faculty since 1967, and served as the Chairman of the Department of Electrical Engineering from 1988 through 1996.

Dr. Goodman's contributions to optics have been recognized in many ways. He has served as President of the International Commission for Optics and of the Optical Society of America (OSA). He received the F.E. **Terman** award of the American Society for Engineering Education (1971), the Max Born Award of the OSA for contributions to physical optics (1983), the Dennis **Gabor** Award of the International Society for Optical Engineering (SPIE, 1987), the Education Medal of the Institute of Electrical and Electronics Engineers (IEEE, 1987), the Frederic Ives Medal of the OSA for overall distinction in optics (1990), and the Esther Hoffman Beller Medal of the OSA for contributions to optics education (1995). He is a Fellow of the OSA, the SPIE, and the IEEE. In 1987 he was elected to the National Academy of Engineering.

In addition to *Introduction to Fourier Optics*, Dr. Goodman is the author of *Statistical Optics* (J. Wiley & Sons, 1985) and the editor of *International Trends in Optics* (Academic Press, 1991). He has authored more than 200 scientific and technical articles in professional journals and books.

*To the memory of my **Mother**, Doris Ryan Goodman,
and my Fathel; Joseph Goodman, Jr.*

CONTENTS

Preface xvii

1 Introduction

- 1.1 Optics, Information, and Communication
- 1.2 The Book

2 Analysis of Two-Dimensional Signals and Systems

- 2.1 Fourier Analysis in Two Dimensions
 - 2.1.1 *Definition and Existence Conditions* / 2.1.2 *The Fourier Transform as a Decomposition* / 2.1.3 *Fourier Transform Theorems* / 2.1.4 *Separable Functions* / 2.1.5 *Functions with Circular Symmetry: Fourier-Bessel Transforms* / 2.1.6 *Some Frequently Used Functions and Some Useful Fourier Transform Pairs*
- 2.2 Local Spatial Frequency and Space-Frequency Localization 16
- 2.3 Linear Systems 19
 - 2.3.1 *Linearity and the Superposition Integral* / 2.3.2 *Invariant Linear Systems: Transfer Functions*
- 2.4 Two-Dimensional Sampling Theory 22
 - 2.4.1 *The Whittaker-Shannon Sampling Theorem* / 2.4.2 *Space-Bandwidth Product*
- Problems—Chapter 2 27

3 Foundations of Scalar Diffraction Theory 32

- 3.1 Historical Introduction 32
- 3.2 From a Vector to a Scalar Theory 36
- 3.3 Some Mathematical Preliminaries 38
 - 3.3.1 *The Helmholtz Equation* / 3.3.2 *Green's Theorem* / 3.3.3 *The Integral Theorem of Helmholtz and Kirchhoff*
- 3.4 The Kirchhoff Formulation of Diffraction by a Planar Screen 42
 - 3.4.1 *Application of the Integral Theorem* / 3.4.2 *The Kirchhoff Boundary Conditions* / 3.4.3 *The Fresnel-Kirchhoff Diffraction Formula*
- 3.5 The Rayleigh-Sommerfeld Formulation of Diffraction
 - 3.5.1 *Choice of Alternative Green's Functions* / 3.5.2 *The Rayleigh-Sommerfeld Diffraction Formula*

3.6	Comparison of the Kirchhoff and Rayleigh-Sommerfeld Theories	50
3.7	Further Discussion of the Huygens-Fresnel Principle	52
3.8	Generalization to Nonmonochromatic Waves	53
3.9	Diffraction at Boundaries	54
3.10	The Angular Spectrum of Plane Waves	55
	3.10.1 <i>The Angular Spectrum and Its Physical Interpretation /</i>	
	3.10.2 <i>Propagation of the Angular Spectrum / 3.10.3 Effects</i>	
	<i>of a Diffracting Aperture on the Angular Spectrum / 3.10.4</i>	
	<i>The Propagation Phenomenon as a Linear Spatial Filter</i>	
	Problems — Chapter 3	61
4	Fresnel and Fraunhofer Diffraction	63
4.1	Background	63
	4.1.1 <i>The Intensity of a Wave Field / 4.1.2 The Huygens-Fresnel</i>	
	<i>Principle in Rectangular Coordinates</i>	
4.2	The Fresnel Approximation	66
	4.2.1 <i>Positive vs. Negative Phases / 4.2.2 Accuracy of the</i>	
	<i>Fresnel Approximation / 4.2.3 The Fresnel Approximation and</i>	
	<i>the Angular Spectrum / 4.2.4 Fresnel Diffraction Between</i>	
	<i>Confocal Spherical Surfaces</i>	
4.3	The Fraunhofer Approximation	73
4.4	Examples of Fraunhofer Diffraction Patterns	75
	4.4.1 <i>Rectangular Aperture / 4.4.2 Circular Aperture /</i>	
	4.4.3 <i>Thin Sinusoidal Amplitude Grating / 4.4.4 Thin</i>	
	<i>Sinusoidal Phase Grating</i>	
4.5	Examples of Fresnel Diffraction Calculations	83
	4.5.1 <i>Fresnel Diffraction by a Square Aperture /</i>	
	4.5.2 <i>Fresnel Diffraction by a Sinusoidal Amplitude</i>	
	<i>Grating — Talbot Images</i>	
	Problems--Chapter 4	90
5	Wave-Optics Analysis of Coherent Optical Systems	96
5.1	A Thin Lens as a Phase Transformation	96
	5.1.1 <i>The Thickness Function / 5.1.2 The Paraxial</i>	
	<i>Approximation / 5.1.3 The Phase Transformation and</i>	
	<i>Its Physical Meaning</i>	
5.2	Fourier Transforming Properties of Lenses	101
	5.2.1 <i>Input Placed Against the Lens / 5.2.2 Input Placed in Front</i>	
	<i>of the Lens / 5.2.3 Input Placed Behind the Lens / 5.2.4 Example</i>	
	<i>of an Optical Fourier Transform</i>	

5.3	Image Formation: Monochromatic Illumination	108
	<i>5.3.1 The Impulse Response of a Positive Lens / 5.3.2 Eliminating Quadratic Phase Factors: The Lens Law / 5.3.3 The Relation Between Object and Image</i>	
5.4	Analysis of Complex Coherent Optical Systems	114
	<i>5.4.1 An Operator Notation / 5.4.2 Application of the Operator Approach to Some Optical Systems</i>	
	Problems--Chapter 5	120
6	Frequency Analysis of Optical Imaging Systems	126
6.1	Generalized Treatment of Imaging Systems	127
	<i>6.1.1 A Generalized Model / 6.1.2 Effects of Diffraction on the Image / 6.1.3 Polychromatic Illumination: The Coherent and Incoherent Cases</i>	
6.2	Frequency Response for Diffraction-Limited Coherent Imaging	134
	<i>6.2.1 The Amplitude Transfer Function / 6.2.2 Examples of Amplitude Transfer Functions</i>	
6.3	Frequency Response for Diffraction-Limited Incoherent Imaging	137
	<i>6.3.1 The Optical Transfer Function / 6.3.2 General Properties of the OTF / 6.3.3 The OTF of an Aberration-Free System / 6.3.4 Examples of Diffraction-Limited OTFs</i>	
6.4	Aberrations and Their Effects on Frequency Response	145
	<i>6.4.1 The Generalized Pupil Function / 6.4.2 Effects of Aberrations on the Amplitude Transfer Function / 6.4.3 Effects of Aberrations on the OTF / 6.4.4 Example of a Simple Aberration: A Focusing Error / 6.4.5 Apodization and Its Effects on Frequency Response</i>	
6.5	Comparison of Coherent and Incoherent Imaging	154
	<i>6.5.1 Frequency Spectrum of the Image Intensity / 6.5.2 Two-Point Resolution / 6.5.3 Other Effects</i>	
6.6	Resolution Beyond the Classical Diffraction Limit	160
	<i>6.6.1 Underlying Mathematical Fundamentals / 6.6.2 Intuitive Explanation of Bandwidth Extrapolation / 6.6.3 An Extrapolation Method Based on the Sampling Theorem / 6.6.4 An Iterative Extrapolation Method / 6.6.5 Practical Limitations</i>	
	Problems--Chapter 6	165
7	Wavefront Modulation	172
7.1	Wavefront Modulation with Photographic Film	173
	<i>7.1.1 The Physical Processes of Exposure, Development, and Fixing / 7.1.2 Definition of Terms / 7.1.3 Film in an Incoherent</i>	

	<i>Optical System / 7.1.4 Film in a Coherent Optical System / 7.1.5 The Modulation Transfer Function / 7.1.6 Bleaching of Photographic Emulsions</i>	
7.2	Spatial Light Modulators	184
	<i>7.2.1 Properties of Liquid Crystals / 7.2.2 Spatial Light Modulators Based on Liquid Crystals / 7.2.3 Magneto-Optic Spatial Light Modulators / 7.2.4 Deformable Mirror Spatial Light Modulators / 7.2.5 Multiple Quantum Well Spatial Light Modulators / 7.2.6 Acousto-Optic Spatial Light Modulators</i>	
7.3	Diffraction Optical Elements	209
	<i>7.3.1 Binary Optics / 7.3.2 Other Types of Diffraction Optics / 7.3.3 A Word of Caution</i>	
	Problems--Chapter 7	215
8	Analog Optical Information Processing	217
8.1	Historical Background	218
	<i>8.1.1 The Abbe-Porter Experiments / 8.1.2 The Zernike Phase-Contrast Microscope / 8.1.3 Improvement of Photographs: Maréchal / 8.1.4 The Emergence of a Communications Viewpoint / 8.1.5 Application of Coherent Optics to More General Data Processing</i>	
8.2	Incoherent Image Processing Systems	224
	<i>8.2.1 Systems Based on Geometrical Optics / 8.2.2 Systems That Incorporate the Effects of Diffraction</i>	
8.3	Coherent Optical Information Processing Systems	232
	<i>8.3.1 Coherent System Architectures / 8.3.2 Constraints on Filter Realization</i>	
8.4	The VanderLugt Filter	237
	<i>8.4.1 Synthesis of the Frequency-Plane Mask / 8.4.2 Processing the Input Data / 8.4.3 Advantages of the VanderLugt Filter</i>	
8.5	The Joint Transform Correlator	243
8.6	Application to Character Recognition	246
	<i>8.6.1 The Matched Filter / 8.6.2 A Character-Recognition Problem / 8.6.3 Optical Synthesis of a Character-Recognition Machine / 8.6.4 Sensitivity to Scale Size and Rotation</i>	
8.7	Optical Approaches to Invariant Pattern Recognition	252
	<i>8.7.1 Mellin Correlators / 8.7.2 Circular Harmonic Correlation / 8.7.3 Synthetic Discriminant Functions</i>	
8.8	Image Restoration	257
	<i>8.8.1 The Inverse Filter / 8.8.2 The Wiener Filter, or the Least-Mean-Square-Error Filter / 8.8.3 Filter Realization</i>	
8.9	Processing Synthetic-Aperture Radar (SAR) Data	264
	<i>8.9.1 Formation of the Synthetic Aperture / 8.9.2 The Collected Data and the Recording Format / 8.9.3 Focal Properties of the</i>	

	<i>Film Transparency / 8.9.4 Forming a Two-Dimensional Image / 8.9.5 The Tilted Plane Processor</i>	
8.10	Acousto-Optic Signal Processing Systems	276
	<i>8.10.1 Bragg Cell Spectrum Analyzer / 8.10.2 Space-Integrating Correlator / 8.10.3 Time-Integrating Correlator / 8.10.4 Other Acousto-Optic Signal Processing Architectures</i>	
8.11	Discrete Analog Optical Processors	282
	<i>8.11.1 Discrete Representation of Signals and Systems / 8.11.2 A Serial Matrix-Vector Multiplier / 8.11.3 A Parallel Incoherent Matrix-Vector Multiplier / 8.11.4 An Outer Product Processor / 8.11.5 Other Discrete Processing Architectures / 8.11.6 Methods for Handling Bipolar and Complex Data</i>	
	Problems — Chapter 8	290
9	Holography	295
9.1	Historical Introduction	295
9.2	The Wavefront Reconstruction Problem	296
	<i>9.2.1 Recording Amplitude and Phase / 9.2.2 The Recording Medium / 9.2.3 Reconstruction of the Original Wavefront / 9.2.4 Linearity of the Holographic Process / 9.2.5 Image Formation by Holography</i>	
9.3	The Gabor Hologram	302
	<i>9.3.1 Origin of the Reference Wave / 9.3.2 The Twin Images / 9.3.3 Limitations of the Gabor Hologram</i>	
9.4	The Leith-Upatnieks Hologram	304
	<i>9.4.1 Recording the Hologram / 9.4.2 Obtaining the Reconstructed Images / 9.4.3 The Minimum Reference Angle / 9.4.4 Holography of Three-Dimensional Scenes / 9.4.5 Practical Problems in Holography</i>	
9.5	Image Locations and Magnification	314
	<i>9.5.1 Image Locations / 9.5.2 Axial and Transverse Magnifications / 9.5.3 An Example</i>	
9.6	Some Different Types of Holograms	319
	<i>9.6.1 Fresnel, Fraunhofer, Image, and Fourier Holograms / 9.6.2 Transmission and Reflection Holograms / 9.6.3 Holographic Stereograms / 9.6.4 Rainbow Holograms / 9.6.5 Multiplex Holograms / 9.6.6 Embossed Holograms</i>	
9.7	Thick Holograms	329
	<i>9.7.1 Recording a Volume Holographic Grating / 9.7.2 Reconstructing Wavefronts from a Volume Grating / 9.7.3 Fringe Orientations for More Complex Recording Geometries / 9.7.4 Gratings of Finite Size / 9.7.5 Diffraction Efficiency-Coupled Mode Theory</i>	

9.8	Recording Materials	346
	<i>9.8.1 Silver Halide Emulsions / 9.8.2 Photopolymer Films / 9.8.3 Dichromated Gelatin / 9.8.4 Photorefractive Materials</i>	
9.9	Computer-Generated Holograms	351
	<i>9.9.1 The Sampling Problem / 9.9.2 The Computational Problem / 9.9.3 The Representational Problem</i>	
9.10	Degradations of Holographic Images	363
	<i>9.10.1 Effects of Film MTF / 9.10.2 Effects of Film Nonlinearities / 9.10.3 Effects of Film-Grain Noise / 9.10.4 Speckle Noise</i>	
9.11	Holography with Spatially Incoherent Light	369
9.12	Applications of Holography	372
	<i>9.12.1 Microscopy and High-Resolution Volume Imagery / 9.12.2 Interferometry / 9.12.3 Imaging Through Distorting Media / 9.12.4 Holographic Data Storage / 9.12.5 Holographic Weights for Artificial Neural Networks / 9.12.6 Other Applications</i>	
	Problems--Chapter 9	388
A	Delta Functions and Fourier Transform Theorems	393
A.1	Delta Functions	393
A.2	Derivation of Fourier Transform Theorems	395
B	Introduction to Paraxial Geometrical Optics	401
B.1	The Domain of Geometrical Optics	401
B.2	Refraction, Snell's Law, and the Paraxial Approximation	403
B.3	The Ray-Transfer Matrix	404
B.4	Conjugate Planes, Focal Planes, and Principal Planes	407
B.5	Entrance and Exit Pupils	411
C	Polarization and Jones Matrices	415
C.1	Definition of the Jones Matrix	415
C.2	Examples of Simple Polarization Transformations	417
C.3	Reflective Polarization Devices	418
	Bibliography	421
	Index	433

PREFACE

Fourier analysis is a ubiquitous tool that has found application to diverse areas of physics and engineering. This book deals with its applications in optics, and in particular with applications to diffraction, imaging, optical data processing, and holography.

Since the subject covered is Fourier Optics, it is natural that the methods of Fourier analysis play a key role as the underlying analytical structure of our treatment. Fourier analysis is a standard part of the background of most physicists and engineers. The theory of linear systems is also familiar, especially to electrical engineers. Chapter 2 reviews the necessary mathematical background. For those not already familiar with Fourier analysis and linear systems theory, it can serve as the outline for a more detailed study that can be made with the help of other textbooks explicitly aimed at this subject. Ample references are given for more detailed treatments of this material. For those who have already been introduced to Fourier analysis and linear systems theory, that experience has usually been with functions of a single independent variable, namely time. The material presented in Chapter 2 deals with the mathematics in two spatial dimensions (as is necessary for most problems in optics), yielding an extra richness not found in the standard treatments of the one-dimensional theory.

The original edition of this book has been considerably expanded in this second edition, an expansion that was needed due to the tremendous amount of progress in the field since 1968 when the first edition was published. The book can be used as a textbook to satisfy the needs of several different types of courses. It is directed towards both physicists and engineers, and the portions of the book used in the course will in general vary depending on the audience. However, by properly selecting the material to be covered, the needs of any of a number of different audiences can be met. This Preface will make several explicit suggestions for the shaping of different kinds of courses.

First a one-quarter or one-semester course on diffraction and image formation can be constructed from the materials covered in Chapters 2 through 6, together with all three appendices. If time is short, the following sections of these chapters can be omitted or left as reading for the advanced student: 3.8, 3.9, 5.4, and 6.6.

A second type of one-quarter or one-semester course would cover the basics of Fourier Optics, but then focus on the application area of analog optical signal processing. For such a course, I would recommend that Chapter 2 be left to the reading of the student, that the material of Chapter 3 be begun with Section 3.7, and followed by Section 3.10, leaving the rest of this chapter to a reading by those students who are curious as to the origins of the Huygens-Fresnel principle. In Chapter 4, Sections 4.2.2 and 4.5.1 can be skipped. Chapter 5 can begin with Eq. (5-10) for the amplitude transmittance function of a thin lens, and can include all the remaining material, with the exception that Section 5.4 can be left as reading for the advanced students. If time is short, Chapter 6 can be skipped entirely. For this course, virtually all of the material presented in Chapter 7 is important, as is much of the material in Chapter 8. If it is necessary to reduce the amount of material, I would recommend that the following sections be omitted: 8.2, 8.8, and 8.9. It is often desirable to include some subset of the material

on holography from Chapter 9 in this course. I would include sections 9.4, 9.6.1, 9.6.2, 9.7.1, 9.7.2, 9.8, 9.9, and 9.12.5. The three appendices should be read by the students but need not be covered in lectures.

A third variation would be a one-quarter or one-semester course that covers the basics of Fourier Optics but focuses on holography as an application. The course can again begin with Section 3.7 and be followed by Section 3.10. The coverage through Chapter 5 can be identical with that outlined above for the course that emphasizes optical signal processing. In this case, the material of Sections 6.1, 6.2, 6.3, and 6.5 can be included. In Chapter 7, only Section 7.1 is needed, although Section 7.3 is a useful addition if there is time. Chapter 8 can now be skipped and Chapter 9 on holography can be the focus of attention. If time is short, Sections 9.10 and 9.11 can be omitted. The first two appendices should be read by the students, and the third can be skipped.

In some universities, more than one quarter or one semester can be devoted to this material. In two quarters or two semesters, most of the material in this book can be covered.

The above suggestions can of course be modified to meet the needs of a particular set of students or to emphasize the material that a particular instructor feels is most appropriate. I hope that these suggestions will at least give some ideas about possibilities.

There are many people to whom I owe a special word of thanks for their help with this new edition of the book. Early versions of the manuscript were used in courses at several different universities. I would in particular like to thank Profs. A.A. Sawchuk, J.F. Walkup, J. Leger, P. Pichon, D. Mehrl, and their many students for catching so many typographical errors and in some cases outright mistakes. Helpful comments were also made by I. Erteza and M. Bashaw, for which I am grateful. Several useful suggestions were also made by anonymous manuscript reviewers engaged by the publisher. A special debt is owed to Prof. Emmett Leith, who provided many helpful suggestions. I would also like to thank the students in my 1995 Fourier Optics class, who competed fiercely to see who could find the most mistakes. Undoubtedly there are others to whom I owe thanks, and I apologize for not mentioning them explicitly here.

Finally, I thank Hon Mai, without whose patience, encouragement and support this book would not have been possible.

Joseph W. Goodman

Introduction to Fourier Optics

CHAPTER 1

Introduction

1.1 OPTICS, INFORMATION, AND COMMUNICATION

Since the late 1930s, the venerable branch of physics known as optics has gradually developed ever-closer ties with the communication and information sciences of electrical engineering. The trend is understandable, for both communication systems and imaging systems are designed to collect or convey information. In the former case, the information is generally of a temporal nature (e.g. a modulated voltage or current waveform), while in the latter case it is of a spatial nature (e.g. a light amplitude or intensity distribution over space), but from an abstract point of view, this difference is a rather superficial one.

Perhaps the strongest tie between the two disciplines lies in the similar mathematics which can be used to describe the respective systems of interest – the mathematics of Fourier analysis and systems theory. The fundamental reason for the similarity is not merely the common subject of "information", but rather certain basic properties which communication systems and imaging systems share. For example, many electronic networks and imaging devices share the properties called *linearity* and *invariance* (for definitions see Chapter 2). Any network or device (electronic, optical, or otherwise) which possesses these two properties can be described mathematically with considerable ease using the techniques of *frequency analysis*. Thus, just as it is convenient to describe an audio amplifier in terms of its (temporal) frequency response, so too it is often convenient to describe an imaging system in terms of its (spatial) frequency response.

The similarities do not end when the linearity and invariance properties are absent. Certain nonlinear optical elements (e.g. photographic film) have input-output relationships which are directly analogous to the corresponding characteristics of nonlinear electronic components (diodes, transistors, etc.), and similar mathematical analysis can be applied in both cases.

2 Introduction to Fourier Optics

It is particularly important to recognize that the similarity of the mathematical structures can be exploited not only for analysis purposes but also for *synthesis* purposes. Thus, just as the spectrum of a temporal function can be intentionally manipulated in a prescribed fashion by filtering, so too can the spectrum of a spatial function be modified in various desired ways. The history of optics is rich with examples of important advances achieved by application of Fourier synthesis techniques – the Zernike phase-contrast microscope is an example that was worthy of a Nobel prize. Many other examples can be found in the fields of signal and image processing.

1.2 THE BOOK

The readers of this book are assumed at the start to have a solid foundation in Fourier analysis and linear systems theory. Chapter 2 reviews the required background; to avoid boring those who are well grounded in the analysis of temporal signals and systems, the review is conducted for functions of two independent variables. Such functions are, of course, of primary concern in optics, and the extension from one to two independent variables provides a new richness to the mathematical theory, introducing many new properties which have no direct counterpart in the theory of temporal signals and systems.

The phenomenon called *diffraction* is of the utmost importance in the theory of optical systems. Chapter 3 treats the foundations of scalar diffraction theory, including the Kirchhoff, Rayleigh-Sommerfeld, and angular spectrum approaches. In Chapter 4, certain approximations to the general results are introduced, namely the Fresnel and Fraunhofer approximations, and examples of diffraction-pattern calculations are presented.

Chapter 5 considers the analysis of coherent optical systems which consist of lenses and free-space propagation. The approach is that of wave optics, rather than the more common geometrical optics method of analysis. A thin lens is modeled as a quadratic phase transformation; the usual lens law is derived from this model, as are certain Fourier transforming properties of lenses.

Chapter 6 considers the application of frequency analysis techniques to both coherent and incoherent imaging systems. Appropriate transfer functions are defined and their properties discussed for systems with and without aberrations. Coherent and incoherent systems are compared from various points of view. The limits to achievable resolution are derived.

In Chapter 7 the subject of wavefront modulation is considered. The properties of photographic film as an input medium for incoherent and coherent optical systems are discussed. Attention is then turned to spatial light modulators, which are devices for entering information into optical systems in real time or near real time. Finally, diffractive optical elements are described in some detail.

Attention is turned to analog optical information processing in Chapter 8. Both continuous and discrete processing systems are considered. Applications to image

enhancement, pattern recognition, and processing of synthetic-aperture radar data are considered.

The final chapter is devoted to the subject of holography. The techniques developed by **Gabor** and by **Leith** and Upatnieks are considered in detail and compared. Both thin and thick holograms are treated. Extensions to three-dimensional imaging are presented. Various applications of holography are described, but emphasis is on the fundamentals.

Analysis of Two-Dimensional Signals and Systems

Many physical phenomena are found experimentally to share the basic property that their response to several stimuli acting simultaneously is identically equal to the sum of the responses that each component stimulus would produce individually. Such phenomena are called *linear*; and the property they share is called *linearity*. Electrical networks composed of resistors, capacitors, and inductors are usually linear over a wide range of inputs. In addition, as we shall soon see, the wave equation describing the propagation of light through most media leads us naturally to regard optical imaging operations as linear mappings of "object" light distributions into "image" light distributions.

The single property of linearity leads to a vast simplification in the mathematical description of such phenomena and represents the foundation of a mathematical structure which we shall refer to here as *linear systems theory*. The great advantage afforded by linearity is the ability to express the response (be it voltage, current, light amplitude, or light intensity) to a complicated stimulus in terms of the responses to certain "elementary" stimuli. Thus if a stimulus is decomposed into a linear combination of elementary stimuli, each of which produces a known response of convenient form, then by virtue of linearity, the total response can be found as a corresponding linear combination of the responses to the elementary stimuli.

In this chapter we review some of the mathematical tools that are useful in describing linear phenomena, and discuss some of the mathematical decompositions that are often employed in their analysis. Throughout the later chapters we shall be concerned with stimuli (system inputs) and responses (system outputs) that may be either of two different physical quantities. If the illumination used in an optical system exhibits a property called *spatial coherence*, then we shall find that it is appropriate to describe the light as a spatial distribution of *complex-valued* field amplitude. When the illumination is totally lacking in spatial coherence, it is appropriate to describe the light as a spatial distribution of *real-valued* intensity. Attention will be focused here on the analysis of linear systems with complex-valued inputs; the results for real-valued inputs are thus included as special cases of the theory.

2.1

FOURIER ANALYSIS IN TWO DIMENSIONS

A mathematical tool of great utility in the analysis of both linear and nonlinear phenomena is Fourier analysis. This tool is widely used in the study of electrical networks and communication systems; it is assumed that the reader has encountered Fourier theory previously, and therefore that he or she is familiar with the analysis of functions of one independent variable (e.g. time). For a review of the fundamental mathematical concepts, see the books by Papoulis [226], Bracewell [32], and Gray and Goodman [131]. A particularly relevant treatment is by Bracewell [33]. Our purpose here is limited to extending the reader's familiarity to the analysis of functions of two independent variables. No attempt at great mathematical rigor will be made, but rather, an operational approach, characteristic of most engineering treatments of the subject, will be adopted.

2.1.1 Definition and Existence Conditions

The Fourier transform (alternatively the Fourier spectrum or frequency spectrum) of a (in general, complex-valued) function g of two independent variables x and y will be represented here by $\mathcal{F}\{g\}$ and is defined by¹

$$\mathcal{F}\{g\} = \iint_{-\infty}^{\infty} g(x, y) \exp[-j2\pi(f_x x + f_y y)] dx dy. \quad (2-1)$$

The transform so defined is itself a complex-valued function of two independent variables f_x and f_y , which we generally refer to as frequencies. Similarly, the inverse Fourier transform of a function $G(f_x, f_y)$ will be represented by $\mathcal{F}^{-1}\{G\}$ and is defined as

$$\mathcal{F}^{-1}\{G\} = \iint_{-\infty}^{\infty} G(f_x, f_y) \exp[j2\pi(f_x x + f_y y)] df_x df_y. \quad (2-2)$$

Note that as mathematical operations the transform and inverse transform are very similar, differing only in the sign of the exponent appearing in the integrand. The inverse Fourier transform is sometimes referred to as the Fourier integral representation of a function $g(x, y)$.

Before discussing the properties of the Fourier transform and its inverse, we must first decide when (2-1) and (2-2) are in fact meaningful. For certain functions, these integrals may not exist in the usual mathematical sense, and therefore this discussion would be incomplete without at least a brief mention of "existence conditions". While a variety of sets of *sufficient* conditions for the existence of (2-1) are possible, perhaps the most common set is the following:

¹When a **single** limit of integration appears above or below a double integral, then that limit applies to both integrations.

6 Introduction to Fourier Optics

1. g must be absolutely integrable over the infinite (x, y) plane.
2. g must have only a finite number of discontinuities and a finite number of maxima and minima in any finite rectangle.
3. g must have no infinite discontinuities.

In general, any one of these conditions can be weakened at the price of strengthening one or both of the companion conditions, but such considerations lead us rather far afield from our purposes here.

As Bracewell [32] has pointed out, "physical possibility is a valid sufficient condition for the existence of a transform." However, it is often convenient in the analysis of systems to represent true physical waveforms by idealized mathematical functions, and for such functions one or more of the above existence conditions may be violated. For example, it is common to represent a strong, **narrow** time pulse by the so-called Dirac delta function² often represented by

$$\delta(t) = \lim_{N \rightarrow \infty} N \exp(-N^2 \pi t^2), \quad (2-3)$$

where the limit operation provides a convenient mental construct but is not meant to be taken literally. See Appendix A for more details. Similarly, an idealized point source of light is often represented by the two-dimensional equivalent,

$$\delta(x, y) = \lim_{N \rightarrow \infty} N^2 \exp[-N^2 \pi (x^2 + y^2)]. \quad (2-4)$$

Such "functions", being infinite at the origin and zero elsewhere, have an infinite discontinuity and therefore fail to satisfy existence condition 3. Other important examples are readily found; for example, the functions

$$f(x, y) = 1 \quad \text{and} \quad f(x, y) = \cos(2\pi f_x x) \quad (2-5)$$

both fail to satisfy existence condition 1.

If the majority of functions of interest are to be included within the framework of Fourier analysis, some generalization of the definition (2-1) is required. Fortunately, it is often possible to find a meaningful transform of functions that do not strictly satisfy the existence conditions, provided those functions can be defined as the limit of a sequence of functions that are transformable. By transforming each member function of the defining sequence, a corresponding sequence of transforms is generated, and we call the limit of this new sequence the generalized Fourier transform of the original function. Generalized transforms can be manipulated in the same manner as conventional transforms, and the distinction between the two cases can generally be ignored, it being understood that when a function fails to satisfy the existence conditions and yet is said to have a transform, then the generalized transform is actually meant. For a more detailed discussion of this generalization of Fourier analysis the reader is referred to the book by Lighthill [194].

To illustrate the calculation of a generalized transform, consider the Dirac delta function, which has been seen to violate existence condition 3. Note that each member function of the defining sequence (2-4) does satisfy the existence requirements and that each, in fact, has a Fourier transform given by (see Table 2.1)

²For a more detailed discussion of the delta function, including definitions, see Appendix A.

$$\mathcal{F}\{N^2 \exp[-N^2 \pi(x^2 + y^2)]\} = \exp\left[-\frac{\pi(f_X^2 + f_Y^2)}{N^2}\right]. \quad (2-6)$$

Accordingly the generalized transform of $\delta(x, y)$ is found to be

$$\mathcal{F}\{\delta(x, y)\} = \lim_{N \rightarrow \infty} \left\{ \exp\left[-\frac{\pi(f_X^2 + f_Y^2)}{N^2}\right] \right\} = 1. \quad (2-7)$$

Note that the spectrum of a delta function extends uniformly over the entire frequency domain.

For other examples of generalized transforms, see Table 2.1.

2.1.2 The Fourier Transform as a Decomposition

As mentioned previously, when dealing with linear systems it is often useful to decompose a complicated input into a number of more simple inputs, to calculate the response of the system to each of these "elementary" functions, and to superimpose the individual responses to find the total response. Fourier analysis provides the basic means of **performing** such a decomposition. Consider the familiar inverse transform relationship

$$g(t) = \int_{-\infty}^{\infty} G(f) \exp(j2\pi f t) df \quad (2-8)$$

expressing the time function g in terms of its frequency spectrum. We may regard this expression as a decomposition of the function $g(t)$ into a linear combination (in this case an integral) of elementary functions, each with a specific form $\exp(j2\pi f t)$. From this it is clear that the complex number $G(f)$ is simply a weighting factor that must be applied to the elementary function of frequency f in order to synthesize the desired $g(t)$.

In a similar fashion, we may regard the *two-dimensional* Fourier transform as a decomposition of a function $g(x, y)$ into a linear combination of elementary functions of the form $\exp[j2\pi(f_X x + f_Y y)]$. Such functions have a number of interesting properties. Note that for any particular frequency pair (f_X, f_Y) the corresponding elementary function has a phase that is zero or an integer multiple of 2π radians along lines described by the equation

$$y = -\frac{f_X}{f_Y} x + \frac{n}{f_Y}, \quad (2-9)$$

where n is an integer. Thus, as indicated in Fig. 2.1, this elementary function may be regarded as being "directed" in the (x, y) plane at an angle θ (with respect to the x axis) given by

$$\theta = \arctan\left(\frac{f_Y}{f_X}\right). \quad (2-10)$$

In addition, the spatial *period* (i.e. the distance between zero-phase lines) is given by

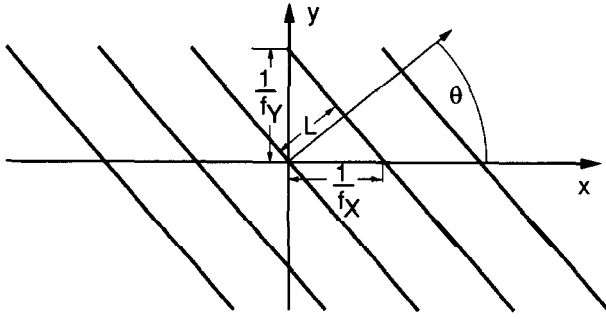


FIGURE 2.1
Lines of zero phase for the function $\exp[j2\pi(f_X x + f_Y y)]$.

$$L = \frac{1}{\sqrt{f_X^2 + f_Y^2}}. \quad (2-11)$$

In conclusion, then, we may again regard the inverse Fourier transform as providing a means for decomposing mathematical functions. The Fourier spectrum G of a function g is simply a description of the weighting factors that must be applied to each elementary function in order to synthesize the desired g . The real advantage obtained from using this decomposition will not be fully evident until our later discussion of invariant linear systems.

2.1.3 Fourier Transform Theorems

The basic definition (2-1) of the Fourier transform leads to a rich mathematical structure associated with the transform operation. We now consider a few of the basic mathematical properties of the transform, properties that will find wide use in later material. These properties are presented as mathematical theorems, followed by brief statements of their physical significance. Since these theorems are direct extensions of the analogous one-dimensional statements, the proofs are deferred to Appendix A.

1. **Linearity theorem.** $\mathcal{F}\{\alpha g + \beta h\} = \alpha \mathcal{F}\{g\} + \beta \mathcal{F}\{h\}$; that is, the transform of a weighted sum of two (or more) functions is simply the identically weighted sum of their individual transforms.
2. **Similarity theorem.** If $\mathcal{F}\{g(x, y)\} = G(f_X, f_Y)$, then

$$\mathcal{F}\{g(ax, by)\} = \frac{1}{|ab|} G\left(\frac{f_X}{a}, \frac{f_Y}{b}\right); \quad (2-12)$$

that is, a "stretch" of the coordinates in the space domain (x, y) results in a contraction of the coordinates in the frequency domain (f_X, f_Y) , plus a change in the overall amplitude of the spectrum.

3. **Shift theorem.** If $\mathcal{F}\{g(x, y)\} = G(f_X, f_Y)$, then

$$\mathcal{F}\{g(x - a, y - b)\} = G(f_X, f_Y) \exp[-j2\pi(f_X a + f_Y b)]; \quad (2-13)$$

that is, translation in the space domain introduces a linear phase shift in the frequency domain.

4. Rayleigh's theorem (Parseval's theorem). If $\mathcal{F}\{g(x, y)\} = G(f_X, f_Y)$, then

$$\iint_{-\infty}^{\infty} |g(x, y)|^2 dx dy = \iint_{-\infty}^{\infty} |G(f_X, f_Y)|^2 df_X df_Y. \quad (2-14)$$

The integral on the left-hand side of this theorem can be interpreted as the energy contained in the waveform $g(x, y)$. This in turn leads us to the idea that the quantity $|G(f_X, f_Y)|^2$ can be interpreted as an energy density in the frequency domain.

5. Convolution theorem. If $\mathcal{F}\{g(x, y)\} = G(f_X, f_Y)$ and $\mathcal{F}\{h(x, y)\} = H(f_X, f_Y)$, then

$$\mathcal{F}\left\{\iint_{-\infty}^{\infty} g(\xi, \eta) h(x - \xi, y - \eta) d\xi d\eta\right\} = G(f_X, f_Y)H(f_X, f_Y). \quad (2-15)$$

The convolution of two functions in the space domain (an operation that will be found to arise frequently in the theory of linear systems) is entirely equivalent to the more simple operation of multiplying their individual transforms and inverse transforming.

6. Autocorrelation theorem. If $\mathcal{F}\{g(x, y)\} = G(f_X, f_Y)$, then

$$\mathcal{F}\left\{\iint_{-\infty}^{\infty} g(\xi, \eta) g^*(\xi - x, \eta - y) d\xi d\eta\right\} = |G(f_X, f_Y)|^2. \quad (2-16)$$

Similarly,

$$\mathcal{F}\{|g(x, y)|^2\} = \iint_{-\infty}^{\infty} G(\xi, \eta) G^*(\xi - f_X, \eta - f_Y) d\xi d\eta. \quad (2-17)$$

This theorem may be regarded as a special case of the convolution theorem in which we convolve $g(x, y)$ with $g^*(-x, -y)$.

7. Fourier integral theorem. At each point of continuity of g ,

$$\mathcal{F}\mathcal{F}^{-1}\{g(x, y)\} = \mathcal{F}^{-1}\mathcal{F}\{g(x, y)\} = g(x, y). \quad (2-18)$$

At each point of discontinuity of g , the two successive transforms yield the angular average of the values of g in a small neighborhood of that point. That is, the successive transformation and inverse transformation of a function yields that function again, except at points of discontinuity.

The above transform theorems are of far more than just theoretical interest. They will be used frequently, since they provide the basic tools for the manipulation of Fourier transforms and can save enormous amounts of work in the solution of Fourier analysis problems.

2.1.4 Separable Functions

A function of two independent variables is called separable with respect to a specific coordinate system if it can be written as a product of two functions, each of which depends on only one of the independent variables. Thus the function g is separable in rectangular coordinates (x, y) if

$$g(x, y) = g_X(x) g_Y(y), \quad (2-19)$$

while it is separable in polar coordinates (r, θ) if

$$g(r, \theta) = g_R(r) g_\Theta(\theta). \quad (2-20)$$

Separable functions are often more convenient to deal with than more general functions, for separability often allows complicated two-dimensional manipulations to be reduced to more simple one-dimensional manipulations. For example, a function separable in rectangular coordinates has the particularly simple property that its two-dimensional Fourier transform can be found as a product of one-dimensional Fourier transforms, as evidenced by the following relation:

$$\begin{aligned} \mathcal{F}\{g(x, y)\} &= \iint_{-\infty}^{\infty} g(x, y) \exp[-j2\pi(f_X x + f_Y y)] dx dy \\ &= \int_{-\infty}^{\infty} g_X(x) \exp[-j2\pi f_X x] dx \int_{-\infty}^{\infty} g_Y(y) \exp[-j2\pi f_Y y] dy \\ &= \mathcal{F}_X\{g_X\} \mathcal{F}_Y\{g_Y\}. \end{aligned} \quad (2-21)$$

Thus the transform of g is itself separable into a product of two factors, one a function of f_X only and the second a function of f_Y only, and the process of two-dimensional transformation simplifies to a succession of more familiar one-dimensional manipulations.

Functions separable in polar coordinates are not so easily handled as those separable in rectangular coordinates, but it is still generally possible to demonstrate that two-dimensional manipulations can be performed by a series of one-dimensional manipulations. For example, the reader is asked to verify in the problems that the Fourier transform of a general function g separable in polar coordinates can be expressed as an infinite sum of weighted *Hankel* transforms

$$\mathcal{F}\{g(r, \theta)\} = \sum_{k=-\infty}^{\infty} c_k (-j)^k \exp(jk\phi) \mathcal{H}_k\{g_R(r)\} \quad (2-22)$$

where

$$c_k = \frac{1}{2\pi} \int_0^{2\pi} g_\Theta(\theta) \exp(-jk\theta) d\theta$$

and $\mathcal{H}_k\{\}$ is the *Hankel* transform operator of order k , defined by

$$\mathcal{H}_k\{g_R(r)\} = 2\pi \int_0^\infty r g_R(r) J_k(2\pi r \rho) dr. \quad (2-23)$$

Here the function J_k is the k th-order Bessel function of the first kind.

2.1.5 Functions with Circular Symmetry: Fourier-Bessel Transforms

Perhaps the simplest class of functions separable in polar coordinates is composed of those possessing circular symmetry. The function g is said to be circularly symmetric if it can be written as a function of r alone, that is,

$$g(r, \theta) = g_R(r). \quad (2-24)$$

Such functions play an important role in the problems of interest here, since most optical systems have precisely this type of symmetry. We accordingly devote special attention to the problem of Fourier transforming a circularly symmetric function.

The Fourier transform of g in a system of rectangular coordinates is, of course, given by

$$G(f_X, f_Y) = \iint_{-\infty}^{\infty} g(x, y) \exp[-j2\pi(f_X x + f_Y y)] dx dy. \quad (2-25)$$

To fully exploit the circular symmetry of g , we make a transformation to polar coordinates in both the (x, y) and the (f_X, f_Y) planes as follows:

$$\begin{aligned} r &= \sqrt{x^2 + y^2} & x &= r \cos \theta \\ \theta &= \arctan\left(\frac{y}{x}\right) & y &= r \sin \theta \\ \rho &= \sqrt{f_X^2 + f_Y^2} & f_X &= \rho \cos \phi \\ \phi &= \arctan\left(\frac{f_Y}{f_X}\right) & f_Y &= \rho \sin \phi. \end{aligned} \quad (2-26)$$

For the present we write the transform as a function of both radius and angle,³

$$\mathcal{F}\{g\} = G_o(\rho, \phi). \quad (2-27)$$

Applying the coordinate transformations (2-26) to **Eq. (2-25)**, the Fourier transform of g can be written

$$G_o(\rho, \phi) = \int_0^{2\pi} d\theta \int_0^\infty dr r g_R(r) \exp[-j2\pi r \rho (\cos \theta \cos \phi + \sin \theta \sin \phi)] \quad (2-28)$$

or equivalently,

$$G_o(\rho, \phi) = \int_0^\infty dr r g_R(r) \int_0^{2\pi} d\theta \exp[-j2\pi r \rho \cos(\theta - \phi)]. \quad (2-29)$$

³Note the subscript in G_o is added simply because the functional form of the expression for the transform in polar coordinates is in general different than the functional form for the same transform in rectangular coordinates.

Finally, we use the Bessel function identity

$$J_0(a) = \frac{1}{2\pi} \int_0^{2\pi} \exp[-ja \cos(\theta - \phi)] d\theta, \quad (2-30)$$

where J_0 is a Bessel function of the first kind, zero order, to simplify the expression for the transform. Substituting (2-30) in (2-29), the dependence of the transform on angle ϕ is seen to disappear, leaving G_0 as the following function of radius ρ ,

$$G_o(\rho, \phi) = G_o(\rho) = 2\pi \int_0^\infty r g_R(r) J_0(2\pi r \rho) dr. \quad (2-31)$$

Thus the Fourier transform of a circularly symmetric function is itself circularly symmetric and can be found by performing the one-dimensional manipulation of (2-31). This particular form of the Fourier transform occurs frequently enough to warrant a special designation; it is accordingly referred to as the Fourier-Bessel *transform*, or alternatively as the *Hankel* transform of zero order (cf. Eq. (2-23)). For brevity, we adopt the former terminology.

By means of arguments identical with those used above, the inverse Fourier transform of a circularly symmetric spectrum $G_o(\rho)$ can be expressed as

$$g_R(r) = 2\pi \int_0^\infty \rho G_o(\rho) J_0(2\pi r \rho) d\rho. \quad (2-32)$$

Thus for circularly symmetric functions there is no difference between the transform and the inverse-transform operations.

Using the notation $\mathcal{B}\{\}$ to represent the Fourier-Bessel transform operation, it follows directly from the Fourier integral theorem that

$$\mathcal{B}\mathcal{B}^{-1}\{g_R(r)\} = \mathcal{B}^{-1}\mathcal{B}\{g_R(r)\} = \mathcal{B}\mathcal{B}\{g_R(r)\} = g_R(r) \quad (2-33)$$

at each value of r where $g_R(r)$ is continuous. In addition, the similarity theorem can be straightforwardly applied (see Prob. 2-6c) to show that

$$\mathcal{B}\{g_R(ar)\} = \frac{1}{a^2} G_o\left(\frac{\rho}{a}\right). \quad (2-34)$$

When using the expression (2-31) for the Fourier-Bessel transform, the reader should remember that it is no more than a special case of the two-dimensional Fourier transform, and therefore any familiar property of the Fourier transform has an entirely equivalent counterpart in the terminology of Fourier-Bessel transforms.

2.1.6 Some Frequently Used Functions and Some Useful Fourier Transform Pairs

A number of mathematical functions will find such extensive use in later material that considerable time and effort can be saved by assigning them special notations of their own. Accordingly, we adopt the following definitions of some frequently used functions:

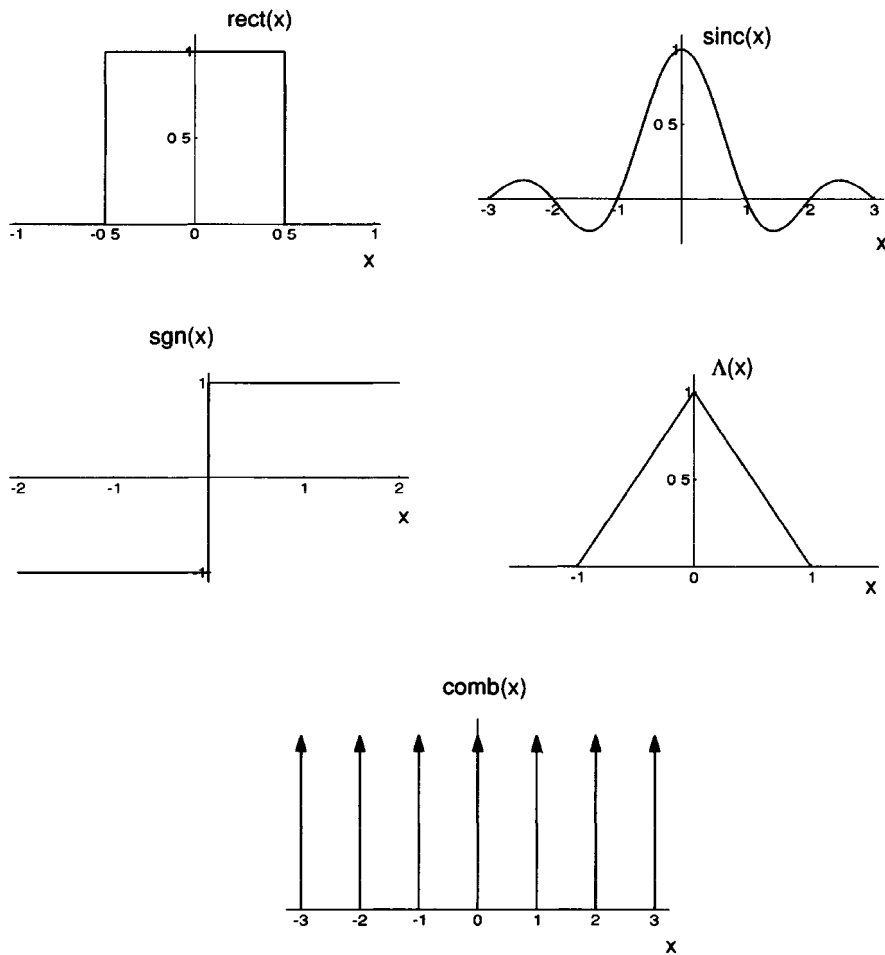


FIGURE 2.2
Special functions.

Rectangle function $\text{rect}(x) = \begin{cases} 1 & |x| < \frac{1}{2} \\ \frac{1}{2} & |x| = \frac{1}{2} \\ 0 & \text{otherwise} \end{cases}$

Sinc function $\text{sinc}(x) = \frac{\sin(\pi x)}{\pi x}$

Signum function $\text{sgn}(x) = \begin{cases} 1 & x > 0 \\ 0 & x = 0 \\ -1 & x < 0 \end{cases}$

Triangle function $\Lambda(x) = \begin{cases} 1 - |x| & |x| \leq 1 \\ 0 & \text{otherwise} \end{cases}$

Comb function $\text{comb}(x) = \sum_{n=-\infty}^{\infty} \delta(x - n)$

TABLE 2.1
Transform pairs for some functions separable in rectangular coordinates.

Function	Transform
$\exp[-\pi(a^2x^2 + b^2y^2)]$	$\frac{1}{ ab } \exp\left[-\pi\left(\frac{f_x^2}{a^2} + \frac{f_y^2}{b^2}\right)\right]$
$\text{rect}(ax) \text{rect}(by)$	$\frac{1}{ ab } \text{sinc}(f_x/a) \text{sinc}(f_y/b)$
$\Lambda(ax) \Lambda(by)$	$\frac{1}{ ab } \text{sinc}^2(f_x/a) \text{sinc}^2(f_y/b)$
$\delta(ax, by)$	$\frac{1}{ ab }$
$\exp[j\pi(ax + by)]$	$\delta(f_x - a/2, f_y - b/2)$
$\text{sgn}(ax) \text{sgn}(by)$	$\frac{ab}{ ab } \frac{1}{j\pi f_x} \frac{1}{j\pi f_y}$
$\text{comb}(ax) \text{comb}(by)$	$\frac{1}{ ab } \text{comb}(f_x/a) \text{comb}(f_y/b)$
$\exp[j\pi(a^2x^2 + b^2y^2)]$	$\frac{j}{ ab } \exp\left[-j\pi\left(\frac{f_x^2}{a^2} + \frac{f_y^2}{b^2}\right)\right]$
$\exp[-(a x + b y)]$	$\frac{1}{ ab } \frac{2}{1 + (2\pi f_x/a)^2} \frac{2}{1 + (2\pi f_y/b)^2}$

$$\text{Circle function} \quad \text{circ}(\sqrt{x^2 + y^2}) = \begin{cases} \sqrt{x^2 + y^2} & \leq 1 \\ 0 & \text{otherwise.} \end{cases}$$

The first five of these functions, depicted in Fig. 2.2, are all functions of only one independent variable; however, a variety of separable functions can be formed in two dimensions by means of products of these functions. The circle function is, of course, unique to the case of two-dimensional variables; see Fig. 2.3 for an illustration of its structure.

We conclude our discussion of Fourier analysis by presenting some specific two-dimensional transform pairs. Table 2.1 lists a number of transforms of functions separable in rectangular coordinates. For the convenience of the reader, the functions are presented with arbitrary scaling constants. Since the transforms of such functions can be found directly from products of familiar one-dimensional transforms, the proofs of these relations are left to the reader (cf. Prob. 2-2).

On the other hand, with a few exceptions (e.g. $\exp[-\pi(x^2 + y^2)]$, which is *both* separable in rectangular coordinates and circularly symmetric), transforms of most circularly symmetric functions cannot be found simply from a knowledge of one-dimensional transforms. The most frequently encountered function with circular symmetry is:

$$\text{circ}(r) = \begin{cases} 1 & r < 1 \\ \frac{1}{2} & r = 1 \\ 0 & \text{otherwise} \end{cases}.$$

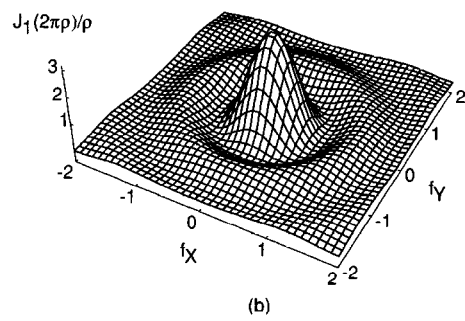
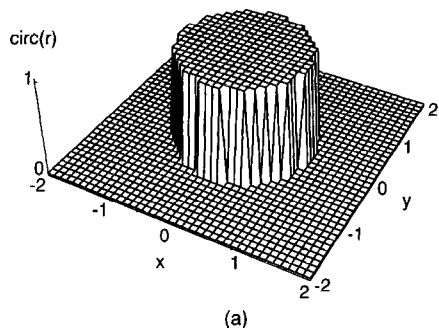


FIGURE 2.3
 (a) The circle function and (b) its transform.

Accordingly, some effort is now devoted to finding the transform of this function. Using the Fourier-Bessel transform expression (2-31), the transform of the circle function can be written

$$\mathcal{B}\{\text{circ}(r)\} = 2\pi \int_0^1 r J_0(2\pi r \rho) dr.$$

Using a change of variables $r' = 2\pi r \rho$ and the identity

$$\int_0^x \xi J_0(\xi) d\xi = x J_1(x),$$

we rewrite the transform as

$$\mathcal{B}\{\text{circ}(r)\} = \frac{1}{2\pi\rho^2} \int_0^{2\pi\rho} r' J_0(r') dr' = \frac{J_1(2\pi\rho)}{\rho} \quad (2-35)$$

where J_1 is a Bessel function of the first kind, order 1. Figure 2.3 illustrates the circle function and its transform. Note that the transform is circularly symmetric, as expected, and consists of a central lobe and a series of concentric rings of diminishing amplitude. Its value at the origin is π . As a matter of curiosity we note that the zeros of this transform are not equally spaced in radius. A convenient normalized version of this function, with value unity at the origin, is $2 \frac{J_1(2\pi\rho)}{2\pi\rho}$. This particular function is called the "besinc" function, or the "jinc" function.

For a number of additional Fourier-Bessel transform pairs, the reader is referred to the problems (see Prob. 2-6).

2.2 LOCAL SPATIAL FREQUENCY AND SPACE-FREQUENCY LOCALIZATION

Each Fourier component of a function is a complex exponential of a unique spatial frequency. As such, every frequency component extends over the entire (x, y) domain. Therefore it is not possible to associate a spatial location with a particular spatial frequency. Nonetheless, we know that in practice certain portions of an image could contain parallel grid lines at a certain fixed spacing, and we are tempted to say that the particular frequency or frequencies represented by these grid lines are localized to certain spatial regions of the image. In this section we introduce the idea of local spatial frequencies and their relation to Fourier components.

For the purpose of this discussion, we consider the general case of complex-valued functions, which we will later see represent the amplitude and phase distributions of monochromatic optical waves. For now, they are just complex functions. Any such function can be represented in the form

$$g(x, y) = a(x, y) \exp[j\phi(x, y)] \quad (2-36)$$

where $a(x, y)$ is a real and nonnegative amplitude distribution, while $\phi(x, y)$ is a real phase distribution. For this discussion we assume that the amplitude distribution $a(x, y)$

is a slowly varying function of (x, y) , so that we can concentrate on the behavior of the phase function $\phi(x, y)$.

We define the local spatial frequency of the function g as a frequency pair (f_{IX}, f_{IY}) given by

$$f_{IX} = \frac{1}{2\pi} \frac{d}{dx} \phi(x, y) \quad f_{IY} = \frac{1}{2\pi} \frac{\partial}{\partial y} \phi(x, y). \quad (2-37)$$

In addition, both f_{IX} and f_{IY} are defined to be zero in regions where the function $g(x, y)$ vanishes.

Consider the result of applying these definitions to the particular complex function

$$g(x, y) = \exp[j2\pi(f_X x + f_Y y)]$$

representing a simple linear-phase exponential of frequencies (f_X, f_Y) . We obtain

$$f_{IX} = \frac{1}{2\pi} \frac{\partial}{\partial x} [2\pi(f_X x + f_Y y)] = f_X \quad f_{IY} = \frac{1}{2\pi} \frac{\partial}{\partial y} [2\pi(f_X x + f_Y y)] = f_Y.$$

Thus we see that for the case of a single Fourier component, the local frequencies do indeed reduce to the frequencies of that component, and those frequencies are constant over the entire (x, y) plane.

Next consider a space-limited version of a quadratic-phase exponential function: which we call a "finite chirp" function,⁵

$$g(x, y) = \exp[j\pi\beta(x^2 + y^2)] \text{rect}\left(\frac{x}{2L_X}\right) \text{rect}\left(\frac{y}{2L_Y}\right). \quad (2-38)$$

Performing the differentiations called for by the definitions of local frequencies, we find that they can be expressed as

$$f_{IX} = \beta x \text{rect}\left(\frac{x}{2L_X}\right) \quad f_{IY} = \beta y \text{rect}\left(\frac{y}{2L_Y}\right). \quad (2-39)$$

We see that in this case the local spatial frequencies do depend on location in the (x, y) plane; within a rectangle of dimensions $2L_X \times 2L_Y$, f_{IX} varies linearly with the x -coordinate while f_{IY} varies linearly with the y -coordinate. Thus for this function (and for most others) there is a dependence of local spatial frequency on position in the (x, y) plane.⁶

Since the local spatial frequencies are bounded to covering a rectangle of dimensions $2L_X \times 2L_Y$, it would be tempting to conclude that the Fourier spectrum of $g(x, y)$ is also limited to the same rectangular region. In fact this is approximately true, but not exactly so. The Fourier transform of this function is given by the expression

⁴For a tutorial discussion of the importance of quadratic-phase functions in various fields of optics, see [229].

⁵The name "chirp function", without the finite length qualifier, will be used for the infinite-length quadratic phase exponential, $\exp[j\pi\beta(x^2 + y^2)]$.

⁶From the definition (2-37) the dimensions of f_{IX} and f_{IY} are both *cycles per meter*; in spite of what might appear to be a contrary implication of Eq. (2-39). The dimensions of β are meters^{-2} .

$$G(f_X, f_Y) = \int_{-L_X}^{L_X} \int_{-L_Y}^{L_Y} e^{j\pi\beta(x^2+y^2)} e^{-j2\pi(f_X x + f_Y y)} dx dy.$$

This expression is separable in rectangular coordinates, so it suffices to find the one-dimensional spectrum

$$G_X(f_X) = \int_{-L_X}^{L_X} e^{j\pi\beta x^2} e^{j2\pi f_X x} dx.$$

Completing the square in the exponent and making a change of variables of integration from x to $t = \sqrt{2\beta}(\mathbf{x} - \frac{f_X}{\beta})$ yields

$$G_X(f_X) = \frac{1}{\sqrt{2\beta}} e^{-j\pi \frac{f_X^2}{\beta}} \int_{-\sqrt{2\beta}(L_X + \frac{f_X}{\beta})}^{\sqrt{2\beta}(L_X - \frac{f_X}{\beta})} \exp\left[j \frac{\pi t^2}{2}\right] dt.$$

This integral can be expressed in terms of tabulated functions, the Fresnel integrals, which are defined by

$$C(z) = \int_0^z \cos\left(\frac{\pi t^2}{2}\right) dt \quad S(z) = \int_0^z \sin\left(\frac{\pi t^2}{2}\right) dt. \quad (2-40)$$

The spectrum G_X can then be expressed as

$$G_X(f_X) = \frac{e^{-j\pi \frac{f_X^2}{\beta}}}{\sqrt{2\beta}} \left\{ C\left[\sqrt{2\beta}\left(L_X - \frac{f_X}{\beta}\right)\right] - C\left[\sqrt{2\beta}\left(-L_X - \frac{f_X}{\beta}\right)\right] \right. \\ \left. + jS\left[\sqrt{2\beta}\left(L_X - \frac{f_X}{\beta}\right)\right] - jS\left[\sqrt{2\beta}\left(-L_X - \frac{f_X}{\beta}\right)\right] \right\}.$$

The expression for G_Y is of course identical, except the Y subscript replaces the X subscript. Figure 2.4 shows a plot of $|G_X(f_X)|$ vs. f_X for the particular case of $L_X = 10$ and $\beta = 1$. As can be seen, the spectrum is almost flat over the region $(-L_X, L_X)$ and

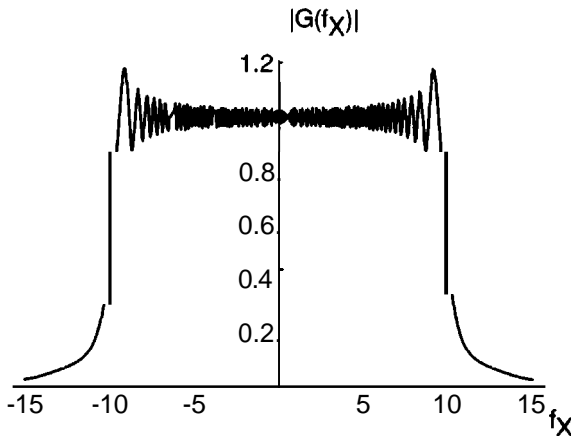


FIGURE 2.4

The spectrum of the finite chirp function, $L_X = 10$, $\beta = 1$.

almost zero outside that region. We conclude that local spatial frequency has provided a good (but not exact) indication of where the significant values of the Fourier spectrum will occur. However, local spatial frequencies are not the same entity as the frequency components of the Fourier spectrum. Examples can be found for which the local spatial frequency distribution and the Fourier spectrum are not in as good agreement as found in the above example. Good agreement can be expected only when the variations of $\phi(x, y)$ are sufficiently "slow" in the (x, y) plane to allow $\phi(x, y)$ to be well approximated by only three terms of its Taylor series expansion about any point (x, y) , i.e. a constant term and two first-partial-derivative terms.

Local spatial frequencies are of special physical significance in optics. When the local spatial frequencies of the complex amplitude of a coherent optical wavefront are found, they correspond to the ray directions of the geometrical optics description of that wavefront. However, we are getting ahead of ourselves; we will return to this idea in later chapters and particularly in Appendix B.

2.3 LINEAR SYSTEMS

For the purposes of discussion here, we seek to define the word system in a way sufficiently general to include both the familiar case of electrical networks and the less-familiar case of optical imaging systems. Accordingly, a system is defined to be a mapping of a set of input functions into a set of output functions. For the case of electrical networks, the inputs and outputs are real-valued functions (voltages or currents) of a one-dimensional independent variable (time); for the case of imaging systems, the inputs and outputs can be real-valued functions (intensity) or complex-valued functions (field amplitude) of a two-dimensional independent variable (space). As mentioned previously, the question of whether intensity or field amplitude should be considered the relevant quantity will be treated at a later time.

If attention is restricted to deterministic (nonrandom) systems, then a specified input must map to a unique output. It is not necessary, however, that each output correspond to a unique input, for as we shall see, a variety of input functions can produce no output. Thus we restrict attention at the outset to systems characterized by many-to-one mappings.

A convenient representation of a system is a mathematical operator, $\mathcal{S}\{\}$, which we imagine to operate on input functions to produce output functions. Thus if the function $g_1(x_1, y_1)$ represents the input to a system, and $g_2(x_2, y_2)$ represents the corresponding output, then by the definition of $\mathcal{S}\{\}$, the two functions are related through

$$g_2(x_2, y_2) = \mathcal{S}\{g_1(x_1, y_1)\}. \quad (2-41)$$

Without specifying more detailed properties of the operator $\mathcal{S}\{\}$, it is difficult to state more specific properties of the general system than those expressed by Eq. (2-41). In the material that follows, we shall be concerned primarily, though not exclusively, with a restricted class of systems that are said to be linear: The assumption of linearity will be found to yield simple and physically meaningful representations of such systems; it will also allow useful relations between inputs and outputs to be developed.

2.3.1 Linearity and the Superposition Integral

A system is said to be linear if the following superposition property is obeyed for all input functions p and q and all complex constants a and b :

$$\mathcal{S}\{ap(x_1, y_1) + bq(x_1, y_1)\} = a\mathcal{S}\{p(x_1, y_1)\} + b\mathcal{S}\{q(x_1, y_1)\}. \quad (2-42)$$

As mentioned previously, the great advantage afforded by linearity is the ability to express the response of a system to an arbitrary input in terms of the responses to certain "elementary" functions into which the input has been decomposed. It is most important, then, to find a simple and convenient means of decomposing the input. Such a decomposition is offered by the so-called sifting property of the δ function (cf. Section 1 of Appendix A), which states that

$$g_1(x_1, y_1) = \iint_{-\infty}^{\infty} g_1(\xi, \eta) \delta(x_1 - \xi, y_1 - \eta) d\xi d\eta. \quad (2-43)$$

This equation may be regarded as expressing g_1 as a linear combination of weighted and displaced δ functions; the elementary functions of the decomposition are, of course, just these δ functions.

To find the response of the system to the input g_1 , substitute (2-43) in (2-41):

$$g_2(x_2, y_2) = \mathcal{S} \left\{ \iint_{-\infty}^{\infty} g_1(\xi, \eta) \delta(x_1 - \xi, y_1 - \eta) d\xi d\eta \right\}. \quad (2-44)$$

Now, regarding the number $g_1(\xi, \eta)$ as simply a weighting factor applied to the elementary function $\delta(x_1 - \xi, y_1 - \eta)$, the linearity property (2-42) is invoked to allow $\mathcal{S}\{\}$ to operate on the individual elementary functions; thus the operator $\mathcal{S}\{\}$ is brought within the integral, yielding

$$g_2(x_2, y_2) = \iint_{-\infty}^{\infty} g_1(\xi, \eta) \mathcal{S}\{\delta(x_1 - \xi, y_1 - \eta)\} d\xi d\eta. \quad (2-45)$$

As a final step we let the symbol $h(x_2, y_2; \xi, \eta)$ denote the response of the system at point (x_2, y_2) of the output space to a δ function input at coordinates (ξ, η) of the input space; that is,

$$h(x_2, y_2; \xi, \eta) = \mathcal{S}\{\delta(x_1 - \xi, y_1 - \eta)\}. \quad (2-46)$$

The function h is called the impulse response (or in optics, the point-spreadfunction) of the system. The system input and output can now be related by the simple equation

$$g_2(x_2, y_2) = \iint_{-\infty}^{\infty} g_1(\xi, \eta) h(x_2, y_2; \xi, \eta) d\xi d\eta. \quad (2-47)$$

This fundamental expression, known as the superposition integral, demonstrates the very important fact that a linear system is completely characterized by its responses

to unit impulses. To completely specify the output, the responses must in general be known for impulses located at all possible points in the input plane. For the case of a linear imaging system, this result has the interesting physical interpretation that the effects of imaging elements (lenses, stops, etc.) can be fully described by specifying the (possibly complex-valued) images of point sources located throughout the object field.

2.3.2 Invariant Linear Systems: Transfer Functions

Having examined the input-output relations for a general linear system, we turn now to an important subclass of linear systems, namely invariant linear systems. An electrical network is said to be time-invariant if its impulse response $h(t; \tau)$ (that is, its response at time t to a unit impulse excitation applied at time τ) depends only on the time difference $(t - \tau)$. Electrical networks composed of fixed resistors, capacitors, and inductors are time-invariant since their characteristics do not change with time.

In a similar fashion, a linear imaging system is space-invariant (or equivalently, isoplanatic) if its impulse response $h(x_2, y_2; \xi, \eta)$ depends only on the distances $(x_2 - \xi)$ and $(y_2 - \eta)$ (i.e. the x and y distances between the excitation point and the response point). For such a system we can, of course, write

$$h(x_2, y_2; \xi, \eta) = h(x_2 - \xi, y_2 - \eta). \quad (2-48)$$

Thus an imaging system is space-invariant if the image of a point source object changes only in location, not in functional form, as the point source explores the object field. In practice, imaging systems are seldom isoplanatic over their entire object field, but it is usually possible to divide that field into small regions (isoplanatic patches), within which the system is approximately invariant. To completely describe the imaging system, the impulse response appropriate for each isoplanatic patch should be specified; but if the particular portion of the object field of interest is sufficiently small, it often suffices to consider only the isoplanatic patch on the optical axis of the system. Note that for an invariant system the superposition integral (2-47) takes on the particularly simple form

$$g_2(x_2, y_2) = \iint_{-\infty}^{\infty} g_1(\xi, \eta) h(x_2 - \xi, y_2 - \eta) d\xi d\eta \quad (2-49)$$

which we recognize as a two-dimensional convolution of the object function with the impulse response of the system. In the future it will be convenient to have a shorthand notation for a convolution relation such as (2-49), and accordingly this equation is written symbolically as

$$g_2 = g_1 \otimes h$$

where a \otimes symbol between any two functions indicates that those functions are to be convolved.

The class of invariant linear systems has associated with it a far more detailed mathematical structure than the more general class of all linear systems, and it is precisely

because of this structure that invariant systems are so easily dealt with. The simplicity of invariant systems begins to be evident when we note that the convolution relation (2-49) takes a particularly simple form after Fourier transformation. Specifically, transforming both sides of (2-49) and invoking the convolution theorem, the spectra $G_2(f_X, f_Y)$ and $G_1(f_X, f_Y)$ of the system output and input are seen to be related by the simple equation

$$G_2(f_X, f_Y) = H(f_X, f_Y) G_1(f_X, f_Y), \quad (2-50)$$

where H is the Fourier transform of the impulse response

$$H(f_X, f_Y) = \iint_{-\infty}^{\infty} h(\xi, \eta) \exp[-j2\pi(f_X\xi + f_Y\eta)] d\xi d\eta. \quad (2-51)$$

The function H , called the transfer function of the system, indicates the effects of the system in the "frequency domain". Note that the relatively tedious convolution operation of (2-49) required to find the system output is replaced in (2-50) by the often more simple sequence of Fourier transformation, multiplication of transforms, and inverse Fourier transformation.

From another point of view, we may regard the relations (2-50) and (2-51) as indicating that, for a linear invariant system, the input can be decomposed into elementary functions that are more convenient than the δ functions of Eq. (2-43). These alternative elementary functions are, of course, the complex-exponential functions of the Fourier integral representation. By transforming g_1 we are simply decomposing the input into complex-exponential functions of various spatial frequencies (f_X, f_Y) . Multiplication of the input spectrum G_1 by the transfer function H then takes into account the effects of the system on each elementary function. Note that these effects are limited to an amplitude change and a phase shift, as evidenced by the fact that we simply multiply the input spectrum by a complex number $H(f_X, f_Y)$ at each (f_X, f_Y) . Inverse transformation of the output spectrum G_2 synthesizes the output g_2 by adding up the modified elementary functions.

The mathematical term eigenfunction is used for a function that retains its original form (up to a multiplicative complex constant) after passage through a system. Thus we see that the complex-exponential functions are the eigenfunctions of linear, invariant systems. The weighting applied by the system to an eigenfunction input is called the eigenvalue corresponding to that input. Hence the transfer function describes the continuum of eigenvalues of the system.

Finally, it should be strongly emphasized that the simplifications afforded by transfer-function theory are only applicable for invariant linear systems. For applications of Fourier theory in the analysis of time-varying electrical networks, the reader may consult Ref. [158]; applications of Fourier analysis to space-variant imaging systems can be found in Ref. [199].

24 TWO-DIMENSIONAL SAMPLING THEORY

It is often convenient, both for data processing and for mathematical analysis purposes, to represent a function $g(x, y)$ by an array of its sampled values taken on a

discrete set of points in the (\mathbf{x}, y) plane. Intuitively, it is clear that if these samples are taken sufficiently close to each other, the sampled data are an accurate representation of the original function, in the sense that g can be reconstructed with considerable accuracy by simple interpolation. It is a less obvious fact that for a particular class of functions (known as bandlimited functions) the reconstruction can be accomplished exactly, provided only that the interval between samples is not greater than a certain limit. This result was originally pointed out by Whittaker [298] and was later popularized by Shannon [259] in his studies of information theory.

The sampling theorem applies to the class of bandlimited functions, by which we mean functions with Fourier transforms that are nonzero over only a finite region \mathcal{R} of the frequency space. We consider first a form of this theorem that is directly analogous to the one-dimensional theorem used by Shannon. Later we very briefly indicate improvements of the theorem that can be made in some two-dimensional cases.

2.4.1 The Whittaker-Shannon Sampling Theorem

To derive what is perhaps the simplest version of the sampling theorem, we consider a rectangular lattice of samples of the function g , as defined by

$$g_s(x, y) = \text{comb}\left(\frac{x}{X}\right) \text{comb}\left(\frac{y}{Y}\right) g(x, y). \quad (2-52)$$

The sampled function g_s thus consists of an array of δ functions, spaced at intervals of width X in the x direction and width Y in the y direction, as illustrated in Fig. 2.5. The area under each δ function is proportional to the value of the function g at that particular point in the rectangular sampling lattice. As implied by the convolution theorem, the spectrum G_s of g_s can be found by convolving the transform of $\text{comb}(x/X) \text{comb}(y/Y)$ with the transform of g , or

$$G_s(f_x, f_y) = \mathcal{F}\left\{\text{comb}\left(\frac{x}{X}\right) \text{comb}\left(\frac{y}{Y}\right)\right\} \otimes G(f_x, f_y)$$

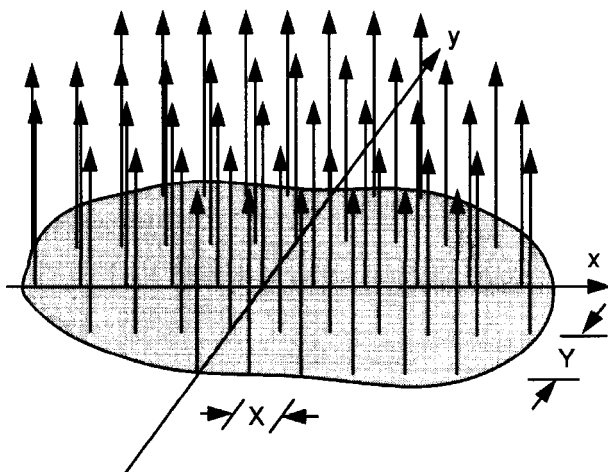


FIGURE 2.5
The sampled function.

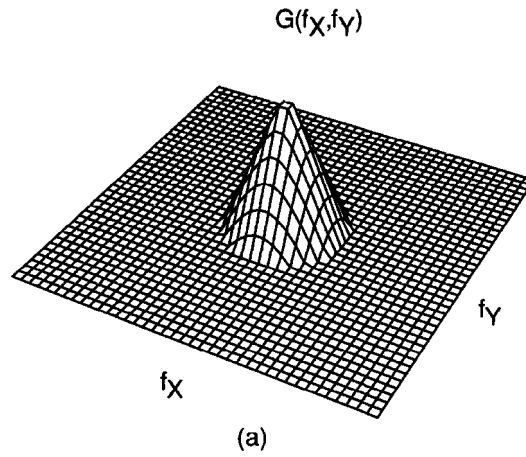


FIGURE 2.6a
Spectrum of the original function.

where the \otimes again indicates that a two-dimensional convolution is to be performed. Now using Table 2.1 we have

$$\mathcal{F} \left\{ \text{comb} \left(\frac{x}{X} \right) \text{comb} \left(\frac{y}{Y} \right) \right\} = XY \text{comb}(X f_x) \text{comb}(Y f_y)$$

while from the results of Prob. 2-1b,

$$XY \text{comb}(X f_x) \text{comb}(Y f_y) = \sum_{n=-\infty}^{\infty} \sum_{m=-\infty}^{\infty} \delta \left(f_x - \frac{n}{X}, f_y - \frac{m}{Y} \right).$$

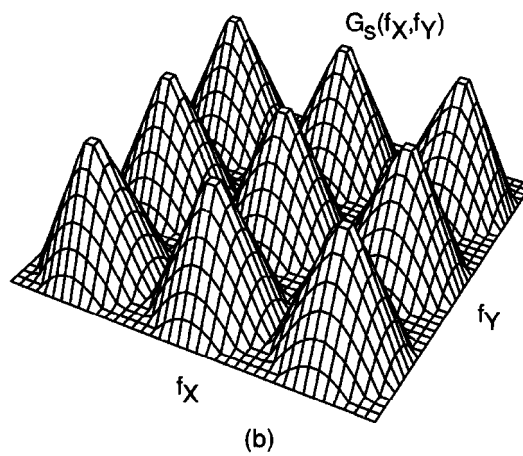


FIGURE 2.6b
Spectrum of the sampled data (only three periods are shown in each direction for this infinitely periodic function).

It follows that

$$G_s(f_x, f_y) = \sum_{n=-\infty}^{\infty} \sum_{m=-\infty}^{\infty} G\left(f_x - \frac{n}{X}, f_y - \frac{m}{Y}\right). \quad (2-53)$$

Evidently the spectrum of g_s can be found simply by erecting the spectrum of g about each point $(n/X, m/Y)$ in the (f_x, f_y) plane as shown in Fig. 2.6b.

Since the function g is assumed to be bandlimited, its spectrum G is nonzero over only a finite region R of the frequency space. As implied by Eq. (2-53), the region over which the spectrum of the sampled function is nonzero can be found by constructing the region R about each point $(n/X, m/Y)$ in the frequency plane. Now it becomes clear that if X and Y are sufficiently small (i.e. the samples are sufficiently close together), then the separations $1/X$ and $1/Y$ of the various spectral islands will be great enough to assure that the adjacent regions do not overlap (see Fig. 2.6b). Thus the recovery of the original spectrum G from G_s can be accomplished exactly by passing the sampled function g_s through a linear invariant filter that transmits the term $(n = 0, m = 0)$ of Eq. (2-53) without distortion, while perfectly excluding all other terms. Thus, at the output of this filter we find an exact replica of the original data $g(x, y)$.

As stated in the above discussion, to successfully recover the original data it is necessary to take samples close enough together to enable separation of the various spectral regions of G_s . To determine the maximum allowable separation between samples, let $2B_X$ and $2B_Y$ represent the widths in the f_x and f_y directions, respectively, of the *smallest* rectangle⁷ that completely encloses the region R . Since the various terms in the spectrum (2-53) of the sampled data are separated by distances $1/X$ and $1/Y$ in the f_x and f_y directions, respectively, separation of the spectral regions is assured if

$$X \geq \frac{1}{2B_X} \quad \text{and} \quad Y \geq \frac{1}{2B_Y}. \quad (2-54)$$

The maximum spacings of the sampling lattice for exact recovery of the original function are thus $(2B_X)^{-1}$ and $(2B_Y)^{-1}$.

Having determined the maximum allowable distances between samples, it remains to specify the exact transfer function of the filter through which the data should be passed. In many cases there is considerable latitude of choice here, since for many possible shapes of the region R there are a multitude of transfer functions that will pass the $(n = 0, m = 0)$ term of G_s and exclude all other terms. For our purposes, however, it suffices to note that if the relations (2-54) are satisfied, there is one transfer function that will always yield the desired result regardless of the shape of R , namely

$$H(f_x, f_y) = \text{rect}\left(\frac{f_x}{2B_X}\right) \text{rect}\left(\frac{f_y}{2B_Y}\right). \quad (2-55)$$

The exact recovery of G from G_s is seen by noting that the spectrum of the output of such a filter is

⁷For simplicity we assume that this rectangle is centered on the origin. If this is not the case, the arguments can be modified in a straightforward manner to yield a somewhat more efficient sampling theorem.

26 Introduction to Fourier Optics

$$G_s(f_x, f_y) \operatorname{rect}\left(\frac{f_x}{2B_x}\right) \operatorname{rect}\left(\frac{f_y}{2B_y}\right) = G(f_x, f_y).$$

The equivalent identity in the space domain is

$$\left[\operatorname{comb}\left(\frac{x}{X}\right) \operatorname{comb}\left(\frac{y}{Y}\right) g(x, y) \right] \otimes h(x, y) = g(x, y) \quad (2-56)$$

where h is the impulse response of the filter,

$$h(x, y) = \mathcal{F}^{-1} \left\{ \operatorname{rect}\left(\frac{f_x}{2B_x}\right) \operatorname{rect}\left(\frac{f_y}{2B_y}\right) \right\} = 4B_x B_y \operatorname{sinc}(2B_x x) \operatorname{sinc}(2B_y y).$$

Noting that

$$\operatorname{comb}\left(\frac{x}{X}\right) \operatorname{comb}\left(\frac{y}{Y}\right) g(x, y) = XY \sum_{n=-\infty}^{\infty} \sum_{m=-\infty}^{\infty} g(nX, mY) \delta(x - nX, y - mY),$$

Eq. (2-56) becomes

$$g(x, y) = 4B_x B_y XY \sum_{n=-\infty}^{\infty} \sum_{m=-\infty}^{\infty} g(nX, mY) \operatorname{sinc}[2B_x(x - nX)] \operatorname{sinc}[2B_y(y - mY)].$$

Finally, when the sampling intervals X and Y are taken to have their maximum allowable values, the identity becomes

$$g(x, y) = \sum_{n=-\infty}^{\infty} \sum_{m=-\infty}^{\infty} g\left(\frac{n}{2B_x}, \frac{m}{2B_y}\right) \operatorname{sinc}\left[2B_x\left(x - \frac{n}{2B_x}\right)\right] \operatorname{sinc}\left[2B_y\left(y - \frac{m}{2B_y}\right)\right]. \quad (2-57)$$

Equation (2-57) represents a fundamental result which we shall refer to as the Whittaker-Shannon sampling theorem. It implies that exact recovery of a bandlimited function can be achieved from an appropriately spaced rectangular array of its sampled values; the recovery is accomplished by injecting, at each sampling point, an interpolation function consisting of a product of **sinc** functions, where each interpolation function is weighted according to the sampled value of g at the corresponding point.

The above result is by no means the only possible sampling theorem. Two rather arbitrary choices were made in the analysis, and alternative choices at these two points will yield alternative sampling theorems. The first arbitrary choice, appearing early in the analysis, was the use of a rectangular sampling lattice. The second, somewhat later in the analysis, was the choice of the particular filter transfer function (2-55). Alternative theorems derived by making different choices at these two points are no less valid than Eq. (2-57); in fact, in some cases alternative theorems are more "efficient" in the sense that fewer samples per unit area are required to assure complete recovery. The reader interested in pursuing this extra richness of multidimensional sampling theory is referred to the works of Bracewell [31] and of Peterson and Middleton [230]. A more modern treatment of multidimensional sampling theory is found in Dudgeon and Mersereau [85].

2.4.2 Space-Bandwidth Product

It is possible to show that no function that is bandlimited can be perfectly space-limited as well. That is, if the spectrum G of a function g is nonzero over only a limited region \mathcal{R} in the (f_X, f_Y) plane, then it is not possible for g to be nonzero over only a finite region in the (x, y) plane simultaneously. Nonetheless, in practice most functions do eventually fall to very small values, and therefore from a practical point-of-view it is usually possible to say that g has *significant* values only in some finite region. Exceptions are functions that do not have Fourier transforms in the usual sense, and have to be dealt with in terms of generalized Fourier transforms (e.g. $g(x, y) = 1$, $g(x, y) = \cos[2\pi(f_X x + f_Y y)]$, etc.).

If $g(x, y)$ is bandlimited and indeed has significant value over only a finite region of the (x, y) plane, then it is possible to represent g with good accuracy by a *finite* number of samples. If g is of significant value only in the region $-L_X \leq x < L_X$, $-L_Y \leq y < L_Y$, and if g is sampled, in accord with the sampling theorem, on a rectangular lattice with spacings $(2B_X)^{-1}$, $(2B_Y)^{-1}$ in the x and y directions, respectively, then the total number of significant samples required to represent $g(x, y)$ is seen to be

$$M = 16L_X L_Y B_X B_Y, \quad (2-58)$$

which we call the space-bandwidth product of the function g . The space-bandwidth product can be regarded as the number of degrees of freedom of the given function.

The concept of space-bandwidth product is also useful for many functions that are not strictly bandlimited. If the function is approximately space-limited and approximately bandlimited, then a rectangle (size $2B_X \times 2B_Y$) within which most of the spectrum is contained can be defined in the frequency domain, and a rectangle (size $2L_X \times 2L_Y$) within which most of the function is contained can be defined in the space domain. The space-bandwidth product of the function is then approximately given by Eq. (2-58).

The space-bandwidth product of a function is a measure of its complexity. The ability of an optical system to accurately handle inputs and outputs having large space-bandwidth products is a measure of performance, and is directly related to the quality of the system.

PROBLEMS-CHAPTER 2

2-1. Prove the following properties of δ functions:

(a) $\delta(ax, by) = \frac{1}{|ab|} \delta(x, y)$.

(b) $\text{comb}(ax)\text{comb}(by) = \frac{1}{|ab|} \sum_{n=-\infty}^m \sum_{m=-\infty}^m \delta\left(x - \frac{n}{a}, y - \frac{m}{b}\right)$.

2-2. Prove the following Fourier transform relations:

(a) $\mathcal{F}\{\text{rect}(x)\text{rect}(y)\} = \text{sinc}(f_X)\text{sinc}(f_Y)$.

$$(b) \mathcal{F}\{\Lambda(x)\Lambda(y)\} = \text{sinc}^2(f_X) \text{sinc}^2(f_Y).$$

Prove the following generalized Fourier transform relations:

$$(c) \mathcal{F}\{1\} = \delta(f_X, f_Y).$$

$$(d) \mathcal{F}\{\text{sgn}(x)\text{sgn}(y)\} = \left(\frac{1}{j\pi f_X}\right)\left(\frac{1}{j\pi f_Y}\right).$$

2-3. Prove the following Fourier transform theorems:

$$(a) \mathcal{F}\mathcal{F}\{g(x, y)\} = \mathcal{F}^{-1}\mathcal{F}^{-1}\{g(x, y)\} = g(-x, -y) \text{ at all points of continuity of } g.$$

$$(b) \mathcal{F}\{g(x, y)h(x, y)\} = \mathcal{F}\{g(x, y)\} \otimes \mathcal{F}\{h(x, y)\}.$$

$$(c) \mathcal{F}\{\nabla^2 g(x, y)\} = -4\pi^2(f_X^2 + f_Y^2)\mathcal{F}\{g(x, y)\} \text{ where } \nabla^2 \text{ is the Laplacian operator}$$

$$\nabla^2 = \frac{\partial^2}{\partial x^2} + \frac{\partial^2}{\partial y^2}.$$

2-4. Let the transform operators $\mathcal{F}_A\{\}$ and $\mathcal{F}_B\{\}$ be defined by

$$\mathcal{F}_A\{g\} = \frac{1}{a} \iint_{-\infty}^{\infty} g(\xi, \eta) \exp\left[-j\frac{2\pi}{a}(f_X\xi + f_Y\eta)\right] d\xi d\eta$$

$$\mathcal{F}_B\{g\} = \frac{1}{b} \iint_{-\infty}^{\infty} g(\xi, \eta) \exp\left[-j\frac{2\pi}{b}(x\xi + y\eta)\right] d\xi d\eta$$

(a) Find a simple interpretation for

$$\mathcal{F}_B\{\mathcal{F}_A\{g(x, y)\}\}.$$

(b) Interpret the result for $a > b$ and $a < b$.

2-5. The "equivalent area" Δ_{XY} of a function $g(x, y)$ can be defined by

$$\Delta_{XY} = \frac{\iint_{-\infty}^{\infty} g(x, y) dx dy}{g(0, 0)},$$

while the "equivalent bandwidth" $\Delta_{f_X f_Y}$ of g is defined in terms of its transform G by

$$\Delta_{f_X f_Y} = \frac{\iint_{-\infty}^{\infty} G(f_X, f_Y) df_X df_Y}{G(0, 0)}$$

Show that $\Delta_{XY} \Delta_{f_X f_Y} = 1$.

2-6. Prove the following Fourier-Bessel transform relations:

(a) If $g_R(r) = \delta(r - r_0)$, then

$$\mathcal{B}\{g_R(r)\} = 2\pi r_0 J_0(2\pi r_0 \rho).$$

(b) If $g_R(r) = 1$ for $a \leq r \leq 1$ and zero otherwise, then

$$\mathcal{B}\{g_R(r)\} = \frac{J_1(2\pi\rho) - aJ_1(2\pi a\rho)}{\rho}$$

(c) If $\mathcal{B}\{g_R(r)\} = G(\rho)$, then

$$\mathcal{B}\{g_R(ar)\} = \frac{1}{a^2} G\left(\frac{\rho}{a}\right).$$

(d) $\mathcal{B}\{\exp(-\pi r^2)\} = \exp(-\pi \rho^2)$.

2-7. Let $g(r, \theta)$ be separable in polar coordinates.

(a) Show that if $g(r, \theta) = g_R(r)e^{jm\theta}$, then

$$\mathcal{F}\{g(r, \theta)\} = (-j)^m e^{jm\phi} \mathcal{H}_m\{g_R(r)\}$$

where $\mathcal{H}_m\{\}$ is the **Hankel** transform of order m ,

$$\mathcal{H}_m\{g_R(r)\} = 2\pi \int_0^\infty r g_R(r) J_m(2\pi r \rho) dr$$

and (ρ, ϕ) are polar coordinates in the frequency space. (Hint: $\exp(ja \sin x) = \sum_{k=-\infty}^{\infty} J_k(a) \exp(jkx)$)

(b) With the help of part (a), prove the general relation presented in Eq. (2-22) for functions separable in polar coordinates.

2-8. Suppose that a sinusoidal input

$$g(x, y) = \cos[2\pi(f_X x + f_Y y)]$$

is applied to a linear system. Under what (sufficient) conditions is the output a real sinusoidal function of the same spatial frequency as the input? Express the amplitude and phase of that output in terms of an appropriate characteristic of the system.

2-9. Show that the zero-order Bessel function $J_0(2\pi\rho_0 r)$ is an eigenfunction of any invariant linear system with a circularly symmetric impulse response. What is the corresponding eigenvalue?

2-10. The Fourier transform operator may be regarded as a mapping of functions into their transforms and therefore satisfies the definition of a system as presented in this chapter.

(a) Is this system linear?

(b) Can you specify a **transfer function** for this system? If yes, what is it? If no, why not?

2-11. The expression

$$p(x, y) = g(x, y) \otimes \left[\text{comb}\left(\frac{x}{X}\right) \text{comb}\left(\frac{y}{Y}\right) \right]$$

defines a periodic function, with period X in the x direction and period Y in the y direction.

(a) Show that the Fourier transform of p can be written

$$P(f_x, f_y) = \sum_{n=-\infty}^{\infty} \sum_{m=-\infty}^{\infty} G\left(\frac{n}{X}, \frac{m}{Y}\right) \delta\left(f_x - \frac{n}{X}, f_y - \frac{m}{Y}\right)$$

where G is the Fourier transform of g .

(b) Sketch the function $p(x, y)$ when

$$g(x, y) = \text{rect}\left(2\frac{x}{X}\right) \text{rect}\left(2\frac{y}{Y}\right)$$

and find the corresponding Fourier transform $P(f_x, f_y)$.

2-12. Show that a function with no nonzero spectral components outside a circle of radius B in the frequency plane obeys the following sampling theorem:

$$g(x, y) = \sum_{n=-\infty}^{\infty} \sum_{m=-\infty}^{\infty} g\left(\frac{n}{2B}, \frac{m}{2B}\right) \frac{\pi}{4} \left\{ 2 \frac{J_1 \left[2\pi B \sqrt{\left(x - \frac{n}{2B}\right)^2 + \left(y - \frac{m}{2B}\right)^2} \right]}{2\pi B \sqrt{\left(x - \frac{n}{2B}\right)^2 + \left(y - \frac{m}{2B}\right)^2}} \right\}.$$

2-13. The input to a certain imaging system is an object complex field distribution $U_o(x, y)$ of unlimited spatial frequency content, while the output of the system is an image field distribution $U_i(x, y)$. The imaging system can be assumed to act as a linear, invariant **lowpass** filter with a transfer function that is identically zero outside the region $|f_x| \leq B_x$, $|f_y| \leq B_y$ in the frequency domain. Show that there exists an "equivalent" object $U'_o(x, y)$ consisting of a rectangular array of point sources that produces exactly the same image U_i as does the true object U_o , and that the field distribution across the equivalent object can be written

$$U'_o(x, y) = \sum_{n=-\infty}^{\infty} \sum_{m=-\infty}^{\infty} \left[\iint_{-\infty}^{\infty} U_o(\xi, \eta) \text{sinc}(n - 2B_x \xi) \text{sinc}(m - 2B_y \eta) d\xi d\eta \right] \times \delta\left(x - \frac{n}{2B_x}, y - \frac{m}{2B_y}\right).$$

2-14. The Wigner distribution function of a one-dimensional function $g(x)$ is defined by

$$W(f, x) = \int_{-\infty}^{\infty} g(\xi + x/2) g^*(\xi - x/2) \exp(-j2\pi f \xi) d\xi$$

and is a description of the simultaneous (one-dimensional) space and spatial-frequency occupancy of a signal.

(a) Find the Wigner distribution function of the infinite-length chirp function by inserting $g(x) = \exp(j\pi\beta x^2)$ in the definition of $W(\mathbf{f}, x)$.

(b) Show that the Wigner distribution function for the one-dimensional finite chirp

$$g(x) = \exp(j\pi\beta x^2) \text{rect}\left(\frac{x}{2L}\right)$$

is given by

$$W(f, x) = (2L - |x|) \text{sinc}[(2L - |x|)(\beta x - f)]$$

for $|x| < 2L$ and zero otherwise.

- (c) If you have access to a computer and appropriate software, plot the Wigner distribution function of the finite-length chirp for $L = 10$ and $\beta = 1$, with x ranging from -10 to 10 and f ranging from -10 to 10 . To make the nature of this function clearer, also plot $W(0, x)$ for $|x| \leq 1$.

Foundations of Scalar Diffraction Theory

The phenomenon known as *diffraction* plays a role of utmost importance in the branches of physics and engineering that deal with wave propagation. In this chapter we consider some of the foundations of scalar diffraction theory. While the theory discussed here is sufficiently general to be applied in other fields, such as acoustic-wave and radio-wave propagation, the applications of primary concern will be in the realm of physical optics. To fully understand the properties of optical imaging and data processing systems, it is essential that diffraction and the limitations it imposes on system performance be appreciated. A variety of references to more comprehensive treatments of diffraction theory will be found in the material that follows.

3.1 HISTORICAL INTRODUCTION

Before beginning a discussion of diffraction, it is first necessary to mention another phenomenon with which diffraction should not be confused — namely refraction. Refraction can be defined as the bending of light rays that takes place when they pass through a region in which there is a gradient of the local velocity of propagation of the wave. The most common example occurs when a light wave encounters a sharp boundary between two regions having different refractive indices. The propagation velocity in the first medium, having refractive index n_1 , is $v_1 = c/n_1$, c being the vacuum velocity of light. The velocity of propagation in the second medium is $v_2 = c/n_2$.

As shown in Fig. 3.1, the incident light rays are bent at the interface. The angles of incidence and refraction are related by *Snell's law*, which is the foundation of geometrical optics,

$$n_1 \sin \theta_1 = n_2 \sin \theta_2, \quad (3-1)$$

where in this example, $n_2 > n_1$ and therefore $\theta_2 < \theta_1$.

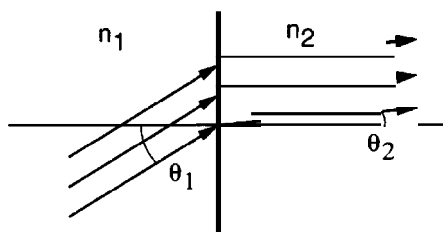


FIGURE 3.1
Snell's law at a sharp boundary.

Light rays are also bent upon reflection, which can occur at a metallic or dielectric interface. The fundamental relation governing this phenomenon is that the angle of reflection is always equal to the angle of incidence.

The term diffraction has been defined by Sommerfeld (Ref. [270]) as "any deviation of light rays from rectilinear paths which cannot be interpreted as reflection or refraction." Diffraction is caused by the confinement of the lateral extent of a wave, and is most appreciable when that confinement is to sizes comparable with a wavelength of the radiation being used. The diffraction phenomenon should also not be confused with the penumbra *effect*, for which the finite extent of a source causes the light transmitted by a small aperture to spread as it propagates away from that aperture (see Fig. 3.2). As can be seen in the figure, the penumbra effect does not involve any bending of the light rays.

There is a fascinating history associated with the discovery and explanation of diffraction effects. The first accurate report and description of such a phenomenon was made by Grimaldi and was published in the year 1665, shortly after his death. The measurements reported were made with an experimental apparatus similar to that shown in Fig. 3.3. An aperture in an opaque screen was illuminated by a light source, chosen small enough to introduce a negligible penumbra effect; the light intensity was observed across a plane some distance behind the screen. The corpuscular theory of light propagation, which was the accepted means of explaining optical phenomena at the time, predicted that the shadow behind the screen should be well defined, with sharp borders. Grimaldi's observations indicated, however, that the transition from light to shadow was gradual rather than abrupt. If the spectral purity of the light source had been better, he might have observed even more striking results, such as the presence of light and dark fringes extending far into the geometrical shadow of the screen. Such effects cannot be explained by a corpuscular theory of light, which requires rectilinear propagation of light rays in the absence of reflection and refraction.

The initial step in the evolution of a theory that would explain such effects was made by the first proponent of the wave theory of light, Christian Huygens, in the year

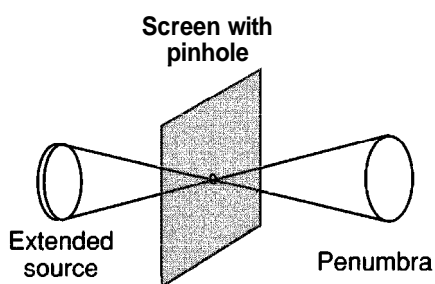


FIGURE 3.2
The penumbra effect.

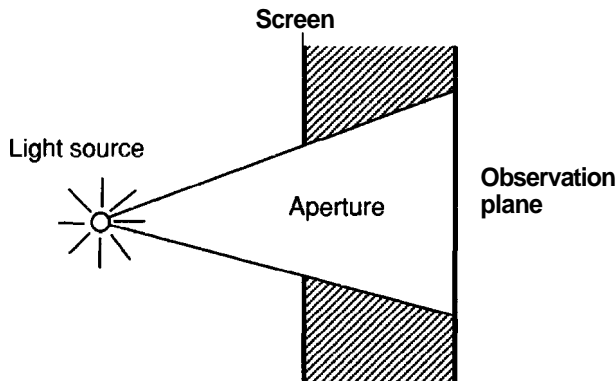


FIGURE 33
Arrangement used for observing
diffraction of light.

1678. Huygens expressed the intuitive conviction that if each point on the wavefront of a disturbance were considered to be a new source of a "secondary" spherical disturbance, then the wavefront at a later instant could be found by constructing the "envelope" of the secondary wavelets, as illustrated in Fig. 3.4.

Progress on further understanding diffraction was impeded throughout the entire 18th century by the fact that Isaac Newton, a scientist with an enormous reputation for his many contributions to physics in general and to optics in particular, favored the corpuscular theory of light as early as 1704. His followers supported this view adamantly. It was not until 1804 that further significant progress occurred. In that year, Thomas Young, an English physician, strengthened the wave theory of light by introducing the critical concept of *interference*. The idea was a radical one at the time, for it stated that under proper conditions, light could be added to light and produce darkness.

The ideas of Huygens and Young were brought together in 1818 in the famous memoir of **Augustin** Jean Fresnel. By making some rather arbitrary assumptions about the amplitudes and phases of Huygens' secondary sources, and by allowing the various wavelets to mutually interfere, Fresnel was able to calculate the distribution of light in diffraction patterns with excellent accuracy.

At Fresnel's presentation of his paper to a prize committee of the French Academy of Sciences, his theory was strongly disputed by the great French mathematician S. Poisson, a member of the committee. He demonstrated the absurdity of the theory

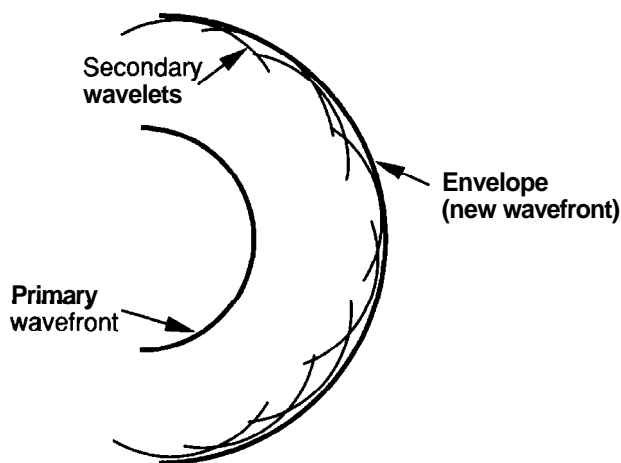


FIGURE 34
Huygens' envelope construction.

by showing that it predicted the existence of a bright spot at the center of the shadow of an opaque disk. F. Arago, who chaired the prize committee, performed such an experiment and found the predicted spot. Fresnel won the prize, and since then the effect has been known as "Poisson's spot".

In 1860 Maxwell identified light as an electromagnetic wave, a step of enormous importance. But it was not until 1882 that the ideas of Huygens and Fresnel were put on a firmer mathematical foundation by Gustav Kirchhoff, who succeeded in showing that the amplitudes and phases ascribed to the secondary sources by Fresnel were indeed logical consequences of the wave nature of light. Kirchhoff based his mathematical formulation upon two assumptions about the boundary values of the light incident on the surface of an obstacle placed in the way of propagation of light. These assumptions were later proved to be inconsistent with each other, by Poincaré in 1892 and by Sommerfeld in 1894.¹ As a consequence of these criticisms, Kirchhoff's formulation of the so-called Huygens-Fresnel principle must be regarded as a first approximation, although under most conditions it yields results that agree amazingly well with experiment. Kottler [174] attempted to resolve the contradictions by reinterpreting Kirchhoff's boundary value problem as a *saltus* problem, where *saltus* is a Latin word signifying a discontinuity or jump. The Kirchhoff theory was also modified by Sommerfeld, who eliminated one of the aforementioned assumptions concerning the light amplitude at the boundary by making use of the theory of Green's functions. This so-called *Rayleigh-Sommerfeld diffraction* theory will be treated in Section 3.5.

It should be emphasized from the start that the Kirchhoff and **Rayleigh-Sommerfeld** theories share certain major simplifications and approximations. Most important, light is treated as a scalar phenomenon, neglecting the fundamentally vectorial nature of the electromagnetic fields. Such an approach neglects the fact that, at boundaries, the various components of the electric and magnetic fields are coupled through Maxwell's equations and cannot be treated independently. Fortunately, experiments in the microwave region of the spectrum [262] have shown that the scalar theory yields very accurate results if two conditions are met: (1) the diffracting aperture must be large compared with a wavelength, and (2) the diffracting fields must not be observed too close to the aperture. These conditions will be well satisfied in the problems treated here. For a more complete discussion of the applicability of scalar theory in instrumental optics the reader may consult Ref. [28] (Section 8.4). Nonetheless, there do exist important problems for which the required conditions are not satisfied, for example in the theory of diffraction from high-resolution gratings and from extremely small pits on optical recording media. Such problems are excluded from consideration here, since the vectorial nature of the fields must be taken into account if reasonably accurate results are to be obtained. Vectorial generalizations of diffraction theory do exist, the first satisfactory treatment being due to Kottler [172].

The first truly rigorous solution of a diffraction problem was given in 1896 by Sommerfeld [268], who treated the two-dimensional case of a plane wave incident on an infinitesimally thin, perfectly conducting half plane. Kottler [173] later compared Sommerfeld's solution with the corresponding results of Kirchhoff's scalar treatment.

¹For a more detailed discussion of these inconsistencies, see Section 3.5.

Needless to say, an historic introduction to a subject so widely mentioned in the literature can hardly be considered complete. The reader is therefore referred to more comprehensive treatments of diffraction theory, for example Refs. [13], [29], and [145].

3.2 FROM A VECTOR TO A SCALAR THEORY

The most fundamental beginning for our analysis is Maxwell's equations. In MKS units and in the absence of free charge, the equations are given by

$$\begin{aligned}\nabla \times \vec{\mathcal{E}} &= -\mu \frac{\partial \vec{\mathcal{H}}}{\partial t} \\ \nabla \times \vec{\mathcal{H}} &= \epsilon \frac{\partial \vec{\mathcal{E}}}{\partial t} \\ \nabla \cdot \epsilon \vec{\mathcal{E}} &= 0 \\ \nabla \cdot \mu \vec{\mathcal{H}} &= 0.\end{aligned}\tag{3-2}$$

Here $\vec{\mathcal{E}}$ is the electric field, with rectilinear components $(\mathcal{E}_X, \mathcal{E}_Y, \mathcal{E}_Z)$, and $\vec{\mathcal{H}}$ is the magnetic field, with components $(\mathcal{H}_X, \mathcal{H}_Y, \mathcal{H}_Z)$. μ and ϵ are the permeability and permittivity, respectively, of the medium in which the wave is propagating. $\vec{\mathcal{E}}$ and $\vec{\mathcal{H}}$ are functions of both position \mathbf{P} and time t . The symbols \times and \cdot represent a vector cross product and a vector dot product, respectively, while $\nabla = \frac{\partial}{\partial x} \hat{i} + \frac{\partial}{\partial y} \hat{j} + \frac{\partial}{\partial z} \hat{k}$, where \hat{i} , \hat{j} and \hat{k} are unit vectors in the x , y , and z directions, respectively.

We assume that the wave is propagating in a dielectric medium. It is important to further specify some properties of that medium. The medium is linear if it satisfies the linearity properties discussed in Chapter 2. The medium is isotropic if its properties are independent of the direction of polarization of the wave (*i.e.* the directions of the $\vec{\mathcal{E}}$ and $\vec{\mathcal{H}}$ vectors). The medium is homogeneous if the permittivity is constant throughout the region of propagation. The medium is nondispersive if the permittivity is independent of wavelength over the wavelength region occupied by the propagating wave. Finally, all media of interest in this book are nonmagnetic, which means that the magnetic permeability is always equal to μ_0 , the vacuum permeability.

Applying the $\nabla \times$ operation to the left and right sides of the first equation for $\vec{\mathcal{E}}$, we make use of the vector identity

$$\nabla \times (\nabla \times \vec{\mathcal{E}}) = \nabla(\nabla \cdot \vec{\mathcal{E}}) - \nabla^2 \vec{\mathcal{E}}.\tag{3-3}$$

If the propagation medium is linear, isotropic, homogeneous (constant ϵ), and **nondispersive**, substitution of the two Maxwell's equations for $\vec{\mathcal{E}}$ in Eq. (3-3) yields

$$\nabla^2 \vec{\mathcal{E}} - \frac{n^2}{c^2} \frac{\partial^2 \vec{\mathcal{E}}}{\partial t^2} = 0\tag{3-4}$$

where n is the refractive index of the medium, defined by

$$n = \left(\frac{\epsilon}{\epsilon_0} \right)^{1/2},\tag{3-5}$$

ϵ_0 is the vacuum permittivity, and c is the velocity of propagation in vacuum, given by

$$c = \frac{1}{\sqrt{\mu_0 \epsilon_0}}. \quad (3-6)$$

The magnetic field satisfies an identical equation,

$$\nabla^2 \vec{\mathcal{H}} - \frac{n^2}{c^2} \frac{\partial^2 \vec{\mathcal{H}}}{\partial t^2} = 0.$$

Since the vector wave equation is obeyed by both $\vec{\mathcal{E}}$ and $\vec{\mathcal{H}}$, an identical scalar wave equation is obeyed by all components of those vectors. Thus, for example, \mathcal{E}_X obeys the equation

$$\nabla^2 \mathcal{E}_X - \frac{n^2}{c^2} \frac{\partial^2 \mathcal{E}_X}{\partial t^2} = 0,$$

and similarly for \mathcal{E}_Y , \mathcal{E}_Z , \mathcal{H}_X , \mathcal{H}_Y , and \mathcal{H}_Z . Therefore it is possible to summarize the behavior of all components of $\vec{\mathcal{E}}$ and $\vec{\mathcal{H}}$ through a single scalar wave equation,

$$\nabla^2 u(\mathbf{P}, t) - \frac{n^2}{c^2} \frac{\partial^2 u(\mathbf{P}, t)}{\partial t^2} = 0, \quad (3-7)$$

where $u(\mathbf{P}, t)$ represents any of the scalar field components, and we have explicitly introduced the dependence of u on both position \mathbf{P} in space and time t .

From above we conclude that in a dielectric medium that is linear, isotropic, homogeneous, and nondispersive, all components of the electric and magnetic field behave identically and their behavior is fully described by a single scalar wave equation. How, then, is the scalar theory only an approximation, rather than exact? The answer becomes clear if we consider situations other than propagation in the uniform dielectric medium hypothesized.

For example, if the medium is inhomogeneous with a permittivity $\epsilon(\mathbf{P})$ that depends on position \mathbf{P} (but not on time t), it is a simple matter to show (see Prob. 3-1) that the wave equation satisfied by $\vec{\mathcal{E}}$ becomes

$$\nabla^2 \vec{\mathcal{E}} + 2\nabla(\vec{\mathcal{E}} \cdot \nabla \ln n) - \frac{n^2}{c^2} \frac{\partial^2 \vec{\mathcal{E}}}{\partial t^2} = 0, \quad (3-8)$$

where n and c are again given by Eqs. (3-5) and (3-6). The new term that has been added to the wave equation will be nonzero for a refractive index that changes over space. More importantly, that term introduces a coupling between the various components of the electric field, with the result that \mathcal{E}_X , \mathcal{E}_Y , and \mathcal{E}_Z may no longer satisfy the same wave equation. This type of coupling is important, for example, when light propagates through a “thick” dielectric diffraction grating.

A similar effect takes place when boundary conditions are imposed on a wave that propagates in a homogeneous medium. At the boundaries, coupling is introduced between $\vec{\mathcal{E}}$ and $\vec{\mathcal{H}}$ as well as between their various scalar components. As a consequence, even when the propagation medium is homogeneous, the use of a scalar theory entails some degree of error. That error will be small provided the boundary conditions have effect over an area that is a small part of the area through which a wave may be passing.

In the case of diffraction of light by an aperture, the $\vec{\mathcal{E}}$ and $\vec{\mathcal{H}}$ fields are modified only at the edges of the aperture where light interacts with the material of which the edges are composed, and the effects extend over only a few wavelengths into the aperture itself. Thus if the aperture has an area that is large compared with a wavelength, the coupling effects of the boundary conditions on the $\vec{\mathcal{E}}$ and $\vec{\mathcal{H}}$ fields will be small. As will be seen, this is equivalent to the requirement that the diffraction angles caused by the aperture are small.

With these discussions as background, we turn away from the vector theory of diffraction to the simpler scalar theory. We close with one final observation. Circuit theory is based on the approximation that circuit elements (resistors, capacitors, and inductors) are small compared to the wavelength of the fields that appear within them, and for this reason can be treated as lumped elements with simple properties. We need not use Maxwell's equations to analyze such elements under these conditions. In a similar vein, the scalar theory of diffraction introduces substantial simplifications compared with a full vectorial theory. The scalar theory is accurate provided that the diffracting structures are large compared with the wavelength of light. Thus the approximation implicit in the scalar theory should be no more disturbing than the approximation used in lumped circuit theory. In both cases it is possible to find situations in which the approximation breaks down, but as long as the simpler theories are used only in cases for which they are expected to be valid, the losses of accuracy will be small and the gain of simplicity will be large.

3.3 SOME MATHEMATICAL PRELIMINARIES

Before embarking on a treatment of diffraction itself, we first consider a number of mathematical preliminaries that form the basis of the later diffraction-theory derivations. These initial discussions will also serve to introduce some of the notation used throughout the book.

3.3.1 The Helmholtz Equation

In accord with the previous introduction of the scalar theory, let the light disturbance at position \mathbf{P} and time t be represented by the scalar function $u(\mathbf{P}, t)$. Attention is now restricted to the case of a purely monochromatic wave, with the generalization to polychromatic waves being deferred to Section 3.8.

For a monochromatic wave, the scalar field may be written explicitly as

$$u(\mathbf{P}, t) = A(\mathbf{P}) \cos[2\pi\nu t + \phi(\mathbf{P})] \quad (3-9)$$

where $A(\mathbf{P})$ and $\phi(\mathbf{P})$ are the amplitude and phase, respectively, of the wave at position \mathbf{P} , while ν is the optical frequency. A more compact form of (3-9) is found by using complex notation, writing

$$u(\mathbf{P}, t) = \text{Re}\{U(\mathbf{P}) \exp(-j2\pi\nu t)\}, \quad (3-10)$$

where $\text{Re}\{\}$ signifies "real part of", and $U(P)$ is a complex function of position (sometimes called a phasor),

$$U(P) = A(P) \exp[-j\phi(P)]. \quad (3-11)$$

If the real disturbance $u(P, t)$ is to represent an optical wave, it must satisfy the scalar wave equation

$$\nabla^2 u - \frac{n^2}{c^2} \frac{\partial^2 u}{\partial t^2} = 0 \quad (3-12)$$

at each source-free point. As before, ∇^2 is the Laplacian operator, n represents the refractive index of the dielectric medium within which light is propagating, and c represents the vacuum velocity of light. The complex function $U(P)$ serves as an adequate description of the disturbance, since the time dependence is known a priori. If (3-10) is substituted in (3-12), it follows that U must obey the time-independent equation

$$(\nabla^2 + k^2)U = 0. \quad (3-13)$$

Here k is termed the wave number and is given by

$$k = 2\pi n \frac{\nu}{c} = \frac{2\pi}{\lambda},$$

and λ is the wavelength in the dielectric medium ($\lambda = c/n\nu$). The relation (3-13) is known as the Helmholtz equation; we may assume in the future that the complex amplitude of any monochromatic optical disturbance propagating in vacuum ($n = 1$) or in a homogeneous dielectric medium ($n > 1$) must obey such a relation.

3.3.2 Green's Theorem

Calculation of the complex disturbance U at an observation point in space can be accomplished with the help of the mathematical relation known as Green's theorem. This theorem, which can be found in most texts on advanced calculus, can be stated as follows:

Let $U(P)$ and $G(P)$ be any two complex-valued functions of position, and let S be a closed surface surrounding a volume V . If U , G , and their first and second partial derivatives are single-valued and continuous within and on S , then we have

$$\iiint_V (U\nabla^2 G - G\nabla^2 U) dv = \iint_S \left(U \frac{\partial G}{\partial n} - G \frac{\partial U}{\partial n} \right) ds \quad (3-14)$$

where $\frac{\partial}{\partial n}$ signifies a partial derivative in the outward normal direction at each point on S .

This theorem is in many respects the prime foundation of scalar diffraction theory. However, only a prudent choice of an auxiliary function G and a closed surface S will allow its direct application to the diffraction problem. We turn now to the former of these problems, considering Kirchhoff's choice of an auxiliary function and the consequent integral theorem that follows.

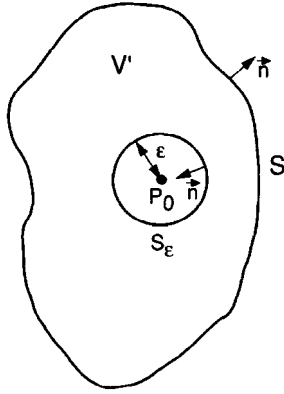


FIGURE 35
Surface of integration.

3.3.3 The Integral Theorem of Helmholtz and Kirchhoff

The Kirchhoff formulation of the diffraction problem is based on a certain integral theorem which expresses the solution of the homogeneous wave equation at an arbitrary point in terms of the values of the solution and its first derivative on an arbitrary closed surface surrounding that point. This theorem had been derived previously in acoustics by H. von Helmholtz.

Let the point of observation be denoted P_0 , and let S denote an arbitrary closed surface surrounding P_0 , as indicated in Fig. 3.5. The problem is to express the optical disturbance at P_0 in terms of its values on the surface S . To solve this problem, we follow Kirchhoff in applying Green's theorem and in choosing as an auxiliary function a unit-amplitude spherical wave expanding about the point P_0 (the so-called *free space* Green's function). Thus the value of Kirchhoff's G at an arbitrary point P_1 is given by²

$$G(P_1) = \frac{\exp(jkr_{01})}{r_{01}}, \quad (3-15)$$

where we adopt the notation that r_{01} is the length of the vector \vec{r}_{01} pointing from P_0 to P_1 .

Before proceeding further, a short diversion regarding Green's functions may be in order. Suppose that we wish to solve an inhomogeneous linear differential equation of the form

$$a_2(x) \frac{d^2 U}{dx^2} + a_1(x) \frac{dU}{dx} + a_0(x)U = V(x) \quad (3-16)$$

where $V(x)$ is a driving function and $U(x)$ satisfies a known set of boundary conditions. We have chosen a one-dimensional variable x but the theory is easily generalized to a multidimensional \vec{x} . It can be shown (see Chapter 1 of [223] and [16]) that if $G(x)$ is the solution to the same differential equation (3-16) when $V(x)$ is replaced by the

²The reader may wish to verify that, for our choice of clockwise rotation of phasors, the description of an expanding wave should have a $+$ sign in the exponential.

impulsive driving function $\delta(\mathbf{x} - \mathbf{x}')$ and with the same boundary conditions applying, then the general solution $U(\mathbf{x})$ can be expressed in terms of the specific solution $G(\mathbf{x})$ through a convolution integral

$$U(\mathbf{x}) = \int G(\mathbf{x} - \mathbf{x}') V(\mathbf{x}') d\mathbf{x}'. \quad (3-17)$$

The function $G(\mathbf{x})$ is known as the Green's *function* of the problem, and is clearly a form of impulse response. Various solutions to the scalar diffraction problem to be discussed in the following sections correspond to results obtained under different assumptions about the Green's function of the problem. The function G appearing in Green's theorem may be regarded either as simply an auxiliary function which we cleverly choose to solve our problem, or it may eventually be related to the Green's function of the problem. Further consideration of the theory of Green's functions is beyond the scope of this treatment.

Returning now to our central discussion, to be legitimately used in Green's theorem, the function G (as well as its first and second partial derivatives) must be continuous within the enclosed volume V . Therefore to exclude the discontinuity at P_0 , a small spherical surface S_ϵ , of radius ϵ , is inserted about the point P_0 . Green's theorem is then applied, the volume of integration V' being that volume lying between S and S_ϵ , and the surface of integration being the composite surface

$$S' = S + S_\epsilon$$

as indicated in Fig. 3.5. Note that the "outward normal to the composite surface points outward in the conventional sense on S , but inward (towards P_0) on S_ϵ .

Within the volume V' , the disturbance G , being simply an expanding spherical wave, satisfies the Helmholtz equation

$$(\nabla^2 + k^2)G = 0. \quad (3-18)$$

Substituting the two Helmholtz equations (3-13) and (3-18) in the left-hand side of Green's theorem, we find

$$\iiint_{V'} (U\nabla^2 G - G\nabla^2 U) dv = - \iiint_{V'} (UGk^2 - GUk^2) dv \equiv 0.$$

Thus the theorem reduces to

$$\begin{aligned} \iint_{S'} \left(U \frac{\partial G}{\partial n} - G \frac{\partial U}{\partial n} \right) ds &= 0 \\ - \iint_{S_\epsilon} \left(U \frac{\partial G}{\partial n} - G \frac{\partial U}{\partial n} \right) ds &= \iint_S \left(U \frac{\partial G}{\partial n} - G \frac{\partial U}{\partial n} \right) ds. \end{aligned} \quad (3-19)$$

Note that, for a general point P_1 on S' , we have

$$G(P_1) = \frac{\exp(jkr_{01})}{r_{01}}$$

and

$$\frac{\partial G(P_1)}{\partial n} = \cos(\vec{n}, \vec{r}_{01}) \left(jk - \frac{1}{r_{01}} \right) \frac{\exp(jkr_{01})}{r_{01}} \quad (3-20)$$

where $\cos(\vec{n}, \vec{r}_{01})$ represents the cosine of the angle between the outward normal \vec{n} and the vector \vec{r}_{01} joining P_0 to P_1 . For the particular case of P_1 on S_ϵ , $\cos(\vec{n}, \vec{r}_{01}) = -1$, and these equations become

$$G(P_1) = \frac{e^{jk\epsilon}}{\epsilon} \quad \text{and} \quad \frac{\partial G(P_1)}{\partial n} = \frac{e^{jk\epsilon}}{\epsilon} \left(\frac{1}{\epsilon} - jk \right).$$

Letting ϵ become arbitrarily small, the continuity of U (and its derivatives) at P_0 allows us to write

$$\begin{aligned} & \lim_{\epsilon \rightarrow 0} \iint_{S_\epsilon} \left(U \frac{\partial G}{\partial n} - G \frac{\partial U}{\partial n} \right) ds \\ &= \lim_{\epsilon \rightarrow 0} 4\pi\epsilon^2 \left[U(P_0) \frac{\exp(jk\epsilon)}{\epsilon} \left(\frac{1}{\epsilon} - jk \right) - \frac{\partial U(P_0) \exp(jk\epsilon)}{\partial n \epsilon} \right] = 4\pi U(P_0). \end{aligned}$$

Substitution of this result in (3-19) (taking account of the negative sign) yields

$$U(P_0) = \frac{1}{4\pi} \iint_S \left\{ \frac{\partial U}{\partial n} \left[\frac{\exp(jkr_{01})}{r_{01}} \right] - U \frac{\partial}{\partial n} \left[\frac{\exp(jkr_{01})}{r_{01}} \right] \right\} ds. \quad (3-21)$$

This result is known as the integral theorem of Helmholtz and *Kirchhoff*; it plays an important role in the development of the scalar theory of diffraction, for it allows the field at any point P_0 to be expressed in terms of the "boundary values" of the wave on any closed surface surrounding that point. As we shall now see, such a relation is instrumental in the further development of scalar diffraction equations.

3.4

THE KIRCHHOFF FORMULATION OF DIFFRACTION BY A PLANAR SCREEN

Consider now the problem of diffraction of light by an aperture in an infinite opaque screen. As illustrated in Fig. 3.6, a wave disturbance is assumed to impinge on the screen and the aperture from the left, and the field at the point P_0 behind the aperture is to be calculated. Again the field is assumed to be monochromatic.

3.4.1 Application of the Integral Theorem

To find the field at the point P_0 , we apply the integral theorem of Helmholtz and *Kirchhoff*, being careful to choose a surface of integration that will allow the calculation to be performed successfully. Following *Kirchhoff*, the closed surface S is chosen to consist of two parts, as shown in Fig. 3.6. Let a plane surface, S_1 , lying directly behind the diffracting screen, be joined and closed by a large spherical cap, S_2 , of radius R and

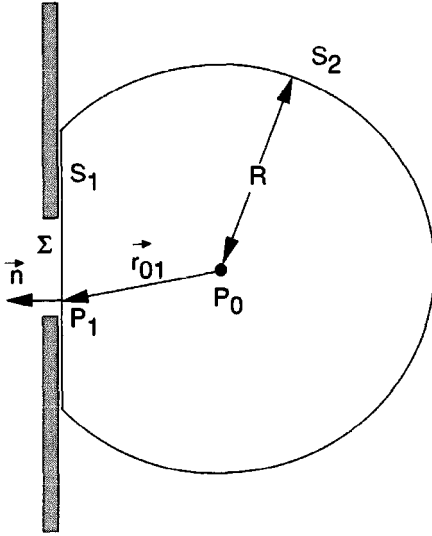


FIGURE 3.6
Kirchhoff formulation of diffraction by a plane screen.

centered at the observation point P_0 . The total closed surface S is simply the sum of S_1 and S_2 . Thus, applying (3-21),

$$U(P_0) = \frac{1}{4\pi} \iint_{S_1+S_2} \left(G \frac{\partial U}{\partial n} - U \frac{\partial G}{\partial n} \right) ds$$

where, as before,

$$G = \frac{\exp(jkr_{01})}{r_{01}}$$

As R increases, S_2 approaches a large hemispherical shell. It is tempting to reason that, since both U and G will fall off as $1/R$, the integrand will ultimately vanish, yielding a contribution of zero from the surface integral over S_2 . However, the area of integration increases as R^2 , so this argument is incomplete. It is also tempting to assume that, since the disturbances are propagating with finite velocity c/n , R will ultimately be so large that the waves have not yet reached S_2 , and the integrand will be zero on that surface. But this argument is incompatible with our assumption of monochromatic disturbances, which must (by definition) have existed for all time. Evidently a more careful investigation is required before the contribution from S_2 can be disposed of.

Examining this problem in more detail, we see that, on S_2 ,

$$G = \frac{\exp(jkR)}{R}$$

and, from (3-20),

$$\frac{\partial G}{\partial n} = \left(jk - \frac{1}{R} \right) \frac{\exp(jkR)}{R} \approx jkG$$

where the last approximation is valid for large R . The integral in question can thus be reduced to

$$\iint_{\Omega} \left[G \frac{\partial U}{\partial n} - U(jkG) \right] ds = \int_{\Omega} G \left(\frac{\partial U}{\partial n} - jkU \right) R^2 d\omega,$$

where Ω is the solid angle subtended by S_2 at P_0 . Now the quantity $|RG|$ is uniformly bounded on S_2 . Therefore the entire integral over S_2 will vanish as R becomes arbitrarily large, provided the disturbance has the property

$$\lim_{R \rightarrow \infty} R \left(\frac{\partial U}{\partial n} - jkU \right) = 0 \quad (3-22)$$

uniformly in angle. This requirement is known as the *Sommerfeld* radiation condition [269] and is satisfied if the disturbance U vanishes at least as fast as a diverging spherical wave (see Prob. 3-2). It guarantees that we are dealing only with outgoing waves on S_2 , rather than incoming waves, for which the integral over S_2 might not vanish as $R \rightarrow \infty$. Since only outgoing waves will fall on S_2 in our problem, the integral over S_2 will yield a contribution of precisely zero.

3.4.2 The Kirchhoff Boundary Conditions

Having disposed of the integration over the surface S_2 , it is now possible to express the disturbance at P_0 in terms of the disturbance and its normal derivative over the infinite plane S_1 immediately behind the screen, that is,

$$U(P_0) = \frac{1}{4\pi} \iint_{S_1} \left(\frac{\partial U}{\partial n} G - U \frac{\partial G}{\partial n} \right) ds. \quad (3-23)$$

The screen is opaque, except for the open aperture which will be denoted Σ . It therefore seems intuitively reasonable that the major contribution to the integral (3-23) arises from the points of S_1 located within the aperture Σ , where we would expect the integrand to be largest. Kirchhoff accordingly adopted the following assumptions [162]:

1. Across the surface Σ , the field distribution U and its derivative $\partial U / \partial n$ are exactly the same as they would be in the absence of the screen.
2. Over the portion of S_1 that lies in the geometrical shadow of the screen, the field distribution U and its derivative $\partial U / \partial n$ are identically zero.

These conditions are commonly known as the *Kirchhoff boundary* conditions. The first allows us to specify the disturbance incident on the aperture by neglecting the presence of the screen. The second allows us to neglect all of the surface of integration except that portion lying directly within the aperture itself. Thus (3-23) is reduced to

$$U(P_0) = \frac{1}{4\pi} \iint_{\Sigma} \left(\frac{\partial U}{\partial n} G - U \frac{\partial G}{\partial n} \right) ds. \quad (3-24)$$

While the Kirchhoff boundary conditions simplify the results considerably, it is important to realize that neither can be exactly true. The presence of the screen will inevitably perturb the fields on Σ to some degree, for along the rim of the aperture certain boundary conditions must be met that would not be required in the absence of the screen. In addition, the shadow behind the screen is never perfect, for fields will inevitably extend behind the screen for a distance of several wavelengths. However, if the dimensions of the aperture are large compared with a wavelength, these fringing

effects can be safely neglected,³ and the two boundary conditions can be used to yield results that agree very well with experiment.

3.4.3 The Fresnel-Kirchhoff Diffraction Formula

A further simplification of the expression for $U(P_0)$ is obtained by noting that the distance r_{01} from the aperture to the observation point is usually many optical wavelengths, and therefore, since $k \gg 1/r_{01}$, Eq. (3-20) becomes

$$\begin{aligned} \frac{\partial G(P_1)}{\partial n} &= \cos(\vec{n}, \vec{r}_{01}) \left(jk - \frac{1}{r_{01}} \right) \frac{\exp(jkr_{01})}{r_{01}} \\ &\approx jk \cos(\vec{n}, \vec{r}_{01}) \frac{\exp(jkr_{01})}{r_{01}}. \end{aligned} \quad (3-25)$$

Substituting this approximation and the expression (3-15) for G in Eq. (3-24), we find

$$U(P_0) = \frac{1}{4\pi} \iint_{\Sigma} \frac{\exp(jkr_{01})}{r_{01}} \left[\frac{\partial U}{\partial n} - jkU \cos(\vec{n}, \vec{r}_{01}) \right] ds. \quad (3-26)$$

Now suppose that the aperture is illuminated by a single spherical wave,

$$U(P_1) = \frac{A \exp(jkr_{21})}{r_{21}}$$

arising from a point source at P_2 , a distance r_{21} from P_1 (see Fig. 3.7). If r_{21} is many optical wavelengths, then (3-26) can be directly reduced (see Prob. 3-3) to

$$U(P_0) = \frac{A}{j\lambda} \iint_{\Sigma} \frac{\exp[jk(r_{21} + r_{01})]}{r_{21}r_{01}} \left[\frac{\cos(\vec{n}, \vec{r}_{01}) - \cos(\vec{n}, \vec{r}_{21})}{2} \right] ds. \quad (3-27)$$

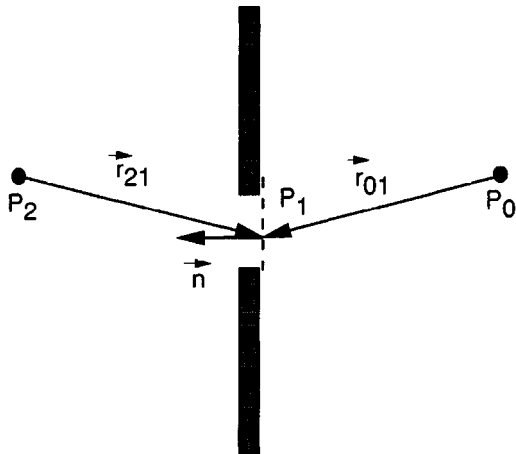


FIGURE 3.7
Point-source illumination of a plane screen.

³As we shall see, objections to the use of the Kirchhoff boundary conditions arise, not because of the fringing effects, but rather because of certain internal inconsistencies.

This result, which holds only for an illumination consisting of a single point source, is known as the *Fresnel-Kirchhoff diffraction formula*.

Note that Eq. (3-27) is symmetrical with respect to the illumination point source at P_2 and the observation point at P_0 . Thus a point source at P_0 will produce at P_2 the same effect that a point source of equal intensity placed at P_2 will produce at P_0 . This result is referred to as the *reciprocity theorem of Helmholtz*.

Finally, we point out an interesting interpretation of the diffraction formula (3-27), to which we will return later for a more detailed discussion. Let that equation be rewritten as follows:

$$U(P_0) = \iint_{\Sigma} U'(P_1) \frac{\exp(jkr_{01})}{r_{01}} ds \quad (3-28)$$

where

$$U'(P_1) = \frac{1}{j\lambda} \left[\frac{A \exp(jkr_{21})}{r_{21}} \right] \left[\frac{\cos(\vec{n}, \vec{r}_{01}) - \cos(\vec{n}, \vec{r}_{21})}{2} \right]. \quad (3-29)$$

Now (3-28) may be interpreted as implying that the field at P_0 arises from an infinity of fictitious "secondary" point sources located within the aperture itself. The secondary sources have certain amplitudes and phases, described by $U'(P_1)$, that are related to the illuminating wavefront and the angles of illumination and observation. Assumptions resembling these were made by Fresnel rather arbitrarily in his combination of Huygens' envelope construction and Young's principle of interference. Fresnel *assumed* these properties to hold in order to obtain accurate results. Kirchhoff showed that such properties are a natural consequence of the wave nature of light.

Note that the above derivation has been restricted to the case of an aperture illumination consisting of a single expanding spherical wave. However, as we shall now see, such a limitation can be removed by the Rayleigh-Sommerfeld theory.

3.5

THE RAYLEIGH-SOMMERFELD FORMULATION OF DIFFRACTION

The Kirchhoff theory has been found experimentally to yield remarkably accurate results and is widely used in practice. However, there are certain internal inconsistencies in the theory which motivated a search for a more satisfactory mathematical development. The difficulties of the Kirchhoff theory stem from the fact that boundary conditions must be imposed on *both* the field strength and its normal derivative. In particular, it is a well-known theorem of potential theory that if a two-dimensional potential function and its normal derivative vanish *together* along any finite curve segment, then that potential function *must vanish over the entire plane*. Similarly, if a solution of the three-dimensional wave equation vanishes on any finite surface element, it must vanish in all space. Thus the two Kirchhoff boundary conditions together imply that the field is zero everywhere behind the aperture, a result which contradicts the known physical situation. A further indication of these inconsistencies is the fact that the Fresnel-Kirchhoff diffraction formula can be shown to fail to reproduce the assumed boundary conditions as the observation point approaches the screen or aperture. In view of these

contradictions, it is indeed remarkable that the Kirchhoff theory has been found to yield such accurate results in practice.⁴

The inconsistencies of the Kirchhoff theory were removed by Sommerfeld, who eliminated the necessity of imposing boundary values on both the disturbance and its normal derivative simultaneously. This so-called Rayleigh-Sommerfeld theory is the subject of this section.

3.5.1 Choice of Alternative Green's Functions

Consider again Eq. (3-23) for the observed field strength in terms of the incident field and its normal derivative across the entire screen:

$$U(P_0) = \frac{1}{4\pi} \iint_{S_1} \left(\frac{\partial U}{\partial n} G - U \frac{\partial G}{\partial n} \right) ds. \quad (3-30)$$

The conditions for validity of this equation are:

1. The scalar theory holds.
2. Both U and G satisfy the homogeneous scalar wave equation.
3. The Sommerfeld radiation condition is satisfied.

Suppose that the Green's function G of the Kirchhoff theory were modified in such a way that, while the development leading to the above equation remains valid, in addition, either G or $\partial G/\partial n$ vanishes over the entire surface S_1 . In either case the necessity of imposing boundary conditions on *both* U and $\partial U/\partial n$ would be removed, and the inconsistencies of the Kirchhoff theory would be eliminated.

Sommerfeld pointed out that Green's functions with the required properties do indeed exist. Suppose G is generated not only by a point source located at P_0 , but also simultaneously by a second point source at a position \tilde{P}_0 which is the mirror image of P_0 on the opposite side of the screen (see Fig. 3.8). Let the source at \tilde{P}_0 be of the same wavelength λ as the source at P_0 , and suppose that the two sources are oscillating with a 180° phase difference. The Green's function in this case is given by

$$G_-(P_1) = \frac{\exp(jkr_{01})}{r_{01}} - \frac{\exp(jk\tilde{r}_{01})}{\tilde{r}_{01}} \quad (3-31)$$

Clearly such a function vanishes on the plane aperture Σ , leaving the following expression for the observed field:

$$U_I(P_0) = \frac{-1}{4\pi} \iint_{\Sigma} U \frac{\partial G_-}{\partial n} ds. \quad (3-32)$$

We refer to this solution as the *first Rayleigh-Sommerfeld solution*.

⁴The fact that one theory is consistent and the other is not does not necessarily mean that the former is *more accurate* than the latter.

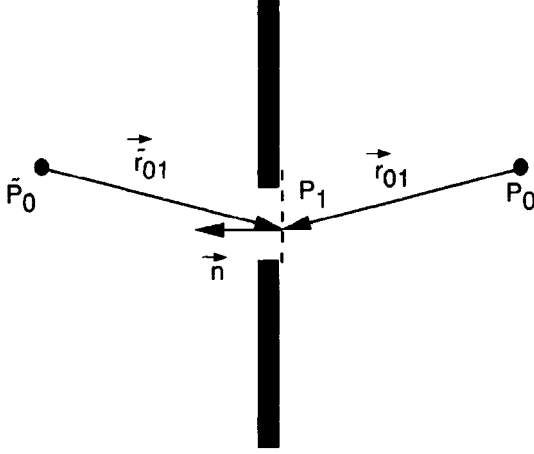


FIGURE 3.8
Rayleigh-Sommerfeld formulation of
diffraction by a plane screen.

To specify this solution further let \tilde{r}_{01} be the distance from \tilde{P}_0 to P_1 . The corresponding normal derivative of G_- is

$$\begin{aligned} \frac{\partial G_-}{\partial n}(P_1) &= \cos(\vec{n}, \vec{r}_{01}) \left(jk - \frac{1}{r_{01}} \right) \frac{\exp(jkr_{01})}{r_{01}} \\ &\quad - \cos(\vec{n}, \vec{\tilde{r}}_{01}) \left(jk - \frac{1}{\tilde{r}_{01}} \right) \frac{\exp(jk\tilde{r}_{01})}{\tilde{r}_{01}}. \end{aligned} \quad (3-33)$$

Now for P_1 on S_1 , we have

$$\begin{aligned} r_{01} &= \tilde{r}_{01} \\ \cos(\vec{n}, \vec{r}_{01}) &= -\cos(\vec{n}, \vec{\tilde{r}}_{01}) \end{aligned}$$

and therefore on that surface

$$\frac{\partial G_-(P_1)}{\partial n} = 2 \cos(\vec{n}, \vec{r}_{01}) \left(jk - \frac{1}{r_{01}} \right) \frac{\exp(jkr_{01})}{r_{01}}. \quad (3-34)$$

For $r_{01} \gg \lambda$, the second term above can be dropped, leaving

$$\frac{\partial G_-(P_1)}{\partial n} = 2jk \cos(\vec{n}, \vec{r}_{01}) \frac{\exp(jkr_{01})}{r_{01}}, \quad (3-35)$$

which is just twice the normal derivative of the Green's function G used in the Kirchhoff analysis, i.e.

$$\frac{\partial G_-(P_1)}{\partial n} = 2 \frac{\partial G(P_1)}{\partial n}.$$

With this result, the first Rayleigh-Sommerfeld solution can be expressed in terms of the more simple Green's function used by Kirchhoff,

$$U_I(P_0) = \frac{-1}{2\pi} \iint_{\Sigma} U \frac{\partial G}{\partial n} ds. \quad (3-36)$$

An alternative and equally valid Green's function is found by allowing the two point sources to oscillate in phase, giving

$$G_+(P_1) = \frac{\exp(jkr_{01})}{r_{01}} + \frac{\exp(jk\tilde{r}_{01})}{\tilde{r}_{01}}. \quad (3-37)$$

It is readily shown (see Prob. 3-4) that the *normal* derivative of this function vanishes across the screen and aperture, leading to the second *Rayleigh-Sommerfeld* solution,

$$U_{II}(P_0) = \frac{1}{4\pi} \iint_{\Sigma} \frac{\partial U}{\partial n} G_+ ds. \quad (3-38)$$

It can be shown that, on Σ and under the condition that $r_{01} \gg A$, G_+ is twice the Kirchhoff Green's function G ,

$$G_+^- = 2G.$$

This leads to an expression for $U(P_0)$ in terms of the Green's function used by Kirchhoff,

$$U_{II}(P_0) = \frac{1}{2\pi} \iint_{\Sigma} \frac{\partial U}{\partial n} G ds. \quad (3-39)$$

3.5.2 The Rayleigh-Sommerfeld Diffraction Formula

Let the Green's function G_- be substituted for G in Eq. (3-23). Using (3-35), it follows directly that

$$U_I(P_0) = \frac{1}{j\lambda} \iint_{S_1} U(P_1) \frac{\exp(jkr_{01})}{r_{01}} \cos(\vec{n}, \vec{r}_{01}) ds \quad (3-40)$$

where it has been assumed that $r_{01} \gg A$. The Kirchhoff boundary conditions may now be applied to U alone, yielding the general result

$$U_I(P_0) = \frac{1}{j\lambda} \iint_{\Sigma} U(P_1) \frac{\exp(jkr_{01})}{r_{01}} \cos(\vec{n}, \vec{r}_{01}) ds \quad (3-41)$$

Since no boundary conditions need be applied to $\partial U/\partial n$, the inconsistencies of the Kirchhoff theory have been removed.

If the alternative Green's function of (3-37) is used, the result can be shown to be

$$U_{II}(P_0) = \frac{1}{2\pi} \iint_{\Sigma} \frac{\partial U(P_1)}{\partial n} \frac{\exp(jkr_{01})}{r_{01}} ds. \quad (3-42)$$

We now specialize Eq. (3-41) and Eq. (3-42) to the case of illumination with a diverging spherical wave, allowing direct comparison with Eq. (3-27) of the Kirchhoff

theory. The illumination of the aperture in all cases is a spherical wave diverging from a point source at position P_2 (see Fig. 3.7 again):

$$U(P_1) = A \frac{\exp(jkr_{21})}{r_{21}}$$

Using G_- we obtain

$$U_I(P_0) = \frac{A}{j\lambda} \iint_{\Sigma} \frac{\exp[jk(r_{21} + r_{01})]}{r_{21}r_{01}} \cos(\vec{n}, \vec{r}_{01}) ds. \quad (3-43)$$

This result is known as the *Rayleigh-Sommerfeld diffraction formula*. Using G_+ , and assuming that $r_{21} \gg \lambda$, the corresponding result is

$$U_{II}(P_0) = -\frac{A}{j\lambda} \iint_{\Sigma} \frac{\exp[jk(r_{21} + r_{01})]}{r_{21}r_{01}} \cos(\vec{n}, \vec{r}_{21}) ds \quad (3-44)$$

where the angle between \vec{n} and \vec{r}_{21} is greater than 90° .

3.6 COMPARISON OF THE KIRCHHOFF AND RAYLEIGH-SOMMERFELD THEORIES

We briefly summarize the similarities and differences of the Kirchhoff and the Rayleigh-Sommerfeld theories. For the purposes of this section, let G_K represent the Green's function for the Kirchhoff theory, while G_- and G_+ are the Green's functions for the two Rayleigh-Sommerfeld formulations. As pointed out earlier, on the surface Σ , $G_+ = 2G_K$ and $\partial G_- / \partial n = 2\partial G_K / \partial n$. Therefore the general results of interest are as follows. For the Kirchhoff theory (cf. Eq. (3-24))

$$U(P_0) = \frac{1}{4\pi} \iint_{\Sigma} \left(\frac{\partial U}{\partial n} G_K - U \frac{\partial G_K}{\partial n} \right) ds, \quad (3-45)$$

for the first Rayleigh-Sommerfeld solution (cf. Eq. (3-36))

$$U_I(P_0) = \frac{-1}{2\pi} \iint_{\Sigma} U \frac{\partial G_K}{\partial n} ds, \quad (3-46)$$

and for the second Rayleigh-Sommerfeld solution (cf. Eq. (3-39))

$$U_{II}(P_0) = \frac{1}{2\pi} \iint_{\Sigma} \frac{\partial U}{\partial n} G_K ds. \quad (3-47)$$

A comparison of the above equations leads us to an interesting and surprising conclusion: *the Kirchhoff solution is the arithmetic average of the two Rayleigh-Sommerfeld solutions!*

Comparing the results of the three approaches for the case of spherical wave illumination, we see that the results derived from the Rayleigh-Sommerfeld theory (i.e. Eqs. (3-43) and (3-44)) differ from the Fresnel-Kirchhoff diffraction formula, Eq. (3-27), only through what is known as the *obliquity factor* ψ , which is the angular dependence introduced by the cosine terms. For *all* cases we can write

$$U(P_0) = \frac{A}{j\lambda} \iint_{\Sigma} \frac{\exp[jk(r_{21} + r_{01})]}{r_{21}r_{01}} \psi ds, \quad (3-48)$$

where

$$\psi = \begin{cases} \frac{1}{2}[\cos(\vec{n}, \vec{r}_{01}) - \cos(\vec{n}, \vec{r}_{21})] & \text{Kirchhoff theory} \\ \cos(\vec{n}, \vec{r}_{01}) & \text{First Rayleigh-Sommerfeld solution} \\ -\cos(\vec{n}, \vec{r}_{21}) & \text{Second Rayleigh-Sommerfeld solution.} \end{cases} \quad (3-49)$$

For the special case of an infinitely distant point source producing normally incident plane wave illumination, the obliquity factors become

$$\psi = \begin{cases} \frac{1}{2}[1 + \cos \theta] & \text{Kirchhoff theory} \\ \cos \theta & \text{First Rayleigh-Sommerfeld solution} \\ 1 & \text{Second Rayleigh-Sommerfeld solution,} \end{cases} \quad (3-50)$$

where θ is the angle between the vectors \vec{n} and \vec{r}_{01} .

Several authors have compared the two formulations of the diffraction problem. We mention in particular Wolf and Marchand [301], who examined differences between the two theories for circular apertures with observation points at a sufficiently great distance from the aperture to be in the "far field" (the meaning of this term will be explained in the chapter to follow). They found the Kirchhoff solution and the two Rayleigh-Sommerfeld solutions to be essentially the same provided the aperture diameter is much greater than a wavelength. Heurtley [143] examined the predictions of the three solutions for observation points on the axis of a circular aperture for all distances behind the aperture, and found differences between the theories only close to the aperture.

When only small angles are involved in the diffraction problem, it is easy to show that all three solutions are identical. In all three cases the obliquity factors approach unity as the angles become small, and the differences between the results vanish. Note that only small angles will be involved if we are far from the diffracting aperture.

In closing it is worth noting that, in spite of its internal inconsistencies, there is one sense in which the Kirchhoff theory is more general than the Rayleigh-Sommerfeld theory. The latter requires that the diffracting screens be *planar*, while the former does not. However, most of the problems of interest here will involve planar diffracting apertures, so this generality will not be particularly significant. In fact, we will generally choose to use the first Rayleigh-Sommerfeld solution because of its simplicity.

3.7

FURTHER DISCUSSION OF THE HUYGENS-FRESNEL PRINCIPLE

The Huygens-Fresnel principle, as predicted by the first Rayleigh-Sommerfeld solution⁵ (see Eq. (3-40)), can be expressed mathematically as follows:

$$U(P_0) = \frac{1}{j\lambda} \iint_{\Sigma} U(P_1) \frac{\exp(jkr_{01})}{r_{01}} \cos \theta \, ds, \quad (3-51)$$

where θ is the angle between the vectors \vec{n} and \vec{r}_{01} . We give a "quasi-physical" interpretation to this integral. It expresses the observed field $U(P_0)$ as a superposition of diverging spherical waves $\exp(jkr_{01})/r_{01}$ originating from secondary sources located at each and every point P_1 within the aperture Σ . The secondary source at P_1 has the following properties:

1. It has a complex amplitude that is proportional to the amplitude of the excitation $U(P_1)$ at the corresponding point.
2. It has an amplitude that is inversely proportional to λ , or equivalently directly proportional to the optical frequency ν .
3. It has a phase that leads the phase of the incident wave by 90° , as indicated by the factor $1/j$.
4. Each secondary source has a directivity pattern $\cos \theta$.

The first of these properties is entirely reasonable. The wave propagation phenomenon is linear, and the wave passed through the aperture should be proportional to the wave incident upon it.

A reasonable explanation of the second and third properties would be as follows. Wave motion from the aperture to the observation point takes place by virtue of changes of the field in the aperture. In the next section we will see more explicitly that the field at P_0 contributed by a secondary source at P_1 depends on the time-rate-of-change of the field at P_1 . Since our basic monochromatic field disturbance is a clockwise rotating phasor of the form $\exp(-j2\pi\nu t)$, the derivative of this function will be proportional to both ν and to $-j = 1/j$.

The last property, namely the obliquity factor, has no simple "quasi-physical" explanation, but arises in slightly different forms in all the theories of diffraction. It is perhaps expecting too much to find such an explanation. After all, there are no material sources within the aperture; rather, they all lie on the rim of the aperture. Therefore the Huygens-Fresnel principle should be regarded as a relatively simple mathematical construct that allows us to solve diffraction problems without paying attention to the physical details of exactly what is happening at the edges of the aperture.

It is important to realize that the Huygens-Fresnel principle, as expressed by Eq. (3-51), is nothing more than a superposition integral of the type discussed in Chapter 2. To emphasize this point of view we rewrite (3-51) as

⁵Hereafter we drop the subscript on the first Rayleigh-Sommerfeld solution, since it will be the solution we use exclusively.

$$U(P_0) = \iint_{\Sigma} h(P_0, P_1) U(P_1) ds, \quad (3-52)$$

where the impulse response $h(P_0, P_1)$ is given explicitly by

$$h(P_0, P_1) = \frac{1}{j\lambda} \frac{\exp(jkr_{01})}{r_{01}} \cos \theta. \quad (3-53)$$

The occurrence of a superposition integral as a result of our diffraction analysis should not be a complete surprise. The primary ingredient required for such a result was previously seen to be linearity, a property that was assumed early in our analysis. When we examine the character of the impulse response $h(P_0, P_1)$ in more detail in Chapter 4, we will find that it is also space-invariant, a consequence of the homogeneity assumed for the dielectric medium. The Huygens-Fresnel principle will then be seen to be a convolution integral.

3.8 GENERALIZATION TO NONMONOCHROMATIC WAVES

The wave disturbances have previously been assumed to be ideally monochromatic in all cases. Such waves can be closely approximated in practice and are particularly easy to analyze. However, the more general case of a nonmonochromatic disturbance will now be considered briefly; attention is restricted to the predictions of the first Rayleigh-Sommerfeld solution, but similar results can be obtained for the other solutions.

Consider the scalar disturbance $u(P_0, t)$ observed behind an aperture Σ in an opaque screen when a disturbance $u(P_1, t)$ is incident on that aperture. The time functions $u(P_0, t)$ and $u(P_1, t)$ may be expressed in terms of their inverse Fourier transforms:

$$u(P_1, t) = \int_{-\infty}^{\infty} U(P_1, \nu) \exp(j2\pi\nu t) d\nu \quad (3-54)$$

$$u(P_0, t) = \int_{-\infty}^{\infty} U(P_0, \nu) \exp(j2\pi\nu t) d\nu,$$

where $U(P_0, \nu)$ and $U(P_1, \nu)$ are the Fourier spectra of $u(P_0, t)$ and $u(P_1, t)$, respectively, and ν represents frequency.

Let Eqs. (3-54) be transformed by the change of variables $\nu' = -\nu$, yielding

$$u(P_1, t) = \int_{-\infty}^{\infty} U(P_1, -\nu') \exp(-j2\pi\nu' t) d\nu' \quad (3-55)$$

$$u(P_0, t) = \int_{-\infty}^{\infty} U(P_0, -\nu') \exp(-j2\pi\nu' t) d\nu'.$$

Now these relations may be regarded as expressing the nonmonochromatic time functions $u(P_1, t)$ and $u(P_0, t)$ as a linear combination of monochromatic time functions of the type represented by Eq. (3-10). The monochromatic elementary functions are of

various frequencies ν' , the complex amplitudes of the disturbance at frequency ν' being simply $U(P_1, -\nu')$ and $U(P_0, -\nu')$. By invoking the linearity of the wave-propagation phenomenon, we use the results of the previous section to find the complex amplitude at P_0 of each monochromatic component of the disturbance, and superimpose these results to yield the general time function $u(P_0, t)$.

To proceed, Eq. (3-51) can be directly used to write

$$U(P_0, -\nu') = -j \frac{\nu'}{v} \iint_{\Sigma} U(P_1, -\nu') \frac{\exp(j2\pi\nu' r_{01}/v)}{r_{01}} \cos(\vec{n}, \vec{r}_{01}) ds, \quad (3-56)$$

where v is the velocity of propagation of the disturbance in a medium of refractive index n ($v = c/n$), and the relation $\nu'\lambda = v$ has been used. Substitution of (3-56) in the second of Eqs. (3-55) and an interchange of the orders of integration give

$$u(P_0, t) = \iint_{\Sigma} \frac{\cos(\vec{n}, \vec{r}_{01})}{2\pi\nu r_{01}} \int_{-\infty}^{\infty} -j2\pi\nu' U(P_1, -\nu') \exp\left[-j2\pi\nu' \left(t - \frac{r_{01}}{v}\right)\right] d\nu' ds.$$

Finally, the identity

$$\begin{aligned} \frac{d}{dt} u(P_1, t) &= \frac{d}{dt} \int_{-\infty}^{\infty} U(P_1, -\nu') \exp(-j2\pi\nu' t) d\nu' \\ &= \int_{-\infty}^{\infty} -j2\pi\nu' U(P_1, -\nu') \exp(-j2\pi\nu' t) d\nu' \end{aligned}$$

can be used to write

$$u(P_0, t) = \iint_{\Sigma} \frac{\cos(\vec{n}, \vec{r}_{01})}{2\pi\nu r_{01}} \frac{d}{dt} u\left(P_1, t - \frac{r_{01}}{v}\right) ds. \quad (3-57)$$

The wave disturbance at point P_0 is seen to be linearly proportional to the time derivative of the disturbance at each point P_1 on the aperture. Since it takes time r_{01}/v for the disturbance to propagate from P_1 to P_0 , the observed wave depends on the derivative of the incident wave at the "retarded" time $t - (r_{01}/v)$.

This more general treatment shows that an understanding of diffraction of monochromatic waves can be used directly to synthesize the results for much more general nonmonochromatic waves. However, the monochromatic results are directly applicable themselves when the optical source has a sufficiently narrow spectrum. See Prob. 3-6 for further elucidation of these points.

3.9 DIFFRACTION AT BOUNDARIES

In the statement of the Huygens-Fresnel principle, we found it convenient to regard each point on the aperture as a new source of spherical waves. It was pointed out that such sources are merely mathematical conveniences and have no real physical significance. A more physical point-of-view, first qualitatively expressed by Thomas Young in 1802,

is to regard the observed field as consisting of a superposition of the incident wave transmitted through the aperture unperturbed, and a diffracted wave originating at the *rim* of the aperture. The possibility of a new wave originating in the material medium of the rim makes this interpretation a more physical one.

Young's qualitative arguments were given added impetus by Sommerfeld's rigorous electromagnetic solution of the problem of diffraction of a plane wave by a semi-infinite, perfectly conducting screen [268]. This rigorous solution showed that the field in the geometrical shadow of the screen has the form of a cylindrical wave originating on the rim of the screen. In the directly illuminated region behind the plane of the screen the field was found to be a superposition of this cylindrical wave with the directly transmitted wave.

The applicability of a boundary diffraction approach in more general diffraction problems was investigated by Maggi [202] and Rubinowicz [249], who showed that the Kirchhoff diffraction formula can indeed be manipulated to yield a form that is equivalent to Young's ideas. More recently, Miyamoto and Wolf [250] have extended the theory of boundary diffraction. For further discussion of these ideas, the reader should consult the references cited.

Another approach closely related to Young's ideas is the geometrical theory of diffraction developed by Keller [161]. In this treatment, the field behind a diffracting obstacle is found by the principles of geometrical optics, modified by the inclusion of "diffracted rays" that originate at certain points on the obstacle itself. New rays are assumed to be generated at edges, corners, tips, and surfaces of the obstacle. This theory can often be applied to calculate the fields diffracted by objects that are too complex to be treated by other methods.

3.10 THE ANGULAR SPECTRUM OF PLANE WAVES

It is also possible to formulate scalar diffraction theory in a framework that closely resembles the theory of linear, invariant systems. As we shall see, if the complex field distribution of a monochromatic disturbance is Fourier-analyzed across any plane, the various spatial Fourier components can be identified as plane waves traveling in different directions away from that plane. The field amplitude at any other point (or across any other parallel plane) can be calculated by adding the contributions of these plane waves, taking due account of the phase shifts they have undergone during propagation. For a detailed treatment of this approach to diffraction theory, as well as its applications in the theory of radio-wave propagation, the reader is referred to the work of Ratcliffe [240].

3.10.1 The Angular Spectrum and Its Physical Interpretation

Suppose that, due to some unspecified system of monochromatic sources, a wave is incident on a transverse (x, y) plane traveling with a component of propagation in the positive z direction. Let the complex field across that $z = 0$ plane be represented by $U(x, y, 0)$; our ultimate objective is to calculate the resulting field $U(x, y, z)$ that appears across a second, parallel plane a distance z to the right of the first plane.

Across the $z = 0$ plane, the function U has a two-dimensional Fourier transform given by

$$A(f_x, f_y; 0) = \iint_{-\infty}^{\infty} U(x, y, 0) \exp[-j2\pi(f_x x + f_y y)] dx dy. \quad (3-58)$$

As pointed out in Chapter 2, the Fourier transform operation may be regarded as a decomposition of a complicated function into a collection of more simple **complex-exponential** functions. To emphasize this point-of-view, we write U as an inverse Fourier transform of its spectrum,

$$U(x, y, 0) = \iint_{-\infty}^{\infty} A(f_x, f_y; 0) \exp[j2\pi(f_x x + f_y y)] df_x df_y. \quad (3-59)$$

To give physical meaning to the functions in the integrand of the above integral, consider the form of a simple plane wave propagating with wave vector \vec{k} , where \vec{k} has magnitude $2\pi/\lambda$ and has direction cosines (α, β, γ) , as illustrated in Fig. 3.9. Such a plane wave has a complex representation of the form

$$p(x, y, z; t) = \exp[j(\vec{k} \cdot \vec{r} - 2\pi\nu t)] \quad (3-60)$$

where $\vec{r} = x\hat{x} + y\hat{y} + z\hat{z}$ is a position vector (the $\hat{}$ symbol signifies a unit vector), while $\vec{k} = \frac{2\pi}{\lambda}(\alpha\hat{x} + \beta\hat{y} + \gamma\hat{z})$. Dropping the time dependence, the complex phasor amplitude of the plane wave across a constant z -plane is

$$P(x, y, z) = \exp(j\vec{k} \cdot \vec{r}) = e^{j\frac{2\pi}{\lambda}(\alpha x + \beta y)} e^{j\frac{2\pi}{\lambda}\gamma z}. \quad (3-61)$$

Note that the direction cosines are interrelated through

$$\gamma = \sqrt{1 - \alpha^2 - \beta^2}.$$

Thus across the plane $z = 0$, a complex-exponential function $\exp[j2\pi(f_x x + f_y y)]$ may be regarded as representing a plane wave propagating with direction cosines

$$\alpha = \lambda f_x \quad \beta = \lambda f_y \quad \gamma = \sqrt{1 - (\lambda f_x)^2 - (\lambda f_y)^2}. \quad (3-62)$$

In the Fourier decomposition of U , the complex amplitude of the plane-wave component with spatial frequencies (f_x, f_y) is simply $A(f_x, f_y; 0) df_x df_y$, evaluated at $(f_x = \alpha/\lambda, f_y = \beta/\lambda)$. For this reason, the function

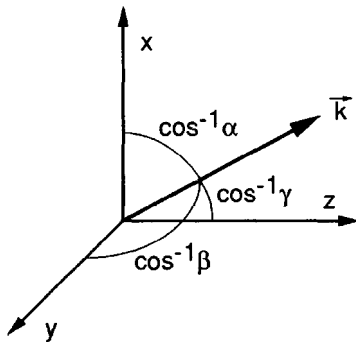


FIGURE 3.9
The wave vector \vec{k} .

$$A\left(\frac{\alpha}{\lambda}, \frac{\beta}{\lambda}; 0\right) = \iint_{-\infty}^{\infty} U(x, y, 0) \exp\left[-j2\pi\left(\frac{\alpha}{\lambda}x + \frac{\beta}{\lambda}y\right)\right] dx dy \quad (3-63)$$

is called the *angular spectrum* of the disturbance $U(x, y, 0)$.

3.10.2 Propagation of the Angular Spectrum

Consider now the angular spectrum of the disturbance U across a plane parallel to the (x, y) plane but at a distance z from it. Let the function $A(\alpha/\lambda, \beta/\lambda; z)$ represent the angular spectrum of $U(x, y, z)$; that is,

$$A\left(\frac{\alpha}{\lambda}, \frac{\beta}{\lambda}; z\right) = \iint_{-\infty}^{\infty} U(x, y, z) \exp\left[-j2\pi\left(\frac{\alpha}{\lambda}x + \frac{\beta}{\lambda}y\right)\right] dx dy. \quad (3-64)$$

Now if the relation between $A(\alpha/\lambda, \beta/\lambda; 0)$ and $A(\alpha/\lambda, \beta/\lambda; z)$ can be found, then the effects of wave propagation on the angular spectrum of the disturbance will be evident.

To find the desired relation, note that U can be written

$$U(x, y, z) = \iint_{-\infty}^{\infty} A\left(\frac{\alpha}{\lambda}, \frac{\beta}{\lambda}; z\right) \exp\left[j2\pi\left(\frac{\alpha}{\lambda}x + \frac{\beta}{\lambda}y\right)\right] d\frac{\alpha}{\lambda} d\frac{\beta}{\lambda}. \quad (3-65)$$

In addition, U must satisfy the Helmholtz equation,

$$\nabla^2 U + k^2 U = 0$$

at all source-free points. Direct application of this requirement to Eq. (3-65) shows that A must satisfy the differential equation

$$\frac{d^2}{dz^2} A\left(\frac{\alpha}{\lambda}, \frac{\beta}{\lambda}; z\right) + \left(\frac{2\pi}{\lambda}\right)^2 [1 - \alpha^2 - \beta^2] A\left(\frac{\alpha}{\lambda}, \frac{\beta}{\lambda}; z\right) = 0.$$

An elementary solution of this equation can be written in the form

$$A\left(\frac{\alpha}{\lambda}, \frac{\beta}{\lambda}; z\right) = A\left(\frac{\alpha}{\lambda}, \frac{\beta}{\lambda}; 0\right) \exp\left(j\frac{2\pi}{\lambda} \sqrt{1 - \alpha^2 - \beta^2} z\right). \quad (3-66)$$

This result demonstrates that when the direction cosines (α, β) satisfy

$$\alpha^2 + \beta^2 < 1, \quad (3-67)$$

as all true direction cosines must, the effect of propagation over distance z is simply a change of the relative phases of the various components of the angular spectrum. Since each plane-wave component propagates at a different angle, each travels a different distance between two parallel planes, and relative phase delays are thus introduced.

However, when (α, β) satisfy

$$\alpha^2 + \beta^2 > 1,$$

a different interpretation is required. Note that since $A(\alpha/\lambda, \beta/\lambda; 0)$ is the Fourier transform of a field distribution on which boundary conditions are imposed in the aperture

plane, it is quite possible that this spectrum will contain components that satisfy the above equation. Under such a condition, α and β are no longer interpretable as direction cosines. Now the square root in Eq. (3-66) is imaginary, and that equation can be rewritten

$$A\left(\frac{\alpha}{\lambda}, \frac{\beta}{\lambda}; z\right) = A\left(\frac{\alpha}{\lambda}, \frac{\beta}{\lambda}; 0\right) \exp(-\mu z) \quad (3-68)$$

where

$$\mu = \frac{2\pi}{\lambda} \sqrt{\alpha^2 + \beta^2 - 1}.$$

Since μ is a positive real number, these wave components are rapidly attenuated by the propagation phenomenon. Such components are called *evanescent waves* and are quite analogous to the waves produced in a microwave waveguide driven below its cutoff frequency. As in the case of the waveguide driven below cutoff, these evanescent waves carry no energy away from the **aperture**.⁶

Finally, we note that the disturbance observed at (x, y, z) can be written in terms of the initial angular spectrum by inverse transforming Eq. (3-66), giving

$$U(x, y, z) = \iint_{-\infty}^{\infty} A\left(\frac{\alpha}{\lambda}, \frac{\beta}{\lambda}; 0\right) \exp\left(j\frac{2\pi}{\lambda} \sqrt{1 - \alpha^2 - \beta^2} z\right) \\ \times \text{circ}(\sqrt{\alpha^2 + \beta^2}) \exp\left[j2\pi\left(\frac{\alpha}{\lambda}x + \frac{\beta}{\lambda}y\right)\right] d\frac{\alpha}{\lambda} d\frac{\beta}{\lambda}, \quad (3-69)$$

where the circ function limits the region of integration to the region within which Eq. (3-67) is **satisfied**.⁷ Note that no angular spectrum components beyond the evanescent wave cutoff contribute to $U(x, y, z)$. This fact is the fundamental reason why no conventional imaging system can resolve a periodic structure with a period that is finer than the wavelength of the radiation used. It is possible, though, to couple to evanescent waves with very fine structures placed in very close proximity to the diffracting object, and thereby recover information that would otherwise be lost. However, we will focus here on conventional optical instruments, for which the evanescent waves are not recoverable.

3.10.3 Effects of a Diffracting Aperture on the Angular Spectrum

Suppose that an infinite opaque screen containing a diffracting structure is introduced in the plane $z = 0$. We now consider the effects of that diffracting screen on the

⁶Note that evanescent waves are predicted only under the very same conditions for which the use of the scalar theory is suspect. Nonetheless, they are a real phenomenon, although perhaps more accurately treated in a full vectorial theory.

⁷We can usually assume that the distance z is larger than a few wavelengths, allowing us to completely drop the evanescent components of the spectrum.

angular spectrum of the disturbance. Define the *amplitude transmittance function* of the aperture as the ratio of the transmitted field amplitude $U_t(x, y; 0)$ to the incident field amplitude $U_i(x, y; 0)$ at each (x, y) in the $z = 0$ plane,

$$t_A(x, y) = \frac{U_t(x, y; 0)}{U_i(x, y; 0)}. \quad (3-70)$$

Then

$$U_t(x, y, 0) = U_i(x, y, 0) t_A(x, y)$$

and the convolution theorem can be used to relate the angular spectrum $A_i(\alpha/\lambda, \beta/\lambda)$ of the incident field and the angular spectrum $A_t(\alpha/\lambda, \beta/\lambda)$ of the transmitted field,

$$A_t\left(\frac{\alpha}{\lambda}, \frac{\beta}{\lambda}\right) = \left[A_i\left(\frac{\alpha}{\lambda}, \frac{\beta}{\lambda}\right) \otimes T\left(\frac{\alpha}{\lambda}, \frac{\beta}{\lambda}\right) \right], \quad (3-71)$$

where

$$T\left(\frac{\alpha}{\lambda}, \frac{\beta}{\lambda}\right) = \iint_{-\infty}^{\infty} t_A(x, y) \exp\left[-j2\pi\left(\frac{\alpha}{\lambda}x + \frac{\beta}{\lambda}y\right)\right] dx dy, \quad (3-72)$$

and \otimes is again the symbol for convolution.

The angular spectrum of the transmitted disturbance is thus seen to be the convolution of the angular spectrum of the incident disturbance with a second angular spectrum that is characteristic of the diffracting structure.

For the case of a unit amplitude plane wave illuminating the diffracting structure normally, the result takes a particularly simple form. In that case

$$A_i\left(\frac{\alpha}{\lambda}, \frac{\beta}{\lambda}\right) = \delta\left(\frac{\alpha}{\lambda}, \frac{\beta}{\lambda}\right)$$

and

$$A_t\left(\frac{\alpha}{\lambda}, \frac{\beta}{\lambda}\right) = \delta\left(\frac{\alpha}{\lambda}, \frac{\beta}{\lambda}\right) \otimes T\left(\frac{\alpha}{\lambda}, \frac{\beta}{\lambda}\right) = T\left(\frac{\alpha}{\lambda}, \frac{\beta}{\lambda}\right).$$

Thus the transmitted angular spectrum is found directly by Fourier transforming the amplitude transmittance function of the aperture.

Note that, if the diffracting structure is an aperture that limits the extent of the field distribution, the result is a broadening of the angular spectrum of the disturbance, from the basic properties of Fourier transforms. The smaller the aperture, the broader the angular spectrum behind the aperture. This effect is entirely analogous to the broadening of the spectrum of an electrical signal as its duration is decreased.

3.10.4 The Propagation Phenomenon as a Linear Spatial Filter

Consider again the propagation of light from the plane $z = 0$ to a parallel plane at nonzero distance z . The disturbance $U(x, y, 0)$ incident on the first plane may be considered to be mapped by the propagation phenomenon into a new field distribution $U(x, y, z)$. Such a mapping satisfies our previous definition of a system. We shall, in

fact, demonstrate that the propagation phenomenon acts as a linear space-invariant system and is characterized by a relatively simple transfer function.

The linearity of the propagation phenomenon has already been discussed; it is directly implied by the linearity of the wave equation, or alternatively, by the superposition integral (3-52). The space-invariance property is most easily demonstrated by actually deriving a transfer function that describes the effects of propagation; if the mapping has a transfer function, then it must be space-invariant.

To find the transfer function, we return to the angular spectrum point-of-view. However, rather than writing the angular spectra as functions of the direction cosines (α, β), it is now more convenient to leave the spectra as functions of spatial frequencies (f_X, f_Y). The spatial frequencies and the direction cosines are related through Eq. (3-62).

Let the spatial spectrum of $U(x, y, z)$ again be represented by $A(f_X, f_Y; z)$, while the spectrum of $U(x, y; 0)$ is again written $A(f_X, f_Y; 0)$. Thus we may express $U(x, y, z)$ as

$$U(x, y, z) = \iint_{-\infty}^{\infty} A(f_X, f_Y; z) \exp[j2\pi(f_X x + f_Y y)] df_X df_Y.$$

But in addition, from Eq.(3-69),

$$U(x, y, z) = \iint_{-\infty}^{\infty} A(f_X, f_Y; 0) \text{circ}(\sqrt{(\lambda f_X)^2 + (\lambda f_Y)^2}) \\ \times \exp\left[j\frac{2\pi}{\lambda} \sqrt{1 - (\lambda f_X)^2 - (\lambda f_Y)^2} z\right] \exp[j2\pi(f_X x + f_Y y)] df_X df_Y,$$

where we have again explicitly introduced the bandwidth limitation associated with evanescent waves through the use of a circ function. A comparison of the above two equations shows that

$$A(f_X, f_Y; z) = A(f_X, f_Y; 0) \text{circ}(\sqrt{(\lambda f_X)^2 + (\lambda f_Y)^2}) \\ \times \exp\left[j2\pi\frac{z}{\lambda} \sqrt{1 - (\lambda f_X)^2 - (\lambda f_Y)^2}\right]. \quad (3-73)$$

Finally, the transfer function of the wave propagation phenomenon is seen to be

$$H(f_X, f_Y) = \begin{cases} \exp\left[j2\pi\frac{z}{\lambda} \sqrt{1 - (\lambda f_X)^2 - (\lambda f_Y)^2}\right] & \sqrt{f_X^2 + f_Y^2} < \frac{1}{\lambda} \\ 0 & \text{otherwise.} \end{cases} \quad (3-74)$$

Thus the propagation phenomenon may be regarded as a linear, dispersive spatial filter with a finite bandwidth. The transmission of the filter is zero outside a circular region of radius λ^{-1} in the frequency plane. Within that circular bandwidth, the modulus of the transfer function is unity but frequency-dependent phase shifts are introduced. The phase dispersion of the system is most significant at high spatial frequencies and

vanishes as both f_X and f_Y approach zero. In addition, for any fixed spatial frequency pair, the phase dispersion increases as the distance of propagation z increases.

In closing we mention the remarkable fact that, despite the apparent differences of their approaches, *the angular spectrum approach and the first Rayleigh-Sommerfeld solution yield identical predictions of diffracted fields!* This has been proved most elegantly by Sherman [260].

PROBLEMS-CHAPTER 3

- 3-1.** Show that in an isotropic, nonmagnetic, and inhomogeneous dielectric medium, Maxwell's equations can be combined to yield Eq. (3-8).
- 3-2.** Show that a diverging spherical wave satisfies the Sommerfeld radiation condition.
- 3-3.** Show that, if $r_{21} \gg \lambda$, Eq. (3-26) can be reduced to Eq. (3-27).
- 3-4.** Show that the normal derivative of Eq. (3-37) for G_+ vanishes across the screen and aperture.
- 3-5.** Assuming unit-amplitude normally incident plane-wave illumination, find the angular spectrum of
- A circular aperture of diameter d .
 - A circular opaque disk of diameter d .
- 3-6.** Consider a real nonmonochromatic disturbance $u(P, t)$ of center frequency ω and bandwidth $\Delta\omega$. Let a related complex-valued disturbance $u_-(P, t)$ be defined as consisting of only the negative-frequency components of $u(P, t)$. Thus

$$u_-(P, t) = \int_{-\infty}^0 U(P, \nu) \exp(j2\pi\nu t) d\nu$$

where $U(P, \nu)$ is the Fourier spectrum of $u(P, t)$. Assuming the geometry of Fig. 3.6 show that if

$$\frac{\Delta\omega}{\bar{\omega}} \ll 1 \text{ and } \frac{1}{\Delta\omega} \gg \frac{nr_{01}}{v}$$

then

$$u_-(P_0, t) = \frac{1}{j\bar{\lambda}} \iint_{-\infty}^{\infty} u_-(P_1, t) \frac{\exp(j\bar{k}r_{01})}{r_{01}} \cos(\bar{n}, \bar{r}_{01}) ds$$

where $\bar{\lambda} = v/\bar{\omega}$ and $\bar{k} = 2\pi/\bar{\lambda}$. In the above equations, n is the refractive index of the medium and v is the velocity of propagation.

- 3-7.** For a wave that travels only in directions that have small angles with respect to the optical axis, the general form of the complex field may be approximated by

$$U(x, y, z) \approx A(x, y, z) \exp(jkz),$$

where $A(x, y, z)$ is a slowly varying function of z .

62 Introduction to Fourier Optics

- (a) Show that for such a wave the Helmholtz equation can be reduced to

$$\nabla_t^2 A + j2k \frac{\partial A}{\partial z} = 0,$$

where $\nabla_t^2 = \partial^2/\partial x^2 + \partial^2/\partial y^2$ is the transverse portion of the Laplacian. This equation is known as the *paraxial* Helmholtz equation.

- (b) Show that a solution to this equation is given by

$$A(x, y, z) = \frac{A_1}{q(z)} \exp \left[jk \frac{x^2 + y^2}{2q(z)} \right]$$

for any complex $q(z)$ having $\frac{d}{dz}q(z) = 1$

- (c) Given

$$\frac{1}{q(z)} = \frac{1}{R(z)} + j \frac{\lambda}{\pi W^2(z)},$$

show that the solution $U(x, y, z)$ takes the form

$$U(x, y, z) = A_1 \frac{W_0}{W(z)} \exp \left[-\frac{\rho^2}{W^2(z)} \right] \exp \left[jkz + jk \frac{\rho^2}{2R(z)} + j\theta(z) \right]$$

where W_0 is a constant (independent of z) and $\theta(z)$ is a phase angle that changes with z . Note that this is a beam with a Gaussian profile and with a quadratic-phase approximation to a spherical wavefront.

Fresnel and Fraunhofer Diffraction

In the preceding chapter the results of scalar diffraction theory were presented in their most general forms. Attention is now turned to certain approximations to the general theory, approximations that will allow diffraction pattern calculations to be reduced to comparatively simple mathematical manipulations. These approximations, which are commonly made in many fields that deal with wave propagation, will be referred to as *Fresnel* and *Fraunhofer* approximations. In accordance with our view of the wave propagation phenomenon as a "system", we shall attempt to find approximations that are valid for a wide class of "input" field distributions.

4.1 BACKGROUND

In this section we prepare the reader for the calculations to follow. The concept of the *intensity* of a wave field is introduced, and the Huygens-Fresnel principle, from which the approximations are derived, is presented in a form that is especially well suited for approximation.

4.1.1 The Intensity of a Wave Field

In the optical region of the spectrum, a photodetector responds directly to the optical power falling on its surface. Thus for a semiconductor detector, if optical power \mathcal{P} is incident on the photosensitive region, absorption of a photon generates an electron in the conduction band and a hole in the valence band. Under the influence of internal and applied fields, the hole and electron move in opposite directions, leading to a **photocurrent** i that is the response to the incident absorbed photon. Under most circumstances

the photocurrent is linearly proportional to the incident power,

$$i = \mathbf{R}P. \quad (4-1)$$

The proportionality constant \mathbf{R} is called the responsivity of the detector and is given by

$$\mathcal{R} = \frac{\eta_{qe}q}{h\nu}, \quad (4-2)$$

where η_{qe} is the quantum *efficiency* of the photodetector (the average number of electron-hole pairs released by the absorption of a photon, a quantity that is less than or equal to unity in the absence of internal gain), q is the electronic charge (1.602×10^{-19} coulombs), h is Planck's constant (6.626196×10^{-34} joule-second), and ν is the optical frequency.¹

Thus in optics the directly measurable quantity is optical power, and it is important to relate such power to the complex scalar fields $u(\mathbf{P}, t)$ and $U(\mathbf{P})$ dealt with in earlier discussions of diffraction theory. To understand this relation requires a return to an electromagnetic description of the problem. We omit the details here, referring the reader to Ref. [253], Sections 5.3 and 5.4, and simply state the major points. Let the medium be isotropic, and the wave monochromatic. Assuming that the wave behaves locally as a transverse electromagnetic plane wave (i.e. $\vec{\mathcal{E}}$, $\vec{\mathcal{H}}$, and \vec{k} form a mutually orthogonal triplet), then the electric and magnetic fields can be expressed locally as

$$\begin{aligned} \vec{\mathcal{E}} &= \text{Re}\{\vec{E}_0 \exp[-j(2\pi\nu t - \vec{k} \cdot \vec{r})]\} \\ \vec{\mathcal{H}} &= \text{Re}\{\vec{H}_0 \exp[-j(2\pi\nu t - \vec{k} \cdot \vec{r})]\}, \end{aligned} \quad (4-3)$$

where \vec{E}_0 and \vec{H}_0 are locally constant and have complex components. The power flows in the direction of the vector \vec{k} and the power density can be expressed as

$$p = \frac{\vec{E}_0 \cdot \vec{E}_0^*}{2\eta} = \frac{E_{0x}^2 + E_{0y}^2 + E_{0z}^2}{2\eta}, \quad (4-4)$$

where η is the characteristic impedance of the medium and is given by

$$\eta = \sqrt{\frac{\mu}{\epsilon}}.$$

In vacuum, η is equal to 377 Ω . The total power incident on a surface of area A is the integral of the power density over A , taking into account that the direction of power flow is in the direction of \vec{k} ,

$$\mathcal{P} = \iint_A p \frac{\vec{k} \cdot \hat{n}}{|\vec{k}|} dx dy.$$

Here \hat{n} is a unit vector pointing into the surface of the detector, while $\vec{k}/|\vec{k}|$ is a unit vector in the direction of power flow. When \vec{k} is nearly normal to the surface, the total power \mathcal{P} is simply the integral of the power density p over the detector area.

¹The reader may wonder why the generation of both an electron and a hole does not lead to a charge $2q$ rather than q in this equation. For an answer, see [253], p. 653.

The proportionality of power density to the squared magnitude of the \vec{E}_0 vector seen in Eq. (4-4) leads us to define the intensity of a scalar monochromatic wave at point P as the squared magnitude of the complex phasor representation $U(P)$ of the disturbance,

$$I(P) = |U(P)|^2. \quad (4-5)$$

Note that power density and intensity are not identical, but the latter quantity is directly proportional to the former. For this reason we regard the intensity as the physically measurable attribute of an optical wavefield.

When a wave is not perfectly monochromatic, but is narrow band, a straightforward generalization of the concept of intensity is given by

$$I(P) = \langle |u(P, t)|^2 \rangle, \quad (4-6)$$

where the angle brackets signify an infinite time average. In some cases, the concept of instantaneous intensity is useful, defined as

$$I(P, t) = |u(P, t)|^2. \quad (4-7)$$

When calculating a diffraction pattern, we will generally regard the intensity of the pattern as the quantity we are seeking.

4.1.2 The Huygens-Fresnel Principle in Rectangular Coordinates

Before introducing a series of approximations to the Huygens-Fresnel principle, it will be helpful to first state the principle in more explicit form for the case of rectangular coordinates. As shown in Fig. 4.1, the diffracting aperture is assumed to lie in the (ξ, η) plane, and is illuminated in the positive z direction. We will calculate the wavefield across the (x, y) plane, which is parallel to the (ξ, η) plane and at normal distance z from it. The z axis pierces both planes at their origins.

According to Eq. (3-41), the Huygens-Fresnel principle can be stated as

$$U(P_0) = \frac{1}{j\lambda} \iint_{\Sigma} U(P_1) \frac{\exp(jkr_{01})}{r_{01}} \cos \theta \, ds, \quad (4-8)$$

where θ is the angle between the outward normal \hat{n} and the vector \vec{r}_{01} pointing from P_0 to P_1 . The term $\cos \theta$ is given exactly by

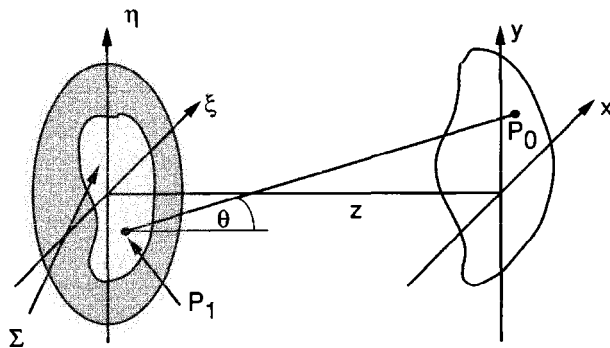


FIGURE 4.1
Diffraction geometry.

$$\cos \theta = \frac{z}{r_{01}},$$

and therefore the Huygens-Fresnel principle can be rewritten

$$U(x, y) = \frac{z}{j\lambda} \iint_{\Sigma} U(\xi, \eta) \frac{\exp(jkr_{01})}{r_{01}^2} d\xi d\eta, \quad (4-9)$$

where the distance r_{01} is given exactly by

$$r_{01} = \sqrt{z^2 + (x - \xi)^2 + (y - \eta)^2} \quad (4-10)$$

There have been only two approximations in reaching this expression. One is the approximation inherent in the scalar theory. The second is the assumption that the observation distance is many wavelengths from the aperture, $r_{01} \gg \lambda$. We now embark on a series of additional approximations.

4.2 THE FRESNEL APPROXIMATION

To reduce the Huygens-Fresnel principle to a more simple and usable expression, we introduce approximations for the distance r_{01} between P_1 and P_0 . The approximations are based on the binomial expansion of the square root in Eq. (4-10). Let b be a number that is less than unity, and consider the expression $\sqrt{1 + b}$. The binomial expansion of the square root is given by

$$\sqrt{1 + b} = 1 + \frac{1}{2}b - \frac{1}{8}b^2 + \cdots, \quad (4-11)$$

where the number of terms needed for a given accuracy depends on the magnitude of b .

To apply the binomial expansion to the problem at hand, factor a z outside the expression for r_{01} , yielding

$$r_{01} = z \sqrt{1 + \left(\frac{x - \xi}{z}\right)^2 + \left(\frac{y - \eta}{z}\right)^2}. \quad (4-12)$$

Let the quantity b in Eq. (4-11) consist of the second and third terms under the square root in (4-12). Then, retaining only the first two terms of the expansion (4-11), we have

$$r_{01} \approx z \left[1 + \frac{1}{2} \left(\frac{x - \xi}{z} \right)^2 + \frac{1}{2} \left(\frac{y - \eta}{z} \right)^2 \right]. \quad (4-13)$$

The question now arises as to whether we need to retain all the terms in the approximation (4-13), or whether only the first term might suffice. The answer to this question depends on which of the several occurrences of r_{01} is being approximated. For the r_{01}^2 appearing in the denominator of Eq. (4-9), the error introduced by dropping **all** terms but z is generally acceptably small. However, for the r_{01} appearing in the exponent,

errors are much more critical. First, they are multiplied by a very large number k , a typical value for which might be greater than 10^7 in the visible region of the spectrum (e.g. $\lambda = 5 \times 10^{-7}$ meters). Second, phase changes of as little as a fraction of a radian can change the value of the exponential significantly. For this reason we retain both terms of the binomial approximation in the exponent. The resulting expression for the field at (x, y) therefore becomes

$$U(x, y) = \frac{e^{jkz}}{j\lambda z} \iint_{-\infty}^{\infty} U(\xi, \eta) \exp\left\{j\frac{k}{2z}[(x - \xi)^2 + (y - \eta)^2]\right\} d\xi d\eta, \quad (4-14)$$

where we have incorporated the finite limits of the aperture in the definition of $U(\xi, \eta)$, in accord with the usual assumed boundary conditions.

Equation (4-14) is readily seen to be a convolution, expressible in the form

$$U(x, y) = \iint_{-\infty}^{\infty} U(\xi, \eta) h(x - \xi, y - \eta) d\xi d\eta \quad (4-15)$$

where the convolution kernel is

$$h(x, y) = \frac{e^{jkz}}{j\lambda z} \exp\left[j\frac{k}{2z}(x^2 + y^2)\right]. \quad (4-16)$$

We will return to this viewpoint a bit later.

Another form of the result (4-14) is found if the term $\exp\left[j\frac{k}{2z}(x^2 + y^2)\right]$ is factored outside the integral signs, yielding

$$U(x, y) = \frac{e^{jkz}}{j\lambda z} e^{j\frac{k}{2z}(x^2 + y^2)} \iint_{-\infty}^{\infty} \left\{ U(\xi, \eta) e^{j\frac{k}{2z}(\xi^2 + \eta^2)} \right\} e^{-j\frac{2\pi}{\lambda z}(x\xi + y\eta)} d\xi d\eta, \quad (4-17)$$

which we recognize (aside from multiplicative factors) to be the *Fourier transform* of the product of the complex field just to the right of the aperture and a quadratic phase exponential.

We refer to both forms of the result, (4-14) and (4-17), as *the Fresnel diffraction integral*. When this approximation is valid, the observer is said to be in the region of Fresnel diffraction, or equivalently in the *near field* of the aperture.²

4.2.1 Positive vs. Negative Phases

We have seen that it is common practice when using the Fresnel approximation to replace expressions for spherical waves by quadratic-phase exponentials. The question often arises as to whether the sign of the phase should be positive or negative in a given

²Recently an interesting relation between the Fresnel diffraction formula and an entity known as the "fractional Fourier transform" has been found. The interested reader can consult Ref. [225] and the references contained therein.

expression. This question is not only pertinent to quadratic-phase exponentials, but also arises when considering exact expressions for spherical waves and when considering plane waves propagating at an angle with respect to the optical axis. We now present the reader with a methodology that will help determine the proper sign of the exponent in all of these cases.

The critical fact to keep in mind is that we have chosen our phasors to rotate in the clockwise direction, i.e. their time dependence is of the form $\exp(-j2\pi\nu t)$. For this reason, if we move in space in such a way as to intercept portions of a wavefield that were emitted later in time, the phasor will have advanced in the clockwise direction, and therefore the phase must become more negative. On the other hand, if we move in space to intercept portions of a wavefield that were emitted earlier in time, the phasor will not have had time to rotate as far in the clockwise direction, and therefore the phase must become more positive.

If we imagine observing a spherical wave that is diverging from a point on the z axis, the observation being in an (x, y) plane that is normal to that axis, then movement away from the origin always results in observation of portions of the wavefront that were emitted earlier in time than that at the origin, since the wave has had to propagate further to reach those points. For that reason the phase must increase in a positive sense as we move away from the origin. Therefore the expressions $\exp(jkr_{01})$ and $\exp[j\frac{k}{2z}(x^2 + y^2)]$ (for positive z) represent a diverging spherical wave and a quadratic-phase approximation to such a wave, respectively. By the same token, $\exp(-jkr_{01})$ and $\exp[-j\frac{k}{2z}(x^2 + y^2)]$ represent a converging spherical wave, again assuming that z is positive. Clearly, if z is a negative number, then the interpretation must be reversed, since a negative sign is hidden in z .

Similar reasoning applies to the expressions for plane waves traveling at an angle with respect to the optical axis. Thus for positive a , the expression $\exp(j2\pi a y)$ represents a plane wave with a wave vector in the (y, z) plane. But does the wave vector point with a positive angle with respect to the z axis or with a negative angle, keeping in mind that a positive angle is one that has rotated counterclockwise with respect to the z axis? If we move in the positive y direction, the argument of the exponential increases in a positive sense, and therefore we are moving to a portion of the wave that was emitted earlier in time. This can only be true if the wave vector points with a positive angle with respect to the z axis, as illustrated in Fig. 4.2.

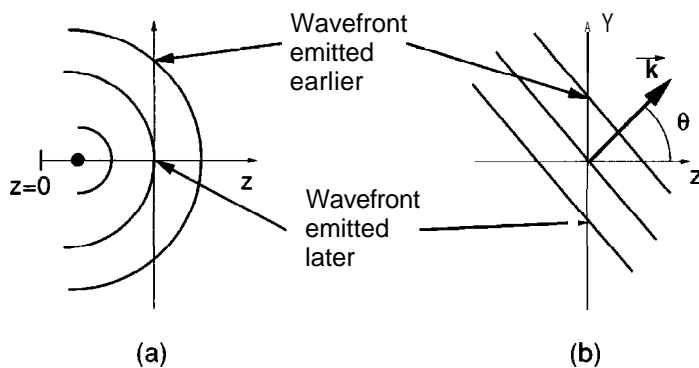


FIGURE 4.2
Determining the sign of the phases of exponential representations of (a) spherical waves and (b) plane waves.

4.2.2 Accuracy of the Fresnel Approximation

Considering the approximation in the exponent, which is the most critical approximation, it can be seen that the *spherical* secondary wavelets of the Huygens-Fresnel principle have been replaced by wavelets with parabolic wavefronts. The accuracy of this approximation is determined by the errors induced when terms higher than first order (linear in b) are dropped in the binomial expansion (4-11). A sufficient condition for accuracy would be that the maximum phase change induced by dropping the $b^2/8$ term be much less than 1 radian. This condition will be met if the distance z satisfies

$$z^3 \gg \frac{\pi}{4\lambda} [(x - \xi)^2 + (y - \eta)^2]_{\max}^2. \quad (4-18)$$

For a circular aperture of size 1 cm, a circular observation region of size 1 cm, and a wavelength of $0.5 \mu\text{m}$, this condition would indicate that the distance z must be $\gg 25$ cm for accuracy. However, as the next comment will explain, this sufficient condition is overly stringent, and accuracy can be expected for much shorter distances.

For the Fresnel approximation to yield accurate results, it is not necessary that the higher-order terms of the expansion be small, only that they not change the value of the Fresnel diffraction integral significantly. Considering the convolution form of the result, Eq. (4-14), if the major contribution to the integral comes from points (ξ, η) for which $\xi \approx x$ and $\eta \approx y$, then the particular values of the higher-order terms of the expansion are unimportant.

To investigate this point more completely, expand the quadratic-phase exponential of Eq. (4-16) into its real and imaginary parts,

$$\frac{1}{j\lambda z} \exp\left[j \frac{\pi}{\lambda z} (x^2 + y^2)\right] = \frac{1}{j\lambda z} \left\{ \cos\left[\frac{\pi}{\lambda z} (x^2 + y^2)\right] + j \sin\left[\frac{\pi}{\lambda z} (x^2 + y^2)\right] \right\}, \quad (4-19)$$

where we have dropped the unit magnitude phasor e^{jkz} simply by redefining the phase reference, and we have replaced k by $2\pi/\lambda$. The volume under this function can readily be shown to be unity (Prob. 4-1). Figure 4.3 shows plots of one-dimensional quadratic-phase cosine and sine functions $\cos(\pi x^2)$ and $\sin(\pi x^2)$. Each of these functions has area $1/\sqrt{2}$. Using this fact it can be shown that all of the unit area under the two-dimensional quadratic-phase exponential is contributed by the two-dimensional sinusoidal term.

Figure 4.4 shows the magnitude of the integral of a quadratic-phase exponential function,

$$\left| \int_{-X}^X \exp(j\pi x^2) dx \right| = \left| \sqrt{2}C(\sqrt{2}X) + j \sqrt{2}S(\sqrt{2}X) \right|$$

which has also been expressed in terms of the Fresnel integrals $C(z)$ and $S(z)$ mentioned in Section 2.2. As can be seen from the figure, the integral grows toward its asymptotic value of unity with increasing X . Note in particular that the integral first reaches unity when $X = 0.5$, and then oscillates about that value with diminishing fluctuations. We conclude that, to a reasonable approximation, the major contributions to a convolution

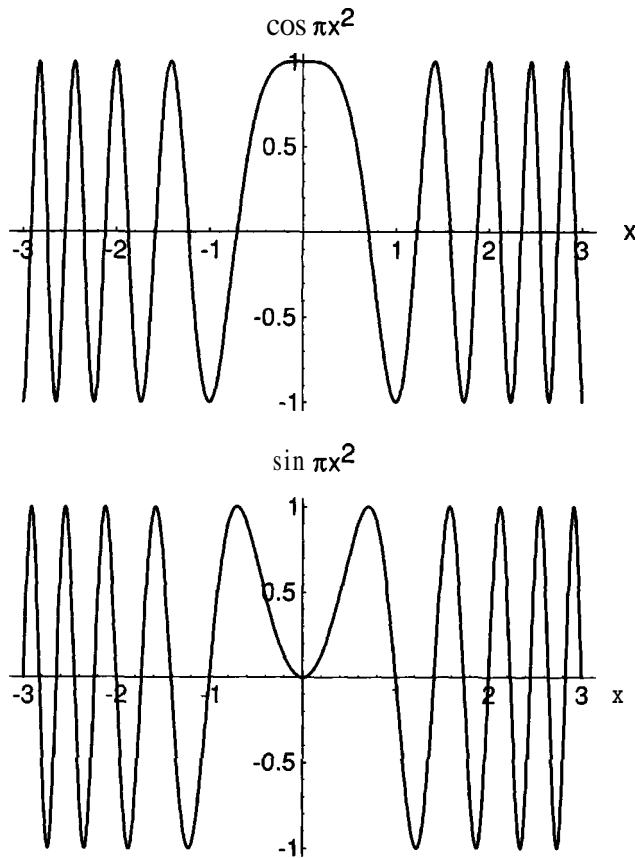


FIGURE 4.3
Quadratic phase cosine and sine functions.

of this function with a second function that is smooth and slowly varying will come from the range $-2 < X < 2$, due to the fact that outside this range the rapid oscillations of the integrand do not yield a significant addition to the total area.

For the *scaled* quadratic-phase exponential of Eqs. (4-14) and (4-16), the corresponding conclusion is that the majority of the contribution to the convolution integral comes from a square in the (ξ, η) plane, with width $4\sqrt{\lambda z}$ and centered on the point $(\xi = x, \eta = y)$. This square grows in size as the distance z behind the aperture increases. In effect, when this square lies entirely within the open portion of the aperture, the field observed at distance z is, to a good approximation, what it would be if the aperture were not present. When the square lies entirely behind the obstruction of the aperture, then the observation point lies in a region that is, to a good approximation, dark

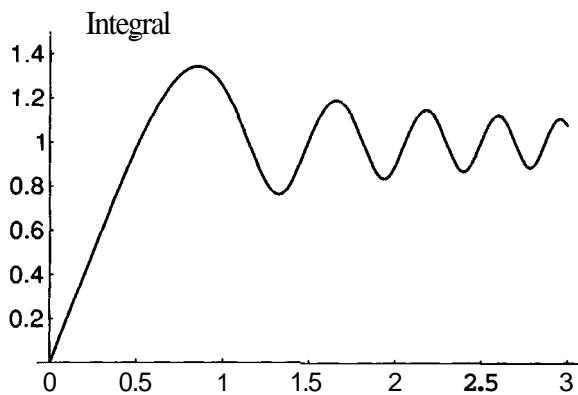


FIGURE 4.4
Magnitude of the integral of the quadratic-phase exponential function.

due to the shadow of the aperture. When the square bridges the open and obstructed parts of the aperture, then the observed field is in the transition region between light and dark. The detailed structure within these regions may be complex, but the general conclusions above are correct. Figure 4.5 illustrates the various regions mentioned. For the case of a one-dimensional rectangular slit, the boundaries between the light region and the transition region, and between the dark region and the transition region, can be shown to be parabolas (see Prob. 4-5).

Note that if the amplitude transmittance **and/or** the illumination of the diffracting aperture is not a relatively smooth and slowly varying function, the above conclusions may not hold. For example, if the amplitude of the field transmitted by the aperture has a high-spatial-frequency sinusoidal component, that component may interact with the high frequencies of the quadratic-phase exponential kernel to produce a nonzero contribution from a location other than the square mentioned above. Thus the restriction of attention to the square of width $4\sqrt{\lambda z}$ must be used with some caution. However, the idea is valid when the diffracting apertures do not contain fine structure and when they are illuminated by uniform plane waves.

If the distance z is allowed to approach zero, *i.e.* the observation point approaches the diffracting aperture, then the two-dimensional quadratic-phase function behaves in the limit like a delta function, producing a field $U(x, y)$ that is identical to the aperture field $U(\xi, \eta)$ in the aperture. In such a case, the predictions of geometrical optics are valid, for such a treatment would predict that the field observed behind the aperture is simply a geometrical projection of the aperture fields onto the plane of observation.

Our discussion above is closely related to the principle of stationary phase, a method for finding the asymptotic values of certain integrals. A good discussion of this method can be found in Appendix III of Ref. [28]. For other examinations of the accuracy of the Fresnel approximation, see Chapter 9 of Ref. [227] and also Ref. [271]. The general conclusions of all of these analyses are similar; namely, the accuracy of the Fresnel approximation is extremely good to distances that are very close to the aperture.

4.2.3 The Fresnel Approximation and the Angular Spectrum

It is of some interest to understand the implications of the Fresnel approximations from the point-of-view of the angular spectrum method of analysis. Such understanding can

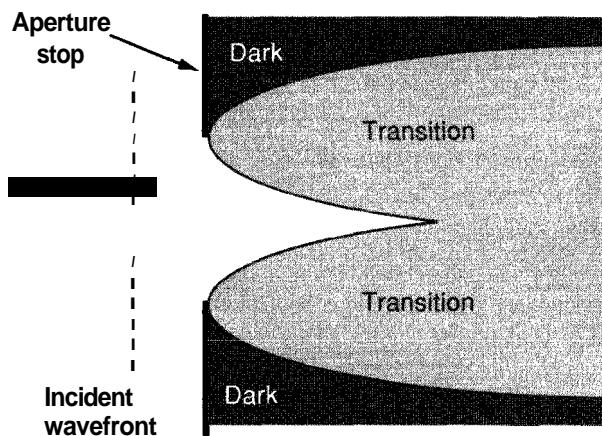


FIGURE 4.5
Light, dark, and transition regions
behind a rectangular slit aperture.

be developed by beginning with Eq. (3-74), which expresses the transfer function of propagation through free space,

$$H(f_X, f_Y) = \begin{cases} \exp \left[j2\pi \frac{z}{\lambda} \sqrt{1 - (\lambda f_X)^2 - (\lambda f_Y)^2} \right] & \sqrt{f_X^2 + f_Y^2} < \frac{1}{\lambda} \\ 0 & \text{otherwise.} \end{cases} \quad (4-20)$$

This result, which is valid subject only to the scalar approximation, can now be compared with the transfer function predicted by the results of the Fresnel analysis. Fourier transforming the Fresnel diffraction impulse response (4-16), we find (with the help of Table 2.1) a transfer function valid for Fresnel diffraction,

$$\begin{aligned} H(f_X, f_Y) &= \mathcal{F} \left\{ \frac{e^{jkz}}{j\lambda z} \exp \left[j \frac{\pi}{\lambda z} (x^2 + y^2) \right] \right\} \\ &= e^{jkz} \exp \left[-j\pi\lambda z (f_X^2 + f_Y^2) \right]. \end{aligned} \quad (4-21)$$

Thus in the Fresnel approximation, the general spatial phase dispersion representing propagation is reduced to a **quadratic** phase dispersion. The factor e^{jkz} on the right of this equation represents a constant phase delay suffered by all plane-wave components traveling between two parallel planes separated by normal distance z . The second term represents the different phase delays suffered by plane-wave components traveling in different directions.

The expression (4-21) is clearly an approximation to the more general transfer function (4-20). We can obtain the approximate result from the general result by applying a binomial expansion to the exponent of (4-20),

$$\sqrt{1 - (\lambda f_X)^2 - (\lambda f_Y)^2} \approx 1 - \frac{(\lambda f_X)^2}{2} - \frac{(\lambda f_Y)^2}{2}, \quad (4-22)$$

which is valid provided $(\lambda f_X) \ll 1$ and $(\lambda f_Y) \ll 1$. Such restrictions on f_X and f_Y are simply restrictions to **small angles**. So we see that, from the perspective of the angular spectrum, the Fresnel approximation is accurate provided only small angles of diffraction are involved. It is for this reason that we often say that the Fresnel approximations and the **paraxial** approximation are equivalent.

4.2.4 Fresnel Diffraction Between Confocal Spherical Surfaces

Until now, attention has been focused on diffraction between two **planes**. An alternative geometry, of more theoretical than practical interest but nonetheless quite instructive, is diffraction between two confocal spherical surfaces (see, for example, [24], [25]). As shown in Fig. 4.6, two spheres are said to be confocal if the center of each lies on the surface of the other. In our case, the two spheres are tangent to the planes previously used, with the points of tangency being the points where the z axis pierces those planes. The distance r_{01} in our previous diffraction analysis is now the distance between the two spherical caps shown.

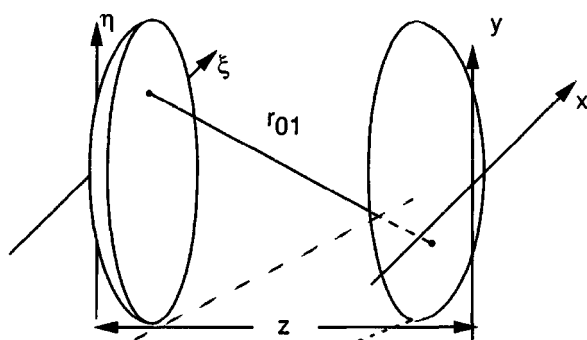


FIGURE 4.6
Confocal spherical surfaces.

A proper analysis would write equations for the left-hand spherical surface and for the right-hand spherical surface, and then use those equations to find the distance r_{01} between the two spherical caps. In the process it would be helpful to simplify certain square roots by using the first two terms of their binomial expansions (i.e. to make *paraxial* approximations to the spherical surfaces). The result of such an analysis is the following simple expression for r_{01} , valid if the extent of the spherical caps about the z -axis is small compared with their radii:

$$r_{01} \approx z - x\xi/z - y\eta/z.$$

The Fresnel diffraction equation now becomes

$$U(x, y) = \frac{e^{jkz}}{j\lambda z} \iint_{-\infty}^{\infty} U(\xi, \eta) e^{-j\frac{2\pi}{\lambda z}(x\xi + y\eta)} d\xi d\eta, \quad (4-23)$$

which, aside from constant multipliers and scale factors, expresses the field observed on the right-hand spherical cap as the *Fourier transform* of the field on the left-hand spherical cap.

Comparison of this result with the previous Fourier-transform version of the Fresnel diffraction integral, Eq. (4-17), shows that the quadratic-phase factors in (x, y) and (ξ, η) have been eliminated by moving from the two planes to the two spherical caps. The two quadratic phase factors in the earlier expression are in fact simply *paraxial* representations of spherical phase surfaces, and it is therefore reasonable that moving to the spheres has eliminated them.

One subtle point worth mention is that, when we analyze diffraction between two spherical caps, it is not really valid to use the Rayleigh-Sommerfeld result as the basis for the calculation, for that result was explicitly valid only for diffraction by a planar aperture. However, the Kirchhoff analysis remains valid, and its predictions are the same as those of the Rayleigh-Sommerfeld approach provided *paraxial* conditions hold.

4.3 THE FRAUNHOFER APPROXIMATION

Before presenting several examples of diffraction pattern calculations, we consider another more stringent approximation which, when valid, greatly simplifies the calculations. It was seen in Eq. (4-17) that, in the region of Fresnel diffraction, the observed

field strength $U(x, y)$ can be found from a Fourier transform of the product of the aperture distribution $U(\xi, \eta)$ and a quadratic phase function $\exp[j(k/2z)(\xi^2 + \eta^2)]$. If in addition to the Fresnel approximation the stronger (Fraunhofer) approximation

$$z \gg \frac{k(\xi^2 + \eta^2)_{\max}}{2} \quad (4-24)$$

is satisfied, then the quadratic phase factor under the integral sign in Eq. (4-17) is approximately unity over the entire aperture, and the observed field strength can be found (up to a multiplicative phase factor in (x, y)) directly from a Fourier transform of the aperture distribution itself. Thus in the region of Fraunhofer *diffraction* (or equivalently, in the *far field*),

$$U(x, y) = \frac{e^{jkz} e^{j\frac{k}{2z}(x^2+y^2)}}{j\lambda z} \iint_{-\infty}^{\infty} U(\xi, \eta) \exp\left[-j\frac{2\pi}{\lambda z}(x\xi + y\eta)\right] d\xi d\eta. \quad (4-25)$$

Aside from multiplicative phase factors preceding the integral, this expression is simply the Fourier transform of the aperture distribution, evaluated at frequencies

$$\begin{aligned} f_x &= x/\lambda z \\ f_y &= y/\lambda z. \end{aligned} \quad (4-26)$$

At optical frequencies, the conditions required for validity of the Fraunhofer approximation can be severe ones. For example, at a wavelength of $0.6 \mu\text{m}$ (red light) and an aperture width of 2.5 cm (1 inch), the observation distance z must satisfy

$$z \gg 1,600 \text{ meters.}$$

An alternative, less stringent condition, known as the "antenna designer's formula", states that for an aperture of linear dimension D , the Fraunhofer approximation will be valid provided

$$z > \frac{2D^2}{\lambda} \quad (4-27)$$

where the inequality is now $>$ rather than \gg . However, for this example the distance z is still required to be larger than 2,000 meters. Nonetheless, the required conditions are met in a number of important problems. In addition, Fraunhofer diffraction patterns can be observed at distances much closer than implied by Eq. (4-24) provided the aperture is illuminated by a spherical wave converging toward the observer (see Prob. 4-16), or if a positive lens is properly situated between the observer and the aperture (see Chapter 5).

Finally, it should be noted that, at first glance, there exists no transfer function that can be associated with Fraunhofer diffraction, for the approximation (4-24) has destroyed the space invariance of the diffraction equation (cf. Prob. 2-10). The secondary wavelets with parabolic surfaces (as implied by the Fresnel approximation) no longer shift laterally in the (x, y) plane with the particular (ξ, η) point under consideration. Rather, when the location of the secondary source shifts, the corresponding quadratic surface tilts in the (x, y) plane by an amount that depends on the location of the

secondary source. Nonetheless, it should not be forgotten that since Fraunhofer diffraction is only a special case of Fresnel diffraction, the transfer function (4-21) remains valid throughout both the Fresnel and the Fraunhofer regimes. That is, it is always possible to calculate diffracted fields in the Fraunhofer region by retaining the full accuracy of the Fresnel approximation.

4.4 EXAMPLES OF FRAUNHOFER DIFFRACTION PATTERNS

We consider next several examples of Fraunhofer diffraction patterns. For additional examples the reader may consult the problems (see Probs. 4-7 through 4-10).

The results of the preceding section can be applied directly to find the complex field distribution across the Fraunhofer diffraction pattern of any given aperture. However, of ultimate interest, for reasons discussed at the beginning of this chapter, is the intensity rather than the complex field strength. The final descriptions of the specific diffraction patterns considered here will therefore be distributions of intensity.

4.4.1 Rectangular Aperture

Consider first a rectangular aperture with an amplitude transmittance given by

$$t_A(\xi, \eta) = \text{rect}\left(\frac{\xi}{2w_X}\right) \text{rect}\left(\frac{\eta}{2w_Y}\right).$$

The constants w_X and w_Y are the half-widths of the aperture in the ξ and η directions. If the aperture is illuminated by a unit-amplitude, normally incident, monochromatic plane wave, then the field distribution across the aperture is equal to the transmittance function t_A . Thus using Eq. (4-25), the Fraunhofer diffraction pattern is seen to be

$$U(x, y) = \frac{e^{jkz} e^{j\frac{k}{2z}(x^2+y^2)}}{jhz} \mathcal{F}\{U(\xi, \eta)\} \Bigg|_{\substack{f_X = x/\lambda z \\ f_Y = y/\lambda z}}.$$

Noting that $\mathcal{F}\{U(\xi, \eta)\} = A \text{sinc}(2w_X f_X) \text{sinc}(2w_Y f_Y)$, where A is the area of the aperture ($A = 4w_X w_Y$), we find

$$U(x, y) = \frac{e^{jkz} e^{j\frac{k}{2z}(x^2+y^2)}}{jhz} A \text{sinc}\left(\frac{2w_X x}{\lambda z}\right) \text{sinc}\left(\frac{2w_Y y}{\lambda z}\right),$$

and

$$I(x, y) = \frac{A^2}{\lambda^2 z^2} \text{sinc}^2\left(\frac{2w_X x}{\lambda z}\right) \text{sinc}^2\left(\frac{2w_Y y}{\lambda z}\right). \quad (4-28)$$

Figure 4.7 shows a cross section of the Fraunhofer intensity pattern along the x axis. Note that the width of the main lobe (i.e. the distance between the first two zeros) is

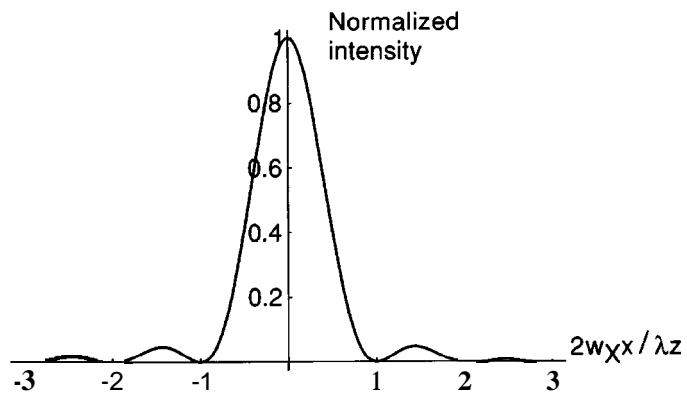


FIGURE 4.7
Cross section of the Fraunhofer diffraction pattern of a rectangular aperture.

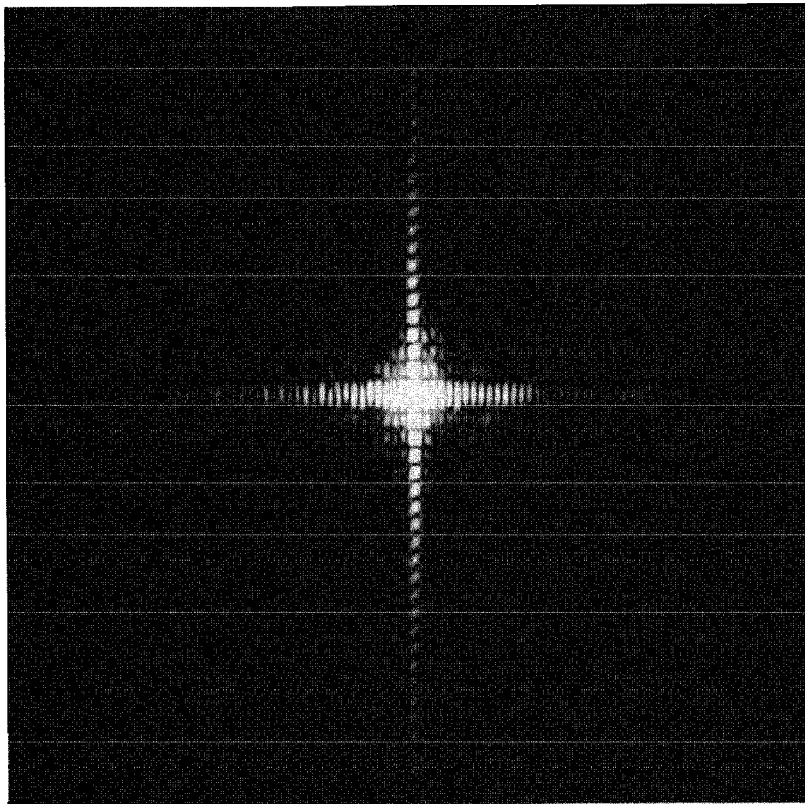


FIGURE 4.8
The Fraunhofer diffraction pattern of a rectangular aperture ($w_x/w_y = 2$).

$$\Delta x = \frac{\lambda z}{w_X}. \quad (4-29)$$

Figure 4.8 shows a photograph of the diffraction pattern produced by a rectangular aperture with a width ratio of $w_X/w_Y = 2$.

4.4.2 Circular Aperture

Consider a diffracting aperture that is circular rather than rectangular, and let the radius of the aperture be w . Thus if q is a radius coordinate in the plane of the aperture, then

$$t_A(q) = \text{circ}\left(\frac{q}{w}\right).$$

The circular symmetry of the problem suggests that the Fourier transform of Eq. (4-25) be rewritten as a Fourier-Bessel transform. Thus if r is the radius coordinate in the observation plane, we have

$$U(r) = \left. j\lambda z \exp\left(j\frac{kr^2}{2z}\right) \mathcal{B}\{U(q)\} \right|_{p=r/\lambda z}, \quad (4-30)$$

where $q = \sqrt{\xi^2 + \eta^2}$ represents radius in the aperture plane, and $p = \sqrt{f_X^2 + f_Y^2}$ represents radius in the spatial frequency domain. For unit-amplitude, normally incident plane-wave illumination, the field transmitted by the aperture is equal to the amplitude transmittance; in addition,

$$\mathcal{B}\left\{\text{circ}\left(\frac{q}{w}\right)\right\} = A \frac{J_1(2\pi w p)}{\pi w p},$$

where $A = \pi w^2$. The amplitude distribution in the Fraunhofer diffraction pattern is seen to be

$$U(r) = e^{jkz} e^{j\frac{kr^2}{2z}} \frac{A}{j\lambda z} \left[2 \frac{J_1(kwr/z)}{kwr/z} \right],$$

and the intensity distribution can be written

$$I(r) = \left(\frac{A}{\lambda z}\right)^2 \left[2 \frac{J_1(kwr/z)}{kwr/z} \right]^2 \quad (4-31)$$

This intensity distribution is referred to as the **Airy pattern**, after **G.B. Airy** who first derived it. Table 4.1 shows the values of the Airy pattern at successive maxima and minima, from which it can be seen that the width of the central lobe, measured along the x or y axis, is given by

$$d = 1.22 \frac{\lambda z}{w}. \quad (4-32)$$

TABLE 4.1
Locations of maxima and minima of the
Airy pattern.

x	$\left[2\frac{J_1(\pi x)}{\pi x}\right]^2$	max, min
0	1	max
1.220	0	min
1.635	0.0175	max
2.233	0	min
2.679	0.0042	max
3.238	0	min
3.699	0.0016	max

Figure 4.9 shows a cross section of the Airy pattern, while Fig. 4.10 is a photograph of the Fraunhofer diffraction pattern of a circular aperture.

4.4.3 Thin Sinusoidal Amplitude Grating

In the previous examples, diffraction was assumed to be caused by apertures in infinite opaque screens. In practice, diffracting objects can be far more complex. In accord with our earlier definition (3-68), the amplitude transmittance $t_A(\xi, \eta)$ of a screen is defined as the ratio of the complex field amplitude immediately behind the screen to the complex amplitude incident on the screen. Until now, our examples have involved only transmittance functions of the form

$$t_A(\xi, \eta) = \begin{cases} 1 & \text{in the aperture} \\ 0 & \text{outside the aperture.} \end{cases}$$

It is possible, however, to introduce a prescribed amplitude transmittance function within a given aperture. Spatial attenuation can be introduced with, for example, an absorbing photographic transparency, thus allowing real values of t_A between zero and unity to be realized. Spatial patterns of phase shift can be introduced by means of transparent plates of varying thickness, thus extending the realizable values of t_A to all points within or on the unit circle in the complex plane.

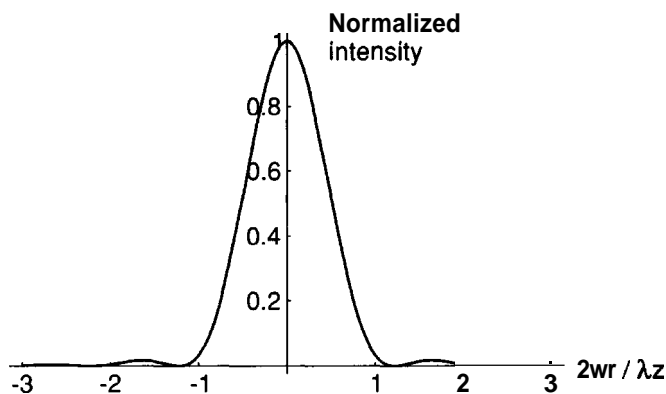


FIGURE 4.9
Cross section of the Fraunhofer
diffraction pattern of a circular
aperture.

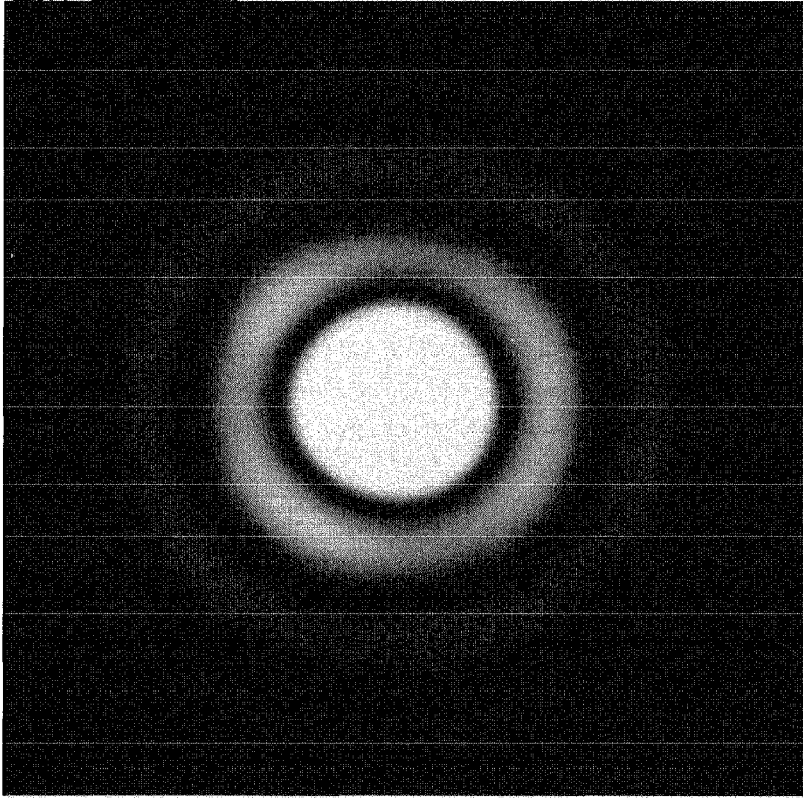


FIGURE 4.10
Fraunhofer diffraction pattern of a circular aperture.

As an example of this more general type of diffracting screen, consider a thin sinusoidal amplitude grating defined by the amplitude transmittance function

$$t_A(\xi, \eta) = \left[\frac{1}{2} + \frac{m}{2} \cos(2\pi f_0 \xi) \right] \text{rect}\left(\frac{\xi}{2w}\right) \text{rect}\left(\frac{\eta}{2w}\right) \quad (4-33)$$

where for simplicity we have assumed that the grating structure is bounded by a square aperture of width $2w$. The parameter m represents the peak-to-peak change of amplitude transmittance across the screen, and f_0 is the spatial frequency of the grating. The term thin in this context means that the structure can indeed be represented by a simple amplitude transmittance. Structures that are not sufficiently thin can not be so represented, a point we shall return to in a later chapter. Figure 4.11 shows a cross section of the grating amplitude transmittance function.

If the screen is normally illuminated by a unit-amplitude plane wave, the field distribution across the aperture is equal simply to t_A . To find the Fraunhofer diffraction pattern, we first Fourier transform

$$\begin{aligned} \mathcal{F}\left\{\frac{1}{2} + \frac{m}{2} \cos(2\pi f_0 \xi)\right\} &= \frac{1}{2} \delta(f_x, f_y) \\ &+ \frac{m}{4} \delta(f_x + f_0, f_y) + \frac{m}{4} \delta(f_x - f_0, f_y) \end{aligned} \quad (4-34)$$

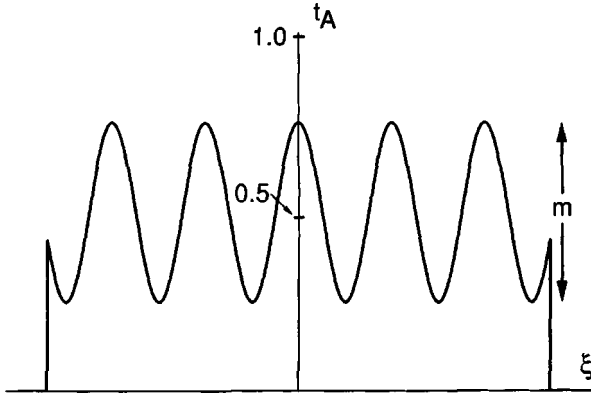


FIGURE 4.11
Amplitude transmittance function of the sinusoidal amplitude grating.

and

$$\mathcal{F}\left\{\text{rect}\left(\frac{\xi}{2w}\right)\text{rect}\left(\frac{\eta}{2w}\right)\right\} = A \text{sinc}(2w f_X) \text{sinc}(2w f_Y),$$

the convolution theorem can be used to write

$$\begin{aligned} \mathcal{F}\{U(\xi, \eta)\} = \frac{A}{2} \text{sinc}(2w f_Y) \left\{ \text{sinc}(2w f_X) + \frac{m}{2} \text{sinc}[2w(f_X + f_0)] \right. \\ \left. + \frac{m}{2} \text{sinc}[2w(f_X - f_0)] \right\}, \end{aligned}$$

where A signifies the area of the aperture bounding the grating. The Fraunhofer diffraction pattern can now be written

$$\begin{aligned} U(x, y) = \frac{A}{j2\lambda z} e^{jkz} e^{j\frac{k}{2z}(x^2+y^2)} \text{sinc}\left(\frac{2wy}{\lambda z}\right) \left\{ \text{sinc}\left(\frac{2wx}{\lambda z}\right) \right. \\ \left. + \frac{m}{2} \text{sinc}\left[\frac{2w}{\lambda z}(x + f_0\lambda z)\right] + \frac{m}{2} \text{sinc}\left[\frac{2w}{\lambda z}(x - f_0\lambda z)\right] \right\}. \quad (4-35) \end{aligned}$$

Finally, the corresponding intensity distribution is found by taking the squared magnitude of Eq. (4-35). Note that if there are many grating periods within the aperture, then $f_0 \gg 1/w$, and there will be negligible overlap of the three sinc functions, allowing the intensity to be calculated as the sum of the squared magnitudes of the three terms in (4-35). The intensity is then given by

$$\begin{aligned} I(x, y) \approx \left[\frac{A}{2\lambda z} \right]^2 \text{sinc}^2\left(\frac{2wy}{\lambda z}\right) \left\{ \text{sinc}^2\left(\frac{2wx}{\lambda z}\right) \right. \\ \left. + \frac{m^2}{4} \text{sinc}^2\left[\frac{2w}{\lambda z}(x + f_0\lambda z)\right] + \frac{m^2}{4} \text{sinc}^2\left[\frac{2w}{\lambda z}(x - f_0\lambda z)\right] \right\}. \quad (4-36) \end{aligned}$$

This intensity pattern is illustrated in Fig. 4.12. Note that some of the incident light is absorbed by the grating, and in addition the sinusoidal transmittance variation across the aperture has deflected some of the energy out of the central diffraction pattern into two additional side patterns. The central diffraction pattern is called the *zero order* of

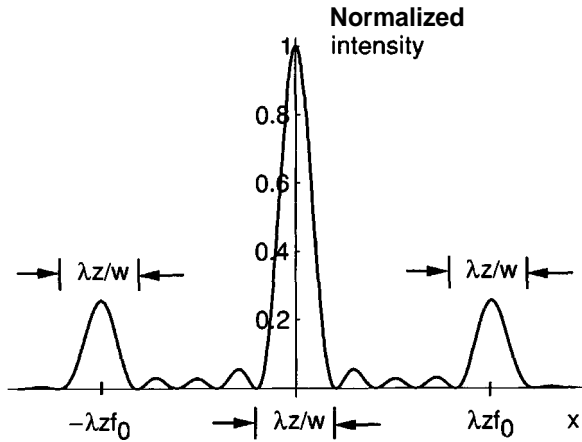


FIGURE 4.12
Fraunhofer diffraction pattern for a thin sinusoidal amplitude grating.

the Fraunhofer pattern, while the two side patterns are called *the first orders*. The spatial separation of the first orders from the zero order is $f_0\lambda z$, while the width of the main lobe of all orders is $\lambda z/w$.

Another quantity of some practical interest in both holography and optical information processing is the *diffraction efficiency* of the grating. The diffraction efficiency is defined as the fraction of the incident optical power that appears in a single diffraction order (usually the +1 order) of the grating. The diffraction efficiency for the grating of interest can be deduced from Eq. (4-34). The fraction of power appearing in each diffraction order can be found by squaring the coefficients of the delta functions in this representation, for it is the delta functions that **determine** the power in each order, not the **sinc** functions that simply spread these impulses. From this equation we conclude that the diffraction efficiencies η_0 , η_{+1} , η_{-1} associated with the three diffraction orders are given by

$$\begin{aligned}\eta_0 &= 0.25 \\ \eta_{+1} &= m^2/16 \\ \eta_{-1} &= m^2/16.\end{aligned}\tag{4-37}$$

Thus a single first diffraction order carries at most $1/16 = 6.25\%$ of the incident power, a rather small fraction. If the efficiencies of the three orders are added up, it will be seen that only $1/4 + m^2/8$ of the total is accounted for. The rest is lost through absorption by the grating.

4.4.4 Thin Sinusoidal Phase Grating

As a final example of Fraunhofer diffraction calculations, consider a thin sinusoidal phase grating defined by the amplitude transmittance function

$$t_A(\xi, \eta) = \exp\left[j\frac{m}{2}\sin(2\pi f_0\xi)\right] \text{rect}\left(\frac{\xi}{2w}\right) \text{rect}\left(\frac{\eta}{2w}\right)\tag{4-38}$$

where, by proper choice of phase reference, we have dropped a factor representing the average phase delay through the grating. The parameter m represents the peak-to-peak excursion of the phase delay.

If the grating is illuminated by a unit-amplitude, normally incident plane wave, then the field distribution immediately behind the screen is given precisely by Eq. (4-38). The analysis is simplified by use of the identity

$$\exp\left[j\frac{m}{2}\sin(2\pi f_0\xi)\right] = \sum_{q=-\infty}^{\infty} J_q\left(\frac{m}{2}\right)\exp(j2\pi q f_0\xi)$$

where J_q is a Bessel function of the first kind, order q . Thus

$$\mathcal{F}\left\{\exp\left[j\frac{m}{2}\sin(2\pi f_0\xi)\right]\right\} = \sum_{q=-\infty}^{\infty} J_q\left(\frac{m}{2}\right)\delta(f_x - qf_0, f_y) \quad (4-39)$$

and

$$\begin{aligned} \mathcal{F}\{U(\xi, \eta)\} &= \mathcal{F}\{t_A(\xi, \eta)\} \\ &= [A \operatorname{sinc}(2wf_x) \operatorname{sinc}(2wf_y)] \otimes \left[\sum_{q=-\infty}^{\infty} J_q\left(\frac{m}{2}\right)\delta(f_x - qf_0, f_y) \right] \\ &= \sum_{q=-\infty}^{\infty} AJ_q\left(\frac{m}{2}\right) \operatorname{sinc}[2w(f_x - qf_0)] \operatorname{sinc}(2wf_y). \end{aligned}$$

Thus the field strength in the Fraunhofer diffraction pattern can be written

$$\begin{aligned} U(x, y) &= \frac{A}{j\lambda z} e^{jkz} e^{j\frac{k}{2z}(x^2+y^2)} \\ &\quad \times \sum_{q=-\infty}^{\infty} J_q\left(\frac{m}{2}\right) \operatorname{sinc}\left[\frac{2w}{\lambda z}(x - qf_0\lambda z)\right] \operatorname{sinc}\left(\frac{2wy}{\lambda z}\right). \end{aligned} \quad (4-40)$$

If we again assume that there are many periods of the grating within the bounding aperture ($f_0 \gg 1/w$), there is negligible overlap of the various diffracted terms, and the corresponding intensity pattern becomes

$$I(x, y) \approx \left(\frac{A}{\lambda z}\right)^2 \sum_{q=-\infty}^{\infty} J_q^2\left(\frac{m}{2}\right) \operatorname{sinc}^2\left[\frac{2w}{\lambda z}(x - qf_0\lambda z)\right] \operatorname{sinc}^2\left(\frac{2wy}{\lambda z}\right). \quad (4-41)$$

The introduction of the sinusoidal phase grating has thus deflected energy out of the zero order into a multitude of higher orders. The peak intensity of the q th order is $[AJ_q(m/2)/\lambda z]^2$, while the displacement of that order from the center of the diffraction pattern is $qf_0\lambda z$. Figure 4.13 shows a cross section of the intensity pattern when the peak-to-peak phase delay m is 8 radians. Note that the strengths of the various orders are symmetric about the zero order.

The diffraction efficiency of the thin sinusoidal phase grating can be found by determining the squared magnitude of the coefficients in Eq. (4-39). Thus the diffraction efficiency of the q th order of this grating is

$$\eta_q = J_q^2(m/2). \quad (4-42)$$

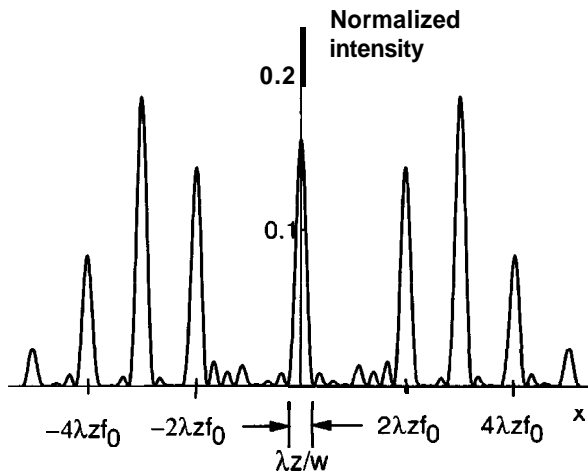


FIGURE 4.13
Fraunhofer diffraction pattern for a thin sinusoidal phase grating. The ± 1 orders have nearly vanished in this example.

Figure 4.14 shows a plot of η_q vs. $m/2$ for various values of q . Note that whenever $m/2$ is a root of J_0 , the central order vanishes entirely! The largest possible diffraction efficiency into one of the ± 1 and -1 diffraction orders is the maximum value of J_1^2 . This maximum is 33.8%, far greater than for the case of a thin sinusoidal amplitude grating. No power is absorbed by this grating, and therefore the sum of the powers appearing in all orders remains constant and equal to the incident power as m is changed.

4.5 EXAMPLES OF FRESNEL DIFFRACTION CALCULATIONS

In a previous section, several different methods for calculating Fresnel diffraction patterns have been introduced. For the beginner, it is difficult to know when one method will be easier than another, and therefore in this section two examples are presented that provide some insight in this regard. The first example, Fresnel diffraction by a square aperture, illustrates the application of the classical approach based on the convolution representation of the diffraction calculation. The second example, Talbot imaging, illustrates a case in which a frequency-domain approach has a large advantage.

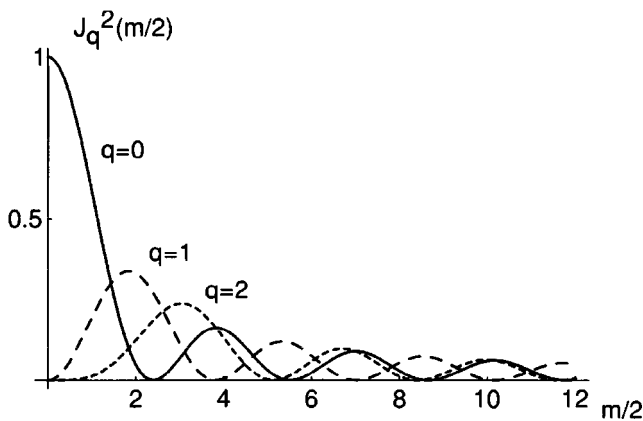


FIGURE 4.14
Diffraction efficiency $J_q^2(m/2)$ vs. $m/2$ for three values of q .

4.5.1 Fresnel Diffraction by a Square Aperture

Suppose that a square aperture of width $2w$ is normally illuminated by a monochromatic plane wave of unit amplitude. The distribution of complex field immediately behind the aperture is

$$U(\xi, \eta) = \text{rect}\left(\frac{\xi}{2w}\right) \text{rect}\left(\frac{\eta}{2w}\right).$$

The convolution form of the Fresnel diffraction equation is most convenient for this problem, yielding

$$U(x, y) = \frac{e^{jkz}}{j\lambda z} \iint_{-w}^w \exp\left\{j\frac{\pi}{\lambda z}[(x - \xi)^2 + (y - \eta)^2]\right\} d\xi d\eta.$$

This expression can be separated into the product of two one-dimensional integrals,

$$U(x, y) = \frac{e^{jkz}}{j} \mathcal{I}(x)\mathcal{I}(y) \quad (4-43)$$

where

$$\mathcal{I}(x) = \frac{1}{\sqrt{\lambda z}} \int_{-w}^w \exp\left[j\frac{\pi}{\lambda z}(\xi - x)^2\right] d\xi$$

$$\mathcal{I}(y) = \frac{1}{\sqrt{\lambda z}} \int_{-w}^w \exp\left[j\frac{\pi}{\lambda z}(\eta - y)^2\right] d\eta.$$

To reduce these integrals to expressions that are related to the Fresnel integrals mentioned on several previous occasions, make the following change of variables:

$$\alpha = \sqrt{\frac{2}{\lambda z}}(\xi - x) \quad \beta = \sqrt{\frac{2}{\lambda z}}(\eta - y),$$

yielding

$$\mathcal{I}(x) = \frac{1}{\sqrt{2}} \int_{\alpha_1}^{\alpha_2} \exp\left(j\frac{\pi}{2}\alpha^2\right) d\alpha$$

$$\mathcal{I}(y) = \frac{1}{\sqrt{2}} \int_{\beta_1}^{\beta_2} \exp\left(j\frac{\pi}{2}\beta^2\right) d\beta,$$

where the limits of integration are

$$\alpha_1 = -\sqrt{\frac{2}{\lambda z}}(w + x) \quad \alpha_2 = \sqrt{\frac{2}{\lambda z}}(w - x)$$

$$\beta_1 = -\sqrt{\frac{2}{\lambda z}}(w + y) \quad \beta_2 = \sqrt{\frac{2}{\lambda z}}(w - y).$$

At this point we define the *Fresnel number*, $N_F = w^2/\lambda z$, and we introduce normalized distance variables in the observation region, $X = x/\sqrt{\lambda z}$ and $Y = y/\sqrt{\lambda z}$, yielding simpler expressions for the limits of integration,

$$\begin{aligned}\alpha_1 &= -\sqrt{2}(\sqrt{N_F} + X) & \alpha_2 &= \sqrt{2}(\sqrt{N_F} - X) \\ \beta_1 &= -\sqrt{2}(\sqrt{N_F} + Y) & \beta_2 &= \sqrt{2}(\sqrt{N_F} - Y).\end{aligned}\tag{4-44}$$

The integrals $\mathcal{I}(x)$ and $\mathcal{I}(y)$ are related to the Fresnel integrals $C(z)$ and $S(z)$ of Sections 2.2 and 4.2. Noting that

$$\int_{\alpha_1}^{\alpha_2} \exp\left(j\frac{\pi}{2}\alpha^2\right) d\alpha = \int_0^{\alpha_2} \exp\left(j\frac{\pi}{2}\alpha^2\right) d\alpha - \int_0^{\alpha_1} \exp\left(j\frac{\pi}{2}\alpha^2\right) d\alpha,$$

we can write

$$\begin{aligned}\mathcal{I}(x) &= \frac{1}{\sqrt{2}} \{ [C(\alpha_2) - C(\alpha_1)] + j[S(\alpha_2) - S(\alpha_1)] \} \\ \mathcal{I}(y) &= \frac{1}{\sqrt{2}} \{ [C(\beta_2) - C(\beta_1)] + j[S(\beta_2) - S(\beta_1)] \}.\end{aligned}\tag{4-45}$$

Finally, substitution of (4-45) in (4-43) yields a complex field distribution

$$\begin{aligned}U(x, y) &= \frac{e^{jkz}}{2j} \{ [C(\alpha_2) - C(\alpha_1)] + j[S(\alpha_2) - S(\alpha_1)] \} \\ &\quad \times \{ [C(\beta_2) - C(\beta_1)] + j[S(\beta_2) - S(\beta_1)] \}.\end{aligned}\tag{4-46}$$

Now recall from Section 4.1 that the measurable physical quantity is the intensity of the wavefield, $I(x, y) = |U(x, y)|^2$, which in this case is given by

$$\begin{aligned}I(x, y) &= \frac{1}{4} \{ [C(\alpha_2) - C(\alpha_1)]^2 + [S(\alpha_2) - S(\alpha_1)]^2 \} \\ &\quad \times \{ [C(\beta_2) - C(\beta_1)]^2 + [S(\beta_2) - S(\beta_1)]^2 \}.\end{aligned}\tag{4-47}$$

The Fresnel integrals are tabulated functions and are available in many mathematical computer programs (e.g. see Ref. [302], p. 576).³ It is therefore a straightforward matter to calculate the above intensity distribution. Note that, for fixed w and A , as z increases the Fresnel number N_F decreases and the normalization increasingly enlarges the true physical distance represented by a fixed distance on the x axis. Figure 4.15 shows a series of graphs of the normalized intensity distribution along the x axis ($y = 0$) for various normalized distances from the aperture, as represented by different Fresnel numbers.

³In the past it has been customary to introduce a graphical aid known as "Cornu's spiral" as a tool for estimating values of Fresnel integrals. Modern computer software packages that contain the Fresnel integrals have made this graphical aid largely obsolete, so we have omitted it here.

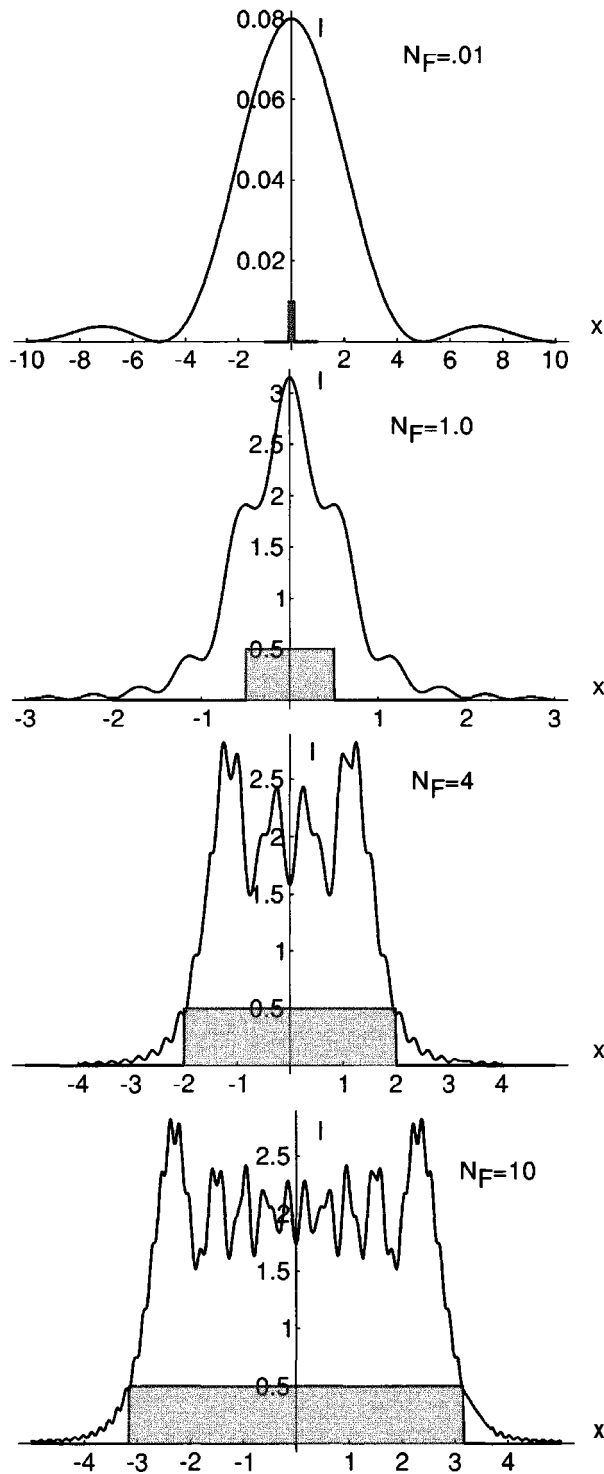


FIGURE 4.15
Fresnel diffraction patterns at different distances from a square aperture. Distance increases as the Fresnel number N_F shrinks. The size of the original rectangular aperture is indicated by the shaded boxes.

Attention is called to the fact that, as the observation plane approaches the plane of the aperture (N_F becomes large), the Fresnel kernel approaches the product of a delta function and a factor e^{jkz} , and the shape of the diffraction pattern approaches the shape of the aperture itself. In fact, the limit of this process is the geometrical optics prediction of the complex field,

$$U(x, y, z) = e^{jkz}U(x, y, 0) = e^{jkz} \text{rect}\left(\frac{x}{w}\right) \text{rect}\left(\frac{y}{2w}\right)$$

where, to avoid confusion, we have explicitly included the z coordinate in the argument of the complex field U .

Note also that, as the distance z becomes large (N_F grows small), the diffraction pattern becomes much wider than the size of the aperture, and comparatively smooth in its structure. In this limit the diffraction pattern is approaching the Fraunhofer limit discussed earlier.

4.5.2 Fresnel Diffraction by a Sinusoidal Amplitude Grating – Talbot Images

Our final example of a diffraction calculation considers again the case of a thin sinusoidal amplitude grating, but this time within the region of Fresnel diffraction rather than Fraunhofer diffraction. For simplicity we neglect the finite extent of the grating and concentrate on the effects of diffraction and propagation on the periodic structure of the fields transmitted by the grating. In effect, we are limiting attention to the central region of the Fresnel diffraction pattern associated with any bounding aperture, between the two transition regions illustrated in Fig. 4.5.

The geometry is illustrated in Fig. 4.16. The grating is modeled as a transmitting structure with amplitude transmittance

$$t_A(\xi, \eta) = \frac{1}{2} [1 + m \cos(2\pi\xi/L)]$$

with period L and with the grating lines running parallel to the η axis. The field and intensity will be calculated some distance z to the right of the grating. The structure is assumed to be illuminated by a unit-amplitude, normally incident plane wave, so the field immediately behind the grating is equal to the amplitude transmittance written above.

There are several possible approaches to calculating the fields behind the grating. We could use the convolution form of the Fresnel diffraction equation, i.e. Eq. (4-14), or the Fourier transform form of Eq. (4-17). Alternatively, we could use the transfer function approach represented by Eq. (4-21), and reproduced here as

$$H(f_X, f_Y) = \exp\{-j\pi\lambda z(f_X^2 + f_Y^2)\}, \quad (4-48)$$

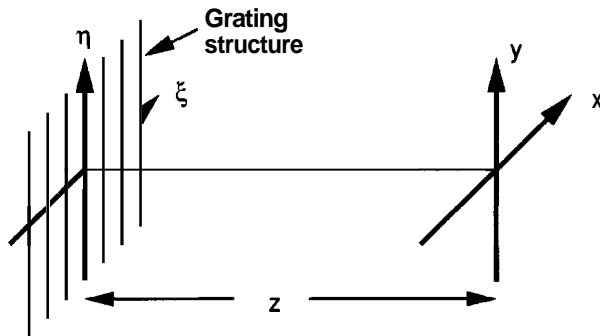


FIGURE 4.16
Geometry for diffraction calculation.

where we have omitted a constant term $\exp(jkz)$. In this problem, and indeed in any problem that deals with a purely periodic structure, the transfer function approach will yield the simplest calculations, and we adopt that approach here.

The solution begins by first finding the spatial frequency spectrum of the field transmitted by the structure. To that end we Fourier transform the amplitude transmittance above, yielding

$$\mathcal{F}\{t_A(\xi, \eta)\} = \frac{1}{2} \delta(f_X, f_Y) + \frac{m}{4} \delta\left(f_X - \frac{1}{L}, f_Y\right) + \frac{m}{4} \delta\left(f_X + \frac{1}{L}, f_Y\right). \quad (4-49)$$

Now the above transfer function has value unity at the origin, and when evaluated at frequencies $(f_X, f_Y) = (\pm \frac{1}{L}, 0)$ yields

$$H\left(\pm \frac{1}{L}, 0\right) = \exp\left\{-j \frac{\pi \lambda z}{L^2}\right\} \quad (4-50)$$

Thus after propagation over distance z behind the grating, the Fourier transform of the field becomes

$$\mathcal{F}\{U(x, y)\} = \frac{1}{2} \delta(f_X, f_Y) + \frac{m}{4} e^{-j \frac{\pi \lambda z}{L^2}} \delta\left(f_X - \frac{1}{L}, f_Y\right) + \frac{m}{4} e^{-j \frac{\pi \lambda z}{L^2}} \delta\left(f_X + \frac{1}{L}, f_Y\right).$$

Inverse **transforming** this spectrum we find the field at distance z from the grating to be given by

$$U(x, y) = \frac{1}{2} + \frac{m}{4} e^{-j \frac{\pi \lambda z}{L^2}} e^{j \frac{2\pi x}{L}} + \frac{m}{4} e^{-j \frac{\pi \lambda z}{L^2}} e^{-j \frac{2\pi x}{L}},$$

which can be simplified to

$$U(x, y) = \frac{1}{2} \left[1 + m e^{-j \frac{\pi \lambda z}{L^2}} \cos\left(\frac{2\pi x}{L}\right) \right] \quad (4-51)$$

Finally, the intensity distribution is given by

$$I(x, y) = \frac{1}{4} \left[1 + 2m \cos\left(\frac{\pi \lambda z}{L^2}\right) \cos\left(\frac{2\pi x}{L}\right) + m^2 \cos^2\left(\frac{2\pi x}{L}\right) \right]. \quad (4-52)$$

We now consider three special cases of this result that have interesting interpretations.

1. Suppose that the distance z behind the grating satisfies $\frac{\pi \lambda z}{L^2} = 2n\pi$ or $z = \frac{2nL^2}{\lambda}$, where n is an integer. Then the intensity observed at this distance behind the grating is

$$I(x, y) = \frac{1}{4} \left[1 + m \cos\left(\frac{2\pi x}{L}\right) \right]^2$$

which can be interpreted as a *perfect image* of the grating. That is, it is an exact replica of the intensity that would be observed just behind the grating. A multiplicity

of such images appear behind the grating, without the help of lenses! Such images are called Talbot images (after the scientist who first observed them), or simply self-images. A good discussion of such images is found in Ref. [280].

- Suppose that the observation distance satisfies $\frac{\pi\lambda z}{L^2} = (2n + 1)\pi$, or $z = \frac{(2n+1)L^2}{\lambda}$. Then

$$I(x, y) = \frac{1}{4} \left[1 - m \cos\left(\frac{2\pi x}{L}\right) \right]^2.$$

This distribution is also an image of the grating, but this time with a 180° spatial phase shift, or equivalently with a contrast reversal. This, too, is called a Talbot image.

- Finally, consider distances satisfying $\frac{\pi\lambda z}{L^2} = (2n - 1)\frac{\pi}{2}$, or $z = \frac{(n-\frac{1}{2})L^2}{\lambda}$. Then $\cos\left(\frac{\pi\lambda z}{L^2}\right) = 0$, and

$$I(x, y) = \frac{1}{4} \left[1 + m^2 \cos^2\left(\frac{2\pi x}{L}\right) \right] = \frac{1}{4} \left[\left(1 + \frac{m^2}{2}\right) + \frac{m^2}{2} \cos\left(\frac{4\pi x}{L}\right) \right].$$

This image has twice the frequency of the original grating and has reduced contrast. Such an image is called a Talbot subimage. Note that if $m \ll 1$, then the periodic image will effectively vanish at the **subimage** planes.

Figure 4.17 shows the locations of the various types of images behind the original grating.

The Talbot image phenomenon is much more general than just the particular case analyzed here. It can be shown to be present for any periodic structure (see Prob. 4-18).

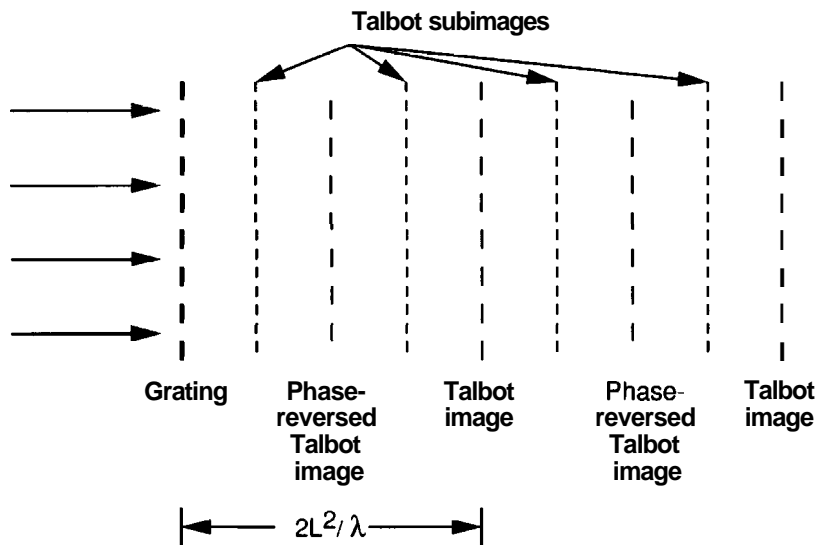


FIGURE 4.17
Locations of Talbot image planes behind the grating.

PROBLEMS-CHAPTER 4

- 4-1.** Consider the quadratic-phase exponential $\frac{1}{j\lambda z} \exp\left[j\frac{\pi}{\lambda z}(x^2 + y^2)\right]$.
- Show that the volume (with respect to x and y) under this function is unity.
 - Show that the two-dimensional quadratic-phase sinusoidal part of this function contributes all of the volume and the two-dimensional quadratic-phase cosinusoidal part contributes none of the volume.
- (Hint: Make use of Table 2.1.)
- 4-2.** Consider a spherical wave expanding about the point $(0, 0, -z_0)$ in a cartesian coordinate system. The wavelength of the light is λ , and $z_0 > 0$.
- Express the phase distribution of the spherical wave across an (x, y) plane located normal to the z axis at coordinate $z = 0$.
 - Using a **paraxial** approximation, express the phase distribution of the parabolic wavefront that approximates this spherical wavefront.
 - Find an exact expression for the phase by which the spherical wavefront *lags* or *leads* the phase of the parabolic wavefront. Does it lag or lead?
- 4-3.** Consider a spherical wave converging towards the point $(0, 0, +z_0)$ in a cartesian coordinate system. The wavelength of the light is λ and $z_0 > 0$.
- Express the phase distribution of the spherical wave across an (x, y) plane located normal to the z axis at coordinate $z = 0$.
 - Using a **paraxial** approximation, express the phase distribution of the parabolic wavefront that approximates this spherical wavefront.
 - Find an exact expression for the phase by which the spherical wavefront lags or leads the phase of the parabolic wavefront. Does it lag or lead?
- 4-4.** Fresnel propagation over a sequence of successive distances z_1, z_2, \dots, z_n must be equivalent to Fresnel propagation over the single distance $z = z_1 + z_2 + \dots + z_n$. Find a simple proof that this is the case.
- 4-5.** Show that the top "transition region" shown in Fig. 4.5 is bounded by the parabola $(w - x)^2 = 4\lambda z$ and the bottom transition region by $(w + x)^2 = 4\lambda z$, where the aperture is $2w$ wide, the origin of the coordinates is at the center of the aperture, z is the distance from the plane of the aperture, and x is the vertical coordinate throughout the figure.
- 4-6.** A spherical wave is converging toward a point $(0, 0, z_0)$ to the right of a circular aperture of radius R , centered on $(0, 0, 0)$. The wavelength of the light is λ . Consider the field observed at an arbitrary point (axial distance z) to the right of the aperture. Show that the wavefront error made in a **paraxial** approximation of the illuminating spherical wave and the error incurred by using a quadratic phase approximation in the Fresnel diffraction equation partially cancel one another. Under what condition does complete cancellation occur?

4-7. Assuming unit-amplitude, normally incident plane-wave illumination:

- (a) Find the intensity distribution in the Fraunhofer diffraction pattern of the double-slit aperture shown in Fig. P4.7.
- (b) Sketch normalized cross sections of this pattern that appear along the x and y axes in the observation plane, assuming $X/\lambda z = 10\text{ m}^{-1}$, $Y/\lambda z = 1\text{ m}^{-1}$, and $A/\lambda z = 312\text{ m}^{-1}$, z being the observation distance and A the wavelength.

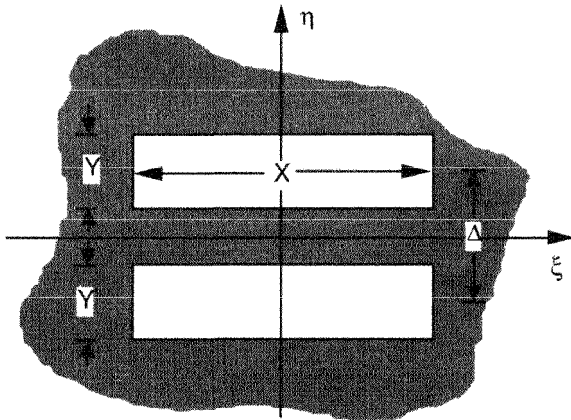


FIGURE P4.7

4-8. (a) Sketch the aperture described by the amplitude transmittance function

$$t_A(\xi, \eta) = \left\{ \left[\text{rect}\left(\frac{\xi}{X}\right) \text{rect}\left(\frac{\eta}{Y}\right) \right] \otimes \left[\frac{1}{\Delta} \text{comb}\left(\frac{\eta}{\Delta}\right) \delta(\xi) \right] \right\} \text{rect}\left(\frac{\eta}{N\Delta}\right)$$

where N is an odd integer and $A > Y$.

- (b) Find an expression for the intensity distribution in the Fraunhofer diffraction pattern of that aperture, assuming illumination by a normally incident plane wave.
 - (c) What relationship between Y and A can be expected to minimize the strength of the even-order diffraction components while leaving the zero-order component approximately unchanged?
- 4-9. Find an expression for the intensity distribution in the Fraunhofer diffraction pattern of the aperture shown in Fig. P4.9. Assume unit-amplitude, normally incident plane-wave illumination. The aperture is square and has a square central obscuration.

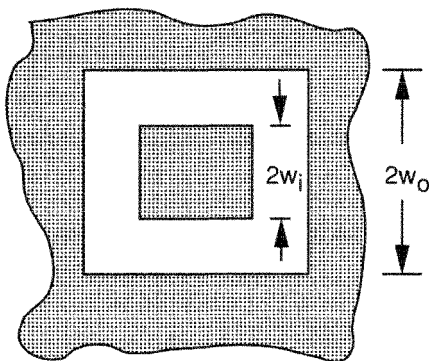


FIGURE P4.9

- 4-10. Find an expression for the intensity distribution in the Fraunhofer diffraction pattern of the aperture shown in Fig. P4.10. Assume unit-amplitude, normally incident plane-wave illumination. The aperture is circular and has a circular central obscuration.

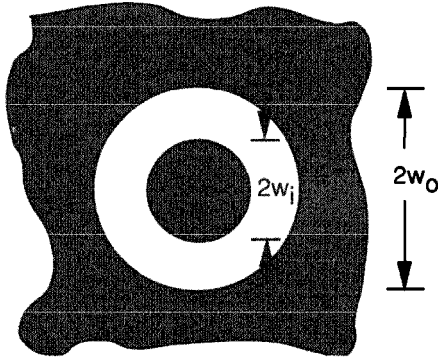


FIGURE P4.10

- 4-11. Two discrete spectral lines of a source are said to be "just resolved" by a diffraction grating if the peak of the 9th-order diffraction component due to source wavelength λ_1 falls exactly on the first zero of the 9th-order diffraction component due to source wavelength λ_2 . The *resolving power* of the grating is defined as the ratio of the mean wavelength λ to the minimum resolvable wavelength difference $\Delta\lambda$. Show that the resolving power of the sinusoidal phase grating discussed in this chapter is

$$\frac{\lambda}{\Delta\lambda} = 2qw f_0 = qM$$

where q is the diffraction order used in the measurement, $2w$ is the width of the square grating, and M is the number of spatial periods of the grating contained in the aperture. What phenomenon limits the use of arbitrarily high diffraction orders?

- 4-12. Consider a thin periodic grating whose amplitude transmittance can be represented by a complex Fourier series,

$$t_A(\xi) = \sum_{k=-\infty}^{\infty} c_k e^{j\frac{2\pi k\xi}{L}}$$

where L is the period of the grating and

$$c_k = \frac{1}{L} \int_{-L/2}^{L/2} t_A(\xi) e^{-j\frac{2\pi k\xi}{L}} d\xi.$$

Neglect the aperture that bounds the grating, since it will not affect the quantities of interest here.

- (a) Show that the diffraction efficiency into the k th order of the grating is simply $\eta_k = |c_k|^2$.
- (b) Calculate the diffraction efficiency into the first diffraction order for a grating with amplitude transmittance given by

$$t_A(\xi) = \left| \cos\left(\frac{\pi\xi}{L}\right) \right|$$

4-13. The amplitude transmittance function of a thin square-wave absorption grating is shown in Fig. P4.13. Find the following properties of this grating:

- The fraction of incident light that is absorbed by the grating.
- The fraction of incident light that is transmitted by the grating.
- The fraction of light that is transmitted into a single first order.

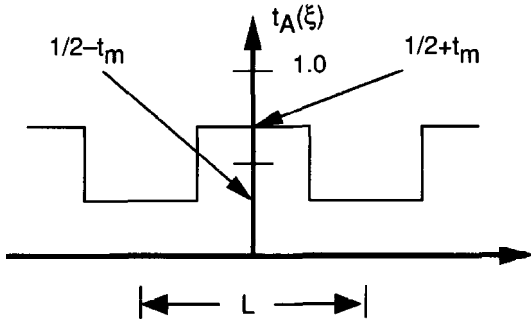


FIGURE P4.13

4-14. A thin square-wave *phase* grating has a thickness that varies periodically (period L) such that the phase of the transmitted light jumps between 0 radians and ϕ radians.

- Find the diffraction efficiency of this grating for the first diffraction orders.
- What value of ϕ yields the maximum diffraction efficiency, and what is the value of that maximum efficiency?

4-15. A "sawtooth" phase grating is periodic with period L and has a distribution of phase within one period from 0 to L given by

$$\phi(\xi) = \frac{2\pi\xi}{L}.$$

- Find the diffraction efficiencies of all of the orders for this grating.
- Suppose that the phase profile of the grating is of the more general form

$$\phi(\xi) = \frac{\phi_0\xi}{L}.$$

Find a general expression for the diffraction efficiency into all the orders of this new grating.

4-16. An aperture Σ in an opaque screen is illuminated by a spherical wave converging towards a point P located in a parallel plane a distance z behind the screen, as shown in Fig. P4.16.

- Find a quadratic-phase approximation to the illuminating wavefront in the plane of the aperture, assuming that the coordinates of P in the (\mathbf{x}, y) plane are $(0, Y)$.
- Assuming *Fresnel* diffraction from the plane of the aperture to the plane containing P , show that in the above case the observed intensity distribution is the *Fraunhofer* diffraction pattern of the aperture, centered on the point P .

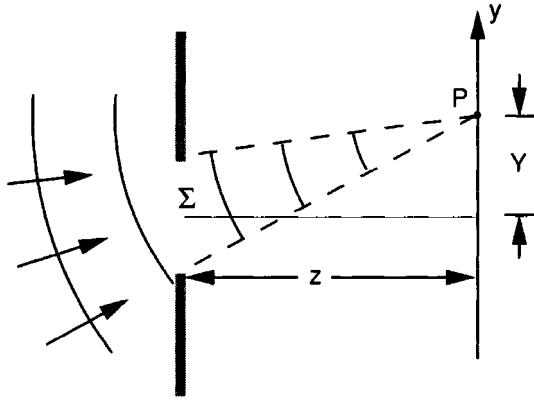


FIGURE P4.16

- 4-17. Find the intensity distribution on the aperture axis in the Fresnel diffraction patterns of apertures with the following transmittance functions (assume normally incident, unit-amplitude, plane-wave illumination):

(a) $t_A(\xi, \eta) = \text{circ } \sqrt{\xi^2 + \eta^2}$.

(b) $t_A(\xi, \eta) = \begin{cases} 1 & a \leq \sqrt{\xi^2 + \eta^2} < b \\ 0 & \text{otherwise} \end{cases}$

where $a < 1$, $b < 1$ and $a < b$.

- 4-18. Consider a one-dimensional periodic object with an amplitude transmittance having an arbitrary periodic profile. Neglect the size of any bounding aperture, ignore the evanescent wave phenomenon, and assume that **paraxial** conditions hold. Show that at certain distances behind this object, perfect images of the amplitude transmittance are found. At what distances do these "self-images" appear?

- 4-19. A certain two-dimensional non-periodic object has the property that all of the frequency components of its amplitude transmittance fall on circles in the frequency plane, the radii of the circles being given by

$$\rho_m = \sqrt{2ma} \quad m = 0, 1, 2, 3, \dots,$$

where a is a constant. Assume uniform plane-wave illumination, neglect the finite size of the object and the evanescent wave phenomenon, and assume that **paraxial** conditions hold. Show that perfect images of the object form at periodic distances behind the object. Find the locations of these images.

- 4-20. A certain circularly symmetric object, infinite in extent, has amplitude transmittance

$$t_A(r) = 2\pi J_0(2\pi r) + 4\pi J_0(4\pi r)$$

where J_0 is a Bessel function of the first kind, zero order, and r is radius in the two-dimensional plane. This object is illuminated by a normally incident, unit-amplitude plane wave. **Paraxial** conditions are assumed to hold. At what distances behind this object will we find a field distribution that is of the same form as that of the object, up to possible complex constants? (Hint: The Fourier transform of the circularly symmetric function $J_0(2\pi r)$ is the circularly symmetric spectrum $\frac{1}{2\pi} \delta(\rho - 1)$.)

- 4-21. An expanding cylindrical wave falls on the "input" plane of an optical system. A **paraxial** approximation to that wave can be written in the form

$$U(y_1) = \exp\left\{j \frac{\pi}{\lambda z_0} [(y_1 - y_0)^2]\right\},$$

where λ is the optical wavelength, while z_0 and y_0 are given constants. The optical system can be represented by a **paraxial** ABCD matrix (see Appendix B, Section B.3) that holds between the input and output planes of the system. Find a **paraxial** expression for the complex amplitude of the field across the "output" plane of the optical system, expressing the results in terms of arbitrary elements of the ray matrix. Assume that the refractive index in the input and output planes is unity. You may treat this problem as one-dimensional.

Wave-Optics Analysis of Coherent Optical Systems

The most important components of optical imaging and data processing systems are lenses. While a thorough discussion of geometrical optics and the properties of lenses would be helpful, such a treatment would require a rather lengthy detour. To provide the most rudimentary background, Appendix B presents a short description of the matrix theory of **paraxial** geometric optics, defining certain quantities that will be important in our purely "wave-optics" approach in this chapter. The reader will be referred to appropriate material in the appendix when needed. However, the philosophy of our approach is to make minimum use of geometrical optics, and instead to develop purely wave-optic analyses of the systems of interest. The results of this approach are entirely consistent with the results of geometrical optics, with the added advantage that diffraction effects are entirely accounted for in the wave-optics approach, but not in the geometrical-optics approach. Our discussions will be limited to the case of monochromatic illumination, with generalization to nonmonochromatic light being deferred to Chapter 6.

5.1

A THIN LENS AS A PHASE TRANSFORMATION

A lens is composed of an optically dense material, usually glass with a refractive index of approximately 1.5, in which the propagation velocity of an optical disturbance is less than the velocity in air. With reference to Appendix B, a lens is said to be a thin lens if a ray entering at coordinates (\mathbf{x}, y) on one face exits at approximately the same coordinates on the opposite face, i.e. if there is negligible translation of a ray within the lens. Thus a thin lens simply delays an incident wavefront by an amount proportional to the thickness of the lens at each point.

Referring to Fig. 5.1, let the maximum thickness of the lens (on its axis) be Δ_0 , and let the thickness at coordinates (x, y) be $\Delta(x, y)$. Then the total phase delay suffered by the wave at coordinates (x, y) in passing through the lens may be written

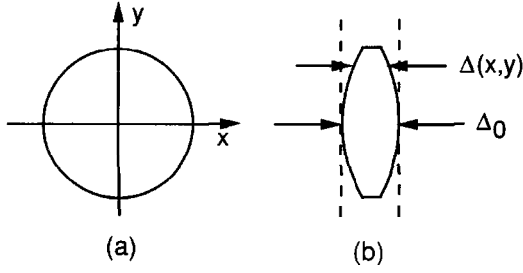


FIGURE 5.1
The thickness function. (a) Front view,
(b) side view

$$\phi(x, y) = kn\Delta(x, y) + k[\Delta_0 - \Delta(x, y)]$$

where n is the refractive index of the lens material, $kn\Delta(x, y)$ is the phase delay introduced by the lens, and $k[\Delta_0 - \Delta(x, y)]$ is the phase delay introduced by the remaining region of free space between the two planes. Equivalently the lens may be represented by a multiplicative phase transformation of the form

$$t_l(x, y) = \exp[jk\Delta_0] \exp[jk(n - 1)\Delta(x, y)]. \quad (5-1)$$

The complex field $U'_l(x, y)$ across a plane immediately behind the lens is then related to the complex field $U_l(x, y)$ incident on a plane immediately in front of the lens by

$$U'_l(x, y) = t_l(x, y) U_l(x, y). \quad (5-2)$$

The problem remains to find the mathematical form of the thickness function $\Delta(x, y)$ in order that the effects of the lens may be understood.

5.1.1 The Thickness Function

In order to specify the forms of the phase transformations introduced by a variety of different types of lenses, we first adopt a sign convention: as rays travel from left to right, each convex surface encountered is taken to have a positive radius of curvature, while each concave surface is taken to have a negative radius of curvature. Thus in Fig. 5.1(b) the radius of curvature of the left-hand surface of the lens is a positive number R_1 , while the radius of curvature of the right-hand surface is a negative number R_2 .

To find the thickness $\Delta(x, y)$, we split the lens into three parts, as shown in Fig. 5.2, and write the total thickness function as the sum of three individual thickness functions,

$$\Delta(x, y) = \Delta_1(x, y) + \Delta_2(x, y) + \Delta_3(x, y). \quad (5-3)$$

Referring to the geometries shown in that figure, the thickness function $\Delta_1(x, y)$ is given by

$$\begin{aligned} \Delta_1(x, y) &= \Delta_{01} - \left(R_1 - \sqrt{R_1^2 - x^2 - y^2} \right) \\ &= \Delta_{01} - R_1 \left(1 - \sqrt{1 - \frac{x^2 + y^2}{R_1^2}} \right). \end{aligned} \quad (5-4)$$

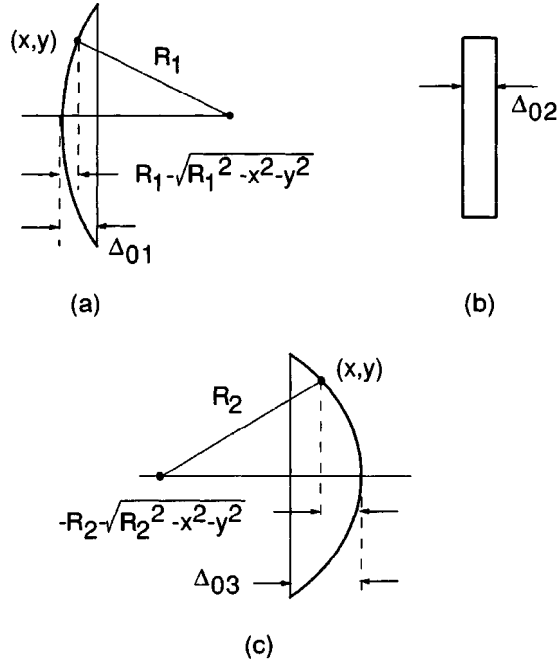


FIGURE 5.2
Calculation of the thickness function.
(a) Geometry for Δ_1 , (b) geometry for Δ_2 ,
and (c) geometry for Δ_3 .

The second component of the thickness function comes from a region of glass of constant thickness Δ_{02} . The third component is given by

$$\begin{aligned} \Delta_3(x, y) &= \Delta_{03} - \left(-R_2 - \sqrt{R_2^2 - x^2 - y^2} \right) \\ &= \Delta_{03} + R_2 \left(1 - \sqrt{1 - \frac{x^2 + y^2}{R_2^2}} \right), \end{aligned} \tag{5-5}$$

where we have factored the positive number $-R_2$ out of the square root. Combining the three expressions for thickness, the total thickness is seen to be

$$\Delta(x, y) = \Delta_0 - R_1 \left(1 - \sqrt{1 - \frac{x^2 + y^2}{R_1^2}} \right) + R_2 \left(1 - \sqrt{1 - \frac{x^2 + y^2}{R_2^2}} \right), \tag{5-6}$$

where $\Delta_0 = \Delta_{01} + \Delta_{02} + \Delta_{03}$.

5.1.2 The Paraxial Approximation

The expression for the thickness function can be substantially simplified if attention is restricted to portions of the wavefront that lie near the lens axis, or equivalently, if only *paraxial* rays are considered. Thus we consider only values of \mathbf{x} and y sufficiently small to allow the following approximations to be accurate:

$$\sqrt{1 - \frac{x^2 + y^2}{R_1^2}} \approx 1 - \frac{x^2 + y^2}{2R_1^2}$$

$$\sqrt{1 - \frac{x^2 + y^2}{R_2^2}} \approx 1 - \frac{x^2 + y^2}{2R_2^2}. \quad (5-7)$$

The resulting phase transformation will, of course, represent the lens accurately over only a limited area, but this limitation is no more restrictive than the usual **paraxial** approximation of geometrical optics. Note that the relations (5-7) amount to approximations of the spherical surfaces of the lens by parabolic surfaces. With the help of these approximations, the thickness function becomes

$$\Delta(x, y) = \Delta_0 - \frac{x^2 + y^2}{2} \left(\frac{1}{R_1} - \frac{1}{R_2} \right). \quad (5-8)$$

5.1.3 The Phase Transformation and Its Physical Meaning

Substitution of Eq. (5-8) into Eq. (5-1) yields the following approximation to the lens transformation:

$$t_l(x, y) = \exp[jkn\Delta_0] \exp \left[-jk(n-1) \frac{x^2 + y^2}{2} \left(\frac{1}{R_1} - \frac{1}{R_2} \right) \right].$$

The physical properties of the lens (that is, n , R_1 , and R_2) can be combined in a single number f called the **focal length**, which is defined by

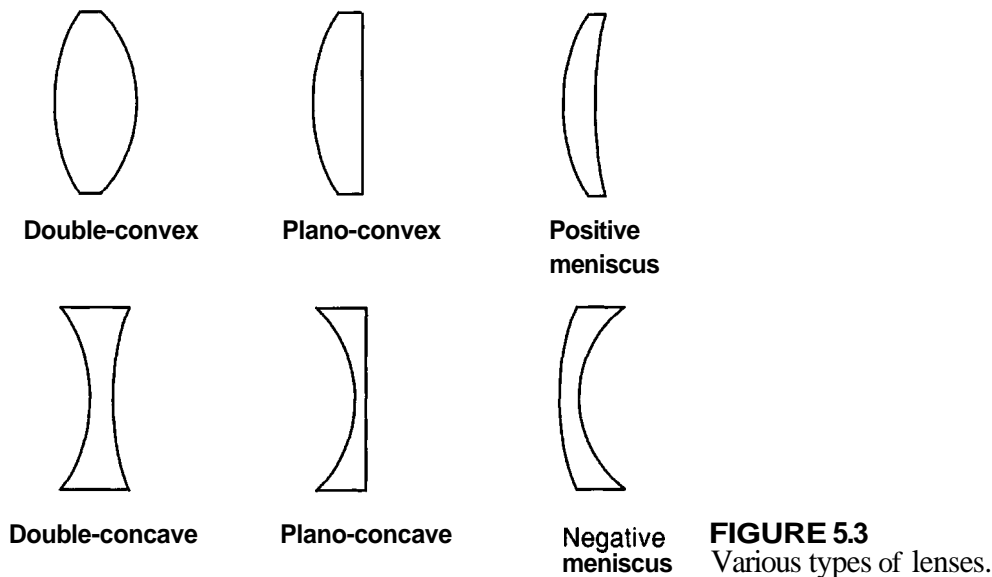
$$\frac{1}{f} \equiv (n-1) \left(\frac{1}{R_1} - \frac{1}{R_2} \right). \quad (5-9)$$

Neglecting the constant phase factor, which we shall drop hereafter, the phase transformation may now be rewritten

$$t_l(x, y) = \exp \left[-j \frac{k}{2f} (x^2 + y^2) \right]. \quad (5-10)$$

This equation will serve as our basic representation of the effects of a thin lens on an incident disturbance. It neglects the finite extent of the lens, which we will account for later.

Note that while our derivation of this expression assumed the specific lens shape shown in Fig. 5.1, the sign convention adopted allows the result to be applied to other types of lenses. Figure 5.3 illustrates several different types of lenses with various combinations of convex and concave surfaces. In Prob. 5-1, the reader is asked to verify that the sign convention adopted implies that the focal length f of a double-convex, plano-convex, or positive meniscus lens is **positive**, while that of a double-concave,



plano-concave, or negative meniscus lens is negative. Thus Eq. (5-10) can be used to represent any of the above lenses, provided the correct sign of the focal length is used.

The physical meaning of the lens transformation can best be understood by considering the effect of the lens on a normally incident, unit-amplitude plane wave. The field distribution U_l in front of the lens is unity, and Eqs. (5-1) and (5-10) yield the following expression for U'_l behind the lens:

$$U'_l(x, y) = \exp\left[-j\frac{k}{2f}(x^2 + y^2)\right].$$

We may interpret this expression as a quadratic approximation to a spherical wave. If the focal length is positive, then the spherical wave is converging towards a point on the lens axis a distance f behind the lens. If f is negative, then the spherical wave is diverging about a point on the lens axis a distance $|f|$ in front of the lens. The two cases are illustrated in Fig. 5.4. Thus a lens with a positive focal length is called a positive or converging lens, while a lens with a negative focal length is a negative or diverging lens.

Our conclusion that a lens composed of spherical surfaces maps an incident plane wave into a spherical wave is very much dependent on the **paraxial** approximation. Under nonparaxial conditions, the emerging wavefront will exhibit departures from perfect sphericity (called aberrations — see Section 6.4), even if the surfaces of the lens are perfectly spherical. In fact, lenses are often "corrected" for aberrations by making their surfaces aspherical in order to improve the sphericity of the emerging wavefront.

We should emphasize, however, that the results which will be derived using the multiplicative phase transformation (5-10) are actually more general than the analysis leading up to that equation might imply. A thorough geometrical-optics analysis of most well-corrected lens systems shows that they behave essentially in the way predicted by our more restrictive theory.

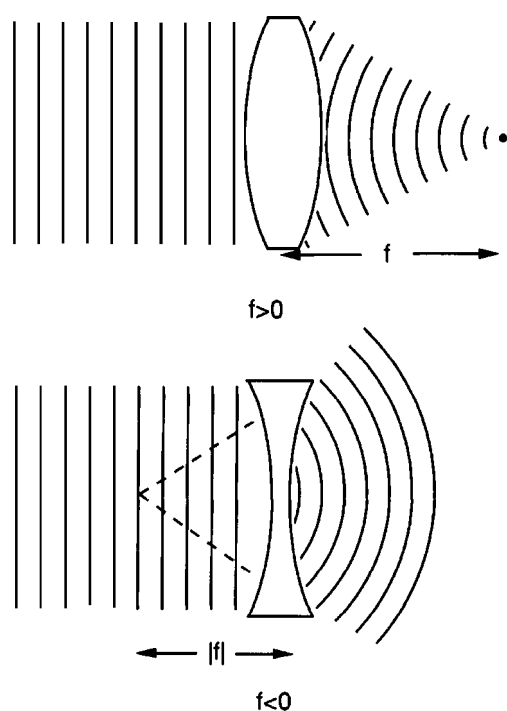


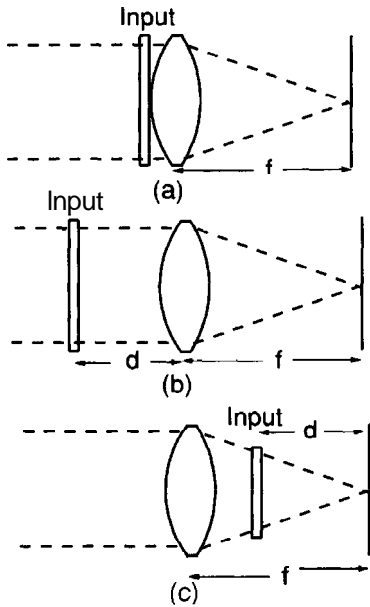
FIGURE 5.4
Effects of a converging lens and a diverging lens on a normally incident plane wave.

5.2 FOURIER TRANSFORMING PROPERTIES OF LENSES

One of the most remarkable and useful properties of a converging lens is its inherent ability to perform two-dimensional Fourier transforms. This complicated analog operation can be performed with extreme simplicity in a coherent optical system, taking advantage of the basic laws of propagation and diffraction of light.

In the material that follows, several different configurations for performing the transform operation are described. In all cases the illumination is assumed to be monochromatic. Under this condition the systems studied are "coherent" systems, which means that they are linear in complex amplitude, and the distribution of light amplitude across a particular plane behind the positive lens is of interest. In some cases this is the *back focal plane* of the lens, which by definition is a plane normal to the lens axis situated a distance f behind the lens (in the direction of propagation of light). The information to be Fourier-transformed is introduced into the optical system by a device with an amplitude transmittance that is proportional to the input function of interest. In some cases this device may consist of a photographic transparency, while in others it may be a nonphotographic *spatial light modulator*, capable of controlling the amplitude transmittance in response to externally supplied electrical or optical information. Such input devices will be discussed in more detail in Chapter 7. We will refer to them as input "transparencies", even though in some cases they may operate by reflection of light rather than transmission of light. We will also often refer to the input as the "object".

Figure 5.5 shows three arrangements that will be considered here. In all cases shown, the illumination is a collimated plane wave which is incident either on the input

**FIGURE 5.5**

Geometries for performing the Fourier transform operation with a positive lens.

transparency or on the lens. In case (a), the input transparency is placed directly against the lens itself. In case (b), the input is placed a distance d in front of the lens. In case (c), the input is placed behind the lens at distance d from the focal plane. An additional, more general case, will be studied in Section 5.4.

For alternative discussions of the Fourier transforming properties of positive lenses, the reader may wish to consult Refs. [243], [73], or [235].

5.2.1 Input Placed Against the Lens

Let a planar input transparency with amplitude transmittance $t_A(x, y)$ be placed immediately in front of a converging lens of focal length f , as shown in Fig. 5.5(a). The input is assumed to be uniformly illuminated by a normally incident, monochromatic plane wave of amplitude A , in which case the disturbance incident on the lens is

$$U_l(x, y) = A t_A(x, y). \quad (5-11)$$

The finite extent of the lens can be accounted for by associating with the lens a pupil function $P(x, y)$ defined by

$$P(x, y) = \begin{cases} 1 & \text{inside the lens aperture} \\ 0 & \text{otherwise.} \end{cases}$$

Thus the amplitude distribution behind the lens becomes, using (5-10),

$$U'_l(x, y) = U_l(x, y) P(x, y) \exp\left[-j\frac{k}{2f}(x^2 + y^2)\right]. \quad (5-12)$$

To find the distribution $U_f(\mathbf{u}, \mathbf{v})$ in the back focal plane of the lens, the Fresnel diffraction formula, Eq. (4-17), is applied. Thus, putting $z = f$,

$$\begin{aligned}
U_f(u, v) &= \frac{\exp\left[j\frac{k}{2f}(u^2 + v^2)\right]}{j\lambda f} \\
&\times \iint_{-\infty}^{\infty} U_l(x, y) \exp\left[j\frac{k}{2f}(x^2 + y^2)\right] \exp\left[-j\frac{2\pi}{\lambda f}(xu + yv)\right] dx dy,
\end{aligned} \tag{5-13}$$

where a constant phase factor has been dropped. Substituting (5-12) in (5-13), the quadratic phase factors within the integrand are seen to exactly cancel, leaving

$$\begin{aligned}
U_f(u, v) &= \frac{\exp\left[j\frac{k}{2f}(u^2 + v^2)\right]}{j\lambda f} \\
&\times \iint_{-\infty}^{\infty} U_l(x, y) P(x, y) \exp\left[-j\frac{2\pi}{\lambda f}(xu + yv)\right] dx dy.
\end{aligned} \tag{5-14}$$

Thus the field distribution U_f is proportional to the two-dimensional Fourier transform of that portion of the incident field subtended by the lens aperture. When the physical extent of the input is smaller than the lens aperture, the factor $P(x, y)$ may be neglected, yielding

$$U_f(u, v) = \frac{\exp\left[j\frac{k}{2f}(u^2 + v^2)\right]}{j\lambda f} \iint_{-\infty}^{\infty} U_l(x, y) \exp\left[-j\frac{2\pi}{\lambda f}(xu + yv)\right] dx dy. \tag{5-15}$$

Thus we see that the complex amplitude distribution of the field in the focal plane of the lens is the Fraunhofer *diffraction pattern* of the field incident on the lens, even though the distance to the observation plane is equal to the focal length of the lens, rather than satisfying the usual distance criterion for observing Fraunhofer diffraction. Note that the amplitude and phase of the light at coordinates (u, v) in the focal plane are determined by the amplitude and phase of the input Fourier component at frequencies $(f_x = u/\lambda f, f_y = v/\lambda f)$.

The Fourier transform relation between the input amplitude transmittance and the focal-plane amplitude distribution is not a complete one, due to the presence of the quadratic phase factor that precedes the integral. While the phase distribution across the focal plane is not the same as the phase distribution across the spectrum of the input, the difference between the two is a simple phase curvature.

In most cases it is the intensity across the focal plane that is of real interest. This phase term is important if the ultimate goal is to calculate another field distribution after further propagation and possibly passage through additional lenses, in which case the complete complex field is needed. In most cases, however, the intensity distribution in the focal plane will be measured, and the phase distribution is of no consequence. Measurement of the intensity distribution yields knowledge of the power spectrum (or more accurately, the energy spectrum) of the input. Thus

$$I_f(u, v) = \frac{A^2}{\lambda^2 f^2} \left| \iint_{-\infty}^{\infty} t_A(x, y) \exp \left[-j \frac{2\pi}{\lambda f} (xu + yv) \right] dx dy \right|^2. \quad (5-16)$$

5.2.2 Input Placed in Front of the Lens

Consider next the more general geometry of Fig. 5.5(b). The input, located a distance d in front of the lens, is illuminated by a normally incident plane wave of amplitude A . The amplitude transmittance of the input is again represented by t_A . In addition, let $F_o(f_X, f_Y)$ represent the Fourier spectrum of the light transmitted by the input transparency, and $F_l(f_X, f_Y)$ the Fourier spectrum of the light incident on the lens; that is,

$$F_o(f_X, f_Y) = \mathcal{F}\{At_A\} \quad F_l(f_X, f_Y) = \mathcal{F}\{U_l\}.$$

Assuming that the Fresnel or **paraxial** approximation is valid for propagation over distance d , then F_o and F_l are related by means of Eq. (4-21), giving

$$F_l(f_X, f_Y) = F_o(f_X, f_Y) \exp \left[-j\pi\lambda d(f_X^2 + f_Y^2) \right], \quad (5-17)$$

where we have dropped a constant phase delay.

For the moment, the finite extent of the lens aperture will be neglected. Thus, letting $P = 1$, Eq. (5-14) can be rewritten

$$U_f(u, v) = \frac{\exp \left[j \frac{k}{2f} (u^2 + v^2) \right]}{j\lambda f} F_l \left(\frac{u}{\lambda f}, \frac{v}{\lambda f} \right). \quad (5-18)$$

Substituting (5-17) into (5-18), we have

$$U_f(u, v) = \frac{\exp \left[j \frac{k}{2f} \left(1 - \frac{d}{f} \right) (u^2 + v^2) \right]}{j\lambda f} F_o \left(\frac{u}{\lambda f}, \frac{v}{\lambda f} \right),$$

$$U_f(u, v) = \frac{A \exp \left[j \frac{k}{2f} \left(1 - \frac{d}{f} \right) (u^2 + v^2) \right]}{j\lambda f}$$

$$\times \iint_{-\infty}^{\infty} t_A(\xi, \eta) \exp \left[-j \frac{2\pi}{\lambda f} (\xi u + \eta v) \right] d\xi d\eta. \quad (5-19)$$

Thus the amplitude and phase of the light at coordinates (u, v) are again related to the amplitude and phase of the input spectrum at frequencies $(u/\lambda f, v/\lambda f)$. Note *that* a quadratic phase factor again precedes the transform integral, but that it vanishes for the very special case $d = f$. ***Evidently when the input is placed in the front focal plane of the lens, the phase curvature disappears, leaving an exact Fourier transform relation!***

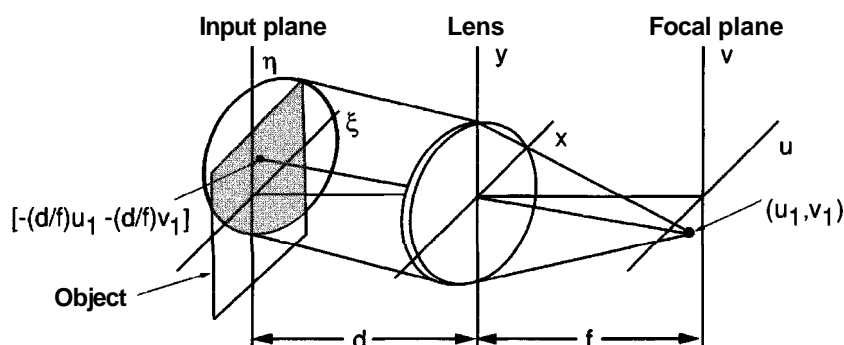


FIGURE 5.6

Vignetting of the input. The shaded area in the input plane represents the portion of the input transparency that contributes to the Fourier transform at (u_1, v_1) .

To this point we have entirely neglected the finite extent of the lens aperture. To include the effects of this aperture, we use a geometrical optics approximation. Such an approximation is accurate if the distance d is sufficiently small to place the input deep within the region of Fresnel diffraction of the lens aperture, if the light were propagating backwards from the focal plane to the plane of the input transparency. This condition is well satisfied in the vast majority of problems of interest. With reference to Fig. 5.6, the light amplitude at coordinates (u_1, v_1) is a summation of all the rays traveling with direction cosines $(\xi \approx u_1/f, \eta \approx v_1/f)$. However, only a finite set of these rays is passed by the lens aperture. Thus the finite extent of the aperture may be accounted for by geometrically projecting that aperture back to the input plane, the projection being centered on a line joining the coordinates (u_1, v_1) with the center of the lens (see Fig. 5.6). The projected lens aperture limits the effective extent of the input, but the particular portion of t that contributes to the field U_f depends on the particular coordinates (u, v) being considered in the back focal plane. As implied by Fig. 5.6, the value of U_f at (u, v) can be found from the Fourier transform of that portion of the input subtended by the projected pupil function P , centered at coordinates $[\xi = -(d/f)u, \eta = -(d/f)v]$. Expressing this fact mathematically,

$$\begin{aligned}
 U_f(u, v) = & \frac{A \exp \left[j \frac{k}{2f} \left(1 - \frac{d}{f} \right) (u^2 + v^2) \right]}{j \lambda f} \\
 & \times \iint_{-\infty}^{\infty} t_A(\xi, \eta) P \left(\xi + \frac{d}{f} u, \eta + \frac{d}{f} v \right) \exp \left[-j \frac{2\pi}{\lambda f} (\xi u + \eta v) \right] d\xi d\eta.
 \end{aligned} \tag{5-20}$$

The limitation of the effective input by the finite lens aperture is known as a *vignetting* effect. Note that for a simple Fourier transforming system, vignetting of the input space is minimized when the input is placed close to the lens and when the lens

aperture is much larger than the input transparency. In practice, when the Fourier transform of the object is of prime interest, it is often preferred to place the input directly against the lens in order to minimize vignetting, although in analysis it is generally convenient to place the input in the front focal plane, where the transform relation is unencumbered with quadratic phase factors.

5.2.3 Input Placed Behind the Lens

Consider next the case of an input that is placed behind the lens, as illustrated in Fig. 5.5(c). The input again has amplitude transmittance t_A , but it is now located a distance d in front of the rear focal plane of the lens. Let the lens be illuminated by a normally incident plane wave of uniform amplitude A . Then incident on the input is a spherical wave converging towards the back focal point of the lens.

In the geometrical optics approximation, the amplitude of the spherical wave impinging on the object is Afd , due to the fact that the linear dimension of the circular converging bundle of rays has been reduced by the factor d/f and energy has been conserved. The particular region of the input that is illuminated is determined by the intersection of the converging cone of rays with the input plane. If the lens is circular and of diameter l , then a circular region of diameter ld/f is illuminated on the input. The finite extent of the illuminating spot can be represented mathematically by projecting the pupil function of the lens down the cone of rays to the intersection with the input plane, yielding an effective illuminated region in that plane described by the pupil function $P[\xi(f/d), \eta(f/d)]$. Note that the input amplitude transmittance t_A will also have a finite aperture associated with it; the effective aperture in the input space is therefore determined by the intersection of the true input aperture with the projected pupil function of the lens. If the finite input transparency is fully illuminated by the converging light, then the projected pupil can be ignored.

Using a paraxial approximation to the spherical wave that illuminates the input, the amplitude of the wave transmitted by the input may be written

$$U_o(\xi, \eta) = \left\{ \frac{Af}{d} P\left(\xi \frac{f}{d}, \eta \frac{f}{d}\right) \exp\left[-j \frac{k}{2d}(\xi^2 + \eta^2)\right] \right\} t_A(\xi, \eta). \quad (5-21)$$

Assuming Fresnel diffraction from the input plane to the focal plane, Eq. (4-17) can be applied to the field transmitted by the input. If this is done it is found that the quadratic phase exponential in (ξ, η) associated with the illuminating wave exactly cancels the similar quadratic phase exponential in the integrand of the Fresnel diffraction integral, with the result

$$U_f(u, v) = \frac{A \exp[j \frac{k}{2d}(u^2 + v^2)] f}{j \lambda d} \frac{f}{d} \times \iint_{-\infty}^{\infty} t_A(\xi, \eta) P\left(\xi \frac{f}{d}, \eta \frac{f}{d}\right) \exp\left[-j \frac{2\pi}{\lambda d}(u\xi + v\eta)\right] d\xi d\eta. \quad (5-22)$$

Thus, up to a quadratic phase factor, the focal-plane amplitude distribution is the Fourier transform of that portion of the input subtended by the projected lens aperture.

The result presented in Eq. (5-22) is essentially the same result obtained when the input was placed directly against the lens itself. However, an extra flexibility has been obtained in the present configuration; namely, the scale of the Fourier transform is under the control of the experimenter. By increasing d , the distance from the focal plane, the size of the transform is made larger, at least until the transparency is directly against the lens (i.e. $d = f$). By decreasing d , the scale of the transform is made smaller. This flexibility can be of utility in spatial filtering applications (see Chapter 8), where some potential adjustment of the size of the transform can be of considerable help.

5.2.4 Example of an Optical Fourier Transform

We illustrate with a typical example the type of two-dimensional Fourier analysis that can be achieved optically with great ease. Figure 5.7 shows a transparent character 3, which is placed in front of a positive lens and illuminated by a plane wave, yielding in the back focal plane the intensity distribution shown in the right-hand part of the figure. Note in particular the high-frequency components introduced by the straight edges in the input.

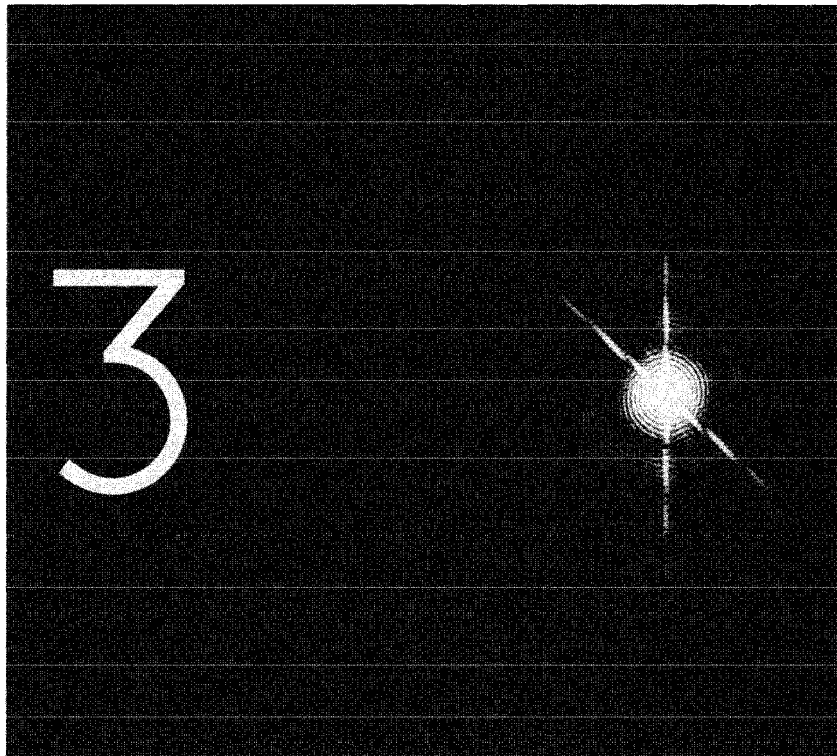


FIGURE 5.7
Optically obtained Fourier transform of the character 3.

5.3

IMAGE FORMATION: MONOCHROMATIC ILLUMINATION

Certainly the most familiar property of lenses is their ability to form images. If an object is placed in front of a lens and illuminated, then under appropriate conditions there will appear across a second plane a distribution of light intensity that closely resembles the object. This distribution of intensity is called an image of the object. The image may be real in the sense that an actual distribution of intensity appears across a plane behind the lens, or it may be virtual in the sense that the light behind the lens appears to originate from an intensity distribution across a new plane in front of the lens.

For the present we consider image formation in only a limited context. First we restrict attention to a positive, aberration-free thin lens that forms a real image. Second, we consider only monochromatic illumination, a restriction implying that the imaging system is linear in complex field amplitude (see Prob. 6-18). Both of these restrictions will be removed in Chapter 6, where the problem of image formation will be treated in a much more general fashion.

5.3.1 The Impulse Response of a Positive Lens

Referring to the geometry of Fig. 5.8, suppose that a planar object is placed a distance z_1 in front of a positive lens and is illuminated by monochromatic light. We represent the complex field immediately behind the object by $U_o(\xi, \eta)$. At a distance z_2 behind the lens there appears a field distribution that we represent by $U_i(u, v)$. Our purpose is to find the conditions under which the field distribution U_i can reasonably be said to be an "image" of the object distribution U_o .

In view of the linearity of the wave propagation phenomenon, we can in all cases express the field U_i by the following superposition integral:

$$U_i(u, v) = \iint_{-\infty}^{\infty} h(u, v; \xi, \eta) U_o(\xi, \eta) d\xi d\eta, \quad (5-23)$$

where $h(u, v; \xi, \eta)$ is the field amplitude produced at coordinates (u, v) by a unit-amplitude point source applied at object coordinates (ξ, η) . Thus the properties of the imaging system will be completely described if the impulse response h can be specified.

If the optical system is to produce high-quality images, then U_i must be as similar as possible to U_o . Equivalently, the impulse response should closely approximate a Dirac delta function,

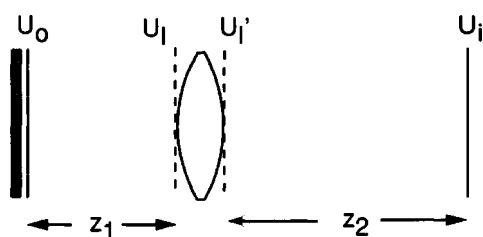


FIGURE 5.8
Geometry for image formation.

$$h(u, v; \xi, \eta) \approx K \delta(u \pm M\xi, v \pm M\eta), \quad (5-24)$$

where K is a complex constant, M represents the system magnification, and the plus and minus signs are included to allow for the absence or presence of image inversion, respectively. We shall therefore specify as the "image plane" that plane where (5-24) is most closely approximated.

To find the impulse response h , let the object be a δ function (point source) at coordinates (ξ, η) . Then incident on the lens will appear a spherical wave diverging from the point (ξ, η) . The **paraxial** approximation to that wave is written

$$U_i(x, y) = \frac{1}{j\lambda z_1} \exp \left\{ j \frac{k}{2z_1} [(x - \xi)^2 + (y - \eta)^2] \right\}. \quad (5-25)$$

After passage through the lens (focal length f), the field distribution becomes

$$U'_i(x, y) = U_i(x, y) P(x, y) \exp \left[-j \frac{k}{2f} (x^2 + y^2) \right] \quad (5-26)$$

Finally, using the Fresnel diffraction equation (4-14) to account for propagation over distance z_2 , we have

$$h(u, v; \xi, \eta) = \frac{1}{j\lambda z_2} \iint_{-\infty}^{\infty} U'_i(x, y) \exp \left\{ j \frac{k}{2z_2} [(u - x)^2 + (v - y)^2] \right\} dx dy \quad (5-27)$$

where constant phase factors have been dropped. Combining (5-25), (5-26), and (5-27), and again neglecting a pure phase factor, yields the formidable result

$$\begin{aligned} h(u, v; \xi, \eta) &= \frac{1}{\lambda^2 z_1 z_2} \exp \left[j \frac{k}{2z_2} (u^2 + v^2) \right] \exp \left[j \frac{k}{2z_1} (\xi^2 + \eta^2) \right] \\ &\times \iint_{-\infty}^{\infty} P(x, y) \exp \left[j \frac{k}{2} \left(\frac{1}{z_1} + \frac{1}{z_2} - \frac{1}{f} \right) (x^2 + y^2) \right] \\ &\times \exp \left\{ -jk \left[\left(\frac{\xi}{z_1} + \frac{u}{z_2} \right) x + \left(\frac{\eta}{z_1} + \frac{v}{z_2} \right) y \right] \right\} dx dy. \end{aligned} \quad (5-28)$$

Equations (5-23) and (5-28) now provide a formal solution specifying the relationship that exists between the object U_o and the image U_i . However, it is difficult to determine the conditions under which U_i can reasonably be called an image of U_o unless further simplifications are adopted.

5.3.2 Eliminating Quadratic Phase Factors: The Lens Law

The most troublesome terms of the impulse response above are those containing quadratic phase factors. Note that two of these terms are independent of the lens coordinates, namely

$$\exp\left[j\frac{k}{2z_2}(u^2 + v^2)\right] \quad \text{and} \quad \exp\left[j\frac{k}{2z_1}(\xi^2 + \eta^2)\right],$$

while one term depends on the lens coordinates (the variables of integration), namely

$$\exp\left[j\frac{k}{2}\left(\frac{1}{z_1} + \frac{1}{z_2} - \frac{1}{f}\right)(x^2 + y^2)\right]$$

We now consider a succession of approximations and restrictions that eliminate these factors. Beginning with the term involving the variables of integration (x, y) first, note that the presence of a quadratic phase factor in what otherwise would be a Fourier transform relationship will generally have the effect of **broadening** the impulse response. For this reason we choose the distance z_2 to the image plane so that this term will identically vanish. This will be true if

$$\frac{1}{z_1} + \frac{1}{z_2} - \frac{1}{f} = 0. \quad (5-29)$$

Note that this relationship is precisely the classical **lens law** of geometrical optics, and must be satisfied for imaging to hold.

Consider next the quadratic phase factor that depends only on image coordinates (u, v). This term can be ignored under either of two conditions:

1. It is the intensity distribution in the image plane that is of interest, in which case the phase distribution associated with the image is of no consequence.
2. The image field distribution is of interest, but the image is measured on a spherical surface, centered at the point where the optical axis pierces the thin lens, and of radius z_2 .

Since it is usually the intensity of the image that is of interest, we will drop this quadratic phase factor in the future.

Finally, consider the quadratic phase factor in the object coordinates (ξ, η). Note that this term depends on the variables over which the convolution operation (5-23) is carried out, and it has the potential to affect the result of that integration significantly. There are three different conditions under which this term can be neglected:

1. The object exists on the surface of a sphere of radius z_1 centered on the point where the optical axis pierces the thin lens.
2. The object is illuminated by a spherical wave that is converging towards the point where the optical axis pierces the lens.
3. The phase of the quadratic phase factor changes by an amount that is only a small fraction of a radian within the region of the object that contributes significantly to the field at the particular image point (u, v).

The first of these conditions rarely occurs in practice. The second can easily be made to occur by proper choice of the illumination, as illustrated in Fig. 5.9. In this case the spherical wave illumination results in the Fourier transform of the object appearing in the pupil plane of the lens. The quadratic phase factor of concern is exactly canceled by this converging spherical wave.

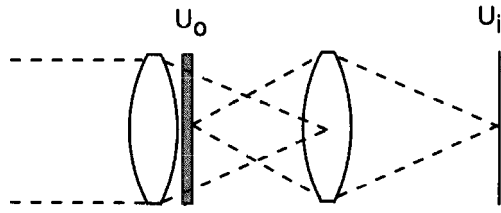


FIGURE 5.9
Converging illumination of the object.

The third possibility for eliminating the effect of the quadratic phase factor in object coordinates requires a more lengthy discussion. In an imaging geometry, the response of the system to an impulse at particular object coordinates should extend over only a small region of image space surrounding the exact image point corresponding to that particular object point. If this were not the case, the system would not be producing an accurate image of the object, or stated another way, it would have an unacceptably large image blur. By the same token, if we view the impulse response for a fixed image point as specifying the weighting function in object space that contributes to that image point, then only a small region on the object should contribute to any given image point. Figure 5.10 illustrates this point-of-view. The gray patch on the left in this figure represents the area from which significant contributions arise for the particular image point on the right. If over this region the factor $\frac{k}{2z_1}(\xi^2 + \eta^2)$ changes by an amount that is only a small fraction of a radian, then the quadratic phase factor in the object plane can be replaced by a single phase that depends on which image point (u, v) is of interest but does not depend on the object coordinates (ξ, η) . The replacement can be stated more precisely as

$$\exp\left[j\frac{k}{2z_1}(\xi^2 + \eta^2)\right] \rightarrow \exp\left[j\frac{k}{2z_1}\left(\frac{u^2 + v^2}{M^2}\right)\right], \quad (5-30)$$

where $M = -z_2/z_1$ is the magnification of the system, to be defined shortly. This new quadratic phase factor in the image space can now be dropped provided that image intensity is the quantity of interest.

Tichenor and Goodman [282] have examined this argument in detail and have found that the approximation stated above is valid provided the size of object is no greater than about $1/4$ the size of the lens aperture. For further consideration of this problem, see Prob. 5-12.

The end result of these arguments is a simplified expression for the impulse response of the imaging system,

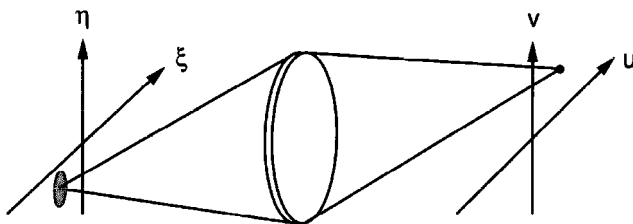


FIGURE 5.10
Region of object space contributing to the field at a particular image point.

$$h(u, v; \xi, \eta) \approx \frac{1}{\lambda^2 z_1 z_2} \iint_{-\infty}^{\infty} P(x, y) \times \exp \left\{ -jk \left[\left(\frac{\xi}{z_1} + \frac{u}{z_2} \right) x + \left(\frac{\eta}{z_1} + \frac{v}{z_2} \right) y \right] \right\} dx dy. \quad (5-31)$$

Defining the *magnification* of the system by

$$M = -\frac{z_2}{z_1}, \quad (5-32)$$

the minus sign being included to remove the effects of image inversion, we find a final simplified form for the impulse response,

$$h(u, v; \xi, \eta) \approx \frac{1}{\lambda^2 z_1 z_2} \iint_{-\infty}^{\infty} P(x, y) \times \exp \left\{ -j \frac{2\pi}{\lambda z_2} [(u - M\xi)x + (v - M\eta)y] \right\} dx dy. \quad (5-33)$$

Thus, if the lens law is satisfied, the impulse response is seen to be given (up to an extra scaling factor $1/\lambda z_1$) by the Fraunhofer diffraction pattern of the lens aperture, centered on image coordinates ($u = M\xi$, $v = M\eta$). The occurrence of a Fraunhofer diffraction formula should not be entirely surprising. By choosing z_2 to satisfy the lens law, we have chosen to examine the plane towards which the spherical wave leaving the lens is converging. From the results of Prob. 4-16, we should expect the distribution of light about this point of convergence to be precisely the Fraunhofer diffraction pattern of the lens aperture that limits the extent of the spherical wave.

5.3.3 The Relation Between Object and Image

Consider first the nature of the image predicted by geometrical optics. If the imaging system is perfect, then the image is simply an inverted and magnified (or demagnified) replication of the object. Thus according to geometrical optics, the image and object would be related by

$$U_i(u, v) = \frac{1}{|M|} U_o \left(\frac{u}{M}, \frac{v}{M} \right). \quad (5-34)$$

Indeed we can show that our wave optics solution reduces to this geometrical optics solution by using the common artifice of allowing the wavelength λ to approach zero, with the result that (see Prob. 5-15)

$$h(u, v; \xi, \eta) \rightarrow \frac{1}{|M|} \delta \left(\xi - \frac{u}{M}, \eta - \frac{v}{M} \right). \quad (5-35)$$

Substitution of this result in the general superposition equation (5-23) yields (5-34).

The predictions of geometrical optics do not include the effects of diffraction. A more complete understanding of the relation between object and image can be obtained only if such effects are included. Towards this end, we return to the expression (5-33) for the impulse response of the imaging system. As it currently stands, the impulse response is that of a linear space-variant system, so the object and image are related by a superposition integral but not by a convolution integral. This space-variant attribute is a direct result of the magnification and image inversion that occur in the imaging operation. To reduce the object-image relation to a convolution equation, we must normalize the object coordinates to remove inversion and magnification. Let the following normalized object-plane variables be introduced:

$$\tilde{\xi} = M\xi \quad \tilde{\eta} = M\eta$$

in which case the impulse response of (5-33) reduces to

$$h(u, v; \tilde{\xi}, \tilde{\eta}) = \frac{1}{\lambda^2 z_1 z_2} \iint_{-\infty}^{\infty} P(x, y) \times \exp\left\{-j\frac{2\pi}{\lambda z_2} \left[(u - \tilde{\xi})x + (v - \tilde{\eta})y\right]\right\} dx dy, \quad (5-36)$$

which depends only on the differences of coordinates $(u - \tilde{\xi}, v - \tilde{\eta})$.

A final set of coordinate normalizations simplifies the results even further. Let

$$\tilde{x} = \frac{x}{\lambda z_2} \quad \tilde{y} = \frac{y}{\lambda z_2} \quad \tilde{h} = \frac{1}{|M|} h.$$

Then the object-image relationship becomes

$$U_i(u, v) = \iint_{-\infty}^{\infty} \tilde{h}(u - \tilde{\xi}, v - \tilde{\eta}) \left[\frac{1}{|M|} U_o\left(\frac{\tilde{\xi}}{M}, \frac{\tilde{\eta}}{M}\right) \right] d\tilde{\xi} d\tilde{\eta}, \quad (5-37)$$

$$U_i(u, v) = \tilde{h}(u, v) \otimes U_g(u, v) \quad (5-38)$$

where

$$U_g(u, v) = \frac{1}{|M|} U_o\left(\frac{u}{M}, \frac{v}{M}\right) \quad (5-39)$$

is the geometrical-optics prediction of the image, and

$$\tilde{h}(u, v) = \iint_{-\infty}^{\infty} P(\lambda z_2 \tilde{x}, \lambda z_2 \tilde{y}) \exp[-j2\pi(u\tilde{x} + v\tilde{y})] d\tilde{x} d\tilde{y} \quad (5-40)$$

is the point-spread function introduced by diffraction.

There are two main conclusions from the analysis and discussion above:

1. The ideal image produced by a diffraction-limited optical system (i.e. a system that is free from aberrations) is a scaled and inverted version of the object.
2. The effect of diffraction is to convolve that ideal image with the Fraunhofer diffraction pattern of the lens pupil.

The smoothing operation associated with the convolution can strongly attenuate the fine details of the object, with a corresponding loss of image fidelity resulting. Similar effects occur in electrical systems when an input with high-frequency components passes through a filter with a limited frequency response. In the case of electrical systems, the loss of signal fidelity is most conveniently described in the frequency domain. The great utility of frequency-analysis concepts in the electrical case suggests that similar concepts might be usefully employed in the study of imaging systems. The application of filtering concepts to imaging systems is a subject of great importance and will be considered in detail in Chapter 6.

5.4 ANALYSIS OF COMPLEX COHERENT OPTICAL SYSTEMS

In the previous sections we have analyzed several different optical systems. These systems involved at most a single thin lens and at most propagation over two regions of free space. More complex optical systems can be analyzed by using the same methods applied above. However, the number of integrations grows as the number of free-space regions grows, and the complexity of the calculations increases as the number of lenses included grows. For these reasons some readers may appreciate the introduction of a certain "operator" notation that is useful in analyzing complex systems. Not all readers will find this approach attractive, and for those the methods already used can simply be extended to the more complex systems.

5.4.1 An Operator Notation

Several different operator methods for analyzing coherent optical systems have been introduced in the literature. The first was that of **VanderLugt** [292] who exploited certain properties of quadratic-phase exponentials to simplify their manipulation. Later papers by Butterweck [46], and Nazarathy and Shamir [219] used what can be called a true operator notation to simplify calculations. We shall follow the approach of Nazarathy and Shamir here.

There are several simplifying assumptions that will be used in the analysis. As has been the case in previous analyses, we restrict attention to monochromatic light, an assumption that will be seen in the next chapter to limit consideration to what we call "coherent" systems. In addition, only **paraxial** conditions will be considered, a limitation also inherent in the usual geometrical optics treatment using ray matrices, as discussed in Appendix B. Finally, for simplicity we will treat the problems in this section as one-dimensional problems rather than two-dimensional problems. For problems

with apertures that are separable in rectangular coordinates, this is not a significant restriction, since the separability of quadratic-phase exponentials allows each of two orthogonal directions to be considered independently. However, if the optical system contains apertures that are not separable in rectangular coordinates, a two-dimensional extension of the treatment is necessary. This extension is not difficult to make, but we will not present it here.

The operator approach is based on several fundamental operations, each of which is represented by an "operator". Most operators have parameters that depend on the geometry of the optical system being analyzed. Parameters are included within square brackets [] following the operator. The operators act on the quantities contained in curly brackets { }.

The basic operators of use to us here are as follows:

Multiplication by a quadratic-phase exponential. The definition of the operator \mathcal{Q} is

$$\mathcal{Q}[c]\{U(x)\} = e^{j\frac{k}{2}cx^2}U(x), \quad (5-41)$$

where $k = 2\pi/\lambda$ and c is an inverse length. The inverse of $\mathcal{Q}[c]$ is $\mathcal{Q}[-c]$.

Scaling by a constant. This operator is represented by the symbol \mathcal{V} and is defined by

$$\mathcal{V}[b]\{U(x)\} = b^{1/2}U(bx), \quad (5-42)$$

where b is dimensionless. The inverse of $\mathcal{V}[b]$ is $\mathcal{V}[1/b]$.

Fourier transformation. This operator is represented by the usual symbol \mathcal{F} and is defined by

$$\mathcal{F}\{U(x)\} = \int_{-\infty}^{\infty} U(x) e^{-j2\pi f x} dx. \quad (5-43)$$

The inverse Fourier transform operator is defined in the usual way, i.e. with a change of the sign of the exponent.

Free-space propagation. Free-space propagation is represented by the operator \mathcal{R} , which is defined by the equation

$$\mathcal{R}[d]\{U(x_1)\} = \frac{1}{\sqrt{j\lambda d}} \int_{-\infty}^{\infty} U(x_1) e^{j\frac{k}{2d}(x_2-x_1)^2} dx_1, \quad (5-44)$$

where d is the distance of propagation and x_2 is the coordinate that applies after propagation. The inverse of $\mathcal{R}[d]$ is $\mathcal{R}[-d]$.

These four operators are sufficient for analyzing most optical systems. Their utility arises from some simple properties and certain relations between them. These properties and relations allow complicated chains of operators to be reduced to simple results, as will shortly be illustrated. Some simple and useful properties are listed below:

$$\mathcal{V}[t_2] \mathcal{V}[t_1] = \mathcal{V}[t_2 t_1] \quad (5-45)$$

$$\mathcal{F} \mathcal{V}[t] = \mathcal{V}\left[\frac{1}{t}\right] \mathcal{F} \quad (5-46)$$

$$\mathcal{F} \mathcal{F} = \mathcal{V}[-1] \quad (5-47)$$

$$\mathcal{Q}[c_2] \mathcal{Q}[c_1] = \mathcal{Q}[c_2 + c_1] \quad (5-48)$$

$$\mathcal{R}[d] = \mathcal{F}^{-1} \mathcal{Q}[-\lambda^2 d] \mathcal{F} \quad (5-49)$$

$$\mathcal{Q}[c] \mathcal{V}[t] = \mathcal{V}[t] \mathcal{Q}\left[\frac{c}{t^2}\right] \quad (5-50)$$

Relations (5-45) and (5-48) are quite obvious and simple to prove. Relation (5-46) is a statement of the similarity theorem of Fourier analysis, while relation (5-47) follows from the Fourier inversion theorem, slightly modified to account for the fact that both transforms are in the forward direction. Relation (5-49) is a statement that free-space propagation over distance d can be analyzed either by a Fresnel diffraction equation or by a sequence of Fourier transformation, multiplication by the transfer function of free space, and inverse Fourier transformation. The left-hand and right-hand sides of relation (5-50) are shown to be equal simply by writing out their definitions.

A slightly more sophisticated relation is

$$\mathcal{R}[d] = \mathcal{Q}\left[\frac{1}{d}\right] \mathcal{V}\left[\frac{1}{\lambda d}\right] \mathcal{F} \mathcal{Q}\left[\frac{1}{d}\right], \quad (5-51)$$

which is a statement that the Fresnel diffraction operation is equivalent to **premultiplication** by a quadratic-phase exponential, a properly scaled Fourier transform, and **postmultiplication** by a quadratic-phase exponential. Another relation of similar complexity is

$$\mathcal{V}\left[\frac{1}{\lambda f}\right] \mathcal{F} = \mathcal{R}[f] \mathcal{Q}\left[-\frac{1}{f}\right] \mathcal{R}[f], \quad (5-52)$$

which is a statement that the fields across the front and back focal planes of a positive lens are related by a properly scaled Fourier transform, with no quadratic-phase exponential multiplier, as proved earlier in this chapter.

Many useful relations between operators are summarized in Table 5.1. With these relations to draw on, we are now ready to apply the operator notation to some simple optical systems.

5.4.2 Application of the Operator Approach to Some Optical Systems

We illustrate the use of the operator notation by analyzing two optical geometries that have not yet been treated. The first is fairly simple, consisting of two spherical lenses, each with the same focal length f , with a separation off between them, as shown in Fig. 5.11. The goal is to determine the relationship between the complex field across a plane \mathcal{S}_1 just to the left of lens L_1 , and the complex field across a plane \mathcal{S}_2 just to the right

TABLE 5.1
Relations between operators.

	\mathcal{V}	\mathcal{F}	\mathcal{Q}	\mathcal{R}
\mathcal{V}	$\mathcal{V}[t_2]\mathcal{V}[t_1] = \mathcal{V}[t_2t_1]$	$\mathcal{V}[t]\mathcal{F} = \mathcal{F}\mathcal{V}\left[\frac{1}{t}\right]$	$\mathcal{V}[t]\mathcal{Q}[c] = \mathcal{Q}[t^2c]\mathcal{V}[t]$	$\mathcal{V}[t]\mathcal{R}[d] = \mathcal{R}\left[\frac{d}{t^2}\right]\mathcal{V}[t]$
\mathcal{F}	$\mathcal{F}\mathcal{V}[t] = \mathcal{V}\left[\frac{1}{t}\right]\mathcal{F}$	$\mathcal{F}\mathcal{F} = \mathcal{V}[-1]$	$\mathcal{F}\mathcal{Q}[c] = \mathcal{R}\left[-\frac{c}{\lambda^2}\right]\mathcal{F}$	$\mathcal{F}\mathcal{R}[d] = \mathcal{Q}[-\lambda^2d]\mathcal{F}$
\mathcal{Q}	$\mathcal{Q}[c]\mathcal{V}[t] = \mathcal{V}[t]\mathcal{Q}\left[\frac{c}{t^2}\right]$	$\mathcal{Q}[c]\mathcal{F} = \mathcal{F}\mathcal{R}\left[-\frac{c}{\lambda^2}\right]$	$\mathcal{Q}[c_2]\mathcal{Q}[c_1] = \mathcal{Q}[c_2 + c_1]$	$\mathcal{Q}[c]\mathcal{R}[d] = \mathcal{R}\left[(d^{-1} + c)^{-1}\right] \cdot \mathcal{V}[1 + cd] \cdot \mathcal{Q}\left[(c^{-1} + d)^{-1}\right]$
\mathcal{R}	$\mathcal{R}[d]\mathcal{V}[t] = \mathcal{V}[t]\mathcal{R}[t^2d]$	$\mathcal{R}[d]\mathcal{F} = \mathcal{F}\mathcal{Q}[-\lambda^2d]$	$\mathcal{R}[d]\mathcal{Q}[c] = \mathcal{Q}\left[(c^{-1} + d)^{-1}\right] \cdot \mathcal{V}[(1 + cd)^{-1}] \cdot \mathcal{R}\left[(d^{-1} + c)^{-1}\right]$	$\mathcal{R}[d_2]\mathcal{R}[d_1] = \mathcal{R}[d_1 + d_2]$

of the lens L_2 . We represent that relationship by a system operator \mathbf{S} . The first operation on the wave takes place as it passes through L_1 , and this operation is represented by the operator $\mathcal{Q}\left[-\frac{1}{f}\right]$. The second operation is propagation through space over distance f , represented by the operator $\mathcal{R}[f]$. The third operation is passage through the lens L_2 , which is represented by the operator $\mathcal{Q}\left[-\frac{1}{f}\right]$. Thus the entire sequence of operations is represented by the operator chain

$$\mathbf{S} = \mathcal{Q}\left[-\frac{1}{f}\right] \mathcal{R}[f] \mathcal{Q}\left[-\frac{1}{f}\right]. \tag{5-53}$$

This set of operators can be simplified by means of Eq. (5-51) applied to $\mathcal{R}[f]$, as now demonstrated,

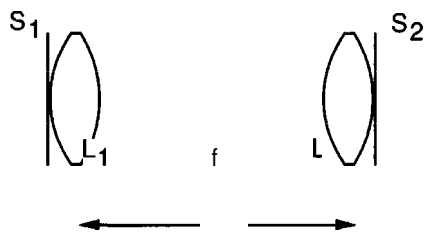


FIGURE 5.11

First problem analyzed.

$$\begin{aligned} \mathcal{S} &= \mathcal{Q}\left[-\frac{1}{f}\right] \mathcal{R}[f] \mathcal{Q}\left[-\frac{1}{f}\right] = \mathcal{Q}\left[-\frac{1}{f}\right] \mathcal{Q}\left[\frac{1}{f}\right] \mathcal{V}\left[\frac{1}{\lambda f}\right] \mathcal{F} \mathcal{Q}\left[\frac{1}{f}\right] \mathcal{Q}\left[-\frac{1}{f}\right] \\ &= \mathcal{V}\left[\frac{1}{\lambda f}\right] \mathcal{F}, \end{aligned}$$

where the relations

$$\mathcal{Q}\left[-\frac{1}{f}\right] \mathcal{Q}\left[\frac{1}{f}\right] = \mathcal{Q}\left[\frac{1}{f}\right] \mathcal{Q}\left[-\frac{1}{f}\right] = 1$$

have been used to simplify the equation. Thus we see that this system of two lenses separated by their common focal length f performs a scaled optical Fourier transform, without quadratic-phase exponentials in the result, similar to the focal-plane-to-focal-plane relationship derived earlier. Stating the result explicitly in terms of the input and output fields,

$$U_f(u) = \frac{1}{\sqrt{\lambda f}} \int_{-\infty}^{\infty} U_o(x) e^{-j\frac{k}{f}xu} dx, \quad (5-54)$$

where U_o is the field just to the left of L_1 and U_f is the field just to the right of L_2 .

The second example we would classify as a complex optical system. Even though its appearance may be simple, the analysis required to find its properties is relatively complex. In addition, the information gleaned from its solution is quite revealing. As shown in Fig. 5.12, the system contains only a single lens. However, the object or the input to the system, located distance d to the left of the lens, is illuminated by a diverging spherical wave, emanating from a point that is distance $z_1 > d$ to the left of the lens. The output of interest here will be in the plane where the point source is imaged, at distance z_2 to the right of the lens, where z_1 , z_2 , and the focal length f of the lens satisfy the usual lens law, $z_1^{-1} + z_2^{-1} - f^{-1} = 0$.

The sequence of operators describing this system is

$$\mathcal{S} = \mathcal{R}[z_2] \mathcal{Q}\left[-\frac{1}{f}\right] \mathcal{R}[d] \mathcal{Q}\left[\frac{1}{z_1 - d}\right], \quad (5-55)$$

where the \mathcal{Q} operator furthest to the right represents the fact that the input is illuminated by a diverging spherical wave, the \mathcal{R} operator second from the right represents propagation over distance d to the lens, the \mathcal{Q} operator next to the left represents the effect of the positive lens, and the operator \mathcal{R} furthest to the left represents the final propagation over distance z_2 . It is simplest to apply the lens law immediately, replacing $\mathcal{Q}[-1/f]$ by $\mathcal{Q}[-1/z_1 - 1/z_2]$.

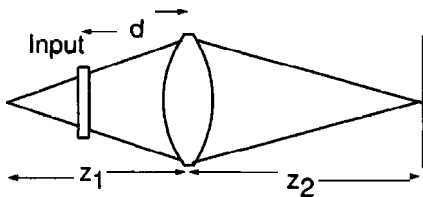


FIGURE 5.12
Second problem analyzed.

There are several different ways to simplify this sequence of operators. Our approach will be to first use the relationship in the 4th row and 3rd column of Table 5.1 to replace the two operators furthest to the left as follows:

$$\mathcal{R}[z_2] \mathcal{Q} \left[-\frac{1}{z_1} - \frac{1}{z_2} \right] = \mathcal{Q} \left[\frac{z_1 + z_2}{z_2^2} \right] \mathcal{V} \left[-\frac{z_1}{z_2} \right] \mathcal{R}[-z_1].$$

The two remaining adjacent \mathcal{R} operators can now be combined using the relation given in the 4th row and 4th column of Table 5.1. The operator sequence is now of the form

$$\mathcal{S} = \mathcal{Q} \left[\frac{z_1 + z_2}{z_2^2} \right] \mathcal{V} \left[-\frac{z_1}{z_2} \right] \mathcal{R}[d - z_1] \mathcal{Q} \left[\frac{1}{z_1 - d} \right].$$

Next Eq. (5-51) is applied to write

$$\mathcal{R}[d - z_1] = \mathcal{Q} \left[\frac{1}{d - z_1} \right] \mathcal{V} \left[\frac{1}{\lambda(d - z_1)} \right] \mathcal{F} \mathcal{Q} \left[\frac{1}{d - z_1} \right].$$

Substitution of this result yields an operator sequence

$$\mathcal{S} = \mathcal{Q} \left[\frac{z_1 + z_2}{z_2^2} \right] \mathcal{V} \left[-\frac{z_1}{z_2} \right] \mathcal{Q} \left[\frac{1}{d - z_1} \right] \mathcal{V} \left[\frac{1}{\lambda(d - z_1)} \right] \mathcal{F},$$

where the last two \mathcal{Q} operators on the right canceled each other. The last steps are to apply the relation (5-50) to invert the order of the \mathcal{V} and \mathcal{Q} operators in the middle of the chain, following which the two adjacent \mathcal{V} operators and the two adjacent \mathcal{Q} operators can be combined. With some algebra the final result becomes

$$\mathcal{S} = \mathcal{Q} \left[\frac{(z_1 + z_2)d - z_1 z_2}{z_2^2(d - z_1)} \right] \mathcal{V} \left[\frac{z_1}{\lambda z_2(z_1 - d)} \right] \mathcal{F}. \quad (5-56)$$

A more conventional statement of the relationship between in the input field $U_1(\xi)$ and the output field $U_2(u)$ is

$$U_2(u) = \frac{\exp \left[j \frac{k}{2} \frac{(z_1 + z_2)d - z_1 z_2}{z_2^2(d - z_1)} u^2 \right]}{\sqrt{\frac{\lambda z_2(z_1 - d)}{z_1}}} \int_{-\infty}^{\infty} U_1(\xi) \exp \left[-j \frac{2\pi z_1}{\lambda z_2(z_1 - d)} u \xi \right] d\xi. \quad (5-57)$$

Thus the field is again seen to be a Fourier transform of the input amplitude distribution. The results of this analysis reveal some important general facts not explicitly evident in our earlier analyses. We emphasize these results because of their generality:

The Fourier transform plane need not be the focal plane of the lens performing the transform! Rather, the Fourier transform always appears in the plane where the source is imaged.

While it is not obvious without some further thought and analysis, our results show that the quadratic-phase factor preceding the Fourier transform operation is always the quadratic-phase factor that would result at the transform plane from a point source of light located on the optical axis in the plane of the input transparency.

The result presented in Eq. (5-57) can be shown to reduce to the results of the previous cases considered if z_1 , z_2 , and d are properly chosen to represent those cases (see Prob. 5-16).

We conclude with a few general comments about the operator method of analysis. Its advantage is that it allows a methodical approach to complex calculations that might otherwise be difficult to treat by the conventional methods. However, the method also has some drawbacks. Being one step more abstract than the diffraction integrals it replaces, the operator method is one step further from the physics of the experiment under analysis. Second, to save time with the operator approach, it is necessary that one be rather familiar with the operator relations of Table 5.1. Good intuition about which operator relations to use on a given problem comes only after experience with the method.

PROBLEMS—CHAPTER 5

- 5-1. Show that the focal lengths of double-convex, plano-convex, and positive meniscus lenses are always positive, while the focal lengths of double-concave, plano-concave, and negative meniscus lenses are always negative.
- 5-2. Consider a thin lens that is composed of a portion of a cylinder, as shown in Fig. P5.2.

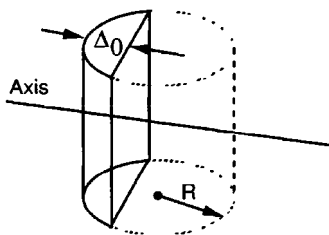


FIGURE P5.2

- (a) Find a paraxial approximation to the phase transformation introduced by a lens of this form.
- (b) What is the effect of such a lens on a plane wave traveling down the optical axis?
- 5-3. A prism (illustrated in Fig. P5.3), which deflects the direction of propagation of a normally incident plane wave to angle θ with respect to the optical axis (the z axis) in the (y, z) plane, can be represented mathematically by an amplitude transmittance

$$t_p(x, y) = \exp\left[-j\frac{2\pi}{\lambda} \sin(\theta)y\right]$$

- (a) Consider a thin transmitting structure with amplitude transmittance given by

$$t_A(x, y) = \exp\{-j\pi[a^2x^2 + (by + c)^2]\},$$

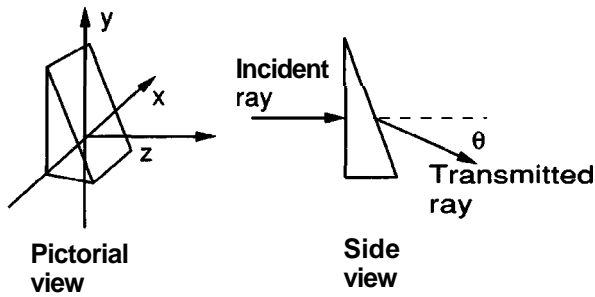


FIGURE P5.3

with a, b, c all real and positive constants. It is claimed that this structure can be considered to consist of a sequence of one spherical lens, one cylindrical lens, and one prism, all placed in contact. Describe such a combination of thin elements that yields this transmittance, specifying the focal lengths of the lenses and the angle of deflection of the prism in terms of a, b, c , and the wavelength λ .

- (b) Can you think of a way to use two cylindrical lenses to achieve an amplitude transmittance

$$t_A(x, y) = \exp(-j\pi dxy)$$

where d is a constant? Explain your conclusion.

- 5-4. Consider a lens that consists of the portion of a cone illustrated in Fig. P5.4.

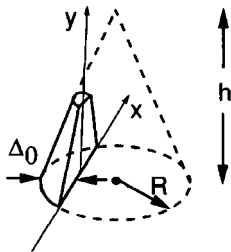


FIGURE P5.4

- (a) Show that a **paraxial** approximation to the phase transformation introduced by such a lens is (under the thin lens assumption)

$$t_l(x, y) = \exp \left\{ jk \left[n\Delta_0 - \frac{(n-1)Ry}{h} - \frac{x^2}{2f(y)} \right] \right\}$$

where

$$f(y) = \frac{R(1 - y/h)}{n - 1}.$$

- (b) What is the effect of such a lens on a plane wave traveling normal to the (x, y) plane?

- 5-5. An input function U_o , bounded by a circular aperture of diameter D and illuminated by a normally incident plane wave, is placed in the front focal plane of a circular positive lens of diameter L . The intensity distribution is measured across the back focal plane of the lens. Assuming $L > D$:

- (a) Find an expression for the maximum spatial frequency of the input for which the measured intensity accurately represents the squared modulus of the input's Fourier spectrum (free from the effects of vignetting).
- (b) What is the numerical value of that spatial frequency (in **cycles/mm**) when $L = 4$ cm, $D = 2$ cm, f (focal length) = 50 cm, and $A = 6 \times 10^{-7}$ meters?
- (c) Above what frequency does the measured spectrum vanish, in spite of the fact that the input may have nonzero Fourier components at such frequencies?

5-6. An array of one-dimensional input functions can be represented by $U_o(\xi, \eta_k)$, where $\eta_1, \eta_2, \dots, \eta_k, \dots, \eta_N$ are N fixed η coordinates in the input plane. It is desired to perform a one-dimensional Fourier transform of all N functions in the ξ direction, yielding an array of transforms

$$G_o(f_x, \eta_k) = \int_{-\infty}^{\infty} U_o(\xi, \eta_k) \exp(-j2\pi f_x \xi) d\xi.$$

Neglecting the finite extent of the lens and object apertures, use the Fourier transforming and imaging properties of lenses derived in this chapter to show how this can be done with

- (a) Two cylindrical lenses of different focal lengths.
- (b) A cylindrical and a spherical lens of the same focal length.

SIMPLIFICATION: You need only display $|G_o|^2$, so phase factors may be dropped.

5-7. A normally incident, unit-amplitude, monochromatic plane wave illuminates a converging lens of 5 cm diameter and 2 meters focal length (see Fig. P5.7). One meter behind the lens and centered on the lens axis is placed an object with amplitude transmittance

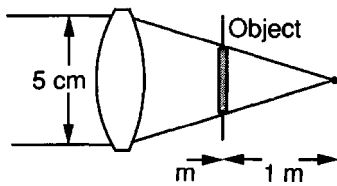


FIGURE P5.7

$$t_A(\xi, \eta) = \frac{1}{2} [1 + \cos(2\pi f_o \xi)] \text{rect}\left(\frac{\xi}{L}\right) \text{rect}\left(\frac{\eta}{L}\right).$$

Assuming $L = 1$ cm, $A = 0.633 \mu\text{m}$, and $f_o = 10$ cycles/mm, sketch the intensity distribution across the u axis of the focal plane, labeling the numerical values of the distance between diffracted components and the width (between first zeros) of the individual components.

5-8. In Fig. P5.8, a monochromatic point source is placed a fixed distance z_1 to the left of a positive lens (focal length f), and a transparent object is placed a variable distance d to the left of the lens. The distance z_1 is greater than f . The Fourier transform and the image of the object appear to the right of the lens.

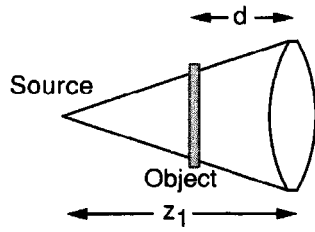


FIGURE P5.8

- (a) How large should the distance d be (in terms of z_1 and f) to assure that the *Fourier plane* and the *object* are equidistant from the lens?
- (b) When the object has the distance found in part (a) above, how far to the right of the lens is its image and what is the magnification of that image?

5-9. A unit-amplitude, normally incident, monochromatic plane wave illuminates an object of maximum linear dimension D , situated immediately in front of a larger positive lens of focal length f (see Fig. P5.9). Due to a positioning error, the intensity distribution is measured across a plane at a distance $f - A$ behind the lens. How small must A be if the measured intensity distribution is to accurately represent the Fraunhofer diffraction pattern of the object?

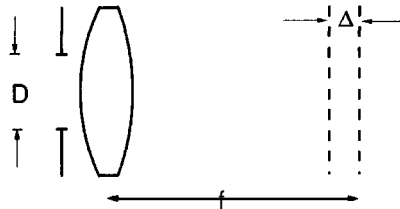


FIGURE P5.9

5-10. Consider the optical system shown in Fig. P5.10. The object on the left is illuminated by a normally incident plane wave. Lens L_1 is a *negative* lens with focal length $-f$, and lens L_2 is a *positive* lens with focal length f . The two lenses are spaced by distance f . Lens L_1 is a distance $2f$ to the right of the object. Use the simplest possible reasoning to predict the distances d and z_2 , respectively, to the Fourier plane and the image plane to the right or left of lens L_2 (specify right or left in the answers).

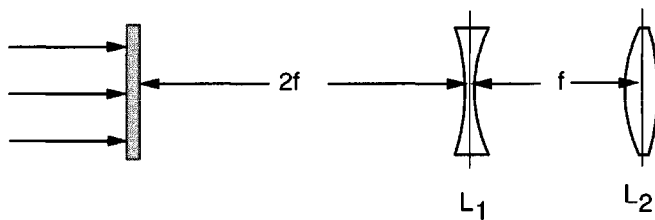


FIGURE P5.10

5-11. In the optical system shown in Fig. P5.11, specify the locations of all Fourier and image planes to the left and right of the lens. The lens shown is positive and has focal length f . The illumination of the object is a converging spherical wave, as indicated.

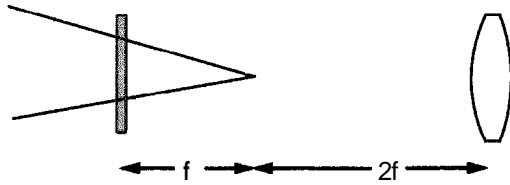


FIGURE P5.11

5-12. With reference to Eq. (5-30):

- (a) At what radius r_0 in the object plane has the phase of $\exp[j\frac{k}{2z_1}(\xi^2 + \eta^2)]$ changed by 1 radian from its value at the origin?
- (b) Assuming a circular pupil function of radius R , what is the radius (in the object plane) to the first zero of the impulse response h , assuming that the observation point in the image space is the origin?
- (c) From the results obtained so far, what relation between R , A , and z_1 will allow the quadratic-phase exponential $\exp[j\frac{k}{2z_1}(\xi^2 + \eta^2)]$ to be replaced by a single complex number, assuming observation near the lens axis?

5-13. A diffracting structure has a circularly symmetric amplitude transmittance function given by

$$t_A(r) = \left(\frac{1}{2} + \frac{1}{2} \cos \gamma r^2\right) \text{circ}\left(\frac{r}{R}\right).$$

- (a) In what way does this screen act like a lens?
- (b) Give an expression for the focal length of the screen.
- (c) What characteristics might seriously limit the use of this screen as an imaging device for polychromatic objects?

5-14. A certain diffracting screen with an amplitude transmittance

$$t_A(r) = \left[\frac{1}{2} + \frac{1}{2} \text{sgn}(\cos \gamma r^2)\right] \text{circ}\left(\frac{r}{R}\right)$$

is normally illuminated by a unit-amplitude, monochromatic plane wave. Show that the screen acts as a lens with multiple focal lengths. Specify the values of these focal lengths and the relative amounts of optical power brought to focus in the corresponding focal planes. (A diffracting structure such as this is known as a *Fresnel* zone plate. Hint: The square wave shown in Fig P5.14 can be represented by the Fourier series

$$f(x) = \sum_{n=-\infty}^{\infty} \left[\frac{\sin(\pi n/2)}{\pi n}\right] \exp\left(j\frac{2\pi nx}{X}\right).$$

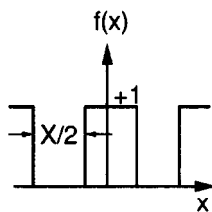


FIGURE P5.14

5-15. Show that in the limit $A \rightarrow 0$, **Eq. (5-33)** approaches the impulse response shown in **Eq. (5-35)**.

5-16. Find the form of the general result of **Eq. (5-57)** under the following limiting conditions:

- (a) $z_1 \rightarrow \infty$ and $d \rightarrow 0$.
- (b) $z_1 \rightarrow \infty$ and $d \rightarrow f$
- (c) $z_1 \rightarrow \infty$, general distance d .

5-17. Consider the simple optical system shown in Fig. P5.17.

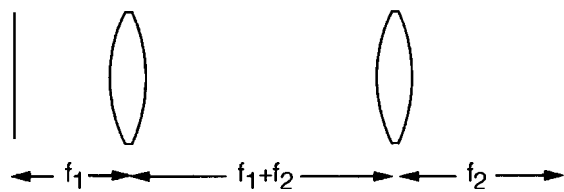


FIGURE P5.17

- (a) Write the operator sequence that describes the successive propagation between planes and through lenses for this system.
- (b) Reduce this operator sequence to a simple scaling operator.

Frequency Analysis of Optical Imaging Systems

Considering the long and rich history of optics, the tools of frequency analysis and linear systems theory have played important roles for only a relatively short period of time. Nevertheless, in this short time these tools have been so widely and successfully used that they now occupy a fundamental place in the theory of imaging systems.

A realization of the utility of Fourier methods in the analysis of optical systems arose rather spontaneously in the late 1930's when a number of workers began to advocate the use of sinusoidal test patterns for system evaluation. Much of the initial stimulus was provided by a French scientist, P.M. Duffieux, whose work culminated in the publication of a book, in 1946, on the use of Fourier methods in optics [86]. This book has recently been translated into English [87]. In the United States, much of the interest in these topics was stimulated by an electrical engineer, Otto Schade, who very successfully employed methods of linear systems theory in the analysis and improvement of television camera lenses [255]. In the United Kingdom, H.H. Hopkins led the way in the use of transfer function methods for the assessment of the quality of optical imaging systems, and was responsible for many of the first calculations of transfer functions in the presence of common aberrations [146]. However, it must be said that the foundations of Fourier optics were laid considerably earlier, particularly in the works of Ernst Abbe (1840-1905) and Lord Rayleigh (1842-1919).

In this chapter we shall consider the role of Fourier analysis in the theory of coherent and incoherent imaging. While historically the case of incoherent imaging has been the more important one, nonetheless the case of coherent imaging has always been important in microscopy, and it gained much additional importance with the advent of the laser. For example, the field of holography is predominantly concerned with coherent imaging.

For additional discussions of various aspects of the subject matter to follow, the reader may wish to consult any of the following references: [223], [103], [196], [76], [300].

6.1 GENERALIZED TREATMENT OF IMAGING SYSTEMS

In the preceding chapter, the imaging properties of a single thin positive lens were studied for the case of monochromatic illumination. In the material to follow, we shall first broaden our discussion beyond a single thin positive lens, finding results applicable to more general systems of lenses, and then remove the restriction to monochromatic light, obtaining results for "quasi-monochromatic" light, both spatially coherent and spatially incoherent. To broaden the perspective, it will be necessary to draw upon some results from the theory of geometrical optics. The necessary concepts are all introduced in Appendix B.

6.1.1 A Generalized Model

Suppose that an imaging system of interest is composed, not of a single thin lens, but perhaps of several lenses, some positive, some negative, with various distances between them. The lenses need not be thin in the sense defined earlier. We shall assume, however, that the system ultimately produces a *real* image in space; this is not a serious restriction, for if the system produces a virtual image, to view that image it must be converted to a real image, perhaps by the lens of the eye.

To specify the properties of the lens system, we adopt the point of view that all imaging elements may be lumped into a single "black box", and that the significant properties of the system can be completely described by specifying only the *terminal properties* of the aggregate. Referring to Fig. 6.1, the "terminals" of this black box consist of the planes containing the entrance and exit pupils (see Appendix B for a discussion of these planes).¹ It is assumed that the passage of light between the entrance pupil and the exit pupil is adequately described by geometrical optics.

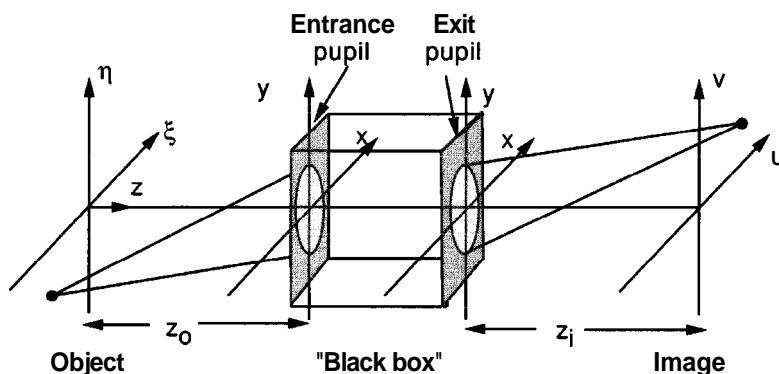


FIGURE 6.1
Generalized model of an imaging system.

¹In general it is not necessary that the entrance pupil lie to the left of the exit pupil as shown in Fig. 6.1. However the conceptual *idea* of a system mapping the light incident on the entrance pupil to the light leaving the exit pupil remains valid, regardless of the order of the two pupils.

The entrance and exit pupils are in fact images of the same limiting aperture within the system. As a consequence there are several different ways to visualize the origin of the spatial limitation of the wavefront that ultimately gives rise to diffraction. It can be viewed as being caused by the physical limiting aperture internal to the system (which is the true physical source of the limitation). Equivalently it can be viewed as arising from the entrance pupil or from the exit pupil of the system.

Throughout this chapter, we shall use the symbol z_o to represent the distance of the plane of the entrance pupil from the object plane, and the symbol z_i to represent the distance of the plane of the exit pupil from the image plane.² The distance z_i is then the distance that will appear in the diffraction equations that represent the effect of diffraction by the exit pupil on the point-spread function of the optical system. We shall refer either to the exit pupil or simply to the "pupil" of the system when discussing these effects.

An imaging system is said to be *diffraction-limited* if a diverging spherical wave, emanating from a point-source object, is converted by the system into a new wave, again perfectly spherical, that converges towards an ideal point in the image plane, where the location of that ideal image point is related to the location of the original object point through a simple scaling factor (the magnification), a factor that must be the same for all points in the image field of interest if the system is to be ideal. Thus the terminal property of a diffraction-limited imaging system is that a diverging spherical wave incident on the entrance pupil is converted by the system into a converging spherical wave at the exit pupil. For any real imaging system, this property will be satisfied, at best, over only finite regions of the object and image planes. If the object of interest is confined to the region for which this property holds, then the system may be regarded as being diffraction-limited.

If in the presence of a point-source object, the wavefront leaving the exit pupil departs significantly from ideal spherical shape, then the imaging system is said to have *aberrations*. Aberrations will be considered in Section 6-4, where it is shown that they lead to defects in the spatial-frequency response of the imaging system.

6.1.2 Effects of Diffraction on the Image

Since geometrical optics adequately describes the passage of light between the entrance and exit pupils of a system, diffraction effects play a role only during passage of light from the object to the entrance pupil, or alternatively and equivalently, from the exit pupil to the image. It is, in fact, possible to associate *all* diffraction limitations with *either* of these two pupils. The two points of view that regard image resolution as being limited by (1) the finite entrance pupil seen from the object space or (2) the finite exit pupil seen from the image space are entirely equivalent, due to the fact that these two pupils are images of each other.

²We reserve the symbols z_1 and z_2 for the distances from the object to the first principal plane and the distance from the second principal plane to the image, respectively.

The view that diffraction effects result from the entrance pupil was first espoused by Ernst Abbe in 1873 [1] in studies of coherent imagery with a microscope. According to the Abbe theory, only a certain portion of the diffracted components generated by a complicated object are intercepted by this finite pupil. The components not intercepted are precisely those generated by the high-frequency components of the object amplitude transmittance. This viewpoint is illustrated in Fig. 6.2 for the case of an object that is a grating with several orders and an imaging system composed of a single positive lens.

A view equivalent to regarding diffraction effects as resulting from the exit pupil was presented by Lord Rayleigh in 1896 [241]. This is the viewpoint that was used in Section 5.3, and we shall adopt it again here.

Again the image amplitude³ is represented by a superposition integral

$$U_i(u, v) = \iint_{-\infty}^{\infty} h(u, v; \xi, \eta) U_o(\xi, \eta) d\xi d\eta, \quad (6-1)$$

where h is the amplitude at image coordinates (u, v) in response to a point-source object at (ξ, η) , and U_o is the amplitude distribution transmitted by the object. In the absence of aberrations, the response h arises from a spherical wave (of limited extent) converging from the exit pupil towards the ideal image point $(u = M\xi, v = M\eta)$. We allow the magnification to be either negative or positive, according to whether the image is inverted or not.

From the result of Prob. 4-16, the discussions of Section 5.3, and in particular, Eq. (5-33), the light amplitude about the ideal image point is simply the Fraunhofer diffraction pattern of the exit pupil, centered on image coordinates $(u = M\xi, v = M\eta)$. Thus

$$h(u, v; \xi, \eta) = \frac{A}{\lambda z_i} \iint_{-\infty}^m P(x, y) \exp \left\{ -j \frac{2\pi}{\lambda z_i} [(u - M\xi)x + (v - M\eta)y] \right\} dx dy, \quad (6-2)$$

where the pupil function P is unity inside and zero outside the projected aperture, A is a constant amplitude, z_i is the distance from the exit pupil to the image plane, and (x, y) are coordinates in the plane of the exit pupil. In writing this equation, we have

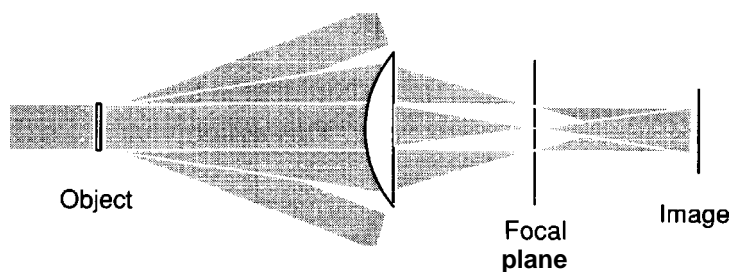


FIGURE 6.2
The Abbe theory of image formation.

³We have retained the assumption of monochromatic illumination but will remove it in the section to follow.

again neglected quadratic phase factors over the object and image planes, as justified in Section 5.3.

In order to achieve space invariance in the imaging operation, it is necessary to remove the effects of magnification and image inversion from the equations. This can be done by defining *reduced coordinates* in the object space⁴ according to

$$\tilde{\xi} = M\xi, \quad \tilde{\eta} = M\eta,$$

in which case the amplitude point-spread function becomes

$$h(u - \tilde{\xi}, v - \tilde{\eta}) = \frac{A}{\lambda z_i} \iint_{-\infty}^{\infty} P(x, y) \exp\left\{-j \frac{2\pi}{\lambda z_i} [(u - \tilde{\xi})x + (v - \tilde{\eta})y]\right\} dx dy.$$

At this point it is convenient to define the *ideal image*, or the geometrical-optics prediction of the image for a perfect imaging system as

$$U_g(\tilde{\xi}, \tilde{\eta}) = \frac{1}{|M|} U_o\left(\frac{\tilde{\xi}}{M}, \frac{\tilde{\eta}}{M}\right), \quad (6-3)$$

yielding a convolution equation for the image,

$$U_i(u, v) = \iint_{-\infty}^{\infty} h(u - \tilde{\xi}, v - \tilde{\eta}) U_g(\tilde{\xi}, \tilde{\eta}) d\tilde{\xi} d\tilde{\eta}, \quad (6-4)$$

where

$$h(u, v) = \frac{A}{\lambda z_i} \iint_{-\infty}^{\infty} P(x, y) \exp\left\{-j \frac{2\pi}{\lambda z_i} (ux + vy)\right\} dx dy. \quad (6-5)$$

Thus in this general case, for a diffraction-limited system we can regard the image as being a convolution of the image predicted by geometrical optics with an impulse response that is the Fraunhofer diffraction pattern of the exit pupil.

6.1.3 Polychromatic Illumination: The Coherent and Incoherent Cases

The assumption of strictly monochromatic illumination has been present in all our discussions of imaging systems up to this point. This assumption is overly restrictive, for the illumination generated by real optical sources, including lasers, is never perfectly monochromatic. The statistical nature of the time variations of illumination amplitude and phase can, in fact, influence the behavior of an imaging system in profound ways.

⁴Often advantages are gained by using much more complex changes of coordinates, particularly when the analysis is nonparaxial. We have chosen to remain with the simplest coordinate system consistent with a paraxially space-invariant system. For discussions of other coordinate mappings (many of which are due to H.H. Hopkins) and their advantages, see Ref. [300].

We therefore digress temporarily to consider the very important effects of polychromaticity.

To treat this subject in a completely satisfactory way, it would be necessary to take a rather long detour through the *theory of partial coherence*. However, for our purposes such a detailed detour would not be practical. We therefore treat the subject from two points of view, one entirely heuristic, and the second more rigorous but not entirely complete. The reader interested in a more complete treatment may wish to consult Refs. [20], [203], [28], or [123].

In the case of monochromatic illumination it was convenient to represent the complex amplitude of the field by a complex phasor U that was a function of space coordinates. When the illumination is polychromatic but narrowband, i.e. occupying a bandwidth that is small compared with its center frequency, this approach can be generalized by representing the field by a *time-varying* phasor that depends on both time and space coordinates. For the narrowband case, the amplitude and phase of the time-varying phasor are readily identified with the envelope and phase of the real optical wave.

Consider the nature of the light that is transmitted by or reflected from an object illuminated by a polychromatic wave. Since the time variations of the phasor amplitude are statistical in nature, only statistical concepts can provide a satisfactory description of the field. As we have seen previously, each object point generates an amplitude impulse response in the image plane. If the amplitude and phase of the light at a particular object point vary randomly with time, then the overall amplitude and phase of the amplitude impulse response will vary in a corresponding fashion. Thus the statistical relationships between the phasor amplitudes at the various points on the object will influence the statistical relationships between the corresponding impulse responses in the image plane. These statistical relationships will greatly affect the result of the time-averaging operation that yields the final image intensity distribution.

We shall consider only two types of illumination here. First, we consider object illumination with the particular property that the phasor amplitudes of the field at all object points vary *in unison*. Thus while any two object points may have different *relative* phases, their absolute phases are varying with time in a perfectly correlated way. Such illumination is called *spatially coherent*. Second, we consider object illumination with the opposite property that the phasor amplitudes at all points on the object are varying in totally uncorrelated fashions. Such illumination is called *spatially incoherent*. (In the future we shall refer to these types of illumination as simply *coherent* or *incoherent*.) Coherent illumination is obtained whenever light appears to originate from a single point.⁵ The most common example of a source of such light is a laser, although more conventional sources (e.g. zirconium arc lamps) can yield coherent light, albeit of weaker brightness than a laser, if their output is first passed through a pinhole. Incoherent light is obtained from diffuse or extended sources, for example gas discharges and the sun.

⁵This is a sufficient but not necessary condition for complete coherence. For example, when light from a point source is passed through a stationary diffuser, the relative phases of the light at any two points behind the diffuser remain correlated. Therefore the transmitted light is still spatially coherent, even though it no longer appears to originate from a point source. Note, however, that before impinging on the diffuser it did originate from a point source.

When the object illumination is coherent, the various impulse responses in the image plane vary in unison, and therefore must be added on a complex amplitude basis. Thus a coherent imaging system is linear in complex amplitude. The results of the monochromatic analysis can therefore be applied directly to such systems, with the understanding that the complex amplitude U is now a time-invariant phasor that depends on the relative phases of the light.

When the object illumination is incoherent, the various impulse responses in the image plane vary in uncorrelated fashions. They must therefore be added on a power or intensity basis. Since the intensity of any given impulse response is proportional to the intensity of the point source that gave rise to it, it follows that an incoherent imaging system is linear in intensity, and the impulse response of such a system is the squared magnitude of the amplitude impulse response.

The preceding arguments have been entirely heuristic, and in fact have certain assumptions and approximations hidden in them. We therefore turn to a more rigorous examination of the problem. To begin, note that in the monochromatic case we obtain the phasor representation of the field by suppressing the positive-frequency component of the cosinusoidal field, and doubling the remaining negative frequency component. To generalize this concept to a polychromatic wave $u(\mathbf{P}, t)$, we suppress all positive-frequency components of its Fourier spectrum, and double its negative-frequency components, yielding a new (complex) function $u_-(\mathbf{P}, t)$. If we further write

$$u_-(\mathbf{P}, t) = U(\mathbf{P}, t) \exp(-j2\pi\bar{\nu}t)$$

where $\bar{\nu}$ represents the mean or center frequency of the optical wave, then the complex function $U(\mathbf{P}, t)$ may be regarded as the time-varying phasor representation of $u(\mathbf{P}, t)$.

Under the narrowband condition assumed above, the amplitude impulse response does not change appreciably for the various frequencies contained within the optical spectrum. Therefore it is possible to express the time-varying phasor representation of the image in terms of the convolution of a wavelength-independent impulse response with the time varying phasor representation of the object (in reduced object coordinates),

$$U_i(u, v; t) = \iint_{-\infty}^{\infty} h(u - \tilde{\xi}, v - \tilde{\eta}) U_g(\tilde{\xi}, \tilde{\eta}; t - \tau) d\tilde{\xi} d\tilde{\eta} \quad (6-6)$$

where τ is a time delay associated with propagation from $(\tilde{\xi}, \tilde{\eta})$ to (u, v) (note that in general, τ is a function of the coordinates involved).

To calculate the image intensity, we must time average the instantaneous intensity represented by $|U_i(u, v; t)|^2$, due to the fact that the detector integration time is usually extremely long compared with the reciprocal of the optical bandwidth, even for narrow-band optical sources. Thus the image intensity is given by $I_i(u, v) = \langle |U_i(u, v; t)|^2 \rangle$, or, after substitution of Eq. (6-6) and interchanging orders of averaging and integration,

$$I_i(u, v) = \iint_{-\infty}^{\infty} d\tilde{\xi}_1 d\tilde{\eta}_1 \iint_{-\infty}^{\infty} d\tilde{\xi}_2 d\tilde{\eta}_2 h(u - \tilde{\xi}_1, v - \tilde{\eta}_1) h^*(u - \tilde{\xi}_2, v - \tilde{\eta}_2) \\ \times \left\langle U_g(\tilde{\xi}_1, \tilde{\eta}_1; t - \tau_1) U_g^*(\tilde{\xi}_2, \tilde{\eta}_2; t - \tau_2) \right\rangle. \quad (6-7)$$

Now for a fixed image point, the impulse response h is nonzero over only a small region about the ideal image point. Therefore the integrand is nonzero only for points $(\tilde{\xi}_1, \tilde{\eta}_1)$ and $(\tilde{\xi}_2, \tilde{\eta}_2)$ that are very close together. Hence we assume that the difference between the time delays τ_1 and τ_2 is negligible under the narrowband assumption, allowing the two delays to be dropped.

The expression for image intensity can now be written

$$I_i(u, v) = \iint_{-\infty}^{\infty} d\tilde{\xi}_1 d\tilde{\eta}_1 \iint_{-\infty}^{\infty} d\tilde{\xi}_2 d\tilde{\eta}_2 h(u - \tilde{\xi}_1, v - \tilde{\eta}_1) h^*(u - \tilde{\xi}_2, v - \tilde{\eta}_2) J_g(\tilde{\xi}_1, \tilde{\eta}_1; \tilde{\xi}_2, \tilde{\eta}_2), \quad (6-8)$$

where

$$J_g(\tilde{\xi}_1, \tilde{\eta}_1; \tilde{\xi}_2, \tilde{\eta}_2) = \left\langle U_g(\tilde{\xi}_1, \tilde{\eta}_1; t) U_g^*(\tilde{\xi}_2, \tilde{\eta}_2; t) \right\rangle \quad (6-9)$$

is known as the *mutual intensity*, and is a measure of the *spatial coherence* of the light at the two object points.

When the illumination is perfectly *coherent*, the time-varying phasor amplitudes across the object plane differ only by complex constants. Equivalently we may write

$$U_g(\tilde{\xi}_1, \tilde{\eta}_1; t) = U_g(\tilde{\xi}_1, \tilde{\eta}_1) \frac{U_g(0, 0; t)}{\langle |U_g(0, 0; t)|^2 \rangle^{\frac{1}{2}}} \quad (6-10)$$

$$U_g(\tilde{\xi}_2, \tilde{\eta}_2; t) = U_g(\tilde{\xi}_2, \tilde{\eta}_2) \frac{U_g(0, 0; t)}{\langle |U_g(0, 0; t)|^2 \rangle^{\frac{1}{2}}}$$

where the phase of the time-varying phasor at the origin has arbitrarily been chosen as a phase reference, the time-independent U_g are phasor amplitudes *relative* to the time varying phasor amplitude at the origin, and the normalizations have been performed to allow the time-independent phasors to retain correct information about the average power or intensity. Substituting these relations in the definition of mutual intensity, Eq. (6-9), for the coherent case we obtain

$$J_g(\tilde{\xi}_1, \tilde{\eta}_1; \tilde{\xi}_2, \tilde{\eta}_2) = U_g(\tilde{\xi}_1, \tilde{\eta}_1) U_g^*(\tilde{\xi}_2, \tilde{\eta}_2). \quad (6-11)$$

When this result is in turn substituted into Eq. (6-8) for the intensity, the result is

$$I_i(u, v) = \left| \iint_{-\infty}^{\infty} h(u - \tilde{\xi}, v - \tilde{\eta}) U_g(\tilde{\xi}, \tilde{\eta}) d\tilde{\xi} d\tilde{\eta} \right|^2. \quad (6-12)$$

Finally, defining a time-invariant phasor amplitude U_i in the image space relative to the corresponding phasor amplitude at the origin, the coherent imaging system is found to be described by an amplitude convolution equation,

$$U_i(u, v) = \iint_{-\infty}^{\infty} h(u - \tilde{\xi}, v - \tilde{\eta}) U_g(\tilde{\xi}, \tilde{\eta}) d\tilde{\xi} d\tilde{\eta}, \quad (6-13)$$

the same result obtained in the monochromatic case. We thus confirm that coherent object illumination yields an imaging system that is linear in *complex amplitude*.

When the object illumination is perfectly *incoherent*, the phasor amplitudes across the object vary in statistically independent fashions. This idealized property may be represented by the equation

$$\left\langle U_g(\tilde{\xi}_1, \tilde{\eta}_1; t) U_g^*(\tilde{\xi}_2, \tilde{\eta}_2; t) \right\rangle = \kappa I_g(\tilde{\xi}_1, \tilde{\eta}_1) \delta(\tilde{\xi}_1 - \tilde{\xi}_2, \tilde{\eta}_1 - \tilde{\eta}_2) \quad (6-14)$$

where κ is a real constant. Such a representation is not exact; in actuality, the minimum distance over which coherence can exist is of the order of one wavelength (see Ref. [20], Section 4.4, for more details). Nonetheless, provided the coherence area on the object is small compared with a resolution cell size in object space, Eq. (6-14) is accurate. When used in Eq. (6-9), the result

$$I_i(u, v) = \kappa \iint_{-\infty}^{\infty} |h(u - \tilde{\xi}, v - \tilde{\eta})|^2 I_g(\tilde{\xi}, \tilde{\eta}) d\tilde{\xi} d\tilde{\eta} \quad (6-15)$$

is obtained. Thus for incoherent illumination, the image intensity is found as a convolution of the *intensity impulse response* $|h|^2$ with the ideal image intensity I . Hence we have confirmed that an incoherent imaging system is linear in *intensity*, rather than amplitude. Furthermore, the impulse response of the incoherent mapping is just the squared modulus of the amplitude impulse response.

When the source of illumination is an extended incoherent source, it is possible to specify the conditions under which the imaging system will behave substantially as an incoherent system and substantially as a coherent system (see Ref [123], page 324). Let θ_s represent the effective angular diameter of the incoherent source that illuminates the object, θ_p the angular diameter of the entrance pupil of the imaging system, and θ_o the angular diameter of the angular spectrum of the object, all angles being measured from the object plane. Then the system can be shown to behave as an incoherent system provided

$$\theta_s \geq \theta_o + \theta_p$$

and will behave a coherent system when

$$\theta_s \ll \theta_p.$$

For conditions between these extremes, the system will behave as a *partially coherent* system, the treatment of which is beyond the scope of this discussion. For information on partially coherent imaging systems, see, for example, [123].

6.2 FREQUENCY RESPONSE FOR DIFFRACTION-LIMITED COHERENT IMAGING

We turn now to the central topic of this chapter, the frequency analysis of imaging systems. Attention in this section is devoted to imaging systems with coherent illumination. Systems with incoherent illumination will be treated in Section 6.3.

As emphasized previously, a coherent imaging system is linear in complex amplitude. This implies, of course, that such a system provides a highly nonlinear intensity mapping. If frequency analysis is to be applied in its usual form, it must be applied to the linear *amplitude* mapping.

6.2.1 The Amplitude Transfer Function

Our analysis of coherent systems has yielded a space-invariant form of the amplitude mapping, as evidenced by the convolution equation (6-13). We would anticipate, then, that transfer-function concepts can be applied directly to this system, provided it is done on an amplitude basis. To do so, define the following frequency spectra⁶ of the input and output, respectively:

$$G_g(f_X, f_Y) = \iint_{-\infty}^{\infty} U_g(u, v) \exp[-j2\pi(f_X u + f_Y v)] du dv \quad (6-16)$$

$$G_i(f_X, f_Y) = \iint_{-\infty}^{\infty} U_i(u, v) \exp[-j2\pi(f_X u + f_Y v)] du dv.$$

In addition, define the *amplitude transfer function H* as the Fourier transform of the space-invariant amplitude impulse response,

$$H(f_X, f_Y) = \iint_{-\infty}^{\infty} h(u, v) \exp[-j2\pi(f_X u + f_Y v)] du dv. \quad (6-17)$$

Now applying the convolution theorem to (6-13), it follows directly that

$$G_i(f_X, f_Y) = H(f_X, f_Y) G_g(f_X, f_Y). \quad (6-18)$$

Thus the effects of the diffraction-limited imaging system have been expressed, at least formally, in the frequency domain. It now remains to relate *H* more directly to the physical characteristics of the imaging system itself.

To this end, note that while Eq. (6-17) defines *H* as the Fourier transform of the amplitude point-spread function *h*, this latter function is itself a Fraunhofer diffraction pattern and can be expressed as a scaled Fourier transform of the pupil function (cf. Eq. 6-5). Thus

$$\begin{aligned} H(f_X, f_Y) &= \mathcal{F} \left\{ \frac{A}{\lambda z_i} \iint_{-\infty}^{\infty} P(x, y) \exp \left\{ -j \frac{2\pi}{\lambda z_i} (ux + vy) \right\} dx dy \right\} \\ &= (A\lambda z_i) P(-\lambda z_i f_X, -\lambda z_i f_Y). \end{aligned} \quad (6-19)$$

⁶Here and throughout, we shall retain the subscripts *X* and *Y* on frequency variables, even though the space variables to which they correspond may have different symbols.

For notational convenience we set the constant $A\lambda z_i$ equal to unity and ignore the negative signs in the arguments of P (almost all applications of interest to us here have pupil functions that are symmetrical in x and y). Thus

$$H(f_X, f_Y) = P(\lambda z_i f_X, \lambda z_i f_Y). \quad (6-20)$$

This relation is of the utmost importance; it supplies very revealing information about the behavior of diffraction-limited coherent imaging systems in the frequency domain. If the pupil function P is indeed unity within some region and zero otherwise, then there exists a finite **passband** in the frequency domain within which the **diffraction-limited** imaging system passes all frequency components without amplitude or phase **distortion**.⁷ At the boundary of this **passband** the frequency response suddenly drops to zero, implying that frequency components outside the **passband** are completely eliminated.

Finally we give some intuitive explanation as to why the scaled pupil function plays the role of the amplitude transfer function. Remember that in order to completely remove the quadratic phase factor across the object, the object should be illuminated with a spherical wave, in this case converging towards the point where the entrance pupil is pierced by the optical axis (cf. discussion leading up to Fig. 5.9). The converging spherical illumination causes the Fourier components of the object amplitude transmittance to appear in the entrance pupil, as well as in the exit pupil, since the latter is the image of the former (see Appendix B). Thus the pupil sharply limits the range of Fourier components passed by the system. If the converging illumination is not present, the same conclusion is approximately true, especially for an object of sufficiently small extent in the object plane, as was discussed in connection with Fig. 5.10.

6.2.2 Examples of Amplitude Transfer Functions

To illustrate the frequency response of diffraction-limited coherent imaging systems, consider the amplitude transfer functions of systems with square (width $2w$) and circular (diameter $2w$) pupils. For these two cases, we have, respectively,

$$P(x, y) = \text{rect}\left(\frac{x}{2w}\right) \text{rect}\left(\frac{y}{2w}\right)$$

$$P(x, y) = \text{circ}\left(\frac{\sqrt{x^2 + y^2}}{w}\right).$$

Thus, from (6-20), the corresponding amplitude transfer functions are

$$H(f_X, f_Y) = \text{rect}\left(\frac{\lambda z_i f_X}{2w}\right) \text{rect}\left(\frac{\lambda z_i f_Y}{2w}\right)$$

$$H(f_X, f_Y) = \text{circ}\left(\frac{\sqrt{f_X^2 + f_Y^2}}{w/\lambda z_i}\right).$$

⁷**Note** that this conclusion has been drawn only for a system free from aberrations. As we shall see in Section 6.4, a system that has aberrations is not free from phase distortion within its passband.

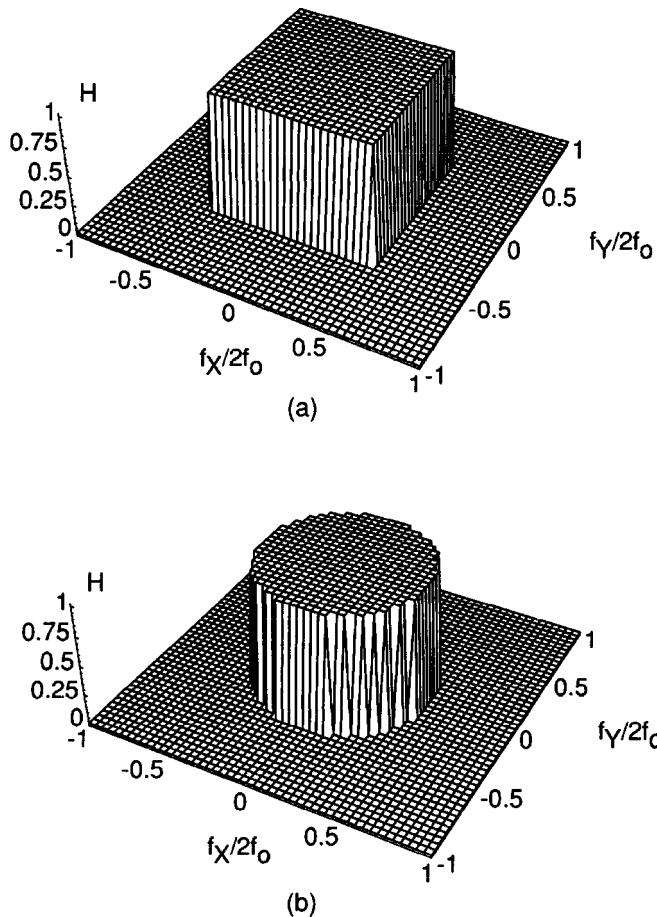


FIGURE 6.3
Amplitude transfer functions for
diffraction-limited systems with (a)
square and (b) circular exit pupils.

These functions are illustrated in Fig. 6.3. Note that a cutoff frequency f_c can be defined in both cases by

$$f_c = \frac{w}{\lambda z_i} \quad (6-21)$$

where in the circular case this cutoff is uniform in all directions in the frequency plane, while in the square case this cutoff applies only along the f_x and f_y axes. To illustrate a particular order-of-magnitude off, suppose that $w = 1$ cm, $z_i = 10$ cm, and $\lambda = 10^{-4}$ cm. Then the cutoff frequency is 100 cycles/mm.

6.3 FREQUENCY RESPONSE FOR DIFFRACTION-LIMITED INCOHERENT IMAGING

In the coherent case, the relation between the pupil and the amplitude transfer function has been seen to be a very direct and simple one. When the object illumination is incoherent, the transfer function of the imaging system will be seen to be determined by the pupil again, but in a less direct and somewhat more interesting way. The theory of imaging with incoherent light has, therefore, a certain extra richness not present in

the coherent case. We turn now to considering this theory; again attention will be centered on *diffraction-limited* systems, although the discussion that immediately follows applies to all incoherent systems, regardless of their aberrations.

6.3.1 The Optical Transfer Function

Imaging systems that use incoherent illumination have been seen to obey the intensity convolution integral

$$I_i(u, v) = \kappa \iint_{-\infty}^m |h(u - \tilde{\xi}, v - \tilde{\eta})|^2 I_g(\tilde{\xi}, \tilde{\eta}) d\tilde{\xi} d\tilde{\eta}. \quad (6-22)$$

Such systems should therefore be frequency-analyzed as linear mappings of intensity distributions. To this end, let the normalized frequency spectra of I_g and I_i be defined by

$$\mathcal{G}_g(f_x, f_y) = \frac{\iint_{-\infty}^{\infty} I_g(u, v) \exp[-j2\pi(f_x u + f_y v)] du dv}{\iint_{-\infty}^{\infty} I_g(u, v) du dv} \quad (6-23)$$

$$\mathcal{G}_i(f_x, f_y) = \frac{\iint_{-\infty}^m I_i(u, v) \exp[-j2\pi(f_x u + f_y v)] du dv}{\iint_{-\infty}^{\infty} I_i(u, v) du dv} \quad (6-24)$$

The normalization of the spectra by their "zero-frequency" values is partly for mathematical convenience, and partly for a more fundamental reason. It can be shown that any real and nonnegative function, such as I_g or I_i , has a Fourier transform which achieves its maximum value at the origin. We choose that maximum value as a normalization constant in defining \mathcal{G}_g and \mathcal{G}_i . Since intensities are nonnegative quantities, they always have a spectrum that is nonzero at the origin. The visual quality of an image depends strongly on the "contrast" of the image, or the relative strengths of the **information-bearing** portions of the image and the ever-present background. Hence the spectra are normalized by that background.

In a similar fashion, the normalized transfer function of the system can be defined by

$$\mathcal{H}(f_x, f_y) = \frac{\iint_{-\infty}^m |h(u, v)|^2 \exp[-j2\pi(f_x u + f_y v)] du dv}{\iint_{-\infty}^{\infty} |h(u, v)|^2 du dv} \quad (6-25)$$

Application of the convolution theorem to Eq. (6-22) then yields the frequency-domain relation

$$\mathcal{G}_i(f_x, f_y) = \mathcal{H}(f_x, f_y) \mathcal{G}_g(f_x, f_y). \quad (6-26)$$

By international agreement, the function 3-1 is known as the optical transfer function (abbreviated OTF) of the system. Its modulus $|H|$ is known as the modulation transfer function (MTF). Note that $\mathcal{H}(f_x, f_y)$ simply specifies the complex weighting factor applied by the system to the frequency component at (f_x, f_y) , relative to the weighting factor applied to the zero-frequency component.

Since the definitions of both the amplitude transfer function and the optical transfer function involve the function h , we might expect some specific relationship between the two. In fact, such a relationship exists and can be readily found with the help of the autocorrelation theorem of Chapter 2. Since

$$H(f_x, f_y) = \mathcal{F}\{h\}$$

and

$$\mathcal{H}(f_x, f_y) = \frac{\mathcal{F}\{|h|^2\}}{\iint_{-\infty}^{\infty} |h(u, v)|^2 du dv},$$

it follows (with the help of Rayleigh's theorem) that

$$\mathcal{H}(f_x, f_y) = \frac{\iint_{-\infty}^{\infty} H(p', q') H^*(p' - f_x, q' - f_y) dp' dq'}{\iint_{-\infty}^{\infty} |H(p', q')|^2 dp' dq'}. \quad (6-27)$$

The simple change of variables

$$p = p' - \frac{f_x}{2} \quad q = q' - \frac{f_y}{2}$$

results in the symmetrical expression

$$\mathcal{H}(f_x, f_y) = \frac{\iint_{-\infty}^{\infty} H\left(p + \frac{f_x}{2}, q + \frac{f_y}{2}\right) H^*\left(p - \frac{f_x}{2}, q - \frac{f_y}{2}\right) dp dq}{\iint_{-\infty}^{\infty} |H(p, q)|^2 dp dq}. \quad (6-28)$$

Thus the OTF is the normalized autocorrelation function of the amplitude transfer function!

Equation (6-28) will serve as our primary link between the properties of coherent and incoherent systems. Note that it is entirely valid for systems both with and without aberrations.

6.3.2 General Properties of the OTF

A number of very simple and elegant properties of the OTF can be stated based only on knowledge that it is a normalized autocorrelation function. The most important of these properties are as follows:

1. $\mathcal{H}(0, 0) = 1$.
2. $\mathcal{H}(-f_X, -f_Y) = \mathcal{H}^*(f_X, f_Y)$.
3. $|\mathcal{H}(f_X, f_Y)| \leq |\mathcal{H}(0, 0)|$.

Property 1 follows directly by substitution of $(f_X = 0, f_Y = 0)$ in Eq. (6-28). The proof of Property 2 is left as an exercise for the reader, it being no more than a statement that the Fourier transform of a real function has Hermitian symmetry.

The proof that the MTF at any frequency is always less than its zero-frequency value of unity requires more effort. To prove Property 3 we use **Schwarz's** inequality ([227], p. 177), which can be stated as follows: If $X(p, q)$ and $Y(p, q)$ are any two complex-valued functions of (p, q) , then

$$\left| \iint XY \, dp \, dq \right|^2 \leq \iint |X|^2 \, dp \, dq \iint |Y|^2 \, dp \, dq \quad (6-29)$$

with equality if and only if $Y = KX^*$ where K is a complex constant. Letting

$$X(p, q) = H\left(p + \frac{f_X}{2}, q + \frac{f_Y}{2}\right) \quad \text{and} \quad Y(p, q) = H^*\left(p - \frac{f_X}{2}, q - \frac{f_Y}{2}\right)$$

we find

$$\begin{aligned} & \left| \iint_{-\infty}^{\infty} H\left(p + \frac{f_X}{2}, q + \frac{f_Y}{2}\right) H^*\left(p - \frac{f_X}{2}, q - \frac{f_Y}{2}\right) dp \, dq \right|^2 \\ & \leq \iint_{-\infty}^{\infty} \left| H\left(p + \frac{f_X}{2}, q + \frac{f_Y}{2}\right) \right|^2 dp \, dq \iint_{-\infty}^{\infty} \left| H\left(p - \frac{f_X}{2}, q - \frac{f_Y}{2}\right) \right|^2 dp \, dq \\ & = \left[\iint_{-\infty}^{\infty} |H(p, q)|^2 dp \, dq \right]^2. \end{aligned}$$

Normalizing by the right-hand side of the inequality, it follows that $|\mathcal{H}(f_X, f_Y)|$ is never greater than unity.

Finally, it should be pointed out that while the OTF is always unity at the zero frequency, this does not imply that the absolute intensity level of the image background is the same as the absolute intensity level of the object background. The normalization used in the definition of the OTF has removed all information about absolute intensity levels.

6.3.3 The OTF of an Aberration-Free System

To this point, our discussions have been equally applicable to systems with and without aberrations. We now consider the special case of a diffraction-limited incoherent system. Recall that for coherent systems we have

$$H(f_X, f_Y) = P(\lambda z_i f_X, \lambda z_i f_Y).$$

For an incoherent system, it follows from Eq. (6-28) (with a simple change of variables) that

$$\mathcal{H}(f_x, f_y) = \frac{\iint_{-\infty}^{\infty} P\left(x + \frac{\lambda z_i f_x}{2}, y + \frac{\lambda z_i f_y}{2}\right) P\left(x - \frac{\lambda z_i f_x}{2}, y - \frac{\lambda z_i f_y}{2}\right) dx dy}{\iint_{-\infty}^{\infty} P(x, y) dx dy}, \quad (6-30)$$

where, in the denominator the fact that P equals either unity or zero has been used to replace P^2 by P .

The expression (6-30) for \mathcal{H} lends itself to an extremely important geometrical interpretation. The numerator represents the area of overlap of two displaced pupil functions, one centered at $(\lambda z_i f_x/2, \lambda z_i f_y/2)$ and the second centered on the diametrically opposite point $(-\lambda z_i f_x/2, -\lambda z_i f_y/2)$. The denominator simply normalizes the area of overlap by the total area of the pupil. Thus

$$\mathcal{H}(f_x, f_y) = \frac{\text{area of overlap}}{\text{total area}}.$$

To calculate the OTF of a diffraction-limited system, the steps indicated by this interpretation can be directly performed, as illustrated in Fig. 6.4. For simple geometrical shapes, closed-form expressions for the normalized overlap area can be found (see examples to follow). Note that this geometrical interpretation of the OTF implies that the OTF of a diffraction-limited system is always *real* and *nonnegative*. It is not necessarily a monotonically decreasing function of frequency, however (see, for example, Prob. 6-3).

For complicated pupils, the OTF can be calculated with the help of a digital computer. A straightforward way to perform such a calculation is to Fourier transform the pupil function (finding the amplitude point-spread function), take the squared

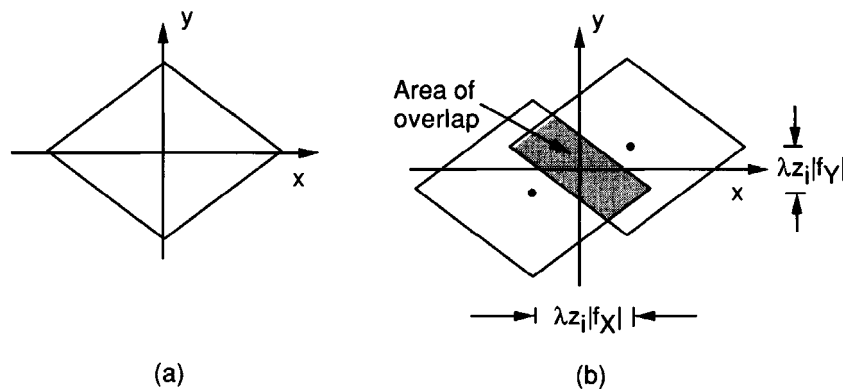


FIGURE 6.4
 Geometrical interpretations of the OTF of a diffraction-limited system. (a) The pupil function—total area is the denominator of the OTF; (b) two displaced pupil functions—the shaded area is the numerator of the OTF.

magnitude of this quantity (thus finding the intensity point-spread function), and inverse Fourier transform the result.

To lend further physical insight into the OTF, consider the ways in which a sinusoidal component of intensity at a particular frequency pair (f_x, f_y) can be generated in the image. We claim that such a fringe can be generated only by interference of light in the image plane from two separate patches on the exit pupil of the system, with a separation between patches that is $(\lambda z_i |f_x|, \lambda z_i |f_y|)$. Only when light contributions from two patches having this particular separation interfere can a fringe with this frequency be generated (cf. Prob. 6-1). However, there are many different pairs of patches of this separation that can be embraced by the pupil of the system. In fact, the relative weight given by the system to this particular frequency pair is determined by how many different ways such a separation can be fit into the pupil. The number of ways a particular separation can be fit into the exit pupil is proportional to the area of overlap of two pupils separated by this particular spacing. See Fig. 6.5.

6.3.4 Examples of Diffraction-Limited OTFs

We consider now as examples the OTFs that correspond to diffraction-limited systems with square (width $2w$) and circular (diameter $2w$) pupils. Figure 6.6 illustrates the calculation for the square case, The area of overlap is evidently

$$\mathcal{A}(f_x, f_y) = \begin{cases} (2w - \lambda z_i |f_x|)(2w - \lambda z_i |f_y|) & |f_x| \leq 2w/\lambda z_i, \\ & |f_y| \leq 2w/\lambda z_i \\ 0 & \text{otherwise.} \end{cases}$$

When this area is normalized by the total area $4w^2$, the result becomes

$$\mathcal{H}(f_x, f_y) = \Lambda\left(\frac{f_x}{2f_o}\right)\Lambda\left(\frac{f_y}{2f_o}\right) \quad (6-31)$$

where Λ is the triangle function of Chapter 2, and f_o is the cutoff frequency of the same system when used with coherent illumination,

$$f_o = \frac{w}{\lambda z_i}.$$

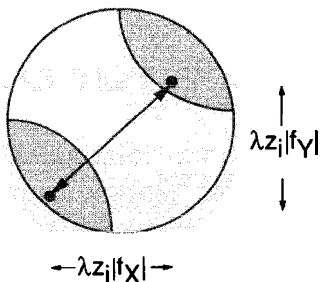


FIGURE 6.5

Light from patches separated by $(\lambda z_i |f_x|, \lambda z_i |f_y|)$ interferes to produce a sinusoidal fringe at frequency (f_x, f_y) . The shaded areas on the pupil are the areas within the light patches can reside while retaining this special separation.

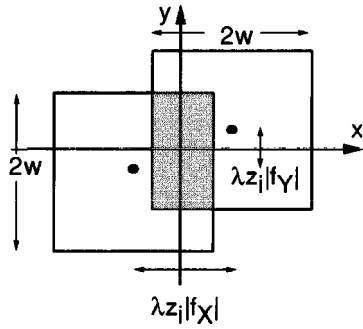


FIGURE 6.6
Calculation of the OTF for a square aperture.

Note that the cutoff frequency of the incoherent system occurs at frequency $2f_o$ along the f_x and f_y axes.⁸ The OTF represented by Eq. (6-31) is illustrated in Fig. 6.7.

When the pupil is circular, the calculation is not quite so straightforward. Since the OTF will clearly be circularly symmetric, it suffices to calculate \mathcal{H} along the positive f_x axis. As illustrated in Fig. 6.8, the area of overlap may be regarded as being equal to four times the shaded area B of the circular sector A + B. But the area of the circular sector is

$$\text{Area}(A + B) = \left[\frac{\theta}{2\pi} \right] (\pi w^2) = \left[\frac{\arccos(\lambda z_i f_x / 2w)}{2\pi} \right] (\pi w^2)$$

while the area of the triangle A is

$$\text{Area}(A) = \frac{1}{2} \left(\frac{\lambda z_i f_x}{2} \right) \sqrt{w^2 - \left(\frac{\lambda z_i f_x}{2} \right)^2}.$$

Finally, we have

$$\mathcal{H}(f_x, 0) = \frac{4[\text{area}(A + B) - \text{area}(A)]}{\pi w^2}$$

or, for a general radial distance ρ in the frequency plane,

$$\mathcal{H}(\rho) = \begin{cases} \frac{2}{\pi} \left[\arccos\left(\frac{\rho}{2\rho_o}\right) - \frac{\rho}{2\rho_o} \sqrt{1 - \left(\frac{\rho}{2\rho_o}\right)^2} \right] & \rho \leq 2\rho_o \\ 0 & \text{otherwise.} \end{cases} \quad (6-32)$$

The quantity ρ_o is the cutoff frequency of the coherent system,

$$\rho_o = \frac{w}{\lambda z_i}.$$

Referring to Fig. 6.9, the OTF is again seen to extend to a frequency that is twice the coherent cutoff frequency.

⁸This should not be taken to imply that the incoherent system has twice the resolving power of the coherent system. See Section 6.5.

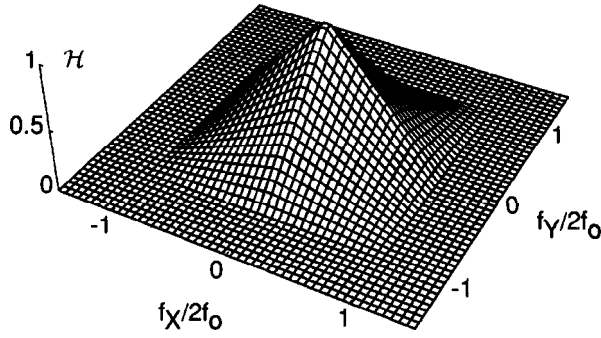


FIGURE 6.7
The optical transfer function of a diffraction-limited system with a square pupil.

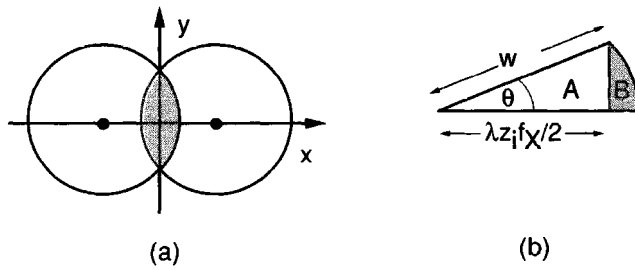


FIGURE 6.8
Calculation of the area of overlap of two displaced circles.
(a) Overlapping circles, (b) geometry of the calculation.

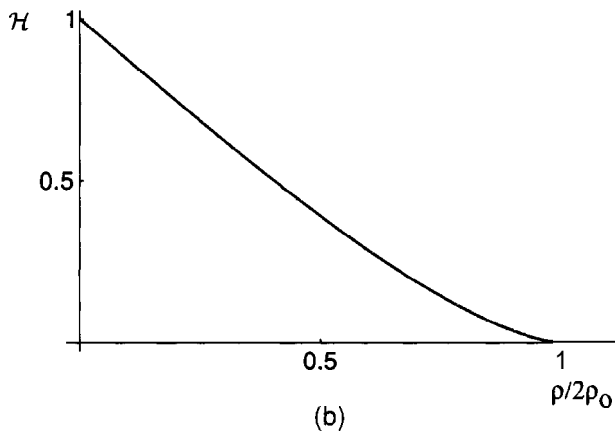
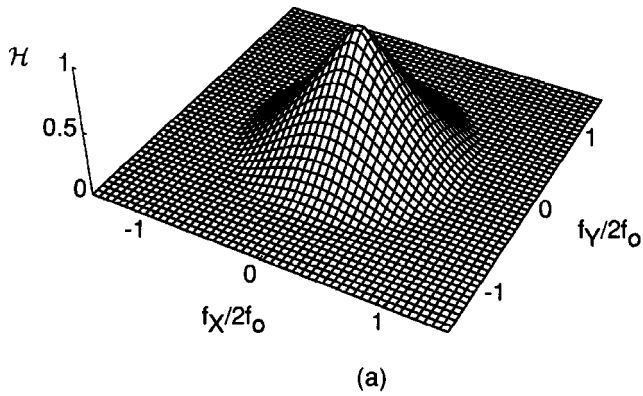


FIGURE 6.9
The optical transfer function of a diffraction-limited system with a circular pupil. (a) Three-dimensional perspective, (b) cross section.

6.4

ABERRATIONS AND THEIR EFFECTS ON FREQUENCY RESPONSE

In the development of a generalized model of an imaging system, it was specifically assumed that the presence of a point-source object yielded at the exit pupil a perfect spherical wave, converging toward the ideal geometrical image point. Such a system was called *diffraction-limited*. We consider now the effects of *aberrations*, or departures of the exit-pupil wavefront from ideal spherical form. Aberrations can arise in a variety of ways, ranging from a defect as simple as a focusing error to inherent properties of perfectly spherical lenses, such as spherical aberration. A complete treatment of aberrations and their detailed effects on frequency response is beyond the scope of this development. Rather we concentrate on very general effects and illustrate with one relatively simple example. For a more complete treatment of various types of aberrations and their effects on frequency response, see, for example, Refs. [300], [146], or [296].

6.4.1 The Generalized Pupil Function

When an imaging system is diffraction limited, the (amplitude) point-spread function has been seen to consist of the Fraunhofer diffraction pattern of the exit pupil, centered on the ideal image point. This fact suggests a convenient artifice which will allow aberrations to be directly included in our previous results. Specifically, when wavefront errors exist, we can imagine that the exit pupil *is* illuminated by a perfect spherical wave, but that a phase-shifting plate exists in the aperture, thus deforming the wavefront that leaves the pupil. If the phase error at the point (x, y) is represented by $kW(x, y)$, where $k = 2\pi/\lambda$ and W is an effective path-length error, then the complex amplitude transmittance $\mathcal{P}(x, y)$ of the imaginary phase-shifting plate is given by

$$\mathcal{P}(x, y) = P(x, y) \exp[jkW(x, y)]. \quad (6-33)$$

The complex function \mathcal{P} may be referred to as the *generalized* pupil function. The amplitude point-spread function of an aberrated coherent system is simply the Fraunhofer diffraction pattern of an aperture with amplitude transmittance \mathcal{P} . The intensity impulse response of an aberrated incoherent system is, of course, the squared magnitude of the amplitude impulse response.

Figure 6.10 shows the geometry that defines the aberration function W . If the system were free from aberrations, the exit pupil would be filled by a perfect spherical wave converging towards the ideal image point. We regard an ideal spherical surface, centered on the ideal image point and passing through the point where the optical axis pierces the exit pupil, as defining a *Gaussian reference sphere* with respect to which the aberration function can be defined. If we trace a ray backward from the ideal image point to the coordinates (x, y) in the exit pupil, the aberration function $W(x, y)$ is the path-length error accumulated by that ray as it passes from the Gaussian reference sphere to the actual wavefront, the latter wavefront also being defined to intercept the optical axis in the exit pupil. The error can be positive or negative, depending on whether the actual wavefront lies to the left or to the right (respectively) of the Gaussian reference sphere.

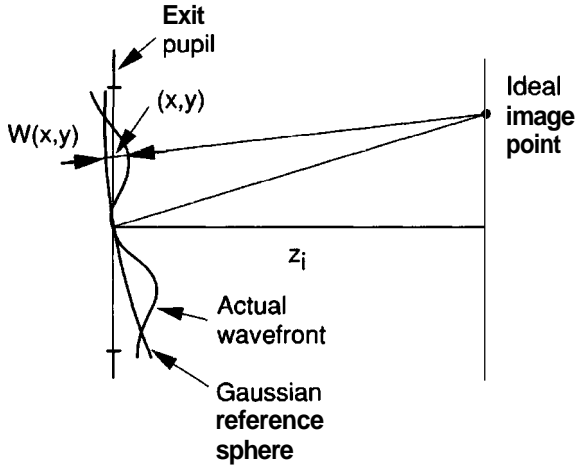


FIGURE 6.10
Geometry for defining the aberration function.

6.4.2 Effects of Aberrations on the Amplitude Transfer Function

When considering a diffraction-limited coherent system, the transfer function was found by noting that (1) the impulse response is the Fourier transform of the pupil function, and (2) the amplitude transfer function is the Fourier transform of the amplitude impulse response. As a consequence of the two Fourier transform relations, the amplitude transfer function was found to be proportional to a scaled pupil function P . Identical reasoning can be used when aberrations are present, provided the generalized pupil function \mathcal{P} replaces P . Thus the amplitude transfer function is written

$$H(f_X, f_Y) = \mathcal{P}(\lambda z_i f_X, \lambda z_i f_Y) = P(\lambda z_i f_X, \lambda z_i f_Y) \exp[jkW(\lambda z_i f_X, \lambda z_i f_Y)]. \quad (6-34)$$

Evidently the band limitation of the amplitude transfer function, as imposed by the finite exit pupil, is unaffected by the presence of aberrations. The sole effect of aberrations is seen to be the introduction of *phase* distortions within the passband. Phase distortions can, of course, have a severe effect on the fidelity of the imaging system.

There is little more of a general nature that can be said about the effects of aberrations on a coherent imaging system. Again the result is a very simple one: as we shall now see, the result for an incoherent system is again more complex and, in many respects, more interesting.

6.4.3 Effects of Aberrations on the OTF

Having found the effects of aberrations on the amplitude transfer function, it is now possible, with the help of Eq. (6-28), to find the effects on the optical transfer function. To simplify the notation, the function $A(f_X, f_Y)$ is defined as the area \mathcal{A} overlap of

$$P\left(x - \frac{\lambda z_i f_X}{2}, y - \frac{\lambda z_i f_Y}{2}\right) \quad \text{and} \quad P\left(x + \frac{\lambda z_i f_X}{2}, y + \frac{\lambda z_i f_Y}{2}\right).$$

Thus the OTF of a diffraction-limited system is given, in this new notation, by

$$\mathcal{H}(f_X, f_Y) = \frac{\iint dx dy}{\iint_{\mathcal{A}(0,0)} dx dy} \mathcal{A}(f_X, f_Y) \quad (6-35)$$

When aberrations are present, substitution of (6-34) into (6-35) yields

$$\mathcal{H}(f_X, f_Y) = \frac{\iint_{\mathcal{A}(f_X, f_Y)} e^{jk \left[W \left(x + \frac{\lambda z_i f_X}{2}, y + \frac{\lambda z_i f_Y}{2} \right) - W \left(x - \frac{\lambda z_i f_X}{2}, y - \frac{\lambda z_i f_Y}{2} \right) \right]} dx dy}{\iint_{\mathcal{A}(0,0)} dx dy} \quad (6-36)$$

This expression allows us, then, to directly relate the wavefront errors and the **OTF**.

As an important general property, it can be shown that aberrations will never increase the **MTF** (the modulus of the **OTF**). To prove this property, Schwarz's inequality (6-29) will be used. Let the functions **X** and **Y** of that equation be defined by

$$X(x, y) = \exp \left[jkW \left(x + \frac{\lambda z_i f_X}{2}, y + \frac{\lambda z_i f_Y}{2} \right) \right]$$

$$Y(x, y) = \exp \left[-jkW \left(x - \frac{\lambda z_i f_X}{2}, y - \frac{\lambda z_i f_Y}{2} \right) \right].$$

Noting that $|X|^2 = |Y|^2 = 1$, it follows that

$$\begin{aligned} & |\mathcal{H}(f_X, f_Y)|_{\text{with aberrations}}^2 \\ &= \left| \frac{\iint_{\mathcal{A}(f_X, f_Y)} e^{jk \left[W \left(x + \frac{\lambda z_i f_X}{2}, y + \frac{\lambda z_i f_Y}{2} \right) - W \left(x - \frac{\lambda z_i f_X}{2}, y - \frac{\lambda z_i f_Y}{2} \right) \right]} dx dy}{\iint_{\mathcal{A}(0,0)} dx dy} \right|^2 \\ &\leq \left[\frac{\iint_{\mathcal{A}(f_X, f_Y)} dx dy}{\iint_{\mathcal{A}(0,0)} dx dy} \right]^2 = |\mathcal{H}(f_X, f_Y)|_{\text{without aberrations}}^2 \end{aligned}$$

Thus aberrations cannot increase the contrast of any spatial-frequency component of the image, and in general will lower the contrast. The absolute cutoff frequency remains unchanged, but severe aberrations can reduce the high-frequency portions of the **OTF** to such an extent that the effective cutoff is much lower than the diffraction-limited cutoff. In addition, aberrations can cause the **OTF** to have negative values in certain bands of frequencies, a result that never occurs for an aberration-free system. When the **OTF** is negative, image components at that frequency undergo a contrast reversal; i.e., intensity peaks become intensity nulls, and vice versa. An example of this effect will be seen in the section that follows.

6.4.4 Example of a Simple Aberration: A Focusing Error

One of the easiest aberrations to deal with mathematically is a simple error of focus. But even in this simple case, the assumption of a *square* aperture (rather than a circular aperture) is needed to keep the mathematics simple.

When a focusing error is present, the center of curvature of the spherical wavefront converging towards the image of an object point-source lies either to the left or to the right of the image plane. Considering an on-axis point for simplicity, this means that the phase distribution across the exit pupil is of the form

$$\phi(x, y) = -\frac{\pi}{\lambda z_a} (x^2 + y^2),$$

where $z_a \neq z_i$. The path-length error $W(x, y)$ can then be determined by subtracting the ideal phase distribution from the actual phase distribution,

$$kW(x, y) = -\frac{\pi}{\lambda z_a} (x^2 + y^2) + \frac{\pi}{\lambda z_i} (x^2 + y^2). \quad (6-37)$$

The path-length error is thus given by

$$W(x, y) = -\frac{1}{2} \left(\frac{1}{z_a} - \frac{1}{z_i} \right) (x^2 + y^2), \quad (6-38)$$

which is seen to depend quadratically on the space variables in the exit pupil.

For a square aperture of width $2w$, the maximum path-length error at the edge of the aperture along the x or y axes, which we represent by W_m , is given by

$$W_m = -\frac{1}{2} \left(\frac{1}{z_a} - \frac{1}{z_i} \right) w^2. \quad (6-39)$$

The number W_m is a convenient indication of the severity of the focusing error. Using the definition of W_m , we can express the path-length error as

$$W(x, y) = W_m \frac{x^2 + y^2}{w^2}. \quad (6-40)$$

If the path-length error W given by (6-40) is substituted in the expression (6-36) for the OTF, a number of straightforward manipulations yield the result

$$\begin{aligned} \mathcal{H}(f_X, f_Y) &= \Lambda\left(\frac{f_X}{2f_o}\right) \Lambda\left(\frac{f_Y}{2f_o}\right) \\ &\times \operatorname{sinc}\left[\frac{8W_m}{\lambda} \left(\frac{f_X}{2f_o}\right) \left(1 - \frac{|f_X|}{2f_o}\right)\right] \operatorname{sinc}\left[\frac{8W_m}{\lambda} \left(\frac{f_Y}{2f_o}\right) \left(1 - \frac{|f_Y|}{2f_o}\right)\right]. \end{aligned} \quad (6-41)$$

Plots of this OTF are shown in Fig. 6.11 for various values of W_m/λ . Note that the diffraction-limited OTF is indeed obtained when $W_m = 0$. Note also that, for values of W_m greater than $\lambda/2$, sign reversals of the OTF occur. These reversals of contrast can readily be observed if the "spoke" target of Fig. 6.12(a) is used as the object. The "local spatial frequency" of this target changes slowly, increasing as the radius from the center is decreased. The local contrast of fringes is thus an indication of the value of

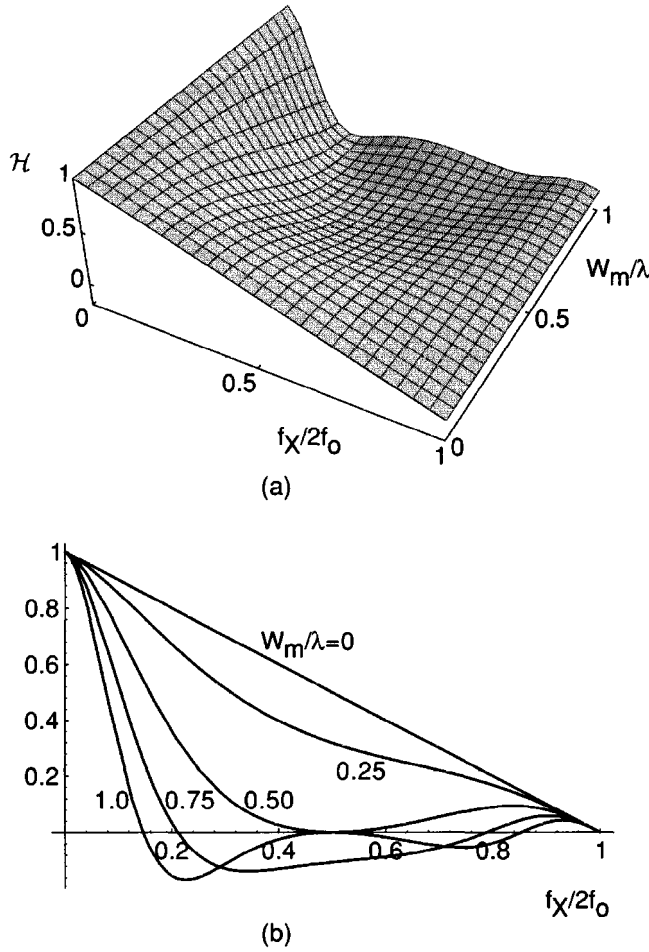


FIGURE 6.11
OTF for a focusing error in a system with a square pupil. (a) Three-dimensional plot with $f_x/2f_o$ along one axis and W_m/λ along the other axis. **(b)** Cross section along the f_x axis with W_m/λ as a parameter.

the MTF at various frequencies. The position of the fringes is determined by the phase associated with the OTF at each frequency. When the system is out of focus, a gradual attenuation of contrast and a number of contrast reversals are obtained for increasing spatial frequency, as illustrated in Fig. 6.12(b).

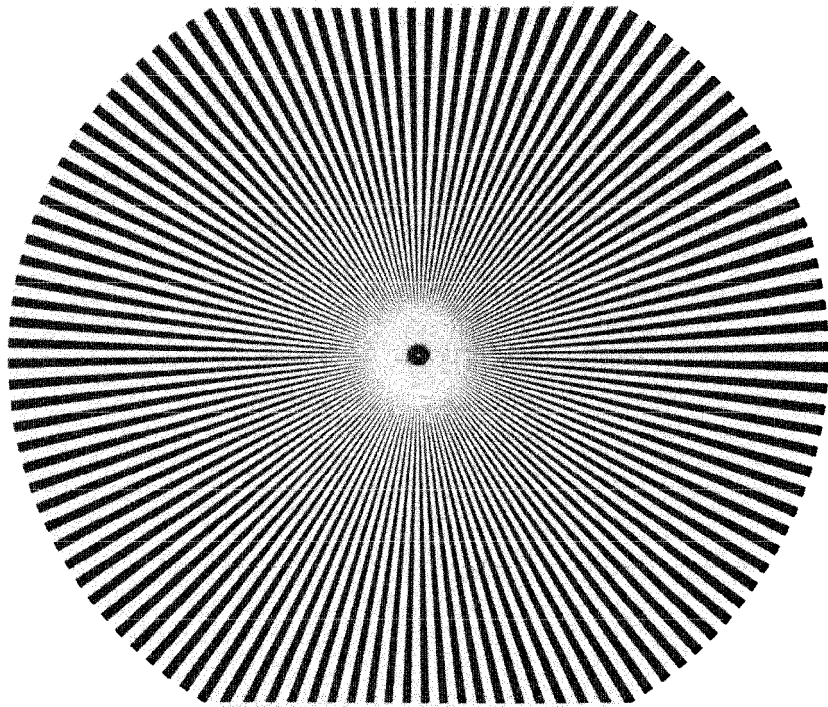
Finally, consider the form of the OTF when the focusing error is very severe (that is, when $W_m \gg \lambda$). In such a case, the frequency response drops towards zero for relatively small values of $f_x/2f_o$ and $f_y/2f_o$. We may therefore write

$$1 - \frac{|f_x|}{2f_o} \approx 1, \quad 1 - \frac{|f_y|}{2f_o} \approx 1,$$

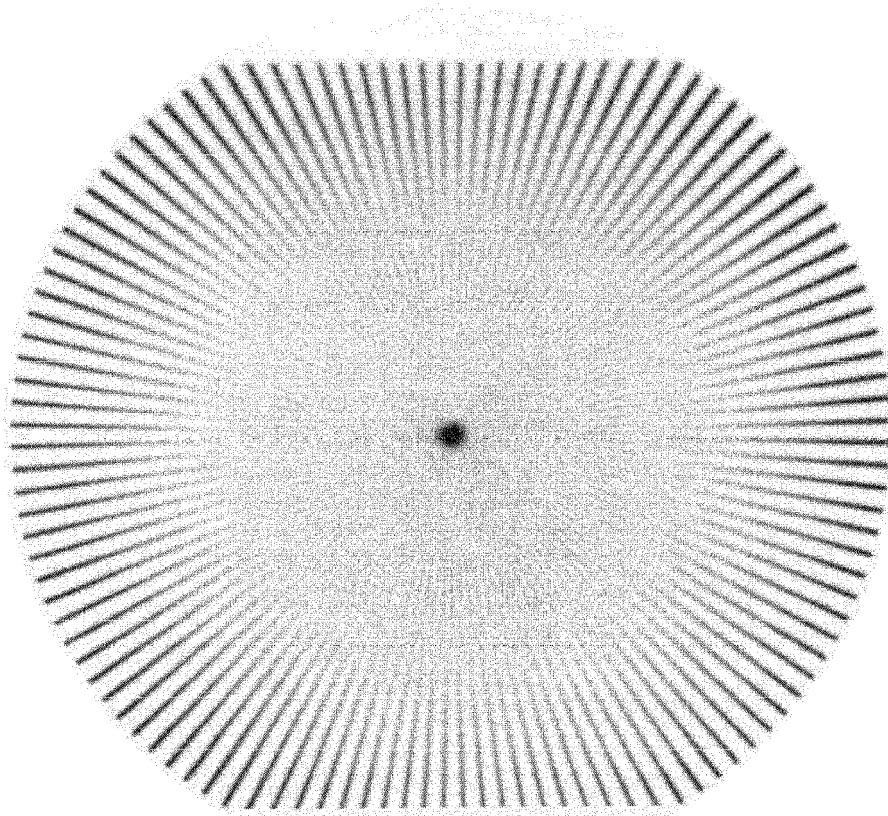
and the OTF reduces to

$$\mathcal{H}(f_x, f_y) = \text{sinc} \left[\frac{8W_m}{\lambda} \left(\frac{f_x}{2f_o} \right) \right] \text{sinc} \left[\frac{8W_m}{\lambda} \left(\frac{f_y}{2f_o} \right) \right]. \quad (6-42)$$

The interested reader can verify that this is precisely the OTF predicted by geometrical optics. Geometrical optics predicts a point-spread function that is the geometrical projection of the exit pupil into the image plane, and therefore the point-spread function should be uniformly bright over a square and zero elsewhere (see Fig. 6.13). The



(a)



(b)

FIGURE 6.12
(a) Focused and (b) misfocused images of a spoke target.

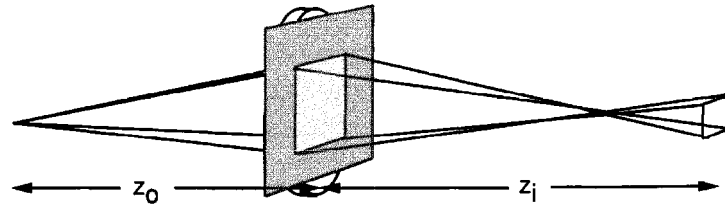


FIGURE 6.13
Geometrical optics prediction of the point-spread function of a system having a square pupil function and a severe focusing error.

Fourier transform of such a spread function yields the OTF of (6-42). More generally, when aberrations of any kind are severe, the geometrical optics predictions of the intensity point-spreadfunction may be Fourier-transformed to yield a good approximation to the OTF of the system. The fundamental reason for this behavior lies in the fact that, when severe aberrations are present, the point-spread function is determined primarily by geometrical-optics effects, and diffraction plays a negligible role in determining its shape.

6.4.5 Apodization and Its Effects on Frequency Response

The point-spread function of a diffraction-limited imaging system generally has side-lobes or side-rings of noticeable strength. While such extraneous responses may be of little concern in many imaging problems, they are of concern in a certain class of situations, such as when we wish to resolve a weak point-source next to a stronger point-source. Such a problem is of considerable importance in astronomy, where the presence or absence of weak companion stars next to a brighter star may often be of interest.

In an attempt to reduce the strength of side-lobes or side-rings, methods known as apodization have been developed. The word apodize is taken from the Greek language, and literally means "to remove the feet". The "feet" being referred to are in fact the side-lobes and side-rings of the diffraction-limited impulse response. Similar techniques are well known in the field of digital signal processing, where they are known by the term windowing (see, for example, [85], Section 3.3).

Generally speaking, apodization amounts to the introduction of attenuation in the exit pupil of an imaging system, attenuation that may be insignificant at the center of the pupil but increases with distance away from the center. Thus it amounts to a "softening" of the edges of the aperture through the introduction of an attenuating mask. Remembering that diffraction by an abrupt aperture can be thought of as coming from edge waves originating around the rim of the aperture, a softening of the edge has the effect of spreading the origin of these diffracted waves over a broader area around the edges of the pupil, thereby suppressing ringing effects caused by edge waves with a highly localized origin. Figure 6.14(a) shows a plot of the unapodized and apodized intensity transmissions through a square pupil with and without a Gaussian intensity apodization that falls to $(1/e)^2$ at the edge of the aperture. Part (b) of the figure shows cross sections of the intensity point-spread functions for the two cases. The logarithm

of intensity is plotted vertically in order to emphasize the side-lobes, and the intensity normalization is proportional to the total integrated intensity passed by the pupil in each case. Note that the side-lobes have been significantly suppressed by the apodization. Also note that the width of the main lobe is increased somewhat by apodization, and that the maximum intensity is also reduced due to extra absorption in the pupil.

The effects of apodization on the frequency response of both coherent and incoherent imaging systems are also of interest. In the coherent case the answer is straightforward due to the direct correspondence between the pupil and the amplitude transfer function. Attenuation that increases with distance from the center of the pupil results in an amplitude transfer function that falls off more rapidly with increasing frequency than it would in the absence of apodization. In the incoherent case, the less direct relationship between the OTF and the pupil makes the effects more subtle. Figure 6.15 shows a plot of cross sections of the apodized and unapodized OTFs of a system with a rectangular pupil, where the apodization is of the Gaussian form described above. As can be seen, the effect of the apodization has been to boost the relative importance of midrange and low frequencies, while diminishing the strength of high frequencies.

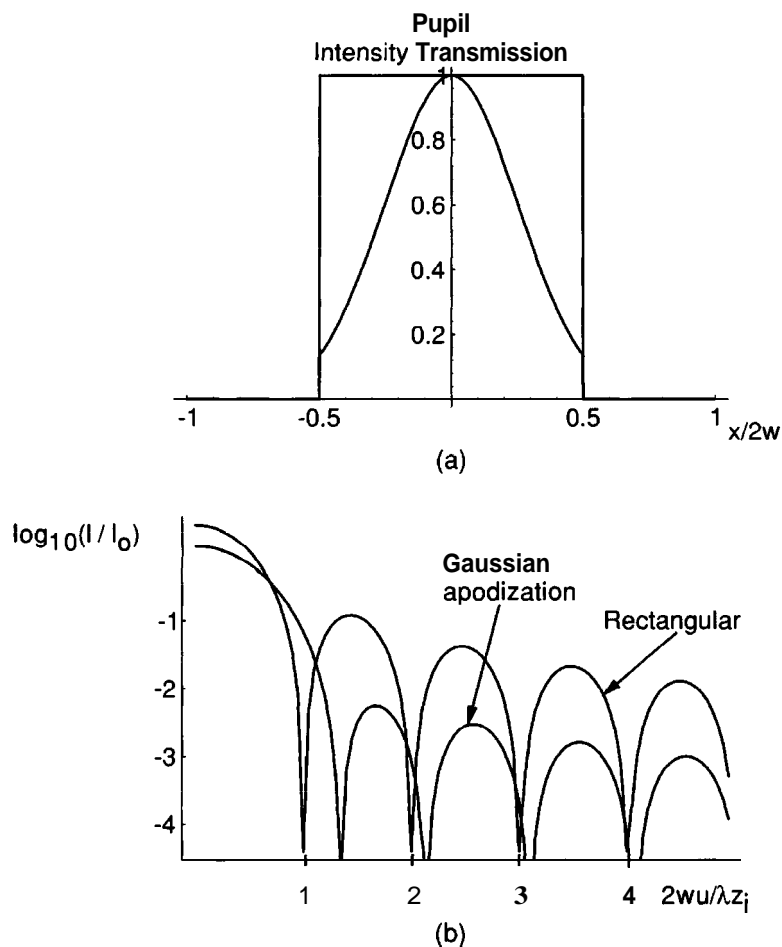


FIGURE 6.14
 Apodization of a rectangular aperture by a Gaussian function.
 (a) Intensity transmissions with and without apodization.
 (b) Point-spread functions with and without apodization.

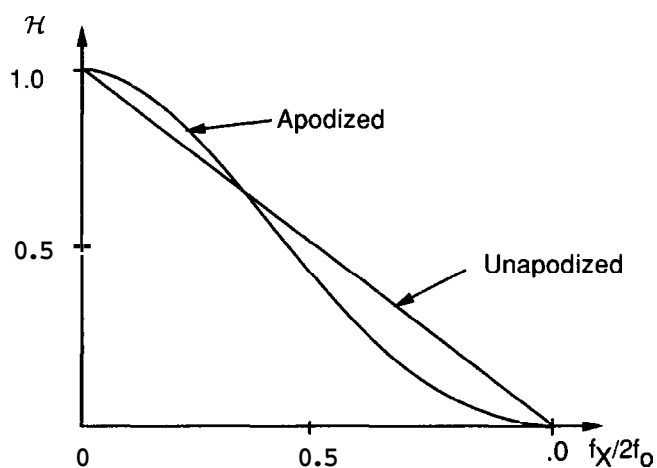


FIGURE 6.15
Optical transfer functions with and without a Gaussian apodization.

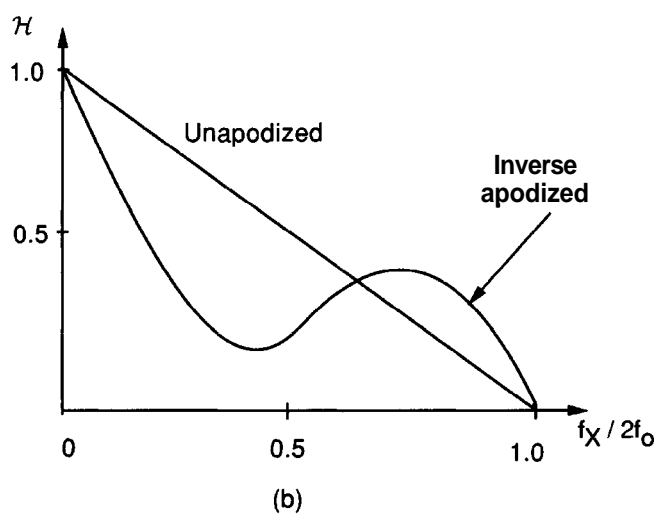
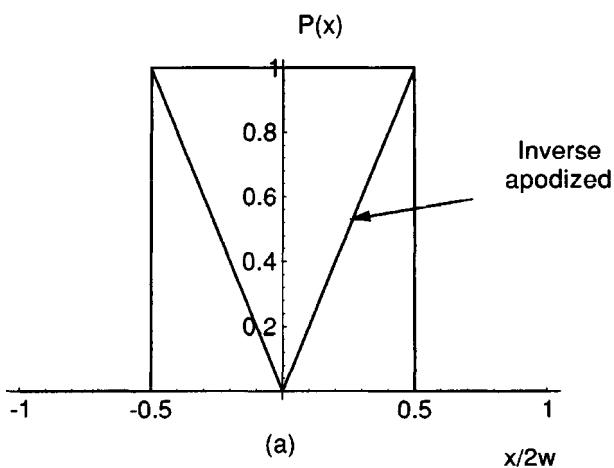


FIGURE 6.16
Pupil amplitude transmittance and the corresponding OTF with and without a particular "inverse" apodization.

While the term *apodization* originally meant a tapering of the transmittance through the pupil near its edges in order to suppress side-lobes of the point-spread function, over time the term has come to be used to describe *any* introduction of absorption into the pupil, whether it lowers or raises the side-lobes. Perhaps a better term for weightings that increase the sidelobes of the point-spread function would be "inverse" apodization. Figure 6.16 shows the amplitude transmittance through the pupil with and without a triangular amplitude weighting that gives extra emphasis to portions of the pupil near the edges, and de-emphasizes the importance of the center of the pupil. Also shown are cross sections of the OTF with and without this weighting. Note that this type of weighting emphasizes the importance of high frequencies relative to low frequencies.

As a closing remark regarding this subject, note that while the OTF of a system with or without apodization always has the value unity at the origin, nonetheless it is *not* true that the amount of light transmitted to the image is the same in the two cases. Naturally the introduction of absorbing material in the pupil diminishes the light that reaches the image, but the normalization of the OTF suppresses this fact. Note also that, unlike the case of aberrations, inverse apodization *can* raise the value of the OTF at certain frequencies, as compared with its unapodized values.

6.5

COMPARISON OF COHERENT AND INCOHERENT IMAGING

As seen in previous sections, the OTF of a diffraction-limited system extends to a frequency that is twice the cutoff frequency of the amplitude transfer function. It is tempting, therefore, to conclude that incoherent illumination will invariably yield "better" resolution than coherent illumination, given that the same imaging system is used in both cases. As we shall now see, this conclusion is in general *not* a valid one; a comparison of the two types of illumination is far more complex than such a superficial examination would suggest.

A major flaw in the above argument lies in the direct comparison of the cutoff frequencies in the two cases. Actually, the two are not directly comparable, since the cutoff of the amplitude transfer function determines the maximum frequency component of the image *amplitude* while the cutoff of the optical transfer function determines the maximum frequency component of image *intensity*. Surely any direct comparison of the two systems must be in terms of the same observable quantity, image intensity.

Even when the quantity to be compared is agreed upon, the comparison remains a difficult one for an additional fundamental reason: the term *better* has not been defined. Thus we have no universal quality criterion upon which to base our conclusions. A number of potential criteria might be considered (e.g., least-mean-square difference between the object and image intensities), but unfortunately the interaction of a human observer is so complex and so little understood that a truly meaningful criterion is difficult to specify.

In the absence of a meaningful quality criterion, we can only examine certain limited aspects of the two types of images, realizing that the comparisons so made will probably bear little direct relation to overall image quality. Nonetheless, such

comparisons are highly instructive, for they point out certain fundamental differences between the two types of illumination.

6.5.1 Frequency Spectrum of the Image Intensity

One simple attribute of the image intensity which can be compared in the two cases is the *frequency spectrum*. Whereas the incoherent system is linear in intensity, the coherent system is highly nonlinear in that quantity. Thus some care must be used in finding the spectrum in the latter case.

In the incoherent case, the image intensity is given by the convolution equation

$$I_i = |h|^2 \otimes I_g = |h|^2 \otimes |U_g|^2.$$

On the other hand, in the coherent case, we have

$$I_i = |h \otimes U_g|^2.$$

Let the symbol \star represent the autocorrelation integral

$$X(f_x, f_y) \star X(f_x, f_y) = \iint_{-\infty}^{\infty} X(p, q) X^*(p - f_x, q - f_y) dp dq. \quad (6-43)$$

Then we can directly write the frequency spectra of the image intensities in the two cases as

$$\begin{aligned} \text{Incoherent:} \quad \mathcal{F}\{I_i\} &= [H \star H][G_g \star G_g] \\ \text{Coherent:} \quad \mathcal{F}\{I_i\} &= HG, \star HG, \end{aligned} \quad (6-44)$$

where G_g is the spectrum of U_g and H is the amplitude transfer function.

The general result (6-44) does not lead to the conclusion that one type of illumination is better than the other in terms of image frequency content. It does, however, illustrate that the frequency content can be quite different in the two cases, and furthermore it shows that the results of any such comparison will depend strongly on both the intensity and *phase* distributions across the object.

To emphasize this latter point, we now consider two objects with the *same* intensity transmittance but different phase distributions, one of which can be said to be imaged better in coherent light and the other better in incoherent light. For simplicity, we suppose that the magnification of the system is unity, so that we may work in either the object or the image space at will without introducing a normalizing factor. Let the intensity transmittance of the ideal image in both cases be

$$\tau(\xi, \eta) = \cos^2 2\pi \tilde{f} \xi$$

where to make our point we will assume that

$$\frac{f_o}{2} < \tilde{f} < f_o,$$

f_o being the cutoff frequency of the amplitude transfer function. The amplitude transmittances of the two objects are taken to be

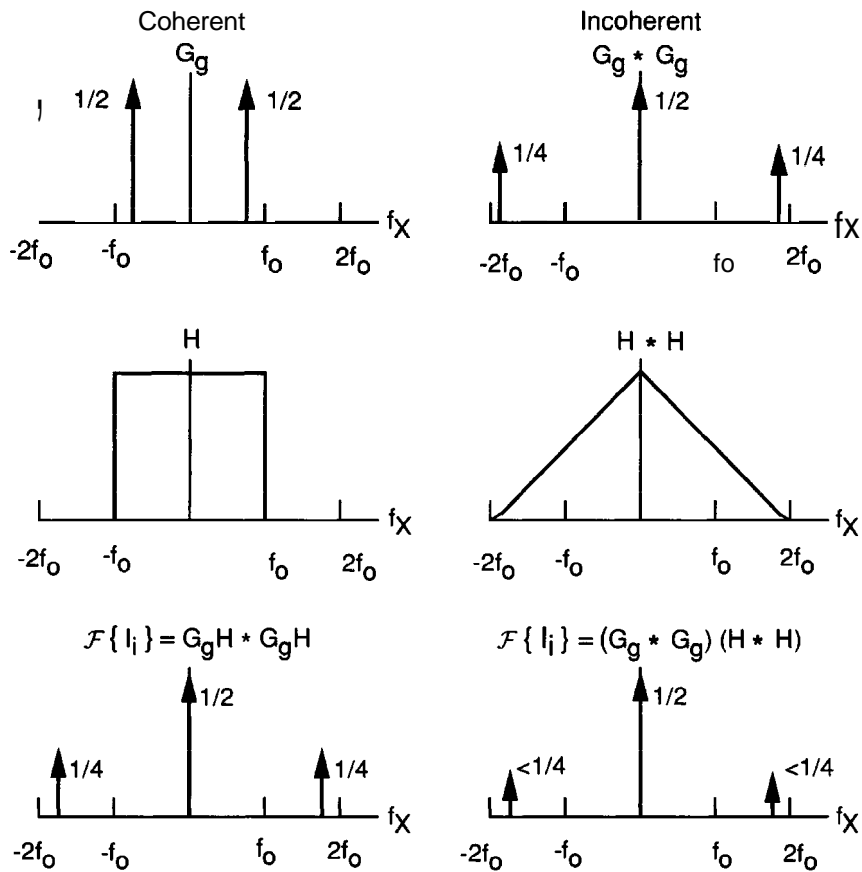


FIGURE 6.17 Calculation of the spectrum of the image intensity for object A.

$$\begin{aligned}
 \mathbf{A}: \quad t_A(\xi, \eta) &= \cos 2\pi \tilde{f} \xi \\
 \mathbf{B}: \quad t_A(\xi, \eta) &= |\cos 2\pi \tilde{f} \xi|.
 \end{aligned}$$

Thus the two objects differ only by a periodic phase distribution.

Figure 6.17 illustrates the various frequency-domain operations that lead to the image spectrum for object A. In all cases the imaging system is assumed to be **diffraction-limited**. Note that the contrast of the image intensity distribution is **poorer** for the incoherent case than for the coherent case. Thus object A is imaged better in coherent light than in incoherent light.

The corresponding comparison for object B requires less detail. The object amplitude distribution is now periodic with fundamental frequency $2\tilde{f}$. But since $2\tilde{f} > f_o$, **no** variations of image intensity will be present for the coherent case, while the incoherent system will form the same image it did for object A. Thus for object B, incoherent illumination must be termed **better** than coherent illumination.

In summary, then, which particular type of illumination is better from the point of view of image spectral content depends very strongly on the detailed structure of the object, and in particular on its phase distribution. It is **not** possible to conclude that one type of illumination is preferred in all cases. The comparison is in general a complex one, although simple cases, such as the one illustrated above, do exist. For a second example, the reader is referred to Prob. 6-10.

6.5.2 Two-Point Resolution

A second possible comparison criterion rests on the ability of the respective systems to resolve two closely spaced point sources. The two-point resolution criterion has long been used as a quality factor for optical systems, particularly in astronomical applications where it has a very real practical significance.

According to the so-called *Rayleigh criterion* of resolution, two incoherent point sources are "barely resolved" by a diffraction-limited system with a circular pupil when the center of the Airy intensity pattern generated by one point source falls exactly on the first zero of the Airy pattern generated by the second. The minimum resolvable separation of the geometrical images is therefore

$$\delta = 0.61 \lambda z_i / w. \quad (6-45)$$

The corresponding result in the nonparaxial case can be shown to be

$$S = 0.61 \frac{\lambda}{\sin \theta} = 0.61 \frac{\lambda}{NA} \quad (6-46)$$

where θ represents the half-angle subtended by the exit pupil when viewed from the image plane, and NA is the *numerical aperture* of the optical system, defined by $NA = \sin \theta$. Figure 6.18 illustrates the intensity distribution in the image of two equally bright incoherent point sources separated by the Rayleigh resolution distance. The central dip is found to fall about 27% below peak intensity.

We can now ask whether the two point-source objects, separated by the same Rayleigh distance S , would be easier or harder to resolve with coherent illumination than with incoherent illumination. This question is academic for astronomical objects, but is quite relevant in microscopy, where the illumination is usually closer to coherent than incoherent, and where in some cases it is possible to control the coherence of the illumination.

As in the previous examples, the answer to this question is found to depend on the *phase distribution* associated with the object. A cross section of the image intensity can be directly written, in normalized image coordinates, as

$$I(x) = \left| 2 \frac{J_1[\pi(x - 0.61)]}{\pi(x - 0.61)} + e^{j\phi} 2 \frac{J_1[\pi(x + 0.61)]}{\pi(x + 0.61)} \right|^2$$

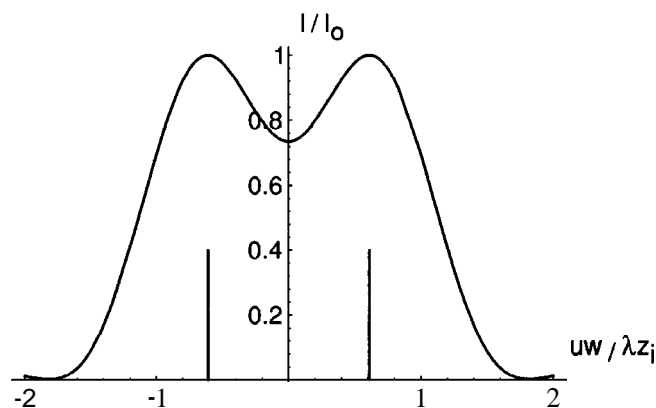


FIGURE 6.18

Image intensity for two equally bright incoherent point sources separated by the Rayleigh resolution distance. The vertical lines show the locations of the two sources.

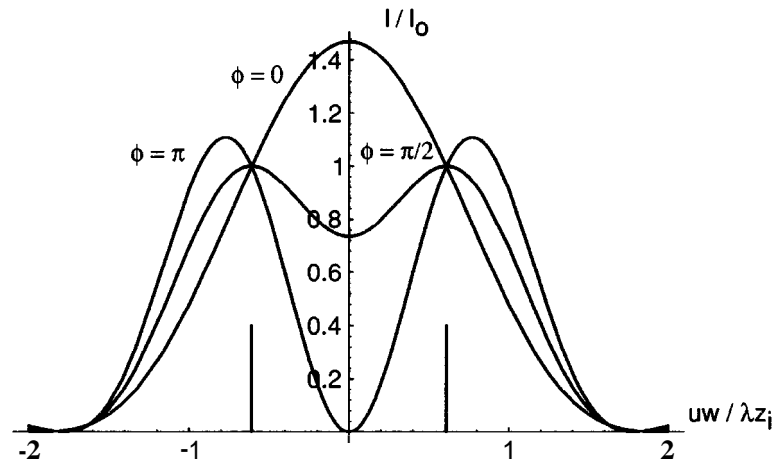


FIGURE 6.19

Image intensities for two equally bright coherent point sources separated by the Rayleigh resolution distance, with the phase difference between the two sources as a parameter. The vertical lines show the locations of the two point sources.

where ϕ is the relative phase between the two point sources. Figure 6.19 shows the distributions of image intensity for point sources in phase ($\phi = 0$ radians), in quadrature ($\phi = \pi/2$ radians), and in phase opposition ($\phi = \pi$ radians). When the sources are in quadrature, the image intensity distribution is identical to that resulting from incoherent point sources. When the sources are in phase, the dip in the image intensity is absent, and therefore the two points are not as well resolved as with incoherent illumination. Finally, when the two objects are in phase opposition, the dip falls all the way to zero intensity (a 100% dip) at the point midway between the locations of the two points, so the two points must be said to be better resolved with coherent illumination than with incoherent illumination. Thus there can again be no generalization as to which type of illumination is preferred for two-point resolution.

6.5.3 Other Effects

There are certain other miscellaneous properties of images formed with coherent light that should be mentioned in any comparison with incoherent images [71]. First, the responses of incoherent and coherent systems to sharp edges are notably different. Figure 6.20 shows the theoretical responses of a system with a circular pupil to a step function object, i.e. an object with amplitude transmittance

$$t_A(\xi, \eta) = \begin{cases} 0 & \xi < 0 \\ 1 & \xi \geq 0. \end{cases}$$

Figure 6.21 shows actual photographs of the image of an edge in the two cases. The coherent system is seen to exhibit rather pronounced "ringing". This property is analogous to the ringing that occurs in video amplifier circuits with transfer functions that fall too abruptly with frequency. The coherent system has a transfer function with sharp discontinuities, while the falloff of the OTF is much more gradual. Another important property of the coherent image is that it crosses the location of the actual edge with

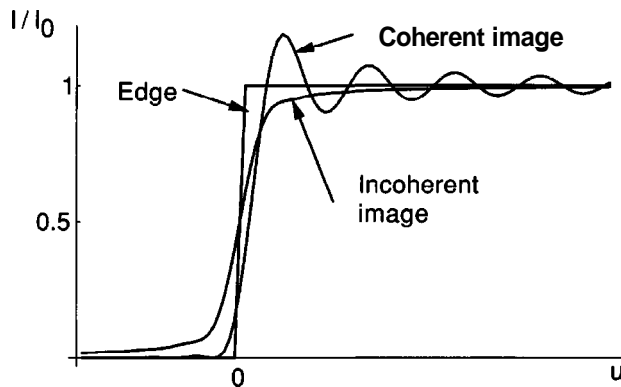


FIGURE 6.20
Images of a step in coherent and
incoherent light.

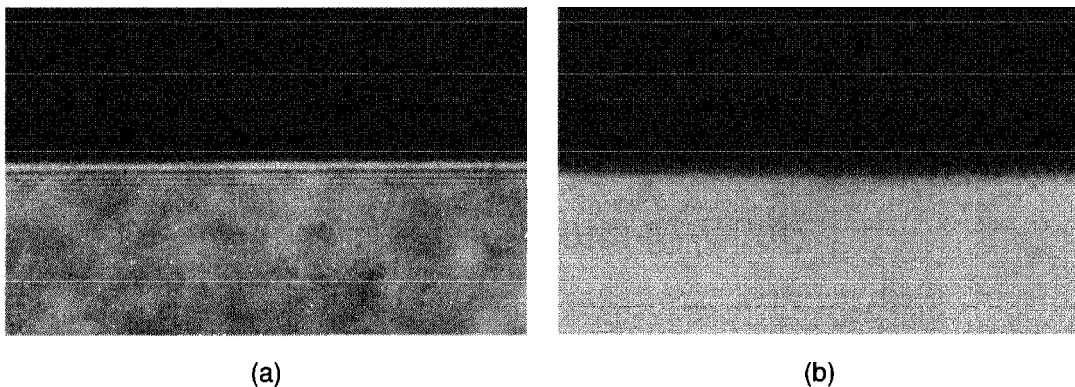


FIGURE 6.21
Photographs of the image of an edge in (a) coherent and (b) incoherent illumination. [By permission of P. S. Considine, Technical Operations, Inc., Burlington, Mass.]

only $1/4$ of its asymptotic value of intensity, whereas the incoherent image crosses with a value of $1/2$ of its asymptotic value. If we were to assume that the actual location of the edge is at the position where the intensity reaches half its asymptotic value, we would arrive at a correct estimate of the position of the edge in the incoherent case, but in the coherent case we would err in the direction of the bright side of the edge. This fact can be important, for example, in estimating the widths of lines on integrated circuit masks.

In addition, we must mention the so-called *speckle effect* that is readily observed with highly coherent illumination. While we shall consider this effect in the context of optical imaging, it has also proven to be a problem in certain other nonoptical imaging modalities, such as microwave side-looking radar and medical ultrasound imaging. Figure 6.22 shows photographs of a transparency object, illuminated through a diffuser (e.g. a piece of ground glass), taken in coherent light and incoherent light. The granular nature of the coherent image is a direct consequence of the complex, random perturbation of the wavefront introduced by the diffuser, together with the coherence of the light. For background on the speckle effect, see, for example, Refs. [221], [120], and [75]. The granularity in the image arises from interference between closely spaced and randomly phased scatterers within the diffuser. The size of the individual *speckles* can be shown [263] to be roughly the size of a *resolution cell* on the image (or object). In the case of incoherent illumination, such interference cannot take place, and speckle is missing from the image. Thus when a particular object of interest is near the resolution



FIGURE 6.22
 (a) Coherent and (b) incoherent images illustrating the speckle effect. (The object is a transparency illuminated through a diffuser.)

limit of an optical system, the speckle effect can be quite bothersome if coherent light is used. Much of this problem can be eliminated by moving the diffuser during the observation, with the result that the coherence of the illumination is at least partially destroyed and the speckles "wash out" during the measurement process. Unfortunately, as we will see in a later chapter, motion of the diffuser is not possible in conventional holography, which by its very nature is almost always a coherent imaging process, so speckle remains a particular problem in holographic imaging. The subject is discussed further in that context in Section 9.10.4.

Finally, highly coherent illumination is particularly sensitive to optical imperfections that may exist along a path to the observer. For example, tiny dust particles on a lens may lead to very pronounced diffraction patterns that will be superimposed on the image. One fundamental reason for the importance of such effects in coherent imaging is the so-called "interference gain" that occurs when a weak undesired signal interferes with a strong desired signal (see Prob. 6-17).

A reasonable conclusion from the above discussion would be that one should choose incoherent illumination whenever possible, to avoid the artifacts associated with coherent illumination. However, there are many situations in which incoherent illumination simply can not be realized or can not be used for a fundamental reason. These situations include high-resolution microscopy, coherent optical information processing, and holography.

6.6 RESOLUTION BEYOND THE CLASSICAL DIFFRACTION LIMIT

The diffraction limit to the resolution attainable by an imaging system is generally regarded to be an accurate estimate of the limits that can actually be reached in practice.

However, it is of some interest to know that *in principle*, for a certain class of objects, resolution beyond the classical diffraction limit is theoretically possible. As we shall show in this section, for the class of *spatially bounded* objects, in the absence of noise it is in principle possible to resolve infinitesimally small object details. Resolution beyond the classical diffraction limit is often referred to as *super-resolution* or *bandwidth extrapolation*.

6.6.1 Underlying Mathematical Fundamentals

There exist very fundamental mathematical reasons why, in the absence of noise and for the cited class of objects, resolution beyond the classical diffraction limit should be possible. These reasons rest on two basic mathematical principles, which we list here as theorems. For proofs of these theorems, see, for example, Ref. [136].

Theorem 1. The two-dimensional Fourier transform of a spatially bounded function is an *analytic* function in the (f_x, f_y) plane.

Theorem 2. If an analytic function in the (f_x, f_y) plane is known exactly in an arbitrarily small (but finite) region of that plane, then the entire function can be found (uniquely) by means of *analytic continuation*.

Now for any imaging system, whether coherent or incoherent, the image information arises from only a finite portion of the object spectrum (i.e. a portion of the spectrum of object amplitude in the coherent case, or a portion of the spectrum of object intensity in the incoherent case), namely, that portion passed by the transfer function of the imaging system. If this finite portion of the object spectrum can be determined exactly from the image, then, for a bounded object, the *entire* object spectrum can be found by analytic continuation. If the entire object spectrum can be found, then the exact object present can be reconstructed with arbitrary precision.

6.6.2 Intuitive Explanation of Bandwidth Extrapolation

A plausibility argument that super-resolution might be possible for a spatially limited object can be presented with the help of a simple example. For this example we assume that the object illumination is incoherent, and for simplicity we argue in one dimension rather than two. Let the object be a sinusoidal intensity distribution of finite extent, with a frequency that exceeds the incoherent cutoff frequency, as illustrated in Fig. 6.23. Note that the sinusoidal intensity necessarily rides on a rectangular background pulse, assuring that intensity remains a positive quantity. The finite-length cosine itself can be expressed as the following intensity distribution:

$$I_g(u) = \frac{1}{2} \left[1 + m \cos(2\pi \tilde{f} u) \right] \text{rect} \left(\frac{u}{L} \right).$$

It follows that the (suitably normalized) spectrum of this intensity distribution is

$$\mathcal{G}_g(f_x) = \text{sinc}(L f_x) + \frac{m}{2} \text{sinc}[L(f_x - \tilde{f})] + \frac{m}{2} \text{sinc}[L(f_x + \tilde{f})],$$

as shown in part (b) of the figure, along with the assumed OTF of the imaging system. Note that the frequency \tilde{f} lies beyond the cutoff of the OTF. The critical point to note from this figure is that the finite width of the cosinusoid has spread its spectral components into sinc functions, and while the frequency \tilde{f} lies beyond the limits of the OTF, nonetheless the tails of the sinc functions centered at $f_X = \pm \tilde{f}$ extend below the cutoff frequency into the observable part of the spectrum. Thus, within the **passband** of the imaging system, there does exist information that originated from the cosinusoidal components that lie outside the passband. To achieve super-resolution, it is necessary to retrieve these extremely weak components and to utilize them in such a way as to recover the signal that gave rise to them.

6.6.3 An Extrapolation Method Based on the Sampling Theorem

While the fundamental mathematical principles are most easily stated in terms of analytic continuation, there are a variety of specialized procedures that have been applied to the problem of bandwidth extrapolation. These include an approach based on the sampling theorem in the frequency domain [140], an approach based on prolate spheroidal wave-function expansions [15], and an iterative approach suitable for digital

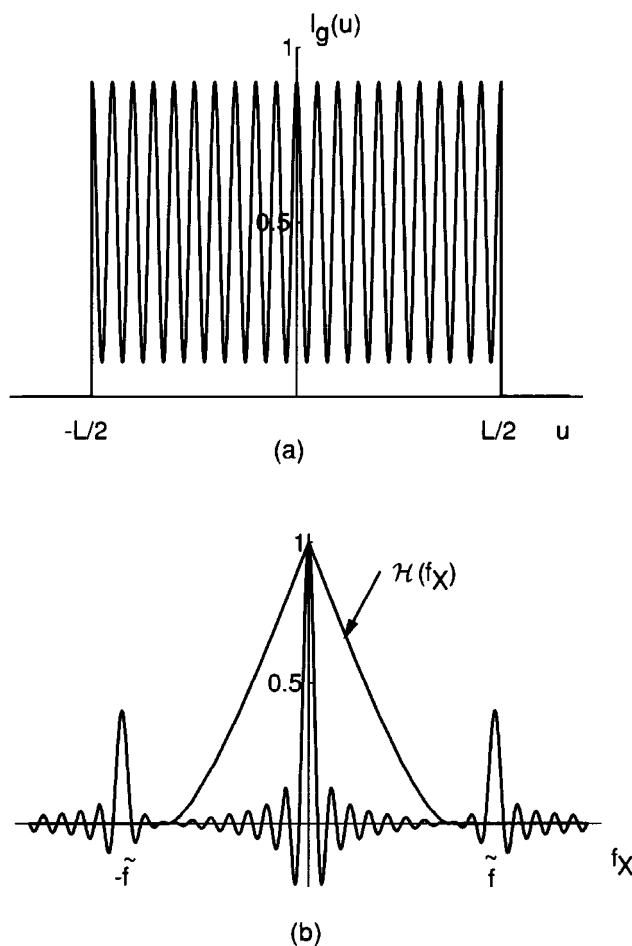


FIGURE 6.23
 (a) Object intensity distribution, and
 (b) object spectrum and the OTF.

implementation that successively reinforces constraints in the space and space-frequency domains [116], [228]. Here we will initially focus on the sampling-theorem approach due to its simplicity.

To make the analysis as simple as possible, we treat only the one-dimensional restoration problem. Extension to two dimensions is straightforward. Suppose that a one-dimensional incoherent object with intensity distribution $I_g(\mathbf{u})$ is bounded to the region $(-L/2, L/2)$ on the \mathbf{u} axis.⁹ By the Whittaker-Shannon sampling theorem, the object spectrum $\mathcal{G}_g(f)$ can be written in terms of its sample values at frequencies n/L :

$$\mathcal{G}_g(f) = \sum_{n=-\infty}^{\infty} \mathcal{G}_g\left(\frac{n}{L}\right) \text{sinc}\left[L\left(f - \frac{n}{L}\right)\right]. \quad (6-47)$$

Now, due to the limited passband of the optical system, values of $\mathcal{G}_g(n/L)$ can be found only for a few low-integer values of n . We would like, of course, to extend our knowledge of the spectrum to larger integer values, say for $-N \leq n \leq N$, so that the approximation

$$\mathcal{G}_g(f) \approx \sum_{n=-N}^N \mathcal{G}_g\left(\frac{n}{L}\right) \text{sinc}\left[L\left(f - \frac{n}{L}\right)\right] \quad (6-48)$$

would be a satisfactory representation of the image, not only within the passband of the imaging system, but also outside that passband over a frequency region of a size that depends on how large an N is chosen. The larger N is, the further beyond the classical diffraction-limited cutoff frequency we will be able to extend our knowledge of the spectrum.

To determine the sample values outside the observable passband, we measure¹⁰ the values of $\mathcal{G}_g(f)$ at any $2N + 1$ distinct frequencies f_k within the passband. The f_k in general will not coincide with the sampling points n/L . (If some of the $\mathcal{G}_g(n/L)$ lie within the observable passband, this makes our job easier, for we can then measure them, rather than find them by manipulating measurements of other quantities.) The value of the object spectrum measured at frequency f_k within the passband is represented by $\hat{\mathcal{G}}_g(f_k)$. Thus the measurements at the $2N + 1$ separate frequencies within the observable passband generate a set of $2N + 1$ simultaneous equations of the form

$$\hat{\mathcal{G}}_g(f_k) = \sum_{n=-N}^N \mathcal{G}_g\left(\frac{n}{L}\right) \text{sinc}\left[L\left(f_k - \frac{n}{L}\right)\right] \quad k = 1, 2, \dots, 2N + 1. \quad (6-49)$$

This is a set of $2N + 1$ linear equations in $2N + 1$ unknowns, the $\mathcal{G}_g(f_k)$.

It is helpful to cast this problem in matrix form. Define a column vector \vec{g} consisting of the $2N + 1$ unknown values of the $\mathcal{G}_g(n/L)$ and a column vector $\vec{\hat{g}}$ consisting of the $2N + 1$ measured values $\hat{\mathcal{G}}_g(f_k)$. In addition define a $(2N + 1) \times (2N + 1)$ matrix D with entry $\text{sinc}[L(f_k - n/L)]$ in the k th row and n th column. Then the set of equations (6-49)

⁹As usual, \mathbf{I} , actually represents the geometrical-optics prediction of the image, or the object geometrically projected into the image plane, but we refer to it as the object.

¹⁰Presumably we know the exact shape of the OTF within the passband, and can compensate for it to determine the actual values of \mathcal{G}_g at each frequency.

can be represented by the simple matrix equation

$$\hat{\vec{g}} = \mathbf{D}\vec{g}.$$

Our goal is to find the vector \vec{g} , from which we can reconstruct a satisfactory extension of the spectrum \mathcal{G}_g beyond the normal cutoff frequency.

Many methods exist for numerically inverting the set of equations of interest. Symbolically, we wish to find the inverse of the matrix \mathbf{D} , allowing us to express the matrix of unknowns \vec{g} through the equation

$$\vec{g} = \mathbf{D}^{-1}\hat{\vec{g}}.$$

It is possible to show that as long as the measurement frequencies f_k are distinct, the determinant of \mathbf{D} is nonzero, and therefore the inverse exists. Thus *in principle* the sample values of the object spectrum outside the **passband** can be determined, and a satisfactory approximation to the object spectrum can be found beyond the cutoff frequency, with the help of the interpolation functions of the sampling theorem.

Before discussing the practical limitations of this and other techniques for extrapolation, we briefly discuss one other approach to the same problem.

6.6.4 An Iterative Extrapolation Method

An iterative method for extrapolation beyond the diffraction limit is especially interesting because this type of method can be applied to many other important problems in optics. The method was applied to the bandwidth extrapolation problem first by Gerchberg [116] and by Papoulis [228]. This method is purely numerical and is readily implemented on a digital computer.

The algorithm is one that iterates between the object domain and the spectral domain, making changes in each domain to reinforce prior knowledge or measured data. Figure 6.24 shows a block diagram illustrating the steps in the algorithm. The original object intensity (the relevant quantity if the system is incoherent) is known to be space-limited and nonnegative. These are the constraints to be reinforced in the object domain. In the spectral domain, we know the object spectrum within the **passband** of the imaging system, for this data was actually measured. This is the constraint that is reinforced in the spectral domain.

Start with the measured image of the object. From a Fourier transform of that image we can discover that part of the spectrum of the object that lies within the **passband** of the imaging system. With these two pieces of data in hand, we begin the algorithm. Due to the finite **passband** of the imaging system, this image is not space-limited (or spatially bounded). We know that the object was space-limited, so we simply truncate the image to the spatial bound of the object (*i.e.* multiply it by a rectangle function of the appropriate width). The effect of the spatial truncation is to change the spectrum of the new image. In particular, spectral components are introduced outside of the **passband** of the imaging system, and in addition the spectral components within the **passband** are changed. The next step is to change the new spectral components within the **passband** to the old values, which were measured and are regarded as prior data. This changes the image again, spreading it beyond the spatial bound. Repeat the spatial bounding

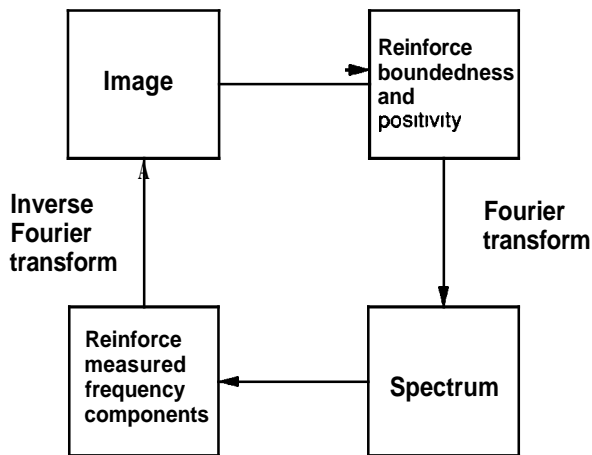


FIGURE 6.24
Block diagram of the iterative extrapolation algorithm.

process, transform again, repeat the reinforcement of the known spectral components, etc. In this way, the spectral components beyond the diffraction limit are introduced and are gradually refined to be consistent with the known information. The algorithm terminates when the image and its spectrum are changing by amounts that are smaller than some threshold. In the absence of noise, this algorithm can be shown to converge.

6.6.5 Practical Limitations

All methods for extrapolating bandwidth beyond the diffraction limit are known to be extremely sensitive to both noise in the measured data and the accuracy of the assumed a priori knowledge. See, for example, [252], [104], [52], and [257]. That this should be so is not entirely a surprise, for as we discussed previously, the information within the **passband** that arose from frequency components outside the **passband** is extremely weak (see Fig. 6.23). This sensitivity also becomes evident in the method based on the sampling theorem when the conditioning of the matrix D is considered. As the spacing between the frequencies f_k shrinks, as it must as we attempt to estimate more and more values of the spectrum on the sampling points outside the observable bandwidth, the matrix becomes more and more ill-conditioned, meaning that the solution vector \vec{g} is ultimately dominated by noise. The growth of noise sensitivity is extremely rapid. Based on these results, it is generally agreed that *the Rayleigh limit to resolution represents a practical limit to the resolution that can be achieved with a conventional imaging system*. Nonetheless, the ideas behind bandwidth extrapolation are important to be aware of and similar methods can be applied to other important problems in optics (see, for example, [273]).

PROBLEMS – CHAPTER 6

- 6-1. The mask shown in Fig. P6.1 is inserted in the exit pupil of an imaging system. Light from the small openings interferes to form a fringe in the image plane.

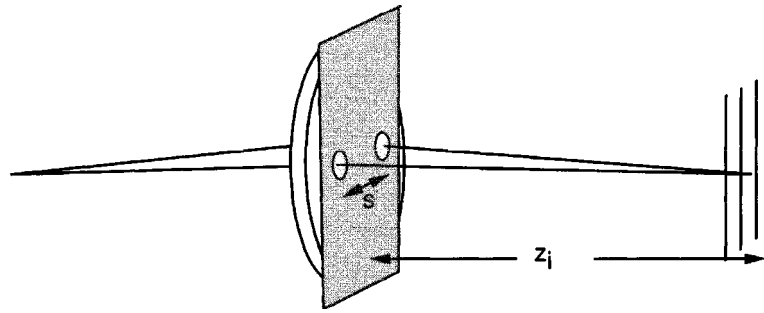


FIGURE P6.1

- (a) Find the spatial frequency of this fringe in terms of the center-to-center spacing s of the two openings, the wavelength λ , and the image distance z_i .
- (b) The openings are circular and have diameter d . Specify the envelope of the fringe pattern caused by the finite openings in the pupil plane.
- 6-2. The line-spread function of a two-dimensional imaging system is defined to be the response of that system to a one-dimensional delta function passing through the origin of the input plane.

- (a) In the case of a line excitation lying along the x axis, show that the line-spread function l and the point-spread function p are related by

$$l(y) = \int_{-\infty}^{\infty} p(x, y) dx,$$

where l and p are to be interpreted as amplitudes or intensities, depending on whether the system is coherent or incoherent, respectively.

- (b) Show that for a line source oriented along the x axis, the (1D) Fourier transform of the line-spread function is equal to a slice through the (2D) Fourier transform of the point-spread function, the slice being along the f_y axis. In other words, if $\mathcal{F}\{l\} = L$ and $\mathcal{F}\{p\} = P$, then $L(f) = P(0, f)$.
- (c) Find the relationship between the line-spread function and the step response of the system, i.e. the response to a unit step excitation oriented parallel to the x axis.
- 6-3. An incoherent imaging system has a square pupil function of width $2w$. A square stop of width w is placed at the center of the pupil, as shown in Fig. P6.3.

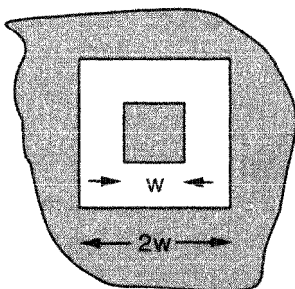


FIGURE P6.3

- (a) Sketch cross sections of the optical transfer function with and without the stop present.
- (b) Sketch the limiting form of the optical transfer function as the size of the stop approaches the size of the full pupil.

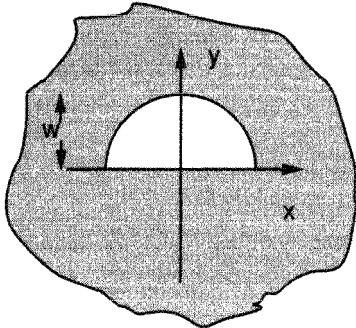


FIGURE P6.4

- 6-4.** An incoherent imaging system has a circular pupil of diameter $2w$. A half-plane stop is inserted in the pupil, yielding the modified pupil shown in Fig. P6.4. Find expressions for the optical transfer function evaluated along the f_x and f_y axes.
- 6-5.** An incoherent imaging system has a pupil consisting of an equilateral triangle, as shown in Fig. P6.5. Find the OTF of this system along the f_x and f_y axes in the spatial frequency domain.

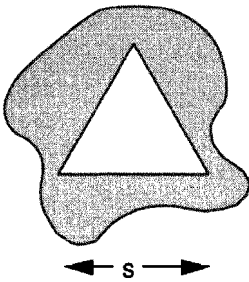


FIGURE P6.5

- 6-6.** Sketch the f_x and f_y cross sections of the optical transfer function of an incoherent imaging system having a pupil function the aperture shown in Fig. P6.6. Be sure to label the various cutoff frequencies and center frequencies on these sketches.

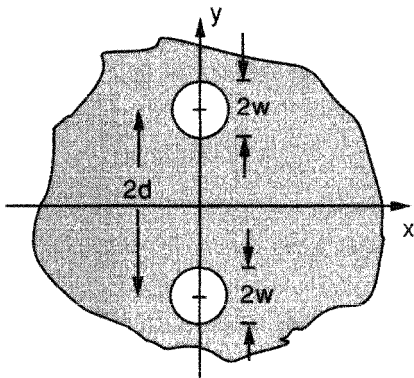


FIGURE P6.6

6-7. Consider a *pinhole camera* shown in Fig. P6.7.

Assume that the object is incoherent and nearly monochromatic, the distance z_o from the object is so large that it can be treated as infinite, and the pinhole is circular with diameter $2w$.

- Under the assumption that the pinhole is large enough to allow a purely geometrical-optics estimation of the point-spread function, find the optical transfer function of this camera. If we define the "cutoff frequency" of the camera to be the frequency where the first zero of the OTF occurs, what is the cutoff frequency under the above geometrical-optics approximation? (Hint: First find the intensity point-spread function, then Fourier transform it. Remember the second approximation above.)
- Again calculate the cutoff frequency, but this time assuming that the pinhole is so small that Fraunhofer diffraction by the pinhole governs the shape of the point-spread function.
- Considering the two expressions for the cutoff frequency that you have found, can you estimate the "optimum" size of the pinhole in terms of the various parameters of the system? Optimum in this case means the size that produces the highest possible cutoff frequency.

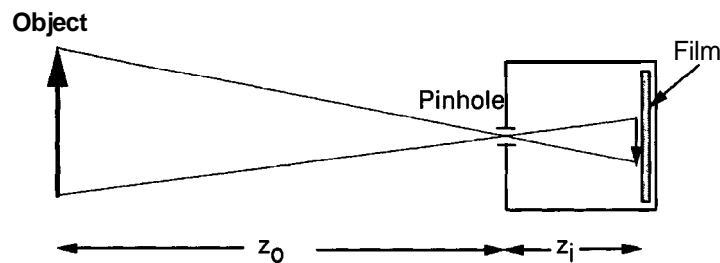


FIGURE P6.7

6-8. Consider the OTF of Eq. (6-41), as predicted for a system having square pupil and a focusing error. It is hypothesized that the point-spread function of this system is the convolution of the diffraction-limited point-spread function with the point-spread function predicted by geometrical optics. Examine the validity of this claim.

6-9. A quantity of considerable utility in determining the seriousness of the aberrations of an optical system is the *Strehl definition* \mathcal{D} , which is defined as the ratio of the light intensity at the maximum of the point-spread function of the system with aberrations to that same maximum for that system in the absence of aberrations. (Both maxima are assumed to exist on the optical axis.) Prove that \mathcal{D} is equal to the normalized volume under the optical transfer function of the aberrated imaging system; that is, prove

$$\mathcal{D} = \frac{\iint_{-\infty}^{\infty} \mathcal{H}(f_x, f_y)_{\text{with}} df_x df_y}{\iint_{-\infty}^{\infty} \mathcal{H}(f_x, f_y)_{\text{without}} df_x df_y},$$

where the notations "with" and "without" refer to the presence or absence of aberrations, respectively.

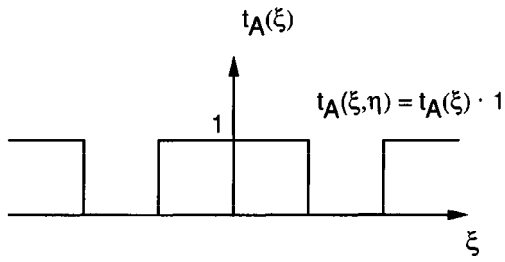


FIGURE P6.10

6-10. An object with a square-wave amplitude transmittance (shown in Fig. P6.10) is imaged by a lens with a circular pupil function. The focal length of the lens is 10 cm, the fundamental frequency of the square wave is 100 cycles/mm, the object distance is 20 cm, and the wavelength is 1 μm. What is the minimum lens diameter that will yield *any variations* of intensity across the image plane for the cases of

- (a) Coherent object illumination?
- (b) Incoherent object illumination?

6-11. An object has an intensity transmittance given by

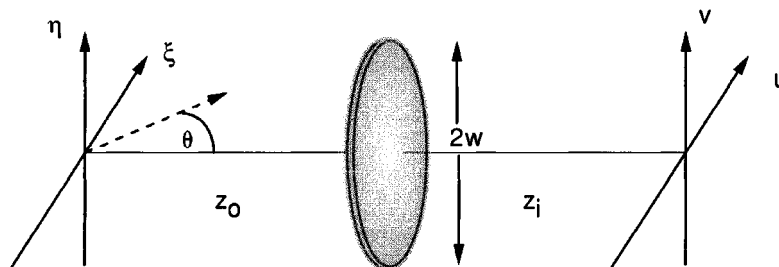
$$\tau(\xi, \eta) = \frac{1}{2} (1 + \cos 2\pi \tilde{f} \xi)$$

and introduces a constant, uniform phase delay across the object plane. This object is placed at distance 2f in front of a positive lens of focal length f, and the image is examined in a plane 2f behind the lens. Compare the maximum frequencies \tilde{f} transmitted by the system for the cases of coherent and incoherent illumination.

6-12. A sinusoidal amplitude grating with transmittance

$$t_A(\xi, \eta) = \frac{1}{2} (1 + \cos 2\pi \tilde{f} \xi)$$

is placed in front of a thin, positive lens (circular with diameter 2w, focal length f) and obliquely illuminated by a monochromatic plane wave traveling at angle θ to the z axis in the (ξ, z) plane, as shown in Fig. P6.12.



- (a) What is the Fourier transform of the amplitude distribution transmitted by the object?
- (b) Assuming $z_i = z_o = 2f$, what is the maximum angle θ for which *any variations* of intensity will appear in the image plane?

- (c) Assuming that this maximum angle is used, what is the intensity in the image plane, and how does it compare with the corresponding intensity distribution for $\theta = 0$?
- (d) Assuming that the maximum angle θ is used, what is the maximum grating frequency \tilde{f} that will yield variations of intensity in the image plane? How does this frequency compare with the cutoff frequency when $\theta = 0$?
- 6-13.** The F-number of a lens with a circular aperture is defined as the ratio of the focal length to the lens diameter. Show that when the object distance is infinite, the cutoff frequency for a coherent imaging system using this lens is given by $f_o = \frac{1}{2\lambda F\#}$, where $F\#$ represents the F-number.
- 6-14.** The Sparrow resolution criterion states that two equally strong incoherent point sources are barely resolved when their separation is the maximum separation for which the image of the pair of points shows no dip at the midpoint. This condition can be equivalently stated as one for which the curvature of the total intensity at the midpoint between the centers of the individual spread functions vanishes.
- (a) Show that, for a spread function that is an even function of u , such a condition occurs when the separation (in the u direction) between the centers of the spread functions is twice the value of u that satisfies the equation

$$\frac{\partial^2 |h(u, 0)|^2}{\partial u^2} = 0$$

where $|h|^2$ is the intensity point-spread function of the system.

- (b) What is the value of the Sparrow separation (in the image space) for a system with a square aperture of width $2w$, where an edge of the aperture runs parallel to the direction of separation of the two sources?
- 6-15.** Consider the step responses of two different imaging systems, one with a circular aperture of diameter $2w$ and the second with a square aperture of width $2w$, with one edge of the aperture parallel with the edge of the step. All other aspects of the two systems are identical.
- (a) Show that, with coherent illumination, the step responses of the two systems are identical.
- (b) Show that, with incoherent illumination, the step responses of the two systems are not identical.
- (c) Describe how you would numerically calculate the step responses in both cases.
- 6-16.** Show that the intensity image of a step-object (edge along the η axis) formed by a coherent imaging system having a square pupil (width $2w$) with edges parallel to and orthogonal to the direction of the step can be expressed as

$$I_i(u, v) = c \left| 1 + \text{Si} \left(\frac{2\pi w u}{\lambda z_i} \right) \right|^2$$

where $\text{Si}(z)$ is defined by

$$\text{Si}(z) = \int_0^z \frac{\sin t}{t} dt$$

and c is a constant. Note: The function $\text{Si}(z)$ is a tabulated function and is known to many mathematical software packages.

- 6-17. Consider the addition of a strong desired field of amplitude A with a weak undesired field of amplitude a . You may assume that $A \gg a$.
- (a) Calculate the relative perturbation $\Delta I/|A|^2$ to the desired intensity caused by the presence of the undesired field when the two fields are mutually coherent.
 - (b) Repeat for the case of mutually incoherent fields.
- 6-18. Using the definition of mutual intensity, show that any purely monochromatic wave is fully coherent spatially, and therefore must be analyzed as a system that is linear in amplitude.

Wavefront Modulation

It is clear from the previous chapter that the tools of linear systems and frequency analysis are useful in the *analysis* of optical systems. However, the theory becomes much more significant if it can be applied to *synthesis* problems as well. In order to synthesize linear optical systems with desired properties, the ability to manipulate light waves is needed. In particular, such an ability is needed to introduce information into an optical system, since the information is carried directly by the optical amplitude in the case of coherent systems, and by the optical intensity in the case of incoherent systems. In addition, for coherent optical information processing systems we require the ability to modify and manipulate the complex optical fields transmitted through the focal plane of a lens, for through such manipulation we are able to filter the input data in various desired ways.

For the above reasons, in this chapter attention is focused on methods for spatially modulating optical wavefields, especially coherent fields. The traditional means of modulation has been through the use of photographic materials, so we consider the properties of such materials in Section 7.1. However, much more powerful optical information processing systems can be realized if photographic film is replaced by *spatial light modulators* capable of changing transmitted light in real time in response to optical or electrical control signals. Many approaches to the construction of spatial light modulators have been studied over the years. In Section 7.2 we focus on just a few of the most important types of such devices.

Finally, in Section 7.3 we consider several approaches to constructing optical elements that control the complex amplitude of transmitted light in fixed but complicated ways, so-called *diffractive optical elements*. As their name implies, these elements control transmitted light through diffraction rather than refraction. Often a computer is employed in the design and construction of these elements, and their properties can be much more complicated than those of refractive elements.

7.1

WAVEFRONT MODULATION WITH PHOTOGRAPHIC FILM

Photographic film is a basic component of optical systems in general and optical information processing systems in particular. Film can play three very fundamental roles in optics. First, it can serve as a detector of optical radiation, a task it performs remarkably efficiently. Second, it can serve as a storage medium for images, capable of retaining information for long periods of time. Third, it can serve as a spatial modulator of transmitted or reflected light, a role of particular importance in optical information processing. All of these functions are achieved at extremely low cost.

Because of the importance of photographic film in optics and in optical information processing, we devote some time here to discussing its properties. For a more comprehensive treatment of the photographic process, see, for example, Ref. [208]. Other useful references include [266] and [22].

7.1.1 The Physical Processes of Exposure, Development, and Fixing

An unexposed photographic film or plate generally consists of a very large number of tiny silver halide (often **AgBr**) grains suspended in a gelatin support, which in turn is attached to a firm "base" consisting of acetate or mylar¹ for films, and glass for plates. The soft emulsion also has a thin layer of a protective overcoating on its exposed surface, as illustrated in the cross section shown in Fig. 7.1. In addition, certain sensitizing agents are added to the gelatin; these agents have a strong influence on the introduction of dislocation centers within the silver halide crystals. Light incident on the emulsion initiates a complex physical process that is outlined as follows:

1. A photon incident on a silver halide grain may or may not be absorbed by that grain. If it is absorbed, an electron-hole pair is released within the grain.
2. The resulting electron is in the conduction band, is mobile within the silver halide crystal, and eventually, with some probability, becomes trapped at a crystal dislocation.

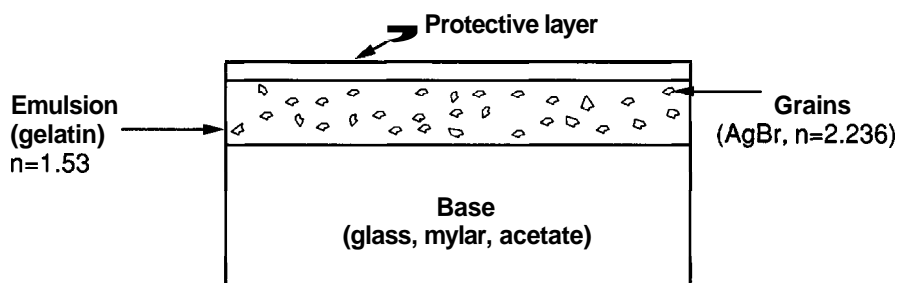


FIGURE 7.1
Structure of a photographic film or plate.

¹Mylar base should be avoided when coherent light is used, due to the fact that it is birefringent and causes unwanted variations of the polarization and phase of the transmitted light.

3. The trapped electron electrostatically attracts a silver ion; such ions are mobile even before exposure by light, a consequence of thermal agitation.
4. The electron and the silver ion combine to form a single atom of metallic silver at the dislocation site. The lifetime of this combination is rather short, of the order of a few seconds.
5. If within the lifetime of this first silver atom, a second silver atom is formed by the same process at the same site, a more stable two-atom unit is formed with a lifetime of at least several days.
6. Typically at least two additional silver atoms must be added to the silver speck in order for it ultimately to be developable. The existence of a threshold, requiring several trapped electrons to activate the development process, is responsible for good stability of unexposed film on the shelf.

The speck of silver formed as above is referred to as a *development speck*, and the collection of development specks present in an exposed emulsion is called the *latent image*. The film is now ready for the development and fixing processes.

The exposed photographic transparency is immersed in a chemical bath, the developer, which acts on silver specks containing more than the threshold number² of silver atoms. For such grains, the developer causes the entire crystal to be entirely reduced to metallic silver. The ratio of the number of silver atoms in a developed grain to the number of photons that must be absorbed to make the grain developable is typically of the order of 10^9 , a number which is often called the "gain"⁷ of the photographic process.

At this point the processed emulsion consists of two types of grains, those that have been turned to silver, and those that did not absorb enough light to form a development center. The latter crystals are still silver halide and, without further processing, will eventually turn to metallic silver themselves simply through thermal processes. Thus in order to assure stability of the image, it is necessary to remove the undeveloped silver halide grains, a process called *fixing* the emulsion. The transparency is immersed in a second chemical bath, which removes the remaining silver halide crystals from the emulsion, leaving only the stable metallic silver.

The processes of exposure, development and fixing are illustrated in Fig. 7.2.

7.1.2 Definition of Terms

The field of photography has developed a certain nomenclature that should be mastered if the properties of photographic emulsions are to be discussed in any detailed way. At this point we introduce the reader to some of these terms.

Exposure. The energy incident per unit area on a photographic emulsion during the exposure process is called the exposure. Represented by the symbol E , it is equal to the product of incident intensity \mathcal{I} at each point and the exposure time T ,

²The threshold is actually not a fixed number, but a statistical one. The assumption that the threshold is four atoms is an approximation.

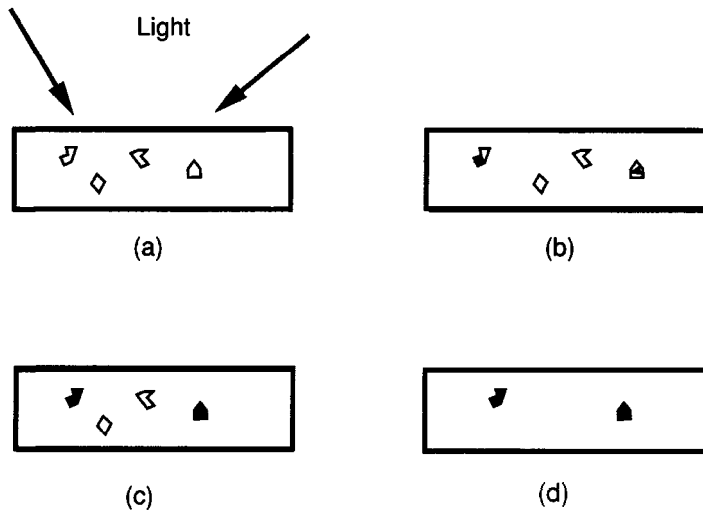


FIGURE 7.2
Pictorial representation of the photographic process. (a) Exposure, (b) latent image, (c) after development, and (d) after fixing. Only the emulsion is shown.

$$E(x, y) = \mathcal{I}(x, y)T.$$

The units for exposure are mJ/cm^2 . Note that the symbol \mathcal{I} is used for intensity incident on the film during exposure; we reserve the symbol \mathbf{I} to represent intensity incident on (or transmitted by) the transparency after development.

Intensity transmittance. The ratio of intensity transmitted by a developed transparency to the intensity incident on that transparency, averaged over a region that is large compared with a single grain but small compared with the finest structure in the original exposure pattern, is called the intensity transmittance. Represented by the symbol τ , it is equal to

$$\tau(x, y) = \frac{\text{local}}{\text{average}} \left\{ \frac{\mathbf{I} \text{ transmitted at } (x, y)}{\mathbf{I} \text{ incident at } (x, y)} \right\}.$$

Photographic density. In the year 1890, F. Hurter and V.C. Driffield published a classic paper in which they showed that the logarithm of the reciprocal of the intensity transmittance of a photographic transparency should be proportional to the silver mass per unit area of that transparency. They accordingly defined the photographic density D as

$$D = \log_{10} \left(\frac{1}{\tau} \right).$$

The corresponding expression for intensity transmittance in terms of density is

$$\tau = 10^{-D}.$$

Hurter-Driffield curve. The most common description of the photometric properties of a photographic emulsion is the Hurter-Driffield curve, or the H&D curve, for short. It is a plot of photographic density D vs. the exposure E that gave rise to that density. A typical H&D curve is shown in Fig. 7.3 for the case of a photographic negative.

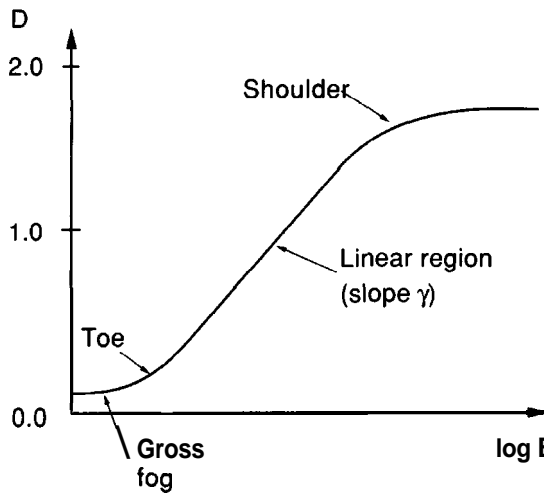


FIGURE 7.3
The Hurter-Driffield curve for a typical emulsion.

Note the various regions of the H&D curve. When the exposure is below a certain level, the density is independent of exposure and equal to a minimum value called gross fog. In the toe of the curve, density begins increasing with exposure. There follows a region of considerable extent in which the density is linearly proportional to the logarithm of exposure—this is the region most commonly used in ordinary photography. The slope of the curve in this linear region is referred to as the gamma of the emulsion and is represented by the symbol γ . Finally the curve saturates in a region called the *shoulder*, beyond which there is no change of density with increasing exposure.

A film with a large value of γ is called a high-contrast film, while a film with a low γ is called a *low-contrast* film. The particular value of γ achieved in any case is influenced by three major factors: (1) the type of emulsion in question (for example, Plus-X and Tri-X are low contrast films, with gammas of 1 or less, while High Contrast Copy has a gamma of 2 or 3); (2) the particular developer used; and (3) the development time. Figure 7.4 illustrates a typical dependence of γ on development time. With a judicious choice of film, developer, and development time, it is possible to achieve a prescribed value of γ with a fair degree of accuracy.

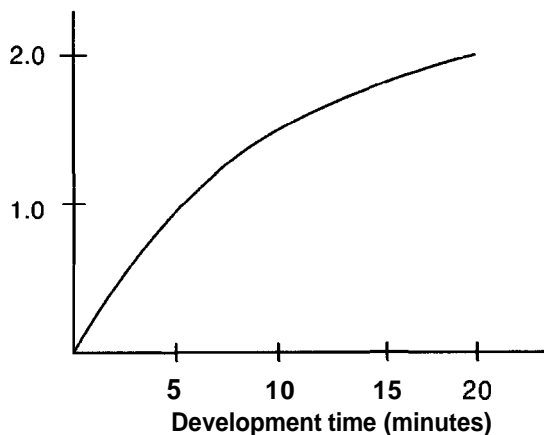


FIGURE 7.4
The dependence of film gamma on development time.

7.1.3 Film in an Incoherent Optical System

In many of its uses, film may be regarded as an element that maps an intensity distribution \mathcal{I} incident during exposure into an intensity distribution \mathbf{I} transmitted after development. Such a point-of-view is particularly appropriate when film is used as an element of an incoherent optical system. We consider now the detailed nature of the mapping so defined.

Assuming that the film is used in the linear region of the **H&D** curve, the density D may be written

$$D = \gamma_n \log_{10} E - D_o = \gamma_n \log_{10}(\mathcal{I} T) - D_o, \quad (7-1)$$

where γ_n is the slope of the linear region of the curve, and $-D_o$ is the value of D where a straight-line approximation would meet the D axis were there no toe. The subscript n on γ is used to indicate that we are dealing with a negative transparency.

The intensity incident during exposure can be related to the intensity transmitted after development by recalling the definition of photographic density,

$$D = \log_{10} \left(\frac{1}{\tau_n} \right).$$

When this definition is substituted into Eq. (7-1), we find

$$\log_{10} \tau_n = -\gamma_n \log_{10}(\mathcal{I} T) + D_o$$

or equivalently

$$\tau_n = 10^{D_o} (\mathcal{I} T)^{-\gamma_n}.$$

Finally,

$$\tau_n = K_n \mathcal{I}^{-\gamma_n} \quad (7-2)$$

where K_n is a positive constant. Note that the intensity mapping defined by this relation is a highly nonlinear one for any positive value of γ_n .

It is also possible to achieve a **positive** power-law relation between intensity transmittance and intensity incident during exposure, although to do so generally requires a two-step photographic process. During the first step, the negative transparency is made in the usual fashion. During the second step, the light transmitted by the negative transparency is used to expose a second emulsion, and the result is a final positive transparency. To understand this process more quantitatively, let the transmittance of the first transparency be written, from (7-2),

$$\tau_n = K_{n1} \mathcal{I}^{-\gamma_{n1}}.$$

If this transparency is placed in contact with a second unexposed emulsion and illuminated with intensity I_0 , then the intensity incident on the second emulsion is simply $\tau_n I_0$, and the resulting intensity transmittance becomes

$$\tau_p = K_{n2} (I_0 \tau_n)^{-\gamma_{n2}} = K_{n2} I_0^{-\gamma_{n2}} K_{n1}^{-\gamma_{n2}} \mathcal{I}^{\gamma_{n1} \gamma_{n2}}$$

or equivalently

$$\tau_p = K_p \mathcal{I}^{-\gamma_p} \quad (7-3)$$

where K_p is a positive constant and by convention $\gamma_p = -\gamma_{n1}\gamma_{n2}$ is a *negative* number, making the overall exponent $-\gamma_p$ a *positive* number. From this result we can see that again the general relation between intensity incident during exposure and intensity transmitted after development is a nonlinear one, but in the specific case of an overall gamma equal to unity, the process becomes a linear one.

While a linear mapping of intensity incident during exposure into intensity transmittance after development has been seen to occur only under very special conditions, nonetheless film can be shown to provide a linear mapping of *incremental changes* of intensity under a much wider class of conditions. For further development of this point, the reader is referred to Prob. 7-1.

7.1.4 Film in a Coherent Optical System

When film is used as an element of a *coherent* optical system, it is more appropriately regarded as providing either (1) a mapping of intensity incident during exposure into complex field transmitted after development, or (2) a mapping of complex amplitude incident during exposure into complex amplitude transmitted after development. The second viewpoint can be used, of course, only when the light that exposes the transparency is itself coherent, and must incorporate the fact that all phase information about the incident complex wavefield is lost upon detection. Only when interferometric detection is used can phase information be captured, and such detection systems will be seen to benefit from the first viewpoint, rather than the second.

Since the complex amplitude of the transmitted light is, from both viewpoints, the important quantity in a coherent system, it is necessary to describe a transparency in terms of its complex *amplitude* transmittance t_A [187]. It is most tempting to define t_A simply as the positive square root of the intensity transmittance τ . However, such a definition neglects the relative phase shifts that can occur as the light passes through the film [152]. Such phase shifts arise as a consequence of variations of the film or plate thickness, which can originate in two distinct ways. First, there are generally random thickness variations across the base of the film, *i.e.* the base is not optically flat. Second, the thickness of the emulsion is often found to vary with the density of the silver in the developed transparency. This latter variation is strongly dependent on the exposure variations to which the film has been subjected. Thus a complete description of the amplitude transmittance of the film must be written

$$t_A(x, y) = \sqrt{\tau(x, y)} \exp[j\phi(x, y)] \quad (7-4)$$

where $\phi(x, y)$ describes the pattern of phase shifts introduced by the transparency.

In most applications, the thickness variations are entirely undesired, for they cannot easily be controlled. It is possible to remove the effects of these variations by means of a device called a *liquid gate*. Such a device consists of two pieces of glass, each ground and polished to be optically flat on one side, between which the transparency and an index matching fluid (often oil) can be sandwiched, as illustrated in Fig. 7.5. The flat surfaces of the glass are, of course, facing the outside, and the index of refraction of the

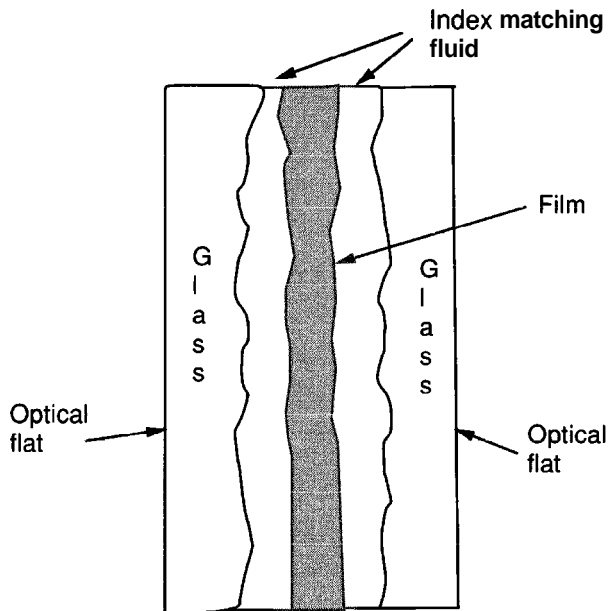


FIGURE 75
A liquid gate for removing film thickness variations. The thickness variations are greatly exaggerated.

fluid must be chosen as a compromise, for it is impossible to match simultaneously the different indices of the base, the emulsion, and the glass. However, with a proper choice of fluid, the optical path length through the liquid gate can be made nearly constant, allowing the amplitude transmittance of the film and gate to be written

$$t_A(x, y) = \sqrt{\tau(x, y)}. \quad (7-5)$$

When the phase shifts have been removed, a combination of Eqs. (7-5), (7-2), and (7-3) allows the amplitude transmittance to be expressed as

$$t_A(x, y) = \kappa I^{-\gamma/2} = \kappa |U|^{-\gamma}, \quad (7-6)$$

where U is the complex amplitude of the field incident during exposure, κ is a constant, and γ is a positive number for a negative transparency and a negative number for a positive transparency.

As will be seen in many of the examples to be discussed in later sections, it is often desirable to have film act as a square-law mapping of complex amplitude. Such behavior can be achieved in a number of ways, one of which is to make a positive transparency with an overall gamma of -2 , as can be seen from (7-6). In order to obtain a maximum dynamic range of exposure over which this relation holds, the first gamma of the two-step process is often chosen less than unity (for example $1/2$), while the second gamma is chosen greater than 2 (for example 4), such that their product remains equal to 2.

It is possible, however, to obtain square-law action over a limited dynamic range with a transparency of any gamma, be it a positive or a negative. This point is most easily seen by abandoning the traditional **H&D** curve description of film and making instead a direct plot of amplitude transmittance vs. exposure (on a linear scale). Such a description was advocated at an early stage by **Maréchal** and was very successfully used by **Kozma** [176] in an analysis of the effects of photographic nonlinearities. Figure 7.6 shows a plot of amplitude transmittance vs. exposure (the t_A - E curve) for a typical

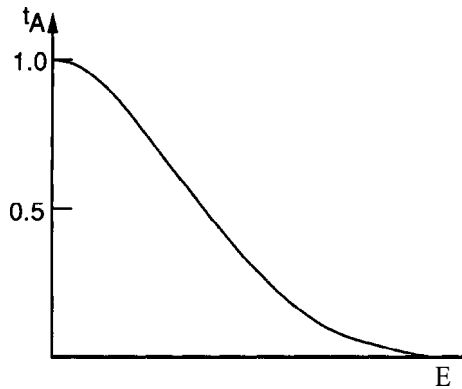


FIGURE 7.6
Typical amplitude transmittance vs. exposure curve.

negative transparency. If the film is "biased" to an operating point that lies within the region of maximum linearity of this curve, then over a certain dynamic range the film will provide a square-law mapping of incremental changes in incident amplitude into incremental changes of amplitude transmittance. Thus if E_b represents the bias exposure and t_b the corresponding bias amplitude transmittance, we may represent the t_A - E curve within its region of linearity by

$$t_A \approx t_b + \beta(E - E_b) = t_b + \beta' |\Delta U|^2 \quad (7-7)$$

where β is the slope of the curve at the bias point, ΔU represents the incremental amplitude changes, and β' is the product of β and the exposure time. Note that β and β' are negative numbers for a negative transparency.

In general, a high-gamma film has a steeper slope to its t_A - E curve than does a low-gamma film and therefore is more efficient in transferring small changes of exposure into changes of amplitude transmittance. However, this increased efficiency is often accompanied by a smaller dynamic range of exposure over which the t_A - E curve remains linear. As an additional point of interest, the bias point at which maximum dynamic range is obtained is found to lie in the toe of the H&D curve.

Before closing this section we note that when thin gratings are recorded by interference in a photographic emulsion, as is often the case in the construction of spatial filters and in recording holograms, it may often be desirable to achieve the highest possible diffraction efficiency, rather than the widest possible dynamic range. It can be shown (see [266], p. 7) that, for small modulations, the maximum diffraction efficiency for a thin sinusoidal grating recorded photographically will occur for a recording made in the region where the magnitude of the slope α of the t_A vs. $\log E$ curve of the emulsion is maximum. This curve is yet another description of the properties of photographic emulsions that is relevant in some applications.

7.1.5 The Modulation Transfer Function

To this point we have tacitly assumed that any variations of exposure, however fine on a spatial scale, will be transferred into corresponding variations of silver density according to the prescription implied by the H&D curve. In practice, one finds that when the spatial scale of exposure variations is too small, the changes of density induced may

be far smaller than would be implied by the **H&D** curve. We can say in very general terms that each given type of film has a limited spatial frequency response.

The spatial frequency response of an emulsion is limited by two separate phenomena:

1. Light scattering within the emulsion during exposure.
2. Chemical diffusion during the development process.

Both of these phenomena are linear ones, although the physical quantities with respect to which they are linear are different. Light scattering is linear in the variable exposure, while chemical diffusion is linear in the variable density. It might be hoped that the linear phenomena that limit spatial frequency response could be separated from the highly nonlinear behavior inherent in the **H&D** curve. This in fact can be done by regarding the photographic process as a cascade of several separate mappings, as illustrated in Fig. 7.7. The first operation in this cascade is a linear, invariant filter representing the effects of light scattering and the resulting spread or blur of the exposure pattern E . The output of this filter, E' , then passes through the **H&D** curve, which is regarded as a *zero-spread nonlinearity*, analogous to the zero-memory nonlinearities often encountered in the analysis of communications systems. The output of the **H&D** curve is a density D' , which is itself subjected to linear spreading and blur by the chemical diffusion process to produce a final density D . This model is often referred to as the "**Kelley model**", after D.H. Kelley who originated it. Often the model is simplified to include only a single linear filter that precedes the nonlinear **H&D** curve, thus ignoring the linear filter associated with diffusion.

The effects of the linear filters are, of course, to limit the spatial frequency response of the emulsion. If the model is simplified to one with a single linear filter preceding

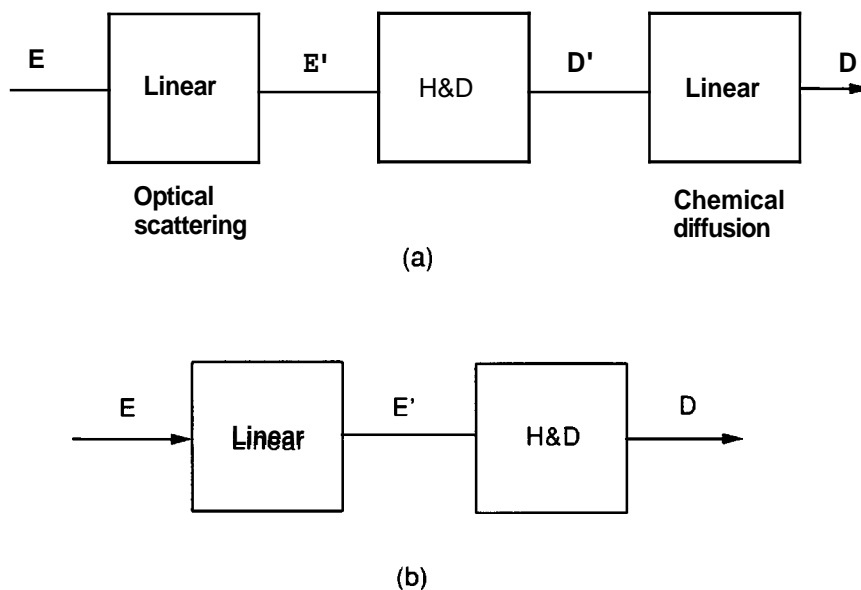


FIGURE 7.7
The Kelley model of the photographic process. (a) Full model;
(b) simplified model.

the nonlinear mapping (Fig. 7.7(b)), then it is of some interest to find the transfer function of the filtering operation, usually referred to as the *modulation transfer function* (*MTF*) of the photographic process. To measure the characteristics of the linear filter, a sinusoidal exposure pattern

$$E = E_0 + E_1 \cos 2\pi f x \quad (7-8)$$

can be applied (such a pattern is easily generated by interference of two mutually coherent plane waves on the emulsion). The "modulation" associated with the exposure is defined as the ratio of the peak variation of exposure to the background exposure level, or

$$M_i = \frac{E_1}{E_0}. \quad (7-9)$$

If the variations of density in the resulting transparency are measured, they can be referred back to the exposure domain through the H&D curve (assumed known) to yield an inferred or "effective" sinusoidal exposure pattern, as indicated in Fig. 7.8. The modulation M_{eff} of the effective exposure distribution will always be less than the modulation M_i of the true exposure distribution. Accordingly the modulation transfer function of the film is defined as

$$M(f) = \frac{M_{\text{eff}}(f)}{M_i(f)}$$

where the dependence on the spatial frequency f of the exposure has been emphasized. In most cases encountered in practice, the form of the point-spread function of the scat-

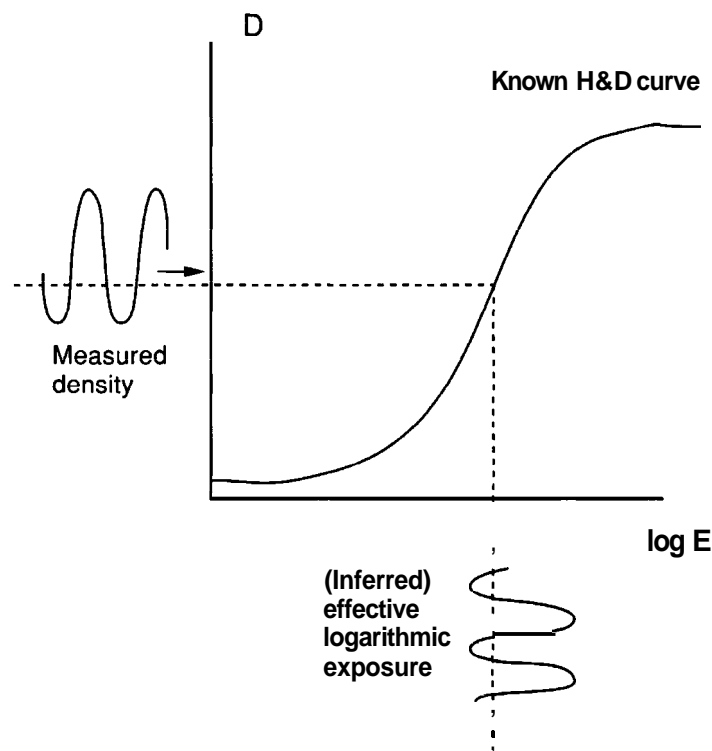


FIGURE 7.8

Measurement of the MTF by projecting back through the H&D curve.

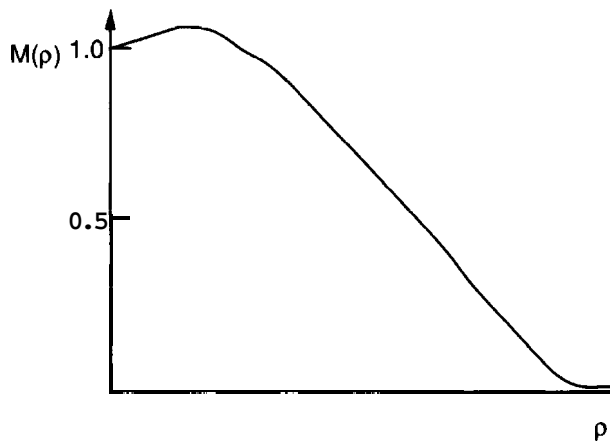


FIGURE 7.9
Typical measured MTF curve.

tering process (approximately Gaussian) is circularly symmetric and there are no phase shifts associated with the transfer function. The effective exposure distribution applied to the nonlinear portion of the film mapping may therefore be written

$$E' = E_0 + M(f)E_1 \cos 2\pi f x. \quad (7-10)$$

Figure 7.9 illustrates the typical measured frequency dependence of the MTF of an emulsion, plotted vs. radial spatial frequency ρ . The small hump rising above unity at low frequencies is caused by chemical diffusion (the final linear filtering box in our model, which was ignored in the procedure for measuring the MTF) and is referred to as arising from the *adjacency effect*.

The range of frequencies over which significant frequency response is obtained varies widely from emulsion to emulsion, depending on grain size, emulsion thickness, and other factors. By way of illustration, Plus-X film has a significant response to about 50 line-pairs (cycles)/mm, while for Kodak 649F spectroscopic plate it extends to beyond 2000 line-pairs/mm.

7.1.6 Bleaching of Photographic Emulsions

Conventional photographic emulsions modulate light primarily through absorption caused by the metallic silver present in the transparency. As a consequence, significant amounts of light are lost when an optical wave passes through such a spatial modulator. In many applications it is desired to have a more efficient modulator, one that can operate primarily through *phase modulation* rather than absorption. Such structures can be realized with photographic materials, provided they are subjected to *chemical bleaching*.

The bleaching process is one that removes metallic silver from the emulsion and leaves in its place either an emulsion thickness variation or a refractive index variation within the emulsion. The chemical processes that lead to these two different phenomena are in general different. A thickness variation results when a so-called *tanning bleach* is used, while a refractive index modulation occurs when a *nontanning bleach* is used.

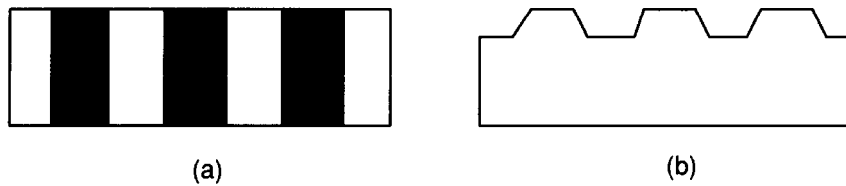


FIGURE 7.10
A relief image produced by a tanning bleach. (a) Original density image. (b) Relief image after bleaching.

Considering first the tanning bleach, the chemical agents used in this type of bleach release certain chemical byproducts as they remove the metallic silver, and these byproducts cause a cross-linking of the gelatin molecules within the emulsion in regions where the silver concentration was high. As the transparency is dried, the hardened areas shrink less than do the unhardened areas, with the result that a relief *image* is formed, with the thickest regions of the emulsion being where the density was highest, and the thinnest regions where the density was lowest. Figure 7.10 illustrates the phenomenon for the case of a square-wave density pattern. This phenomenon is found to depend strongly on the spatial frequency content of the density pattern, and to act as a **bandpass** filter, with no relief induced at very low spatial frequencies and at very high spatial frequencies. For a **15- μm -thick** emulsion, the peak thickness variations are found to occur at a spatial frequency of about **10 cycles/mm**, with a maximum relief height in the **1- to 2- μm** range. Using such a bleach it is possible, for example, to make an approximately sinusoidal relief grating, which will exhibit diffraction efficiencies typical of sinusoidal phase gratings, considerably higher than those of sinusoidal amplitude gratings.

Nontanning bleaches, on the other hand, produce internal refractive index changes within the emulsion, rather than relief images. For such bleaches, the metallic silver within the developed transparency is changed back by the chemical bleach to a transparent silver halide crystal, with a refractive index considerably larger than that of the surrounding gelatin. In addition, the bleach must remove the sensitizing agents found in unexposed silver halide crystals to prevent them from turning to metallic silver due to thermal effects and additional exposure to light. The resulting refractive index structures constitute a pure phase image. The spatial frequency response of this kind of bleached transparency is not a **bandpass** response, but rather is similar to that of the original silver image. Very high-frequency phase structures can be recorded using this method. Phase shifts of the order of 2π radians can be induced in a wavefront passing through the bleached emulsion, although this number obviously depends on the emulsion thickness.

7.2 SPATIAL LIGHT MODULATORS

The technology of photographic emulsions has a long history and is extremely well developed. However, such materials have one distinct disadvantage when image or signal

processing is of concern, namely the long time delays required for chemical processing. In the event that the data to be processed is originally in photographic form, this may not pose a significant problem. However, if information is being rapidly gathered, perhaps by some electronic means, one would prefer a more direct interface between the electronic information and the data processing system. For this reason those working in the field of optical **information** processing have explored a large number of devices capable of converting data in electronic form (or sometimes in incoherent optical form) into spatially modulated coherent optical signals. Such a device is called a *spatial light modulator*, a term that is abbreviated by SLM.

There is a broad categorization of SLMs into two classes that can be made at the start: (1) electrically written SLMs and (2) optically written SLMs. In the former case, electrical signals representing the information to be input to the system (perhaps in raster format) directly drive a device in such a way as to control its spatial distribution of absorption or phase shift. In the latter case, the information may be input to the SLM in the form of an optical image at the start (for example from a CRT display), rather than in electrical form. In this case the function of the SLM may be, for example, to convert an incoherent image into a coherent image for subsequent processing by a coherent optical system. Often a given SLM technology may have two different forms, one suitable for electrical addressing and one suitable for optical addressing.

Optically addressed SLMs have several key properties besides their fast temporal response that are very useful for optical processing systems. First, they can convert incoherent images into coherent images, as alluded to above. Second, they can provide image amplification: a weak incoherent image input to an optically addressed SLM can be read out with an intense coherent source. Third, they can provide wavelength conversion: e.g. an incoherent image in the infrared could be used to control the amplitude transmittance of a device in the visible.

SLMs are used not only to input data to be processed, but also to create spatial filters that can be modified in real time. In such a case the SLM is placed in the back focal plane of a Fourier transforming lens, where it modifies the transmitted amplitude of the fields in accord with a desired complex spatial filter.

Over the history of optical information processing, a great many different SLM technologies have been explored. Books have been written on this subject (see, for example, [91]). For a review article covering the properties of more types of SLMs than will be discussed here, the reader may wish to consult Ref. [220] and its associated references. In addition, a series of meeting proceedings on the subject provides valuable information [88], [89], [90]. Here we limit ourselves to presenting the barest outlines of the principles of operation of what are currently regarded as the most important SLM technologies. These include (1) liquid crystal SLMs, (2) magneto-optic SLMs, (3) deformable mirror SLMs, (4) multiple-quantum-well (MQW) SLMs, and (5) acousto-optic Bragg cells.

7.2.1 Properties of Liquid Crystals

The use of liquid crystals in low-cost displays is commonplace. Examples include watch displays and screens for laptop computers. In such applications voltages applied to

pixelated electrodes cause a change in the intensity of the light transmitted by or reflected from the display. Similar principles can be used to construct a spatial light modulator for input to an optical information processing system.

Background on the optics of liquid crystals can be found in [253], Sections 6.5 and 18.3. For an additional reference that covers liquid crystal displays in detail, the reader can consult [159]. See also Ref. [91], Chapters 1 and 2.

Mechanical properties of liquid crystals

Liquid crystal materials are interesting from a physical point-of-view because they share some of the properties of both solids and liquids. The molecules composing such materials can be visualized as ellipsoids, with a single long axis about which there is circular symmetry in any transverse plane. These ellipsoidal molecules can stack next to one another in various ways, with different geometrical configurations defining different general types of liquid crystals. Adjacent molecules are not rigidly bound to one another, and can rotate or slide with respect to one another under the application of mechanical or electrical forces, thus exhibiting some of the properties of a liquid. However, there are constraints on the geometrical organization of collections of molecules, and these constraints introduce some properties normally associated with solids.

There are three different general classes (or phases) of liquid crystals that are of general interest in optics: (1) nematic, (2) smectic, and (3) cholesteric. The classes are differentiated by the different molecular orders or organizational constraints, as illustrated in Fig. 7.11. For nematic liquid crystals (NLC), the molecules throughout the entire volume of the material favor a parallel orientation, with randomly located centers within that volume. For smectic liquid crystals, the molecules again favor parallel alignment, but their centers lie in parallel layers, with randomness of location only within a layer. Finally, a cholesteric liquid crystal is a distorted form of a smectic liquid crystal in which, from layer to layer, the alignment of molecules undergoes helical rotation about an axis. Spatial light modulators are based primarily on nematic liquid crystals and on a special class of smectic liquid crystals (the so-called smectic-C* phase) called ferroelectric liquid crystals (FLC), so our discussions will focus on these types primarily.

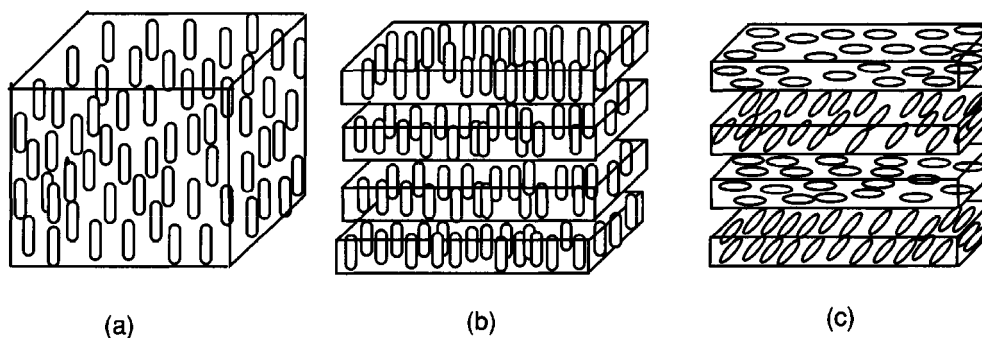
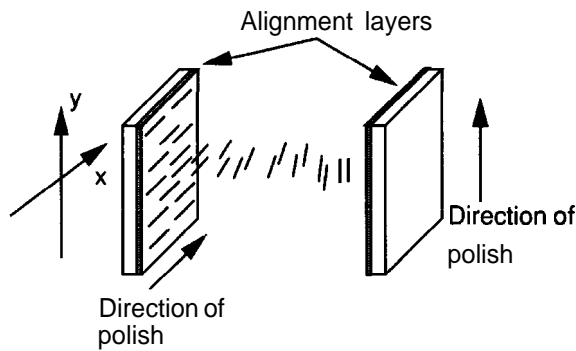


FIGURE 7.11

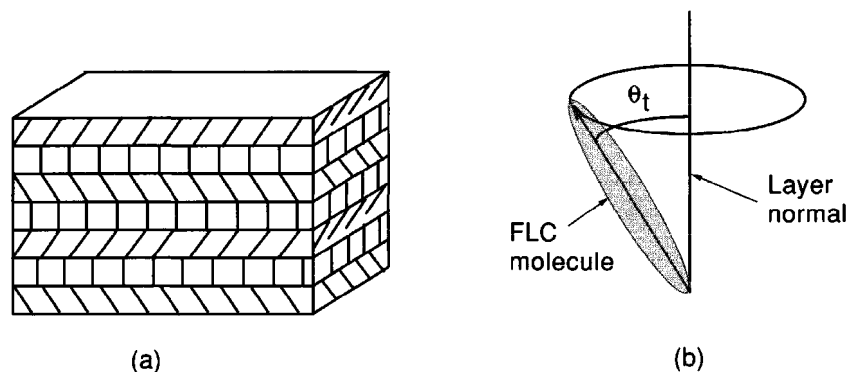
Molecular arrangements for different types of liquid crystals. (a) Nematic liquid crystal, (b) smectic liquid crystal, and (c) cholesteric liquid crystal. The layers in (b) and (c) have been separated for clarity.

**FIGURE 7.12**

Molecular arrangements in a twisted nematic liquid crystal. The lines between the alignment layers indicate the direction of molecular alignment at various depths within the cell. Only a small column of molecules is shown.

It is possible to impose boundary conditions on the alignment of nematic liquid crystal molecules contained between two glass plates by polishing soft alignment layers coated on those plates with strokes in the desired alignment direction. The small scratches associated with the polishing operation establish a preferred direction of alignment for the molecules that are in contact with the plate, with their long direction parallel with the scratches. If the two alignment layers are polished in different directions (for example, in orthogonal directions), then the tendency of the molecules to remain aligned with one another (characteristic of the nematic liquid crystal phase) and the alignment of the molecules with the direction of polish at the glass plates combine to create a *twisted* nematic liquid crystal, as illustrated in Fig. 7.12. Thus as we move between the two plates, the directions of the long axes of the various molecules remain parallel to one another in planes parallel to the glass plates, but gradually rotate between those planes to match the boundary conditions at the alignment layers.

The structure of ferroelectric liquid crystals is more complex. Since they are of the smectic type, their molecules are arranged in layers. Within a given layer, the molecules are aligned in the same direction. For smectic-C* materials, the angle of the molecules within a single layer is constrained to lie at a specific declination angle θ_t with respect to the layer normal, and thus there is a cone of possible orientations for any given layer. Figure 7.13 illustrates the structure of the surface stabilized FLC for large cell thickness. The directions of orientation between layers form a helical spiral.

**FIGURE 7.13**

Ferroelectric liquid crystal (a) smectic-C* layered structure, and (b) allowed molecular orientations.

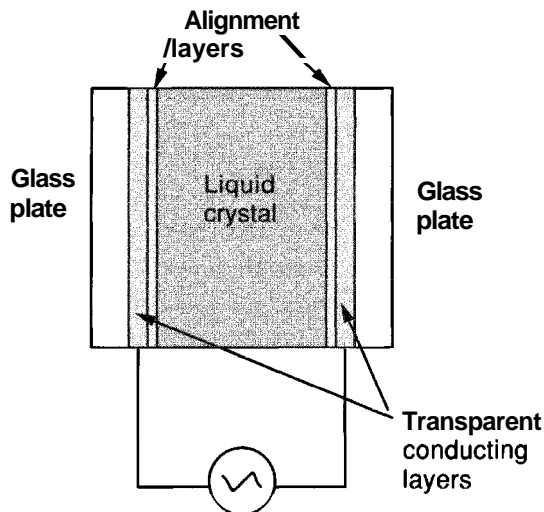


FIGURE 7.14
Structure of an electrically controlled liquid crystal cell.

The angular directions of the two layers at the interfaces with the glass plates can be stabilized by aligned polishing³ [64]. In practice, the cells are made sufficiently thin (typically only a very few microns of thickness) to eliminate the possibility that different layers will be in different allowed states.

Electrical properties of liquid crystals

Both displays and SLMs exploit the ability to change the transmittance of a liquid crystal by means of applied electric fields. Usually those fields are applied between the glass plates that contain the liquid crystal material using transparent conductive layers (indium tin oxide films) coated on the inside of the glass plates. In order to achieve alignment of the liquid crystal at the interface, the conductive layer is covered with a thin alignment layer (often polyimide) which is subjected to polishing, as shown in Fig. 7.14.

The application of an electric field across such a device can induce an electric dipole in each liquid crystal molecule, and can interact with any permanent electric dipoles that may be present. If, as is usually the case, the dielectric constant of a molecule is larger in the direction of the long axis of the molecule than normal to that axis, the induced dipoles have charge at opposite ends of the long direction of the molecule. Under the influence of the applied fields, the torques exerted on these dipoles can cause the liquid crystal molecules to change their natural spatial orientation.

For nematic liquid crystals, which do not have the extra constraints of smectic and cholesteric materials, a sufficiently large applied voltage will cause the molecules that are not in close proximity to the alignment layers to rotate freely and to align their long axes with the applied field. Thus the arrangement of the molecules within the twisted nematic liquid crystal cell shown previously in Fig. 7.12 will change under sufficient applied field to the arrangement shown in Fig. 7.15, in which the vast majority

³The liquid crystal cell is filled with material at an elevated temperature, where the phase of the liquid crystal is smectic-A. Such a phase has no tilt angle, and therefore the molecules align with their long direction parallel to the alignment grooves. When the material cools, it is transformed to the smectic-C* state, which has the tilt mentioned above.

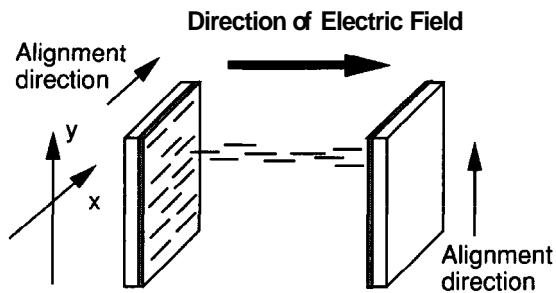


FIGURE 7.15 Twisted nematic liquid crystal with a voltage applied. Only a small column of molecules is shown.

of the molecules have their long axis aligned with the field, i.e. pointing in a direction normal to the glass plates. As we shall discuss shortly, the change in the orientation of the molecules changes the optical properties of the cell as well. To avoid permanent chemical changes to the NLC material, cells of this type are driven by AC voltages, typically with frequencies in the range of 1 kHz to 10 kHz and with voltages of the order of 5 volts. Note that because the dipole moment of a nematic liquid crystal is an induced moment rather than a permanent moment, the direction of the moment reverses when the applied field reverses in polarity. Thus the direction of the torque exerted by the field on the molecules is independent of the polarity of the applied voltage, and they align in the same direction with respect to the applied field, regardless of polarity.

In the case of the ferroelectric liquid crystal cell, the molecules can be shown to have a permanent electric dipole (with an orientation normal to the long dimension of the molecules), which enhances their interaction with the applied fields, and leads to only two allowable orientation states, one for each possible direction of the applied field. Figure 7.16 shows the molecules oriented at angle $+\theta_f$ to the surface normal for one direction of the applied field and $-\theta_f$ to the surface normal for the other direction of applied field. Because of the permanent dipole moment of the FLC molecules, the current state is retained by the material even after the applied field is removed. The FLC cell is thus bistable and has memory. It is because of the permanent dipole moment that

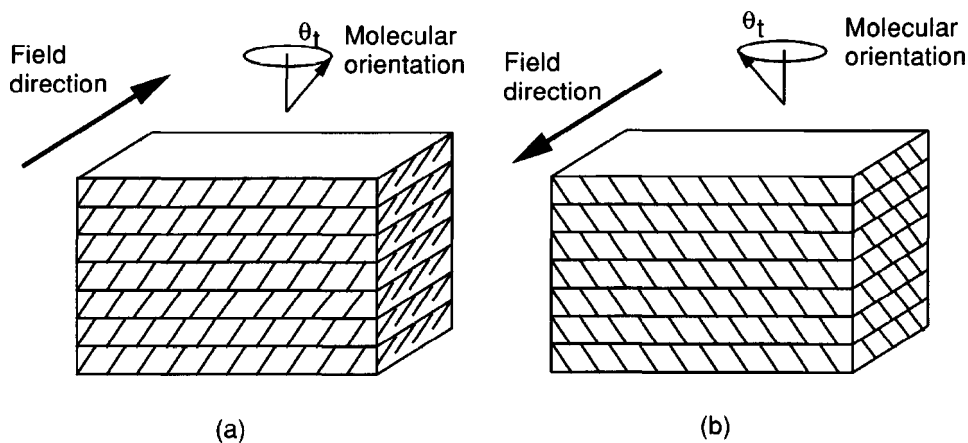


FIGURE 7.16 Ferroelectric liquid crystal molecules align in one of two allowed directions, depending on the direction of the field. The angles of orientation in the two states are separated by $2\theta_f$.

the direction of the applied field matters. Unlike the case of nematic liquid crystals, DC fields of opposite polarity must be applied to the ferroelectric liquid crystal in order to switch between states.

Liquid crystals have high resistivity and therefore act primarily as an electrical dielectric material. The electrical response of a liquid crystal cell is predominantly that of a simple RC circuit, where the resistance arises from the finite resistivity of the transparent electrodes and the capacitance is that of a parallel plate capacitor (the NLC cell is typically 5 to 10 μm thick). For sufficiently small cells, or sufficiently small pixels on a large array, the electrical time constant is small by comparison with the time constant associated with the mechanical rotation of the molecules. Typical time constants for NLC materials are approximately 100 μs for the molecules to align with an applied field, and 20 ms for the molecules to relax back to their original state. The permanent dipole moment of the FLC materials makes them considerably faster; cell thicknesses are typically in the 1- to 2- μm range, applied voltages are typically in the 5- to 10-volt range, and switching times of the order of 50 μs . In some cases even submicrosecond response times are observed [65].

Optical properties of nematic and ferroelectric liquid crystals

A quantitative understanding of the behavior of SLMs based on liquid crystals, as well as many other types of SLMs that operate by means of polarization effects, requires the use of a mathematical formalism known as the Jones calculus. This formalism is outlined in Appendix C, to which the reader is referred. The state of polarization of a monochromatic wave with X and Y components of polarization expressed in terms of complex phasors U_X and U_Y is represented by a polarization vector \vec{U} with components U_X and U_Y ,

$$\vec{U} = \begin{bmatrix} U_X \\ U_Y \end{bmatrix} \quad (7-11)$$

The passage of light through a linear polarization-sensitive device is described by a 2×2 Jones matrix, such that the new polarization vector \vec{U}' is related to the old polarization vector \vec{U} through the matrix equation

$$\vec{U}' = \mathbf{L} \vec{U} = \begin{bmatrix} l_{11} & l_{12} \\ l_{21} & l_{22} \end{bmatrix} \vec{U}. \quad (7-12)$$

If we can characterize a given device by specifying its Jones matrix, we will then be able to understand completely the effect of that device on the state of polarization of an incident wave.

The elongated structure of liquid crystal molecules causes such materials to be anisotropic in their optical behavior, and in particular to exhibit birefringence, or to have different refractive indices for light polarized in different directions. The most common materials have larger refractive indices for optical polarization parallel to the long axis of the crystal (the extraordinary refractive index, n_e), and smaller uniform refractive index for all polarization directions normal to the long axis (the ordinary refractive index, n_o). One of the highly useful properties of these materials is the very large difference between the extraordinary and ordinary refractive indices they exhibit, often in the range of 0.2 or more.

It can be shown (see [253], pp. 228-230) that for a twisted nematic liquid crystal with no voltage applied, having a helical twist of α radians per meter in the right-hand sense along the direction of wave propagation and introducing a relative retardation β radians per meter between the extraordinary and ordinary polarization components, a wave polarized initially in the direction of the long molecular axis at the entrance surface of the cell will undergo polarization rotation as the light propagates through the cell, with the direction of polarization closely tracking the direction of the long crystal axis, provided only that $\beta \gg \alpha$. The Jones matrix describing such a transformation can be shown to be the product of a coordinate rotation matrix $\mathbf{L}_{\text{rotate}}(-\alpha d)$ and a wave retarder $\mathbf{L}_{\text{retard}}(\beta d)$,

$$\mathbf{L} = \mathbf{L}_{\text{rotate}}(-\alpha d)\mathbf{L}_{\text{retard}}(\beta d), \quad (7-13)$$

where the coordinate rotation matrix is given by

$$\mathbf{L}_{\text{rotate}}(-\alpha d) = \begin{bmatrix} \cos \alpha d & -\sin \alpha d \\ \sin \alpha d & \cos \alpha d \end{bmatrix}, \quad (7-14)$$

and the retardation matrix is (neglecting constant phase multipliers)

$$\mathbf{L}_{\text{retard}}(\beta d) = \begin{bmatrix} 1 & 0 \\ 0 & e^{-j\beta d} \end{bmatrix}, \quad (7-15)$$

where β is given by

$$\beta = \frac{2\pi(n_e - n_o)}{\lambda_o} \quad (7-16)$$

Here λ_o is the vacuum wavelength of light and d is the cell thickness. With the help of this Jones matrix, the effects of the twisted nematic cell with no voltage applied can be found for any initial state of polarization.

When voltage is applied to an NLC cell along the direction of wave propagation, the molecules rotate so that the long axis coincides with that direction, and no polarization rotation occurs. Under this condition both α and β go to zero, and the cell has no effect on the incident polarization state. Thus an NLC can be used as a changeable polarization rotator, with rotation experienced in the unexcited state (no voltage applied) by an amount determined by the orientation of the alignment layers on the two glass plates as well as the thickness of the cell, and no rotation experienced in the excited state (voltage applied).

To consider the case of an FLC cell, a bit of further background is needed. When a liquid crystal cell of thickness d has all of its molecules tilted such that the long dimension of the molecule lies in the (x, y) plane, but tilted at angle $+\theta_t$ to the y (vertical) axis, the effects of the cell on incident light can be represented by a Jones matrix that is the sequence of a coordinate rotation with angle θ_t , which aligns the direction of the y axis with the long axis of the molecules, a retardation matrix representing the phase shift experienced by polarization components oriented parallel to the long and short axes of the liquid crystal molecule, followed by a second rotation matrix with angle $-\theta_t$, which returns the y axis to its original orientation at angle $-\theta_t$ to the long axis of the molecule. Taking account of the proper ordering of the matrix product,

$$\begin{aligned} \mathbf{L} &= \mathbf{L}_{\text{rotate}}(-\theta_t) \mathbf{L}_{\text{retard}}(\beta d) \mathbf{L}_{\text{rotate}}(\theta_t) \\ &= \begin{bmatrix} \cos \theta_t & -\sin \theta_t \\ \sin \theta_t & \cos \theta_t \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & e^{-j\beta d} \end{bmatrix} \begin{bmatrix} \cos \theta_t & \sin \theta_t \\ -\sin \theta_t & \cos \theta_t \end{bmatrix}, \end{aligned} \quad (7-17)$$

where β is again given by Eq. (7-16).

The Jones matrix for such an FLC cell has two possible forms, one for each direction of the applied field. When the applied field switches the direction of alignment so that the long molecular axis is at angle $+\theta_t$ to the y axis, then from Eq. (7-17) the Jones matrix is of the form

$$\mathbf{L}_+ = \mathbf{L}_{\text{rotate}}(-\theta_t) \mathbf{L}_{\text{retard}}(\beta d) \mathbf{L}_{\text{rotate}}(\theta_t), \quad (7-18)$$

whereas for the field in the opposite direction we have

$$\mathbf{L}_- = \mathbf{L}_{\text{rotate}}(\theta_t) \mathbf{L}_{\text{retard}}(\beta d) \mathbf{L}_{\text{rotate}}(-\theta_t). \quad (7-19)$$

A case of special interest is that of a cell thickness d such that the retardation satisfies $\beta d = \pi$ (i.e. the cell is a half-wave plate). The reader is asked to verify (see Prob. 7-3) that the two Jones matrices above can be reduced to the forms

$$\begin{aligned} \mathbf{L}_+ &= \begin{bmatrix} \cos 2\theta_t & \sin 2\theta_t \\ \sin 2\theta_t & -\cos 2\theta_t \end{bmatrix} \\ \mathbf{L}_- &= \begin{bmatrix} \cos 2\theta_t & -\sin 2\theta_t \\ -\sin 2\theta_t & -\cos 2\theta_t \end{bmatrix}. \end{aligned} \quad (7-20)$$

Furthermore, when the input to the FLC cell is a linearly polarized wave with polarization vector inclined at angle $+\theta_t$ to the y axis, the output polarization vectors in the two respective cases are found to be

$$\begin{aligned} \vec{U}'_+ &= \begin{bmatrix} \sin \theta_t \\ -\cos \theta_t \end{bmatrix} \\ \vec{U}'_- &= \begin{bmatrix} -\sin 3\theta_t \\ -\cos 3\theta_t \end{bmatrix}. \end{aligned} \quad (7-21)$$

Finally we note that, if the tilt angle of the liquid crystal is 22.5° , the two vectors above are orthogonal, aside from a sign change indicating a 180° phase shift. Thus for this particular tilt angle, a wave with linear polarization coincident with the long molecular axis in one state of the device is rotated by 90° when the device is switched to the opposite state. Such a device is therefore a 90° rotator for this particular direction of input polarization.

Liquid crystal cells are often used to construct intensity modulators, and indeed such modulation is important for several different types of **SLMs**. Consider first the case of nematic liquid crystals. If the NLC cell has a polarizer on its front surface and a polarization analyzer on its rear surface, it can modulate the intensity of the light it transmits. For example, in the case of a 90° twist illustrated previously in Fig. 7.12, with a polarizer oriented parallel to the front-surface alignment and an analyzer oriented parallel to the rear-surface alignment, light will pass through the exit analyzer when no voltage is applied to the cell (a consequence of rotation), but will be blocked due to the absence of rotation when the full extinction voltage is applied to the cell. If less than the

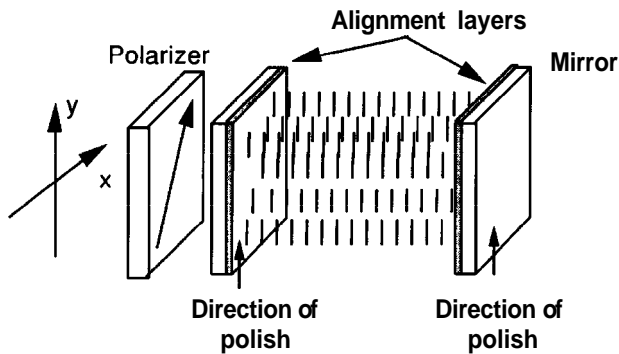


FIGURE 7.17
Intensity modulation with a reflective
NLC cell.

full extinction voltage is applied, then over a certain range of voltage, partial intensity transmission will occur, with a limited dynamic range of analog operation. Similarly, an FLC can act as a 90° polarization rotator (as explained above) and therefore can act as a binary intensity modulator.

It is also possible to make a reflection modulator using a liquid crystal cell, as illustrated in Fig. 7.17. For NLC materials, an untwisted cell is simplest. Consider a cell with the long molecular axis (the "slow" axis) aligned parallel to the y axis throughout the cell. Let the thickness of the cell be chosen to assure a 90° relative retardation of polarization components oriented along and orthogonal to the slow axis after one pass through the cell (i.e. the cell is a quarter-wave plate). The output glass plate on the cell is replaced by a mirror, and a polarizer oriented at 45° to the x axis is inserted at the front of the cell.

The operation of this cell can be understood intuitively as follows. The light incident on the cell is linearly polarized at $+45^\circ$ to the x axis due to the presence of the polarizer. When no voltage is applied across the cell, there is no molecular rotation. After the first passage through the cell, the incident linear polarization has been converted to circular polarization. Reflection from the mirror reverses the sense of circular polarization, and a second pass back through the quarter-wave plate results in a linear polarization that is orthogonal to the original polarization. Thus the reflected light is blocked by the polarizer.

On the other hand, in the presence of a sufficiently large applied voltage, the long axes of the molecules all rotate to alignment with the direction of the applied field, which coincides with the direction of propagation of the wave, eliminating the birefringence of the cell. The direction of linear polarization is therefore maintained after passage through the cell, is unchanged after reflection from the mirror, and is unchanged after the second passage through the cell. The reflected light is therefore transmitted by the polarizer.

Application of a voltage that is less than that required to fully rotate the molecules will result in partial transmission of the reflected light.

In a similar fashion, it is possible to show that an FLC cell with tilt angle 22.5° will act as a binary reflection intensity modulator if the input polarizer is aligned along one of the long molecular orientation axes and the cell thickness is chosen to realize a quarter-wave plate.

This completes the background on liquid crystal cells, and we now turn attention to specific spatial light modulators based on these materials.

7.2.2 Spatial Light Modulators Based on Liquid Crystals

Of the SLM technologies that have been explored over a period of many years, the liquid crystal devices have survived the longest and remain important devices in practice. There are many variants of these devices, some using nematic liquid crystals, and others using ferroelectric liquid crystals. We present a brief overview of the most important types.

The Hughes liquid crystal light valve

The most widely used liquid crystal SLM in optical information processing is the Hughes liquid crystal light valve. Unlike the devices discussed in the previous sections, which had their states changed by application of an electric field, this device is written optically, rather than electrically. However, optical writing results in the establishment of certain internal electric fields, and therefore the functioning of this device can be understood based on the previous background. A complete description of this rather complex device can be found in Ref. [133]. Our description will be somewhat simplified.

A diagram of the structure of the device is shown in Fig. 7.18. The device can be written with incoherent or coherent light of any state of polarization, and it is read with polarized coherent light. A polarizer and analyzer are external to the device, as will be discussed. To understand the operation of the device, we begin with the "write" side shown on the right of Fig. 7.18.

Let an optical image be cast on the right-hand entrance of the device, which can consist of a glass plate or, for better preservation of resolution, a fiber-optic faceplate. The light passes through a transparent conducting electrode and is detected by a photosensor, which in the most common version of the device is cadmium sulfide (CdS).

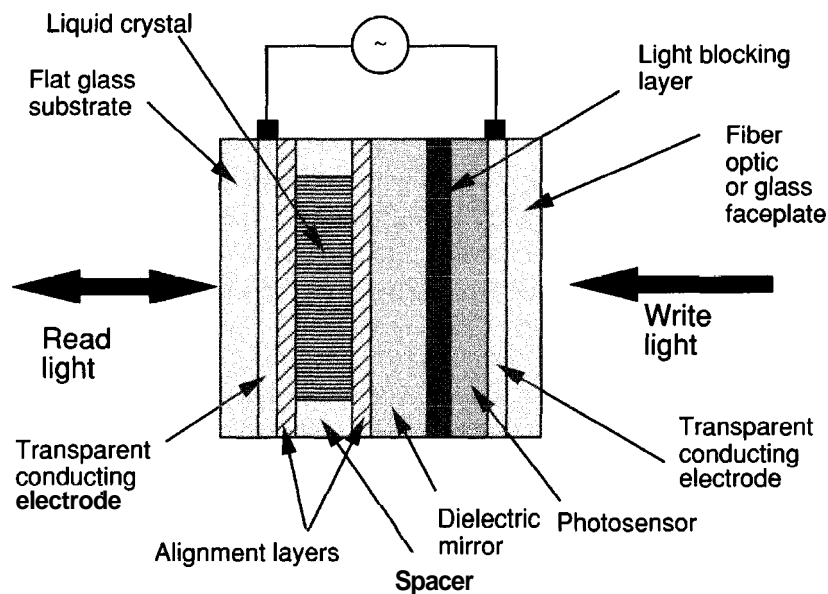


FIGURE 7.18
Hughes liquid crystal SLM.

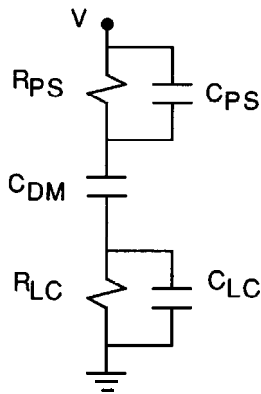


FIGURE 7.19 Electrical model for the optically written SLM. R_{PS} and C_{PS} are the resistance and capacitance of the photosensor, C_{DM} is the capacitance of the dielectric mirror, and R_{LC} and C_{LC} are the resistance and capacitance of the liquid crystal layer.

The photoconductor should have the highest possible resistivity in the absence of write light, and the lowest possible resistivity in the presence of strong write light. Thus light absorbed by the photoconductor increases its local electrical conductivity in proportion to the incident optical intensity. To the left of the photoconductor is a light-blocking layer composed of cadmium telluride (CdTe), which optically isolates the write side of the device from the read side. An audio frequency AC voltage, with an rms voltage in the 5- to 10-volt range, is applied across the electrodes of the device.

On the read side of the device, an optically flat glass faceplate is followed to the right by a transparent conducting electrode, to the right of which is a thin NLC cell with alignment layers on both sides. The alignment layers are oriented at 45° to one another, so that with no applied voltage the liquid crystal molecules undergo a 45° twist. Following the liquid crystal is a dielectric mirror which reflects incident read light back through the device a second time. The dielectric mirror also prevents DC currents from flowing through the device, which extends its lifetime.

From the electrical point-of-view, it is the rms AC voltage applied across the liquid crystal layer that determines the optical state of the read side of the device. A simplified electrical model [14] for the device is shown in Fig. 7.19. In the off state (no write light applied), the two resistances are sufficiently large that they can be neglected, and the values of the capacitances of the photosensor and the dielectric stack must be sufficiently small (i.e. their impedances at the drive frequency must be sufficiently high) compared with the capacitance of the liquid crystal layer that the rms voltage across the liquid crystal layer is too small to cause the molecules to depart from their original twisted state. In the on state, ideally there is no voltage drop across the photosensor, and the fraction of the applied rms voltage appearing across the liquid crystal must be large enough to cause significant rotation of the molecules. The capacitances involved can be controlled in the design of the device, through appropriate choice of layer thicknesses, to satisfy these requirements.⁴

Figure 7.20 illustrates the write and read operations. The liquid crystal layer is operated in a so-called "hybrid-field-effect" mode, which is explained as follows. The polarization of the incident read light is chosen to be in a direction parallel to the long

⁴In the real device, operation is complicated by the fact that the photosensor and the light-blocking layer together form an electrical diode with asymmetric $I - V$ properties.

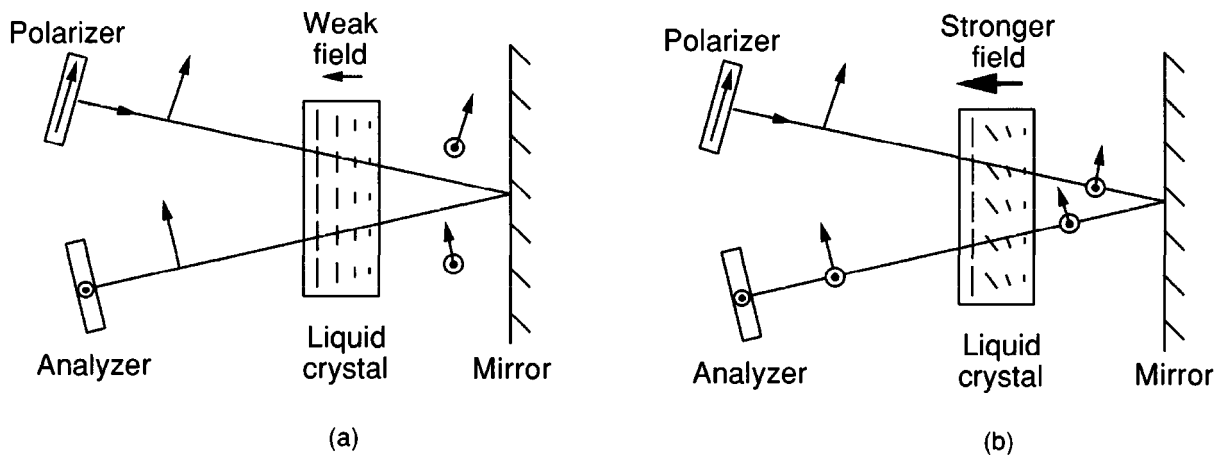


FIGURE 7.20 Readout of the Hughes liquid crystal SLM with (a) no write light present, and (b) write light present.

axis of the aligned liquid-crystal molecules at the left-hand alignment layer. Thus as light passes through the liquid crystal layer, the direction of polarization follows the twisted direction of the liquid crystal molecules, arriving at the dielectric mirror with a 45° polarization rotation. After reflection, the light propagates back through the liquid crystal a second time, with the direction of polarization again following the alignment of the molecules, thus returning to its original state. A polarization analyzer oriented at 90° to the direction of incident polarization then blocks the reflected light, yielding a uniformly dark output image when there is no write light. If write light is applied to the device, a spatially varying AC electric field is established across the liquid crystal layer, and the long axis of the liquid crystal molecules begins to tilt away from the plane of the electrode. If the electric field were strong enough to fully rotate the molecules, then the birefringence of the material would vanish, the device would not change the direction of polarization, and again the reflected light would be completely blocked by the output analyzer. However, the fields are not sufficient to fully rotate the molecules, and hence they only partially tip away from the transverse plane, with an amount of tip that is proportional to the strength of the field (and therefore the strength of the write image). The partially tipped molecules retain some birefringent effect, and therefore the linearly polarized input light is transformed into elliptically polarized light, with a degree of ellipticity that depends on the strength of the applied field. The elliptically polarized field has a component that is parallel to the direction of the output analyzer, and therefore some of the reflected light passes that analyzer.

Contrast ratios of the order of 100 : 1 can be achieved with this device, and its resolution is several tens of line pairs per mm. The write time is of the order of 10 msec and the erase time about 15 msec. Due to the optically flat faceplate on the read side, the wavefront exiting the device is of good optical quality and the device is therefore suitable for use within a coherent optical data processing system. The nonmonotonic dependence of reflectance on applied voltage (both no voltage and a very high voltage result in the analyzer blocking all or most of the light) allows the device to be operated in several different linear and nonlinear modes, depending on that voltage.

Liquid crystal TV displays

The use of liquid crystal displays in small, light-weight portable televisions is widespread, and the technology of such displays has advanced rapidly in recent years. While displays of this type are not made for use in coherent light, nonetheless they can be adapted for use in a coherent optical system [132].

TV displays of this type are, of course, electrically addressed, and they display on the order of 100 to 200 pixels in both the horizontal and vertical dimensions. They are made from nematic liquid crystals, usually with either a 90° or a 270° twist. To use them in a coherent optical processing system, it is first necessary to remove polarizers attached to the display, and to remove any attached diffusing screen. In general the quality of the liquid crystal displays manufactured for projection TV are superior to those used in small TV sets.

Displays of this kind have not been manufactured with attention to their optical flatness, since the TV display application does not require it. As a consequence their optical quality is not outstanding, and they are useful mainly for rudimentary demonstrations, rather than as the basis for a system of very high performance. Their most important attribute is their extremely low cost, as compared with other SLM technologies.

Ferroelectric liquid crystal spatial light modulators

Ferroelectric liquid crystals provide the basis for several different approaches to the construction of spatial light modulators. SLMs based on these materials are inherently binary in nature, but gray scales can be created with the help of half-tone techniques. Both optically addressed and electrically addressed FLC SLMs have been reported. An excellent reference can be found in [91], Chapter 6.

Optically addressed FLC SLMs embody some of the same principles used in the Hughes liquid crystal light valve, but they also have some significant differences. Their general structure is similar to that of Fig. 7.18, but different materials are used and different conditions must be satisfied. Unlike NLC based devices, the FLC device must operate by reversal of the direction of the electric field across the liquid crystal layer. A different photoconductor, hydrogenated amorphous silicon, which has a faster response time than CdS, has been used. These devices are driven with audio-frequency square waves. The layer thicknesses (and therefore the capacitances in Fig. 7.19) are chosen so that the voltages appearing across the liquid crystal layer always remain sufficiently negative or sufficiently positive (depending on whether write light is or is not present) to drive the FLC material into its appropriate state. The tilt angles of the FLC molecules are again 45° apart and the FLC layer thickness is chosen for quarter-wave retardation, appropriate for a reflective modulator operating by polarization rotation.

Unlike the optically addressed SLMs, electrically addressed FLC SLMs are discrete pixelated devices, i.e. they display sampled images rather than continuous images. The FLC SLM is a pixelated version of the FLC intensity modulator described in a previous section. Particularly interesting SLMs can be realized when silicon is used as the substrate on which the pixelated cells are fabricated. Pixelated metallic electrodes can be deposited on the silicon surface, and will also serve as mirror surfaces. These electrodes can be matrix addressed on the chip. A variety of electronic devices can also be integrated in the silicon substrate. For example, drive electronics and various

electronic logic functions can be associated with each pixel of the device. Such an approach is often referred to as providing "smart pixels" [157]. FLC-on-silicon electrically driven light modulators with as many as 256×256 pixels have been reported [206]. This technology remains in a stage of rapid development.

7.2.3 Magneto-Optic Spatial Light Modulators

The SLMs discussed up to this point operate by means of the electrooptic effect, with polarization rotation being induced by a changing electric field across the device. We turn attention now to a different type of device, one that operates by means of polarization rotation under the application of a magnetic field, or the Faraday effect. For an alternative reference, see [91], Chapter 7.

SLMs of this type were developed by Litton [248] under the name "Light-Mod" and for some time were marketed by Semetex Corporation under the name "Sight-Mod". We shall use the abbreviation "MOSLM" here, standing for Magneto-Optic Spatial Light Modulator.

The MOSLM device consists of a two-dimensional array of magneto-optic elements in the form of individually isolated mesas on an epitaxially grown magnetic garnet film, mounted on a transparent nonmagnetic backing substrate. The garnet mesas are largely transparent to light, but when fully magnetized, they rotate the polarization of incident light as a consequence of the Faraday effect. The direction of rotation depends on the direction of magnetization of a mesa. When the magnetization direction coincides with the direction of propagation of the light, linearly polarized light will be rotated in a right-hand screw sense, by an angle $+\theta_f$ that depends on the thickness of the garnet film, and when the magnetization is opposite to the direction of propagation, the rotation of polarization is by angle $-\theta_f$. Thus, like the FLC SLMs, the MOSLM is a binary device with memory.⁵

The magnetization directions of the pixels are controlled by a combination of an external magnetic field, supplied by a bias coil, and a magnetic field introduced at the corner of each pixel by means of row and column metallic electrodes. Figure 7.21 illus-

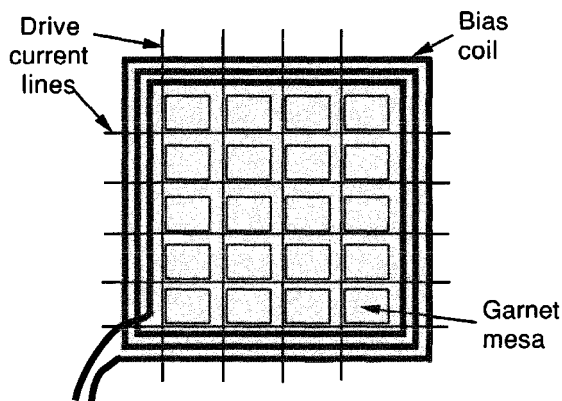


FIGURE 7.21
MOSLM device with bias coil and row-column address lines.

⁵There also exists a third, intermediate state of magnetization, in which the pixel consists of a multitude of small randomly oriented domains. This state is usually not used.

trates the geometry. Surrounding the pixelated garnet film is a bias coil which can be driven with current in either of two directions, thereby establishing a strong magnetic field in either of two directions, *i.e.* either parallel to the direction of light propagation, or anti-parallel to that direction. In addition, a row-column matrix of metallic conductors is deposited by photolithographic techniques such that a row electrode and a column electrode pass over one another at the corner of each pixel. The row and column electrodes are separated in the vertical direction by an insulating film.

To change the state of an individual pixel, the following sequence of operations must take place. First, the bias coil must be driven with current such that it establishes a strong magnetic field in the direction of the desired magnetization. Second, current pulses must be injected into the row and column electrodes that intersect at the pixel of interest, with the direction of those currents being such as to establish a small magnetic field which nucleates a magnetic domain with magnetization in the desired direction at the corner of the pixel. While all the pixels in the selected row and column experience some magnetic field from the current pulses, only where the two electrodes overlap is the magnetic field strong enough to nucleate a change of state of the pixel. With the arrival of the nucleating field, a change of state of magnetization is initiated at the corner of the pixel. The presence of the strong bias field causes this change to propagate at high speed across the entire pixel, thus changing the magnetization state of that mesa.

Pixels are written one at a time. Note that if two pixels must be changed to states of magnetization that are opposite from one another, two write cycles must be used, with a reversal of the current in the bias coil taking place between changes.

Quantitative analysis of the MOSLM device is aided by use of Jones matrices. It can be shown that the origin of Faraday rotation lies in different refractive indices experienced by the left-hand and right-hand circularly polarized components of a propagating wave (see [134], pp. 590-596). For the magnetic field oriented in one direction, the left-hand circularly polarized component experiences n_1 and the right-hand circularly polarized component experiences n_2 , whereas, when the direction of the magnetic field reverses, the refractive indices also reverse. From this fact it can be shown (see Prob. 7-5) that, aside from a phase factor that is common to both, the Jones matrices for the two directions of the magnetic field are simply rotation matrices,

$$\begin{aligned} \mathbf{L}_+ &= \begin{bmatrix} \cos \theta_f & -\sin \theta_f \\ \sin \theta_f & \cos \theta_f \end{bmatrix} \\ \mathbf{L}_- &= \begin{bmatrix} \cos \theta_f & \sin \theta_f \\ -\sin \theta_f & \cos \theta_f \end{bmatrix}, \end{aligned} \quad (7-22)$$

where for a film of thickness d , the rotation angle is given by

$$\theta_f = \frac{\pi(n_2 - n_1)d}{\lambda_o} \quad (7-23)$$

The Faraday rotation angle θ_f is in general quite small and therefore the total amount of polarization rotation that takes place between states of the device is much less than 90° . As a consequence, to use the device as an intensity modulator, the output polarization analyzer should be oriented orthogonal to the direction of polarization of

the light in one of the two states of rotation. One state of the device will then be entirely off, and the other will be partially on. When light polarized along the y axis is used for illumination and the analyzer is oriented at angle $+\theta_f$ to the x axis, the intensity transmission of the pixel in the "off" state is zero and in the "on" state can be shown to be

$$\tau = \eta_P e^{-\alpha d} \sin^2(2\beta d), \quad (7-24)$$

where η_P is the combined efficiency of the polarizer-analyzer combination, d is the film thickness (μm), α is the loss per μm of film thickness, and β is the rotation per μm of film thickness ($\beta d = \theta_f$). The Faraday rotation angle θ_f thus increases with the thickness of the garnet film, but at the same time the attenuation of the device, due to absorption, also increases with that thickness. Therefore for any given film there is an optimum thickness that maximizes the intensity transmission in the "on" state.

Typical parameters for spatial light modulators of this type operating at 633-nm wavelength are [80] [220]:

- Array sizes from 128×128 pixels to 256×256 pixels
- Faraday rotation parameter β as high as 1.46" per μm
- Absorption coefficient α of $0.086 \mu\text{m}^{-1}$
- Film thickness of 6 μm
- Optical efficiency in the "on" state of a few percent
- Frame rate of approximately 1 kHz.

This technology is relatively mature, but improvements, including the construction of high-performance reflective devices, are still taking place [247].

7.2.4 Deformable Mirror Spatial Light Modulators

A variety of devices have been reported that use electrostatically induced mechanical deformation to modulate a reflected light wave. Such devices are usually referred to as "deformable mirror devices", or **DMDs**. The most advanced SLMs of this type have been developed by Texas Instruments, Inc. Early devices utilized continuous membranes which deformed under the fields exerted by pixelated driving electrodes. These SLMs gradually evolved into deformable mirror devices, in which discrete cantilevered mirrors were individually addressed via voltages set on floating MOS (metal oxide semiconductor) sources, the entire device being integrated on silicon. The most recent versions have used mirror elements with two points of support, which twist under the application of an applied field. An excellent discussion of all of these approaches is found in [147].

Figure 7.22 shows the structures for a membrane device and for a cantilever beam device. For the membrane device, a metallized polymer membrane is stretched over a spacer grid, the spacers forming an air gap between the membrane and the underlying address electrodes. A negative bias voltage is applied to the metallized membrane. When a positive voltage is applied to the address electrode under the membrane, it deflects downward under the influence of the electrostatic forces. When the address voltage is removed, the membrane moves upward to its original position. In this way a

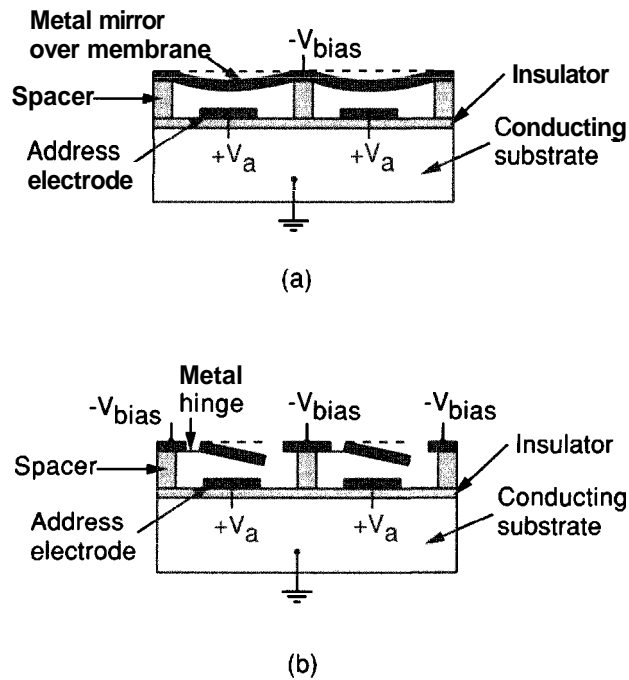


FIGURE 7.2
Deformable mirror pixel structures for (a) a membrane SLM and (b) a cantilever beam SLM.

phase modulation is introduced, but that phase modulation can be converted to an intensity modulation by appropriate optics following the mirror (cf. Probs. 8-2 and 8-3).

For the cantilever beam device, the structure is quite different. The metallized beam, which is biased to a negative voltage, is attached to a spacer post through a thin metal hinge. When the underlying address electrode is activated with a positive voltage, the cantilever rotates downward, although not far enough to touch the address electrode. An incident optical beam is thus deflected by the tilted pixel, and will not be collected by an optical system that follows. In this way an intensity modulation is induced at each pixel.

The most advanced DMD structures are based on a geometry related to that of the cantilever beam, but instead use a torsion beam which is connected at two points rather than through a single metal hinge. Figure 7.23 shows a top view of the metallized pixel. As shown in part (a) of the figure, the torsion rod connects the mirror to supports at the ends of one diagonal. Again the mirror is metallized and connected to a negative bias voltage.

As shown in part (b) of the figure, two address electrodes exist for each such pixel, one on either side of the rotation axis. When one address electrode is activated with a positive voltage, the mirror twists in one direction, and when the other electrode is **activated**, the mirror twists in the opposite direction. Under each mirror element are two landing electrodes, held at the bias voltage, so that when the mirror tip twists so far as to hit the **underlying** landing electrode, there is no electrical discharge. The light incident on each **pixel** is deflected in either of two directions by the mirror when it is activated, and is **not** deflected when it is not activated. The device can be operated in either an analog mode, in which twist is a continuous function of applied address voltage, or in a digital mode, in which the device has either two stable states or three stable states, depending on the bias voltage applied [147].

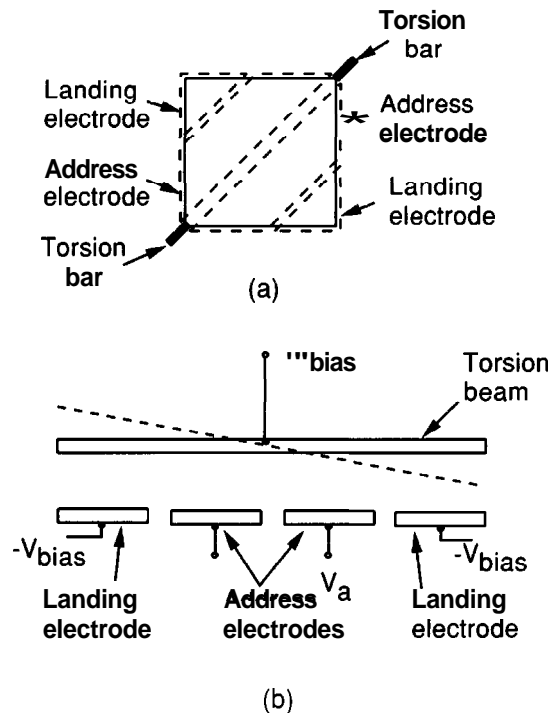


FIGURE 7.23
Torsion beam DMD: (a) Top view, and
(b) side view.

A major advantage of this type of SLM technology is that it is silicon-based and compatible with the use of CMOS (complementary metal oxide silicon) drivers on the same substrate used for the SLM pixels. Both line-addressed DMDs and frame-addressed DMDs have been reported in sizes 128×128 and above. Devices of this type with as many as 1152×2048 pixels have been reported for use as high-definition TV (HDTV) displays. A second advantage is the ability of the device to operate at any optical wavelength where good mirrors can be fabricated in integrated form.

Measurements of the electrical and optical properties of this type of DMD have been reported in the literature [278]. Maximum deflection angles approaching 10° are measured with applied voltages of about 16 volts. Deflection times of about $28 \mu\text{sec}$ were measured for an individual pixel, but this number depends on pixel size and can be shorter for smaller pixels. The resonant frequency of a pixel was found to be of the order of 10 kHz.

7.2.5 Multiple Quantum Well Spatial Light Modulators

The use of molecular beam epitaxy to fabricate sophisticated electronic and optoelectronic devices consisting of large numbers of extremely thin layers of different semiconductor materials has led to new approaches to the construction of spatial light modulators. A typical material system for such structures would involve alternating layers of GaAs and AlGaAs with thicknesses of the order of 10 nm. The small thickness of these layers, known as quantum wells, results in certain quantum-mechanical effects, in particular new absorption peaks associated with structures known as excitons. Such structures consist of an electron-hole pair for which the electron and hole are normally separated by a distance that is larger than the thickness of a single layer, but which are

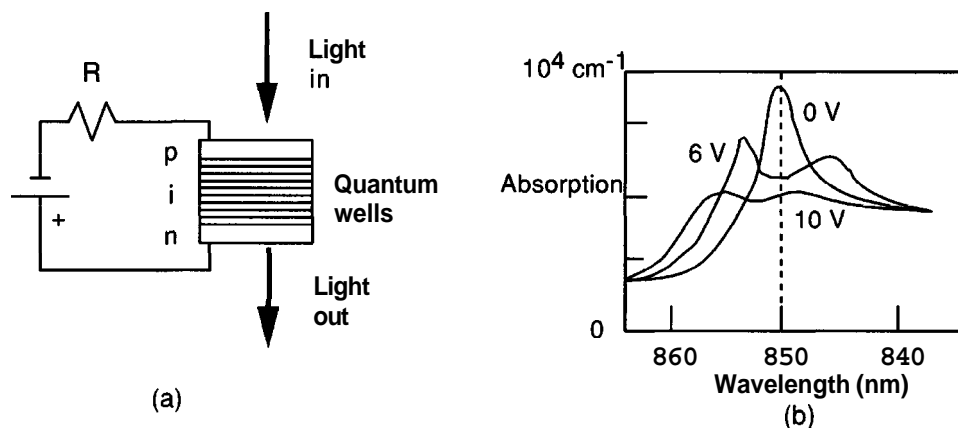


FIGURE 7.24 Shift of the exciton resonance in a multiple quantum well p-i-n diode. (a) Circuit schematic, and (b) absorption spectra at various applied voltages.

brought closer together by the constraints imposed by the thin layers of the device. The structure introduces constraints on the movement of the excitons (known as "quantum confinement") and the emergence of sharp absorption peaks that would not normally be observed in structures with thicker layers. Structures of this type are usually referred to as *multiple quantum well (MQW)* structures. For a survey of the use of such structures in spatial light modulators, see [91], Chapter 5.

A further consequence of quantum confinement is a dependence of the spectral location of the absorption peak on an electric field applied normal to the quantum wells in the structure, which is explained by a mechanism known as the quantum confined Stark effect (QCSE) [214]. When an electric field is applied across the structure, the exciton resonances move to lower photon energies, or to longer wavelengths. As a consequence, if an MQW device is illuminated by light with a properly chosen wavelength, the application of an applied field to the device can change the absorption experienced by that light as it passes through the structure, an effect that can serve as the basis for realizing an optical intensity modulator⁶ [303]. Any method for modulating the intensity of light is also a candidate technology for the construction of a spatial light modulator. Figure 7.24 illustrates typical absorption curves for a quantum well p-i-n diode [212]. Note the shift of the absorption peak as the applied voltage increases, as well as the gradual reduction in the strength of that peak. Note also that a source having wavelength 850 nm, as indicated by the vertical dashed line, will experience decreasing absorption as the voltage is increased.

Modulators of the type above can be fabricated in modest size arrays and addressed electrically to produce a discrete, pixelated SLM. The size of an individual active pixels can be in the range of **10 μm** to **100 μm** on a side. Modest contrast ratios can be achieved (e.g. **3:1** on to off). The speeds of the modulators can be quite fast, depending on the size of the pixels and their related capacitance, with modulation bandwidths

⁶The electroabsorption effect (i.e. the change of absorption with an applied electric field) in MQW devices is approximately 50 times greater than the same effect in bulk GaAs.

of tens of GHz having been demonstrated in waveguide devices. The inclusion of a dielectric mirror stack on the back of the device, made during the same molecular beam epitaxy (MBE) process used for the rest of the device, leads to a double pass of light through the device and improvement of on/off contrast. Devices of this kind can also be made within a Fabry-Perot Ctalon structure, yielding even better contrast at the price of narrower optical bandwidth. Other approaches to MQW SLM construction are also possible (cf. [12], [92], and [30]).

The self-electro-optic effect device

The MQW modulator arrays discussed above are electrically addressed. It is also possible to utilize the optical absorption and current generation associated with back-biased p-i-n structures to create pixelated arrays of devices that can have their states changed by the application of optical signals. The most common type of such a device is known as the self-electro-optic effect device (the SEED). As will be seen, such devices typically exhibit bistability.

The simplest SEED structure is actually that shown in Fig. 7.24(a), which is known as the resistor-biased SEED, or R-SEED. In this case the diode structure is used simultaneously as a detector and a modulator. Suppose that initially there is no light incident on the diode, and as a result there is no current flowing in the circuit. In such a case, all of the applied voltage must appear across the MQW device, and none across the resistor. As Fig. 7.24(b) indicates, when the full voltage exists across the device, the absorption is low but still significant. If light at a wavelength indicated by the dashed line is now incident on the device, some of it will be absorbed, and the back-biased diode, acting as a photodetector, will generate current. That current will cause a portion of the applied voltage to appear across the resistor, and less of the voltage to fall across the diode. Referring again to Fig. 7.24(b), less voltage across the diode results in higher absorption by the device, more current generation, and an even lower voltage across the diode. This positive feedback action results in the MQW modulator switching to its highest possible absorption under the application of light to the device (i.e. its "absorptive" state). If the incident optical power is now decreased, the current will decrease, the voltage across the diode increases, the absorption drops, and again a positive feedback mechanism switches the device back to its "transparent" initial state.

The action of the R-SEED described above does not lead to a very useful device in itself. Each pixel that is written with light becomes maximally absorbing, thereby transmitting little or no light, and each pixel that is not written by light remains maximally transparent. More complex structures are needed to produce a useful SLM.

A wide variety of different SEED structures have been conceived of. We will explain the operation of one more of these structures, the symmetric SEED, or S-SEED, even though it is primarily of interest for digital logic rather than analog processing. A typical transmissive S-SEED structure is shown in Fig. 7.25. In this case, two MQW diodes are integrated into a single pixel and electrically interconnected with one another, as shown. The device pair operates with a complementary set of inputs and produces a complementary set of outputs. It is the ratio of the intensities in the two complementary light beams that carries the information of interest.

Suppose that initially, in the absence of light, the voltage V is equally divided across the two diodes. Imagine now that a complementary pair of beams is applied to the inputs

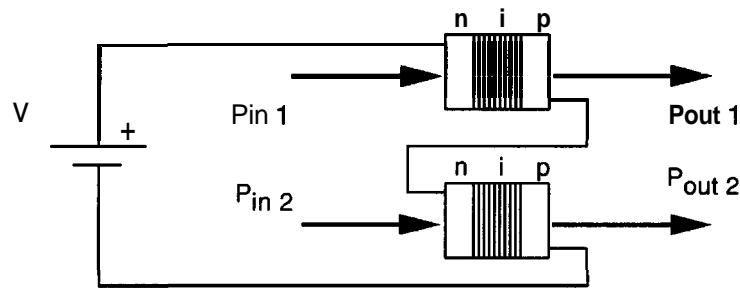


FIGURE 7.25
Symmetric SEED or S-SEED.

of the two diodes, with bright on top and dark on the bottom. The top diode will absorb light, generate current, and the voltage will drop across that diode and rise across the lower diode, which was not illuminated. A rise in voltage across the lower diode makes it less absorptive, which means that its voltage will rise further, while at the same time a drop in voltage across the upper diode will make it more absorptive, generate more current, and so forth. Thus the diode pair switches into a stable state in which the top diode is maximally absorptive and the bottom diode is maximally transmissive.

If the beams applied to the diode pair had been the complement of that assumed, i.e. had been dark on top and bright on the bottom, the device pair would have gone to a complementary state, namely maximally absorptive on the bottom and maximally transmissive on the top. Thus the two stable states of the device are represented by the two possible combinations of transmissions for the top and bottom diodes.

Once the state of the device has been set by the applied optical beams, it is possible to read out that state nondestructively and pass it on to a subsequent S-SEED pair. Let the diodes in the pair be illuminated by a pair of equally bright read-out spots. Since the illuminations on the two diodes are equal, there is no imbalance that would cause the diode pair to change its state. Thus the pair of beams transmitted by the diode pair will carry with it the state of the device, but with a brightness inversion compared with the pair of beams that set the state of the diode pair.

In practice, S-SEED arrays are made with reflective devices such that the beams traverse a given device twice, with the result that the contrast between the two beams is increased. Arrays with as many as 512×256 S-SEED pairs have been made by AT&T.

Research continues at a rapid pace on this technology. One new device type is the FET-SEED, in which field effect transistor (FET) devices are integrated with the diode detector/modulators, allowing logic of considerable complexity to be performed electrically on the chip at each pixel [192], another example of "smart pixel" technology. An additional development of considerable interest is the report of arrays of SEED devices that are not bistable and can operate as analog modulators, suitable for continuous gray-scale SLMs [213].

7.2.6 Acousto-Optic Spatial Light Modulators

The SLMs considered in the above sections are capable of modulating a two-dimensional wavefront, either in a continuous fashion or with a discrete two-dimensional

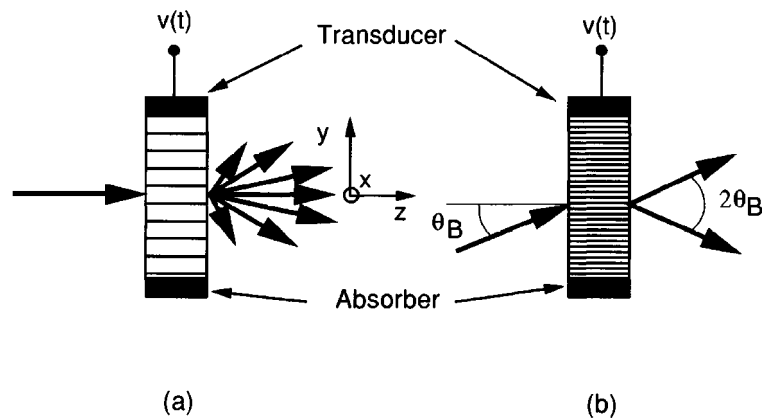


FIGURE 7.26
Acousto-optic cells operating in the (a) Raman-Nath regime
and the (b) Bragg regime.

array of elements. We turn now to an SLM technology that is most commonly one-dimensional, but which has been developed over a period of many years into a highly mature technology. This approach to wavefront modulation uses the interaction of a column of traveling acoustic waves with an incident coherent optical beam to modulate the properties of the transmitted optical wavefront. For alternative references that treat acousto-optic interactions and their applications in coherent optical systems, see, for example, [171], [21], and [293].

Figure 7.26 illustrates two versions of acousto-optic SLMs, each operating in a different physical regime. In both cases, the acousto-optic cell consists of a transparent medium (e.g. a liquid or a transparent crystal) into which acoustic waves can be launched by a piezoelectric transducer. The transducer is driven by an RF voltage source and launches a compressional wave (or, in some cases, a shear wave) into the acoustic medium. The acoustic wave propagates in the medium through small local displacements of molecules (strain). Associated with these strains are small changes of the local refractive index, a phenomenon known as the acousto-optic or the photo-elastic effect. The driving voltage has an RF spectrum that is centered at some center frequency f_c with a bandwidth B about that center frequency.

A CW drive voltage

For a perfectly sinusoidal drive voltage of frequency f_c (i.e. a CW voltage), the transducer launches a sinusoidal traveling acoustic wave in the cell, which moves with the acoustic velocity V characteristic of the medium. This traveling wave induces a moving sinusoidal phase grating with period $\Lambda = V/f_c$, and interacts with the incident optical wavefront to produce various diffraction orders (cf. Section 4.4). However, there are two different regimes within which the acousto-optic interaction exhibits different properties, the Raman-Nath regime and the Bragg regime.

In the Raman-Nath regime, which is typically encountered for center frequencies in the range of several tens of MHz in cells that use liquid as the acoustic medium, the moving grating acts as a thin phase sinusoidal grating exactly as described in the example of Section 4.4, with the one exception that, as a consequence of the grating mo-

tion through the cell, the various diffraction orders emerge from the cell with different optical frequencies. If the cell is illuminated normal to the direction of acoustic wave propagation, as shown in Fig. 7.26(a), the zero-order component remains centered at frequency ν_o of the incident light, but higher-order components suffer frequency translations, which can be interpreted as Doppler shifts due to the motion of the grating. Since the period of the grating is Λ , the q th diffraction order leaves the cell with angle θ_q with respect to the incident wave, where

$$\sin \theta_q = q \frac{\lambda}{\Lambda}, \quad (7-25)$$

λ being the optical wavelength within the acousto-optic medium. The optical frequency of the q th diffraction order can be determined from the Doppler-shift relation

$$\nu_q = \nu_o \left(1 + \frac{V}{c} \right) \sin \theta_q \approx \nu_o + q f_c. \quad (7-26)$$

Thus the optical frequency of the q th diffraction order is translated by q times the RF frequency, where q can be a positive or a negative number. q is a positive integer for diffraction orders with components of direction parallel to the direction of motion of the acoustic wave (i.e. downwards in Fig. 7.26), and negative for diffraction orders with components of direction opposite to that of the motion of the acoustic wave (i.e. upwards in Fig. 7.26). As for any thin sinusoidal phase grating, the intensities associated with the various diffraction orders are proportional to the squares of the Bessel functions of the first kind, $J_q^2(\Delta\phi)$, where $\Delta\phi$ is the peak-to-peak phase modulation, as shown in Fig. 4.13.

For RF frequencies in the hundreds of MHz to the GHz range, and in acoustic media consisting of crystals, the thickness of the acousto-optic column compared with the acoustic wavelength introduces a preferential weighting for certain diffraction orders, and suppresses others. This effect is known as the *Bragg effect* and will be discussed at greater length in Chapter 9. For the moment it suffices to point out that in this regime the dominant diffraction orders are the zero order and a single first order. Strong diffraction into a first diffraction order occurs only when the angle of the incident beam, with respect to plane of the acoustic wavefronts, has the particular value θ_B satisfying

$$\sin \theta_B = \pm \frac{\lambda}{2\Lambda} \quad (7-27)$$

(cf. Fig. 7.26(b)), where again λ is the optical wavelength within the acoustic medium. An angle satisfying the above relation is known as a *Bragg angle*. Equivalently, if \vec{k}_i is the wave vector of the incident optical wave ($|\vec{k}_i| = 2\pi/\lambda$) and \vec{K} is the wave vector of the acoustic wave ($|\vec{K}| = 2\pi/\Lambda$), then

$$\sin \theta_B = \pm \frac{|\vec{K}|}{2|\vec{k}_i|}. \quad (7-28)$$

The frequency of the first-order diffracted component is $\nu_o + f_c$ for the geometry shown in Fig. 7.26(b). The strength of the first-order component can be far greater than in the Raman-Nath regime, as discussed in more detail in Chapter 9.

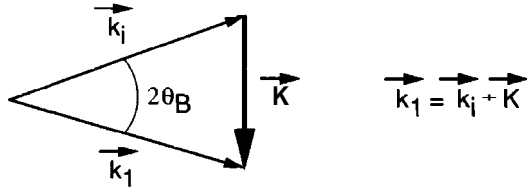


FIGURE 7.27

Wave vector diagram for Bragg interaction. \vec{k}_i is the incident optical wave vector, \vec{k}_1 is the optical wave vector of the component diffracted into the first diffraction order, and \vec{K} is the acoustical wave vector.

An aid for visualizing the relations between the optical and acoustical wave vectors is a wave vector diagram, as shown in Fig. 7.27. For strong Bragg diffraction, the wave vector diagram must close as shown, a property that can be viewed as a statement of conservation of momentum.

The boundary between the Raman-Nath regime and the Bragg regime is not a sharp one, but is often described in terms of the so-called Q factor given by

$$Q = \frac{2\pi\lambda_0 d}{n\Lambda^2} \quad (7-29)$$

where d is the thickness of the acoustic column in the z direction, n is the refractive index of the acousto-optic cell, and Λ is the vacuum wavelength of the light. If $Q < 2\pi$, operation is in the Raman-Nath regime, while if $Q > 2\pi$, operation is in the Bragg regime.

A modulated drive voltage

Until now the voltage driving the acousto-optic cell has been assumed to be a perfect CW signal. We now generalize by allowing the voltage to be an amplitude and phase-modulated CW signal, of the form

$$v(t) = A(t) \sin[2\pi f_c t - \psi(t)], \quad (7-30)$$

where $A(t)$ and $\psi(t)$ are the amplitude and phase modulations, respectively. The refractive index disturbance generated by this applied voltage then propagates through the cell with velocity V . With reference to Fig. 7.26, if y is a coordinate running opposite to the direction of travel of the acoustic wave and is centered in the middle of the cell (as indicated in Fig. 7.26), and x is normal to the page of that figure, then at any instant of time t the distribution of refractive index perturbation in the cell can be written

$$\Delta n(y; t) = \sigma v\left(\frac{y}{V} + t - \tau_o\right), \quad (7-31)$$

where σ is a proportionality constant, $\tau_o = L/2V$ is the time delay required for acoustic propagation over half the length L of the cell, and we neglect the x dependence because it plays no role here or in what follows.

In the Raman-Nath regime, the optical wavefront is simply phase modulated by the moving refractive index grating, yielding a complex amplitude of the transmitted signal given by

$$U(y; t) = U_o \exp\left\{j \frac{2\pi\sigma d}{\lambda_o} A\left(\frac{y}{V} + t - \tau_o\right) \sin\left[2\pi f_c \left(\frac{y}{V} + t - \tau_o\right) - \psi\left(\frac{y}{V} + t - \tau_o\right)\right]\right\} \text{rect} \frac{y}{L}, \quad (7-32)$$

where U_o is the complex amplitude of the incident monochromatic optical wave. Now the expansion

$$\exp[j\phi \sin \beta] = \sum_{q=-\infty}^{\infty} J_q(\phi) \exp(jq\beta) \quad (7-33)$$

can be applied to the expression for $U(y; t)$. In addition, the peak phase modulation suffered by the optical wave as it passes through the cell is usually quite small, with the result that the approximation

$$J_{\pm 1}(\phi) \approx \pm \phi/2$$

holds for the first diffraction orders, which are the orders of main interest to us. As a result the complex amplitudes transmitted into the two first orders, represented by $U_{\pm 1}$, are given approximately by

$$U_{\pm 1} \approx \pm \frac{\pi \sigma d}{\lambda_o} U_o A \left(\frac{y}{V} \pm t - \tau_o \right) e^{\mp j\psi(y/V \pm t - \tau_o)} e^{\pm j2\pi y/\Lambda} e^{\pm j2\pi f_c(t - \tau_o)} \text{rect} \frac{y}{L}, \quad (7-34)$$

where the top sign corresponds to what we will call the “+1” diffraction order (diffracted downwards in Fig. 7.26) and the bottom sign corresponds to the “-1” order (diffracted upwards).

From Eq. (7-34) we see that the **+1** diffracted order consists of a wavefront that is proportional to a moving version of the complex representation $A(y/V)e^{j\psi(y/V)}$ of applied voltage, while the **-1** diffracted order contains the complex conjugate of this representation. The spatial argument of the moving field is scaled by the acoustic velocity V . A simple spatial filtering operation (see Chapter 8) can eliminate the unwanted diffraction orders and pass only the desired order. Thus the acousto-optic cell has acted as a one-dimensional spatial light modulator, transforming an electrical voltage modulation applied to the cell into an optical wavefront exiting the cell.

The discussion above has been framed in terms of Raman-Nath diffraction, but a similar expression for the **+1** diffraction order is found in the case of Bragg diffraction, the primary difference lying in the strengths of the various orders. As mentioned earlier, the diffraction efficiency into one first order is generally considerably larger in the Bragg regime than in the Raman-Nath regime, and other orders are generally strongly suppressed by the diffraction process itself. Thus an acousto-optic cell operating in the Bragg regime again acts as a one-dimensional spatial light modulator, translating the applied voltage modulation into a spatial wavefront, albeit more efficiently than in the case of Raman-Nath diffraction.

7.3 DIFFRACTIVE OPTICAL ELEMENTS

The vast majority of optical instruments in use today use *refractive* or *reflective* optical elements (e.g. lenses, mirrors, prisms, etc.) for controlling the distribution of light. In some cases it is possible to replace refractive or reflective elements with *diffractive* elements, a change that can lead to some significant benefits in certain applications.

Diffractive optics can be made to perform functions that would be difficult or impossible to achieve with more conventional optics (e.g., a single diffractive optical element can have several or many different focal points simultaneously). Diffractive optical elements also generally have much less weight and occupy less volume than their refractive or reflective counterparts. They may also be less expensive to manufacture and in some cases may have superior optical performance (e.g. a wider field of view). Examples of applications of such components include optical heads for compact disks, beam shaping for lasers, grating beamsplitters, and reference elements in **interferometric** testing.

Along with these several advantages comes one significant difficulty with **diffractive** optical components: because they are based on diffraction, they are highly dispersive (i.e. wavelength sensitive). For this reason they are best applied in problems for which the light is highly monochromatic. Such is the case for most coherent optical systems. However, diffractive optics can be used together with either refractive optics or additional diffractive elements in such a way that their dispersive properties partially cancel (cf. [274], [217], [99]), allowing their use in systems for which the light is not highly monochromatic.

For additional background on diffractive optics, the reader may wish to consult review articles [279],[98], and Vol. 2, Chapter 8 of [17].

Our discussion of the subject will consider in detail only one approach to the construction of diffractive optics, known as binary optics. This approach is well developed and applicable to a broad range of different applications.

7.3.1 Binary Optics

The term **binary** optics has come to have different meanings to different people, but there are certain threads that are common and which can serve to define the field. First and foremost is the fact that binary optical elements are manufactured using VLSI fabrication techniques, namely photolithography and micromachining. Second, binary optical elements depend solely on the surface relief profile of the optical element. They are usually thin structures, with relief patterns on the order of one to several microns in depth, and as such they can be inexpensively replicated using well-established methods of embossing. Surprisingly, the relief patterns utilized are often not binary at all, and therefore in a certain sense these elements are misnamed. However, such elements are usually defined through a series of binary exposure steps, and this fact has provided the rationale for retention of the name.

Approximation by a stepped thickness function

Binary optical elements have stepped approximations to ideal continuous phase distributions. We briefly discuss the approximation process here, and then turn to the most common methods of fabrication.

We suppose that a certain thickness function $\Delta(x, y)$ is desired for the element (as usual, x and y are the transverse coordinates on the face of the element). Presumably this function has been derived from a design process, which may have been simple or may have been quite complex itself. As an example of a simple case, the element may be a

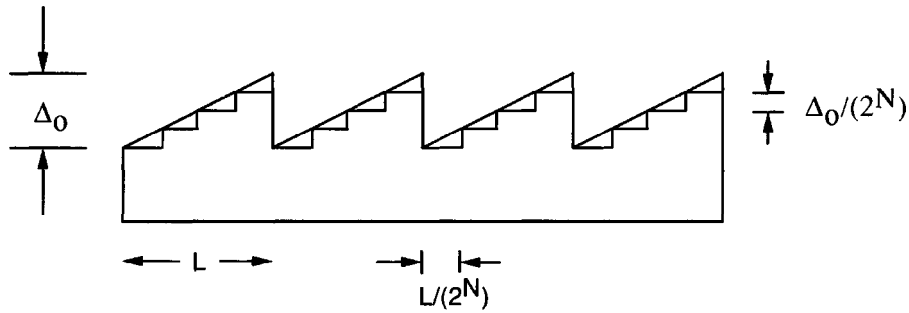


FIGURE 7.28
Ideal sawtooth thickness profile for a blazed grating, and binary optic approximation to that profile ($N = 2$).

grating of constant spatial frequency, the purpose of which is to deflect the incident light through a certain angle with the highest possible optical efficiency. An example of a more complex case might be a focusing element which generates an aspheric wavefront such that certain aberrations are reduced or eliminated. We shall assume that the desired thickness function $\Delta(x, y)$ is known and that the problem at hand is how to fabricate a thin relief element that closely approximates this desired thickness function.

An approximation to the desired thickness function is made by quantizing that function to a set of 2^N discrete levels (usually equally spaced). Figure 7.28 shows an ideal phase grating profile with a perfect sawtooth period, and a quantized version of that grating with 2^N levels. The continuous blazed grating has the property that, if the peak-to-peak phase variation it introduces is exactly 2π radians, 100% of the incident light will be diffracted into a single first diffraction order (cf. Prob. 4-15). The binary optic approximation to the grating is a quantized version with 4 discrete levels. More generally 2^N quantization levels can be realized through a series of N exposure and micromachining operations, as described below. The peak-to-peak thickness change of the quantized element is $\frac{2^N-1}{2^N}$ times the peak-to-peak thickness of the unquantized element.⁷

The diffraction efficiency of the step approximation to the sawtooth grating can be obtained by expanding its periodic amplitude transmittance in a Fourier series. A straightforward but tedious calculation shows that the diffraction efficiency of the q th diffraction order can be expressed by [79]

$$\eta_q = \text{sinc}^2\left(\frac{q}{2^N}\right) \frac{\text{sinc}^2\left(q - \frac{\phi_o}{2\pi}\right)}{\text{sinc}^2\left(\frac{q - \frac{\phi_o}{2\pi}}{2^N}\right)}, \tag{7-35}$$

where ϕ_o is the peak-to-peak phase difference of the continuous sawtooth grating, and is related to the peak-to-peak thickness variation (again, of the continuous grating) through

⁷These ideal and quantized gratings may be considered to be local approximations to more general gratings for which the local period, and therefore the angle of deflection, change across the grating.

$$\phi_o = 2\pi \frac{\Delta_o(n_2 - n_1)}{\lambda_o}, \tag{7-36}$$

n_2 being the refractive index of the substrate and n_1 that of the surround, and λ_o being the vacuum wavelength of the light.

Of special interest is the case of a quantized approximation to the blazed grating with a peak-to-peak phase difference of $\phi_o = 2\pi$. Substitution in Eq. (7-35) yields

$$\eta_q = \text{sinc}^2\left(\frac{q}{2^N}\right) \frac{\text{sinc}^2(q - 1)}{\text{sinc}^2\left(\frac{q-1}{2^N}\right)}. \tag{7-37}$$

Consider for the moment only the last factor, consisting of the ratio of two sinc functions. The numerator is zero for all integer q except $q = 1$, when it is unity. The denominator is also unity for $q = 1$, and is nonzero except when

$$q - 1 = p2^N,$$

where p is any integer other than zero. For values of q for which the numerator and denominator vanish simultaneously, l'Hôpital's rule can be used to show that the ratio of the two factors is unity. Thus the factor in question will be zero except when

$$q = p2^N + 1,$$

in which case it is unity. The diffraction efficiency therefore is given by

$$\eta_{(p2^N+1)} = \text{sinc}^2\left(p + \frac{1}{2^N}\right). \tag{7-38}$$

As the number, 2^N , of phase levels used increases, the angular separation between nonzero diffraction orders increases as well, since it is proportional to 2^N . The primary order of interest is the $+1$ order ($p = 0$), for which the diffraction efficiency is

$$\eta_1 = \text{sinc}^2\left(\frac{1}{2^N}\right). \tag{7-39}$$

Figure 7.29 shows the diffraction efficiencies of various nonzero orders as a function of the number of levels. It can be seen that, as $N \rightarrow \infty$, all diffraction orders except the

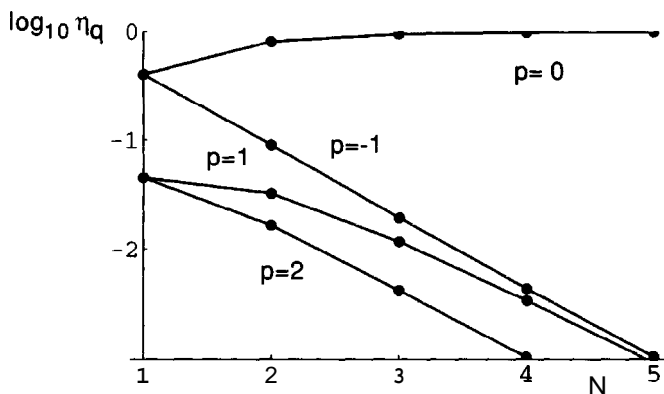


FIGURE 7.29 Diffraction efficiencies of various orders of a stepped approximation to a sawtooth grating. The parameter p determines the particular diffraction order, with the order number given by $p2^N + 1$, and the number of discrete levels is 2^N .

$+1$ order vanish, and the diffraction efficiency of that nonvanishing order approaches 100%, identical with the case of a continuous blazed grating with the same peak-to-peak phase shift. Thus the properties of stepped approximation to the continuous blazed grating do indeed approach those of the continuous grating as the number of steps increases.

The fabrication process

Figure 7.30 illustrates the process by which a four-level binary optic approximation to a sawtooth thickness function is generated. The process consists of a number of discrete steps, each of which consists of photoresist application, exposure through one of several binary masks, photoresist removal, and etching. Masks are usually made by electron-beam writing. For a binary optic element with 2^N levels, N separate masks are required. **Part (a)** of the figure shows a substrate overcoated with photoresist, which is exposed through the first binary mask, having transparent cells of width equal to $1/2^N$ th of the period of the desired final structure. After exposure, the photoresist is developed. For a positive photoresist, the development process removes the exposed areas and leaves the unexposed areas, while for a negative photoresist the opposite is true. We will assume a positive photoresist here. Following the photoresist development process, micromachining is applied to remove material from the uncovered portions of the substrate, as illustrated in part (b) of the figure. The two most common micromachining methods are reactive ion etching and ion milling. This first micromachining step removes substrate material to a depth of $1/2^N$ th of the desired peak-to-peak depth of the grating. Now photoresist is spun onto the substrate a second time and is exposed through a second mask which has openings of width equal to $1/2^{N-1}$ th of the desired

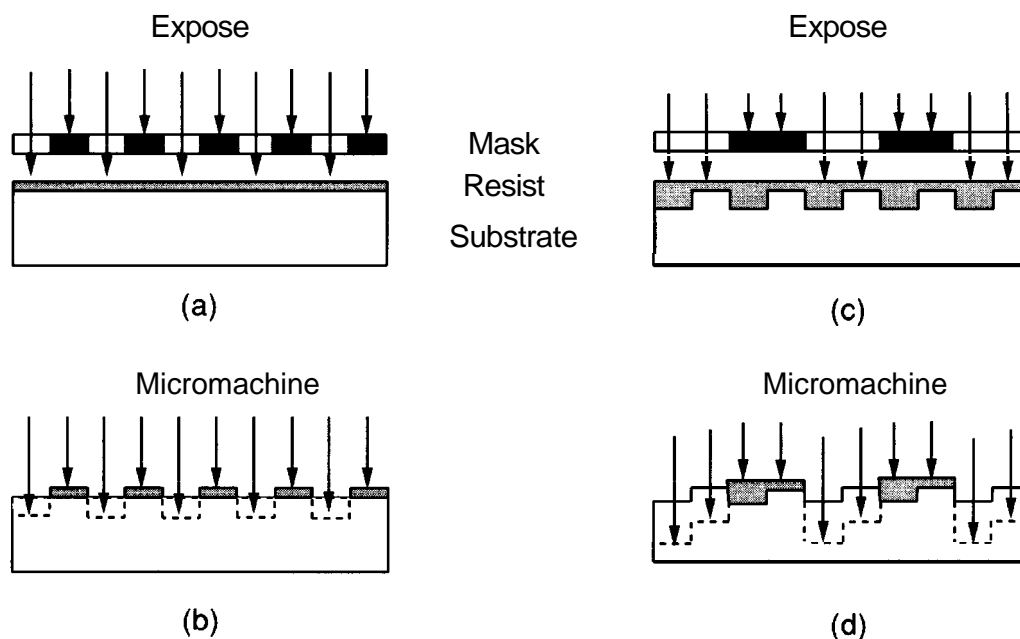


FIGURE 7.30 Steps in the fabrication of a four-level binary optic element.

final period, as shown in part (c) of the figure. Micromachining again removes the exposed portions of the substrate, this time with an etch depth $1/2^{N-1}$ th of the final desired maximum depth, as illustrated in part (d) of the figure. For the case of a 4-level element, the fabrication process now terminates. If there are 2^N levels desired, N different masks, exposures, development, and etching processes are required. The last etch process must be to a depth that is $1/2$ of the total desired peak-to-peak depth. A variety of different materials can be used for the substrate of such elements including silicon and glass. It is also possible to make reflective optical devices by overcoating the etched profile with a thin layer of metal. With the use of electron beam writing, it is possible to control the accuracy of the masks to about one-tenth of a μm . When the profile is more complex than binary, alignment of several masks is required, and the accuracy is reduced.

Diffraction efficiencies of 80 to 90 percent are quite common for these types of elements.

7.3.2 Other Types of Diffractive Optics

Attention has been focused above on binary diffractive optics, which are fabricated by the techniques widely used in the semiconductor industry. Many other approaches to fabricating diffractive optical elements exist. Some methods use similar substrates to those mentioned above, but use different methods of micromachining, for example diamond turning or laser ablation. Some differ through their use of photographic film, rather than **etchable** substrates, as the means for creating the element. **Computer-generated** holographic optical elements are an example that will be discussed in more detail in Chapter 9. Some depend on more conventional methods for recording holograms.

For an overview of the field, including examples of many different approaches, the reader is referred to the proceedings of a series of meetings held on this general subject [59], [60], [61], [62], [63].

7.3.3 A Word of Caution

The capability of semiconductor fabrication techniques to make structures of ever smaller physical size has led already to the construction of diffractive optical elements with individual feature sizes that are comparable with and even smaller than the size of a wavelength of the light with which the element will be used. Such small structures lie in the domain where the use of a scalar theory to predict the properties of these optical elements is known to yield results with significant inaccuracies. It is therefore important to use some caution when approaching the analysis of the properties of diffractive optical elements. If the minimum scale size in the optical element is smaller than a few optical wavelengths, then a more rigorous approach to diffraction calculations will probably be needed, depending on the accuracy desired from the computation. For a discussion of such issues, see, for example, Ref. [232].

PROBLEMS-CHAPTER 7

7-1. A low-contrast intensity distribution

$$\mathcal{I}(x, y) = \mathcal{I}_0 + \Delta\mathcal{I}(x, y) \quad |\Delta\mathcal{I}| \ll \mathcal{I}_0$$

exposes a photographic plate, and a negative transparency is made. Assuming that \mathcal{I}_0 is fixed and biases the film in the linear region of the H&D curve, show that, when the contrast $\Delta\mathcal{I}/\mathcal{I}_0$ is sufficiently low, the contrast distribution transmitted by the transparency is linearly related to the exposing contrast distribution.

7-2. The interference between two plane waves of the form

$$U_1(x, y) = A \exp(j2\pi\beta_1 y)$$

$$U_2(x, y) = B \exp(j2\pi\beta_2 y)$$

is recorded on a photographic film. The film has an MTF of known form $M(f)$, and it is processed to produce a **positive** transparency with a gamma of -2 . This transparency (dimensions $L \times L$) is then placed in front of a positive lens with focal length f , is illuminated by a normally incident plane wave, and the distribution of intensity across the back focal plane is measured. The wavelength of the light is λ . Assuming that the entire range of exposure experiences the same photographic gamma, plot the distribution of light intensity in the rear focal plane, labeling in particular the relative strengths and locations of the various frequency components present.

7-3. Show that, for a retardation of $\beta d = \pi$, the Jones matrices of Eqs. (7-18) and (7-19) reduce to those of Eq. (7-20).

7-4. A ferroelectric liquid crystal cell has a tilt angle of 22.5° . The input of the cell has a polarizer oriented parallel to the long molecular axis when the cell is in one of its two states, and the rear of the cell is a mirror. Using Jones matrices, show that if the retardation of the cell is one-quarter of a wave, the FLC cell can be used as a binary intensity modulator.

7-5. Consider a linear polarized wave with polarization direction at angle $+\theta$ to the x axis.

(a) Show that such a wave can be expressed as a linear combination of a left-hand circularly polarized wave and a right-hand circularly polarized wave, and find the coefficients of that expansion.

(b) Given that, for Faraday rotation, with the magnetic field pointing in the direction of wave propagation, the left-hand circularly polarized component experiences a refractive index n_1 and the right-hand circularly polarized component experiences refractive index n_2 , show that the Jones matrix describing this polarization transformation is given by

$$L_+ = \begin{bmatrix} \cos \Delta/2 & -\sin \Delta/2 \\ \sin \Delta/2 & \cos \Delta/2 \end{bmatrix}$$

where d is the thickness of the magnetic film, λ_0 is the vacuum wavelength, and

$$\Delta = \frac{2\pi(n_2 - n_1)d}{\lambda_0}$$

- (c) Given that when the magnetic field reverses, the roles of n_1 and n_2 reverse, show that the Jones matrix for the device when the magnetic field points in the direction opposite to the direction of wave propagation is

$$\mathbf{L}_- = \begin{bmatrix} \cos \Delta/2 & \sin \Delta/2 \\ -\sin \Delta/2 & \cos \Delta/2 \end{bmatrix}.$$

- 7-6. Show that if a Faraday-rotating magnetic film is illuminated by light polarized in the y direction and the film is followed by an analyzer oriented in the x direction, a reversal of the direction of the magnetic field results in a change of the phase of the transmitted light by 180° , with no change in the transmitted intensity. Thus in this configuration the MOSLM can be used as a binary phase SLM.
- 7-7. A magneto-optic film has a Faraday rotation coefficient of $1.46^\circ/\mu\text{m}$ and an absorption coefficient of $0.086 \mu\text{m}^{-1}$. Find the thickness of the film that will maximize the light efficiency of the device in the "on" state, given that the polarization analyzer is oriented to assure complete extinction in the "off" state.
- 7-8. An ideal grating has a profile that is illustrated by the triangular curve in Fig. **W.8** This ideal profile is approximated by a four-level quantized grating profile also shown in the figure. The peak-to-peak phase difference introduced by the continuous grating is exactly 2π radians.

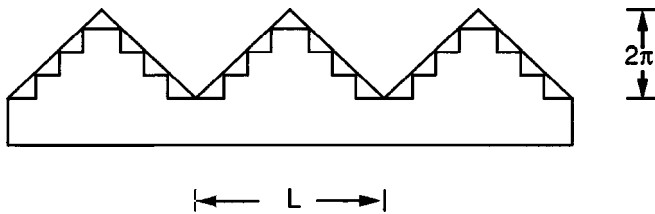


FIGURE P7.8
Profiles of ideal and quantized gratings.

- (a) Find the diffraction efficiencies of the $+4$, $+3$, $+2$, ± 1 , and 0 orders of the continuous grating.
- (b) Find the diffraction efficiencies of the same orders for the quantized grating.

Analog Optical Information Processing

The broad utility of linear systems concepts in the analysis of imaging systems is evident from the preceding chapters. However, if these concepts were useful *only* for analysis purposes, they would occupy a far less important position in modern optics than they in fact enjoy today. Their true importance comes into full perspective only when the exciting possibilities of system *synthesis* are considered.

There exist many examples of the benefits reaped by the application of linear systems concepts to the synthesis of optical systems. One class of such benefits has arisen from the application of frequency-domain reasoning to the improvement of various types of imaging instruments. Examples of this type of problem are discussed in their historical perspective in Section 8.1.

There are equally important applications that do not fall in the realm of imaging as such, but rather are more properly considered in the general domain of *information processing*. Such applications rest on the ability to perform general linear transformations of input data. In some cases, a vast amount of data may, by its sheer quantity, overpower the effectiveness of the human observer. A linear transformation can then play a crucial role in the *reduction* of large quantities of data, yielding indications of the particular portions of the data that warrant the attention of the observer. An example of this type of application is found in the discussion of *character recognition* (Section 8.6). In other cases, a body of data may simply not be in a form compatible with a human observer, and a linear transformation of the data may place it in a compatible form. An example of this type of application is found in the discussion of processing synthetic-aperture radar data (Section 8.9).

The entire subject of optical information processing is too broad to be fully treated in any single chapter; in fact, many books devoted exclusively to the subject already exist (e.g. see Refs. [231], [284], [182], [47], [21], [148], and [293]). We shall limit our goals here to a presentation of the most important and widely used analog optical information processing architectures and applications. We explicitly exclude from consideration the subject of "digital" or "numerical" optical computing, since this field

is not yet well developed. The reader interested in the digital domain can consult, for example, Refs. [100], [297], [218], [204], [160], or [156].

8.1 HISTORICAL BACKGROUND

The history of Fourier synthesis techniques can be said to have begun with the first intentional manipulations of the spectrum of an image. Experiments of this type were first reported by Abbe in 1873 [1] and later by Porter in 1906 [233]. In both cases the express purposes of the experiments were verification of Abbe's theory of image formation in the microscope and an investigation of its implications. Because of the beauty and simplicity of these experiments, we discuss them briefly here.

8.1.1 The Abbe-Porter Experiments

The experiments performed by Abbe and Porter provide a powerful demonstration of the detailed mechanism by which coherent images are formed, and indeed the most basic principles of Fourier analysis itself. The general nature of these experiments is illustrated in Fig. 8.1. An object consisting of a fine wire mesh is illuminated by collimated, coherent light. In the back focal plane of the imaging lens appears the Fourier spectrum of the periodic mesh, and finally in the image plane the various Fourier components passed by the lens are recombined to form a replica of the mesh. By placing various obstructions (e.g. an iris, a slit, or a small stop) in the focal plane, it is possible to directly manipulate the spectrum of the image in a variety of ways.

Figure 8.2(a) shows a photograph of the spectrum of the mesh; Fig. 8.2(b) is the full image of the original mesh. The periodic nature of the object generates in the focal plane a series of isolated spectral components, each spread somewhat by the finite extent of the circular aperture within which the mesh is confined. Bright spots along the horizontal axis in the focal plane arise from complex-exponential components of the object that are directed horizontally (cf. Fig. 2.1); bright spots along the vertical axis correspond to vertically directed complex-exponential components. Off-axis spots correspond to components directed at corresponding angles in the object plane.

The power of *spatial filtering* techniques is well illustrated by inserting a narrow slit in the focal plane to pass only a single row of spectral components. Figure 8.3(a)

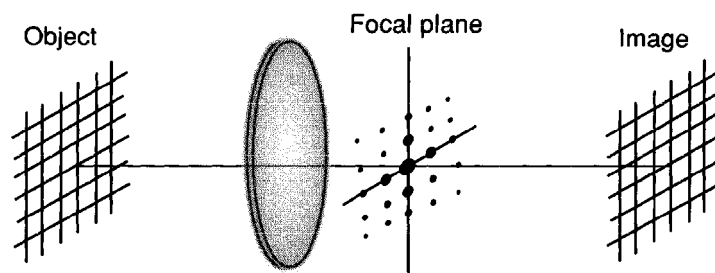
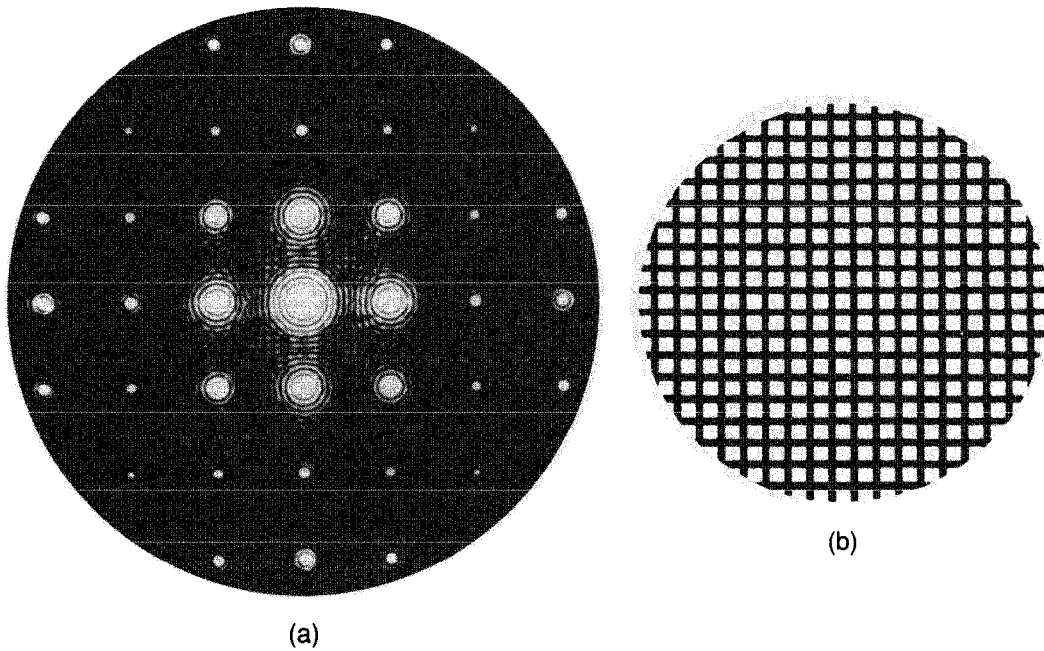


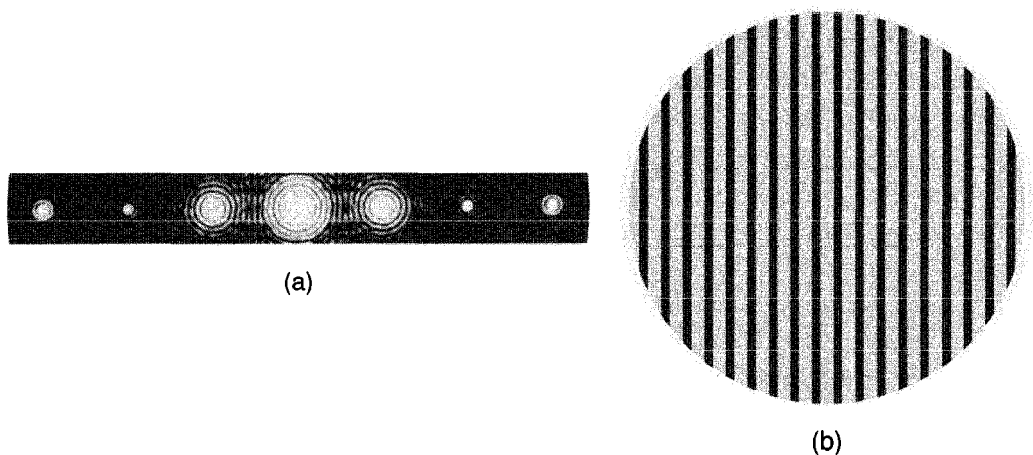
FIGURE 8.1
The Abbe-Porter experiment.

**FIGURE 8.2**

Photograph of (a) the spectrum of mesh and (b) the original mesh.

shows the transmitted spectrum when a horizontal slit is used. The corresponding image, seen in Fig. 8.3(b), contains only the vertical structure of the mesh; it is precisely the horizontally directed complex-exponential components that contribute to the structure in the image that is uniform vertically. The suppression of the horizontal structure is quite complete.

When the slit is rotated by 90° to pass only the spectral column of Fig. 8.4(a), the image in part (b) of the figure is seen to contain only horizontal structure. Other interesting effects can also be readily observed. For example, if an iris is placed in the focal plane and stopped down to pass only the on-axis Fourier component, then with a gradual expansion of the iris the Fourier synthesis of the mesh can be watched step by step. In addition, if the iris is removed and a small central stop is placed on the optical

**FIGURE 8.3**

Mesh filtered with a horizontal slit in the focal plane. (a) Spectrum, (b) image.

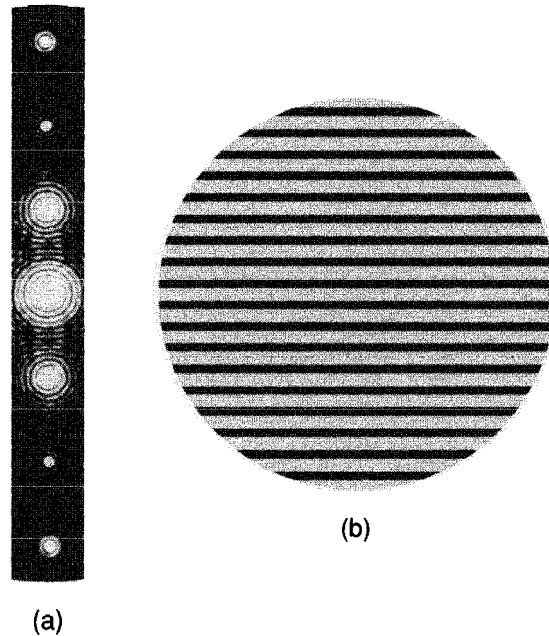


FIGURE 8.4
Mesh filtered with a vertical slit in the focal
plane. (a) Spectrum, (b) image.

axis in the focal plane to block only the central order or "zero-frequency" component, then a contrast reversal can be seen in the image of the mesh (see Prob. 8-1).

8.1.2 The Zernike Phase-Contrast Microscope

Many objects of interest in microscopy are largely transparent, thus absorbing little or no light (e.g. an unstained bacterium). When light passes through such an object, the predominant effect is the generation of a spatially varying phase shift; this effect is not directly observable with a conventional microscope and a sensor that responds to light intensity. A number of techniques for viewing such objects have been known for many years; these include interferometric techniques, the *central dark ground method* in which a small stop is used on the optical axis in the focal plane to block only the zero-frequency spectral component (see Prob. 8-2), and the *schlieren method* in which all spectral components to one side of the zero-frequency component are excluded (see Prob. 8-3). All these techniques suffer from a similar defect—the observed intensity variations are *not* linearly related to the phase shift and therefore cannot be taken as directly indicative of the thickness variations of the object.

In 1935, Frits Zernike [305] proposed a new *phase contrast* technique which rests on spatial filtering principles and has the advantage that the observed intensity *is* (under certain conditions to be discussed) linearly related to the phase shift introduced by the object.¹ This development represents an early success of synthesis ideas and therefore will be treated in some additional detail.

¹For a discussion of the history of the phase contrast technique, as well as the scientific life of Frits Zernike, see [101].

Suppose that a transparent object, with amplitude transmittance

$$t_A(\xi, \eta) = \exp[j\phi(\xi, \eta)] \quad (8-1)$$

is coherently illuminated in an image-forming system. For mathematical simplicity we assume a magnification of unity and neglect the finite extent of the entrance and exit pupils of the system. In addition, a necessary condition to achieve linearity between phase shift and intensity is that the variable part of the object-induced phase shift, $\Delta\phi$, be small compared with 2π radians, in which case the crudest approximation to amplitude transmittance might be

$$t_A(\xi, \eta) = e^{j\phi_o} e^{j\Delta\phi} \approx e^{j\phi_o} [1 + j\Delta\phi(\xi, \eta)]. \quad (8-2)$$

In this equation we have neglected terms in $(\Delta\phi)^2$ and higher powers, assuming them to be zero in our approximation, and the quantity ϕ_o represents the average phase shift through the object, so $\Delta\phi(\xi, \eta)$ by definition has no zero-frequency spectral component. Note that the first term on the right of Eq. (8-2) represents a strong wave component that passes through the sample suffering a uniform phase shift ϕ_o , while the second term generates weaker diffracted light that is deflected away from the optical axis.

The image produced by a conventional microscope could be written, in our approximation, as

$$I_i \approx |1 + j\Delta\phi|^2 \approx 1$$

where, to remain consistent with our approximation, the term $\Delta\phi^2$ has been replaced by zero. Zernike realized that the diffracted light arising from the phase structure is not observable in the image plane because it is in *phase quadrature* with the strong background, and that if this phase-quadrature relation could be modified, the two terms might interfere more directly to produce observable variations of image intensity. Recognizing that the background is brought to focus on the optical axis in the focal plane while the diffracted light, arising from higher spatial frequencies, is spread away from the optical axis, he proposed that a phase-changing plate be inserted in the focal plane to modify the phase relation between the focused and diffracted light.

The phase-changing plate can consist of a glass substrate on which a small transparent dielectric dot has been **deposited**.² The dot is centered on the optical axis in the focal plane and has a thickness and index of refraction such that it retards the phase of the focused light by either $\pi/2$ radians or $3\pi/2$ radians relative to the phase retardation of the diffracted light. In the former case the intensity in the image plane becomes

$$I_i = |\exp[j(\pi/2)] + j\Delta\phi|^2 = |j(1 + \Delta\phi)|^2 \approx 1 + 2\Delta\phi \quad (8-3)$$

while in the latter case we have

$$I_i = |\exp[j(3\pi/2)] + j\Delta\phi|^2 = |-j(1 - \Delta\phi)|^2 \approx 1 - 2\Delta\phi. \quad (8-4)$$

²In practice, phase-contrast microscopes usually have a source that is a circular ring and a phase-shifting structure that is also a circular ring, placed over the image of the source in the focal plane. However, the explanation based on the assumption of point-source illumination is somewhat simpler to explain and to understand.

Thus the image intensity has become linearly related to the variations of phase shift $\Delta\phi$. The case of Eq. (8-3) is referred to as *positive phase contrast* while the case of Eq. (8-4) is referred to as *negative phase contrast*. It is also possible to improve the contrast of the phase-induced variations of intensity in the image by making the phase-shifting dot partially absorbing (see Prob. 8-4).

The phase-contrast method is one technique for converting a spatial phase modulation into a spatial intensity modulation. The reader with a background in communications may be interested to note that one year after Zernike's invention a remarkably similar technique was proposed by E.H. Armstrong [8] for converting amplitude-modulated electrical signals into phase-modulated signals. As we have seen in Chapter 6 and will continue to see in this chapter, the disciplines of optics and electrical engineering were to develop even closer ties in the years to follow.

8.1.3 Improvement of Photographs: Maréchal

In the early 1950s, workers at the Institut d'Optique, Université de Paris, became actively engaged in the use of coherent optical filtering techniques to improve the quality of photographs. Most notable was the work of A. Maréchal, whose success with these techniques was to provide a strong motivation for future expansion of interest in the optical information processing field.

Maréchal regarded undesired defects in photographs as arising from corresponding defects in the optical transfer function of the incoherent imaging system that produced them. He further reasoned that if the photographic transparencies were placed in a coherent optical system, then by insertion of appropriate attenuating and phase-shifting plates in the focal plane, a *compensating filter* could be synthesized to at least partially remove the undesired defects. While the optical transfer function of the initial imaging system might be poor, the product of that transfer function with the (amplitude) transfer function of the compensating system would hopefully yield an overall frequency response that was more satisfactory.

A variety of types of improvements to photographs were successfully demonstrated by Maréchal and his co-workers. For example, it was shown that small details in the image could be strongly emphasized if the low-frequency components of the object spectrum were simply attenuated. Considerable success was also demonstrated in the removal of image blur. In the latter case, the original imaging system was badly defocused, producing an impulse response which (in the geometrical-optics approximation) consisted of a uniform circle of light. The corresponding optical transfer function was therefore of the form

$$\mathcal{H}(\rho) \approx 2 \frac{J_1(\pi a \rho)}{\pi a \rho},$$

where a is a constant and $\rho = \sqrt{f_x^2 + f_y^2}$. The compensating filter was synthesized by placing both an absorbing plate and a phase-shifting plate in the focal plane of the coherent filtering system, as shown in Fig. 8.5(a). The absorbing plate attenuated the large low-frequency peak of \mathcal{H} , while the phase-shifting plate shifted the phase of

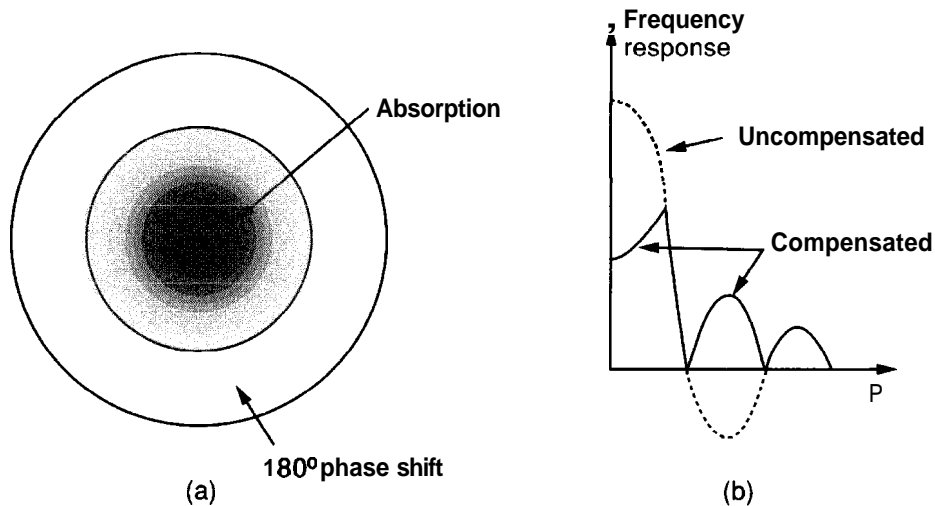


FIGURE 8.5 Compensation for image blur. (a) Focal-plane filter; (b) transfer functions.

the first negative lobe of \mathcal{H} by 180°. The original and compensated transfer functions are illustrated in Fig. 8.5(b).

As an additional example, it was shown that the periodic structure associated with the halftone process used in printing photographs, for example, in newspapers, could be suppressed by a simple spatial filter. The halftone process is, in many respects, similar to the periodic sampling procedures discussed in Section 2.4. The spectrum of a picture printed in this fashion has a periodic structure much like that illustrated in Fig. 2.5. By inserting an iris in the focal plane of the filtering system, it is possible to pass only the harmonic zone centered on zero frequency, thereby removing the periodic structure of the picture while passing all of the desired image data.

Notice a common requirement in all the applications mentioned above: a picture or photograph taken in incoherent light is filtered in a system that uses coherent light. To assure that linear systems are used, and therefore that transfer function concepts remain valid, it is necessary that *the amplitude introduced into the coherent system be proportional to the intensity of the image we wish to filter.*

8.1.4 The Emergence of a Communications Viewpoint

In the early 1950s it became evident that an exchange between the disciplines of communications and optics could reap high profits. Many of the problems facing those working in optics bore strong resemblances to the optimum filtering, detection, and estimation problems of communications theory. Much initial stimulus toward an exchange was provided by a communication theorist, Peter Elias, and his associates D.S. Gray and D.Z. Robinson, with the publication of a paper in 1952 entitled "Fourier treatment of optical processes" [95], and again by Elias with the publication of the paper "Optics and communication theory" [94] in 1953. However, the most complete wedding between the two viewpoints was provided by a physicist, E.L. O'Neill, with the

publication of his paper "Spatial filtering in optics" in 1956 [222], and more generally through the great impact of his research and teaching.

Since the time of this early work, the merger of the two points of view has become so complete that it is sometimes difficult to judge whether a particular piece of work should be published in an optics journal or an electrical engineering journal.

8.1.5 Application of Coherent Optics to More General Data Processing

While the early 1950s were characterized by a growing realization on the part of physicists that certain aspects of electrical engineering were of particular relevance to optics, the late fifties and early sixties saw a gradual realization on the part of electrical engineers that spatial filtering systems might be usefully employed in their more general data-processing problems. The potentials of coherent filtering were particularly evident in the field of radar signal processing and were exploited at an early stage by L.J. Cutrona and his associates at the University of Michigan Radar Laboratory. The publication of the paper "Optical data processing and filtering systems" [73] by the Michigan group in 1960 stimulated much interest in these techniques among electrical engineers and physicists alike. One of the most successful early applications of coherent filtering in the radar realm has been to the processing of data collected by synthetic aperture radar systems [74], a subject that will be briefly treated in Section 8.9. A survey of the literature from the mid-1960s shows application of coherent processing techniques in such widely diverse fields as, for example, Fourier spectroscopy [276] and seismic-wave analysis [153].

8.2 INCOHERENT OPTICAL INFORMATION PROCESSING SYSTEMS

The use of spatially incoherent light in optical information processing provides certain advantages, but also certain disadvantages. Important advantages include the general freedom of incoherent systems from coherent artifacts, for example, those associated with dust specks on the optical components and those that arise from the speckle phenomenon. These advantages can be attributed to a certain redundancy found in incoherent systems, namely the fact that light from a single pixel or resolvable spot of an input passes through the system via many spatially separate channels, due to the extended nature of the incoherent source. In addition, incoherent systems allow the introduction of data into the system by means of light-emitting diode arrays or cathode-ray tube (CRT) displays, and do not require the more complex and expensive **SLMs** discussed in the previous chapter. Generally speaking, incoherent systems are somewhat more simple than coherent systems in their physical realization. For a historical perspective on the field of incoherent optical processing, see Ref. [246].

However, the above advantages are accompanied by some serious disadvantages as well. An incoherent optical processing system has no "frequency plane", as is found in the focal plane of a coherent optical system, and the manipulation of the spectrum of an input must therefore resort to less direct methods than simply modifying the fields

in the Fourier plane in proportion to a desired transfer function. Second, the intensity of light is fundamentally a nonnegative and real physical quantity, and the representation of data by intensity places limits on the type of data manipulations that can be carried out in purely optical form. For example, there is no natural optical way to subtract two intensity patterns, whereas complex amplitude patterns can in principle be subtracted by adding them with a π radian phase shift between them. The fact that an incoherent image always has its maximum spectral content at the origin often leads to problems of low contrast at the output of incoherent processing systems. As a consequence, incoherent systems often must have a heavy intrusion of electronics at their output in order to achieve a flexibility comparable with that of coherent systems.

Incoherent data processing systems can be broadly divided into three separate categories: (1) systems based on geometrical optics, (2) systems based on diffraction, and (3) discrete systems. The first two categories of systems are designed to accommodate spatially continuous inputs; both will be discussed here. The third category, discrete systems, will be deferred to a subsequent section of this chapter.

8.2.1 Systems Based on Geometrical Optics

A variety of methods are known for designing optical information processing systems based purely on geometrical optics. Such approaches ignore the diffraction phenomenon, and as will be pointed out later, suffer from limitations on the achievable space-bandwidth product.

Systems based on image casting

Systems based on geometrical optics almost invariably use one form or another of what could be called "image casting" or "shadow casting", namely the geometrical projection of one image onto another. Such a system was proposed as early as 1927 by Emanuel **Goldberg** of Dresden, Germany, in a patent application. Goldberg, who was granted a U.S. patent in 1931 [118], fully recognized the potential application of his invention to the field of character recognition.

The principles underlying the most simple image-casting system, namely a system that performs a spatial integration of the product of two functions, are straightforward. If a transparency with intensity transmittance τ_1 is imaged onto a second transparency with intensity transmittance τ_2 , then according to geometrical optics the intensity at each point immediately behind the second transparency is $\tau_1\tau_2$. A photodetector can be used to measure the *total* intensity transmitted through the pair, yielding a photocurrent I given by³

$$I = k \iint_{-\infty}^{\infty} \tau_1(x, y) \tau_2(x, y) dx dy. \quad (8-5)$$

³As usual, in writing infinite limits of integration, we have assumed that the finite sizes of the transparencies are incorporated in the functions τ_1 and τ_2 .

Two means of achieving this operation are illustrated in Fig. 8.6. For the technique (a), the lens L_1 casts a magnified image of the uniform incoherent source onto the two transparencies which are placed in direct contact. The lens L_2 then casts a demagnified image of the light transmitted by τ_2 onto the photodetector D. The photocurrent is then given by Eq. (8-5).

In some cases it may be desired to change one of the inputs rapidly, in which case physical separation of the two transparencies may be advantageous. Such separation can be achieved with the geometry shown in Fig. 8.6(b). The lens L_1 again casts a magnified image of the source onto τ_1 . Lens L_2 images τ_1 onto τ_2 , and lens L_3 casts a demagnified image of the light transmitted by τ_2 onto the detector. Note that the transparency τ_1 must be inserted in an inverted geometry to compensate for the inversion introduced by the imaging operation performed by L_2 . Again the photocurrent is given by Eq. (8-5).

While the operation described above is a useful one in a number of applications, including character recognition, it is often desired to realize the related but more general operation of convolution. A one-dimensional convolution of two functions can be realized with either of the above systems by moving one of the transparencies with

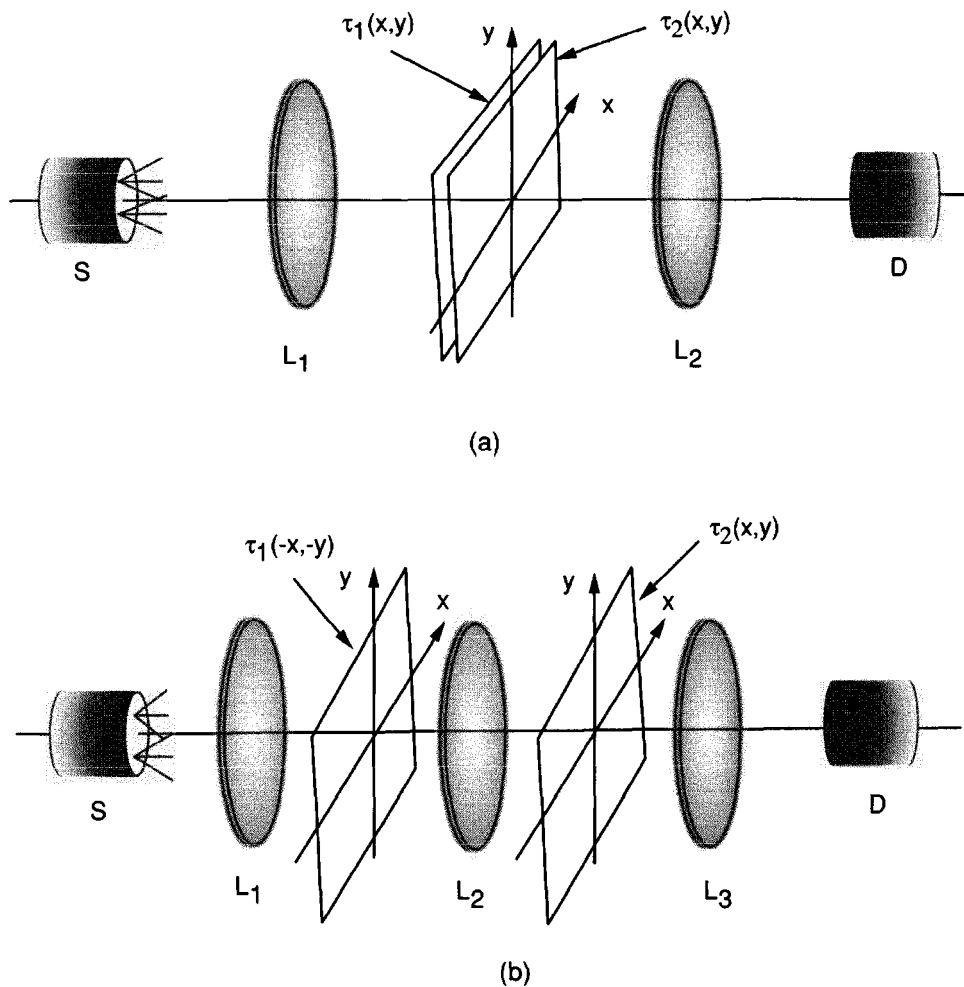


FIGURE 8.6
Systems for realizing the integral of a product of two functions.

uniform velocity and measuring the photodetector response as a function of time. More specifically, with reference to Fig. 8.6(b), let the transparency τ_1 be introduced without the inversion referred to earlier, so that the operation (8-5) becomes

$$I = k \iint_{-\infty}^{\infty} \tau_1(-x, -y) \tau_2(x, y) dx dy.$$

If the transparency τ_1 is moved in the negative x direction with speed v , the detector response as a function of time will be given by

$$I(t) = k \iint_{-\infty}^{\infty} \tau_1(vt - x, -y) \tau_2(x, y) dx dy.$$

If the scans are repeated sequentially, each for a different y displacement $-y_m$, then the detector responses will be

$$I_m(t) = k \iint_{-\infty}^{\infty} \tau_1(vt - x, y_m - y) \tau_2(x, y) dx dy. \quad (8-6)$$

The array of functions $I_m(t)$ represents a full two-dimensional convolution, albeit sampled in the y displacement.

Convolution without motion

The preceding technique for performing convolutions is extremely awkward and time-consuming due to the mechanical scanning required. It is possible to perform the same operation *without relative motions* if the optical configuration is modified [175]. Referring to Fig. 8.7, let the distributed incoherent source S be placed in the front focal plane of the lens L_1 . Immediately behind L_1 is placed a transparency with intensity transmittance $\tau_1(-x, -y)$. At a distance d from τ_1 and immediately in front of lens L_2 the transparency $\tau_2(x, y)$ appears. The intensity distribution across the back focal

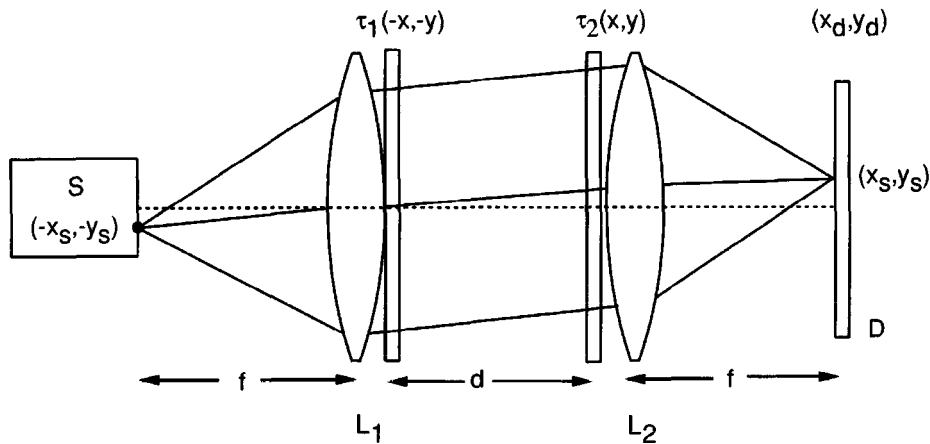


FIGURE 8.7
Systems for performing convolution without motion.

plane of L_2 is then measured, perhaps with film, although the use of a two-dimensional electronic detector, such as a vidicon, is also possible.

To understand the operation of this system, consider first the light generated by a particular point with coordinates $(-x_s, -y_s)$ on the source. The rays from that point emerge from L_1 (and from τ_1) parallel with each other and illuminate τ_2 with an intensity distribution proportional to $\tau_1[-x + (d/f)x_s, -y + (d/f)y_s]$. After passing through τ_2 the rays are focused onto the detector at coordinates (x_d, y_d) , where we have assumed that the two lenses have identical focal lengths. Thus the intensity distribution across the detector may be written

$$I(x_d = x_s, y_d = y_s) = k \iint_{-\infty}^{\infty} \tau_1 \left(\frac{d}{f} x_s - x, \frac{d}{f} y_s - y \right) \tau_2(x, y) dx dy, \quad (8-7)$$

which is the desired convolution.

Impulse response synthesis with a misfocused system

Direct synthesis of a desired impulse is possible, within the confines of geometrical optics, by means of the "misfocused" system illustrated in Fig. 8.8. While this system is nearly identical with that of Fig. 8.7, the point of view is sufficiently different to warrant a separate discussion. The lens L_1 again serves to illuminate the "input" transparency (intensity transmittance τ_1) with uniform light from the extended source S . Lens L_2 forms an image of τ_1 across the plane P' . For simplicity we assume that τ_1 and P' are each at distance $2f$ from the lens L_2 , thus yielding magnification unity in the proper image plane. The transparency τ_2 , having an intensity transmittance equal in form to that of the desired impulse response, is inserted directly against lens L_2 ; the system output is found across the plane P , located distance A from the ideal image plane P' .

The operation of this system is most easily understood by applying a unit-intensity point source at coordinates (x, y) on τ_1 and finding the resulting intensity distribution across P . In the geometrical-optics approximation, the rays passing through τ_2 converge

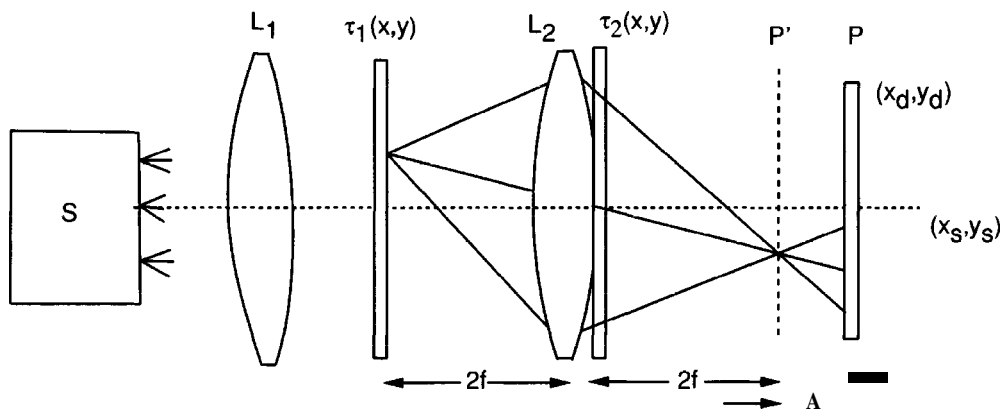


FIGURE 8.8
Impulse response synthesis with a misfocused system.

to an ideal point in plane \mathbf{P}' , and then diverge to form a demagnified projection of τ_2 in plane \mathbf{P} . The projection is centered at coordinates

$$x_d = -\left(1 + \frac{\Delta}{2f}\right)x$$

$$y_d = -\left(1 + \frac{\Delta}{2f}\right)y,$$

and the demagnification of τ_2 is $\Delta/2f$. Taking into account the reflection of τ_2 when projected, the response to the point source may be written

$$|h(x_d, y_d; x, y)|^2 = k\tau_2 \left\{ -\frac{2f}{\Delta} \left[x_d + \left(1 + \frac{\Delta}{2f}\right)x \right], -\frac{2f}{\Delta} \left[y_d + \left(1 + \frac{\Delta}{2f}\right)y \right] \right\}. \quad (8-8)$$

The intensity at output coordinates $(-x_d, -y_d)$ can then be written as the convolution integral

$$I(-x_d, -y_d) = k' \iint_{-\infty}^{\infty} \tau_1(x, y) |h(-x_d, -y_d; x, y)|^2 dx dy \quad (8-9)$$

where Eq. (8-8) must be substituted after the change to arguments $(-x_d, -y_d)$. Except for the scaling factor $\Delta/2f$ that precedes x and y , the resulting form is a convolution. By properly scaling τ_2 this scaling factor can be removed and a true convolution can be obtained.

Limitations

All systems designed on the basis of geometrical optics must satisfy a common constraint: the geometry of the system must be chosen in such a way that diffraction effects are entirely negligible. This requirement is readily satisfied with the system of Fig. 8.6(a), but is difficult to meet in all the other systems presented, to varying degrees.

A measure of the power of an optical information processing system is the space-bandwidth product of the input function that it will accept. To maximize the input space-bandwidth product, we would attempt to place as many independent data points as possible on the transparencies. But as the structure on the input transparencies gets finer and finer, more and more of the light passing through them will be diffracted, with less and less of the light obeying the laws of geometrical optics. Thus the methods used in the analyses of these systems become less and less accurate, and the system outputs will depart more and more severely from their predicted forms.

While we have considered only a few specific examples of systems based on geometrical optics, a fundamental fact is clear: if large quantities of data are to be squeezed into an aperture of a given size, ultimately diffraction effects must be taken into account. It is therefore extremely important to be sure that when a system is designed on the basis of geometrical optics, it is used in a way that assures accuracy of the laws of geometrical optics.

8.2.2 Systems That Incorporate the Effects of Diffraction

It is possible to design incoherent optical information processing systems that take full account of the laws of diffraction. When considering such systems, some of the other difficulties associated with the use of incoherent light become more evident. The two major difficulties encountered in attempting to perform general filtering operations with incoherent light are (1) the point-spread functions that can be synthesized must be non-negative and real, since they are intensity distributions—a constraint that restricts the generality of the operations that can be directly achieved, and (2) there are many different pupil-plane masks that will generate the same intensity point-spread function, but no known method for finding the simplest such mask.

Nonetheless, interesting operations can be performed even under the constraints mentioned. We illustrate here with one particular approach introduced by Rhodes [244] called two-pupil OTF synthesis. This approach provides a means for performing bandpass filtering using incoherent light for at least part of the processing operation. For an alternative example of a system that uses incoherent processing in a portion of its operation and is consistent with the laws of diffraction, see Prob. 8-10 [7].

Bandpass filtering is an operation that fundamentally requires subtraction, for the large low-frequency components that are always present in incoherent images must be removed by the processing operations. While some nonlinear optical phenomena allow the light from one strong incoherent image to suppress the strength of a second, weaker incoherent image, there is no optical operation, linear or nonlinear, that can produce a negative intensity, so true subtraction is not possible with purely optical operations. For this reason, incoherent processing must be supplemented with some other form of processing, either electronic processing or coherent optical processing, to achieve the desired **bandpass** operation.

As the name "two-pupil OTF synthesis" implies, this method accounts for the laws of diffraction by manipulating the OTFs of the optical systems used. The OTF of an optical system is found from the (normalized) Fourier transform of the point-spread function of the system, and the calculation of a point-spread function is a calculation that is based on the laws of diffraction.

If we collect an incoherent image using an optical system with the pupil shown in Fig. 8.9(a), the resulting OTF of the system is as shown in part (b) of the same figure. We have eliminated some midfrequency components and have emphasized the spatial frequencies present in the **passband** of interest, but the low frequencies remain very prominent. Suppose that we now place a phase plate over one of the two apertures in the pupil, a phase plate that introduces a 180° phase shift, which we assume to be approximately constant over the narrow band of optical wavelengths used in this experiment. The autocorrelation function of this new pupil yields an OTF as shown in part (c) of the figure, with the sign of the **bandpass** regions of the OTF reversed, but with the low-frequency part of the OTF left unchanged. Finally, suppose that the two image intensities collected with the OTFs shown in parts (b) and (c) of the figure are subtracted, perhaps by an electronic system. The effective transfer function for the difference image is the difference of the two OTFs used in collecting those images, or the transfer function shown in part (d) of the figure, which indeed provides a true **bandpass** filter.

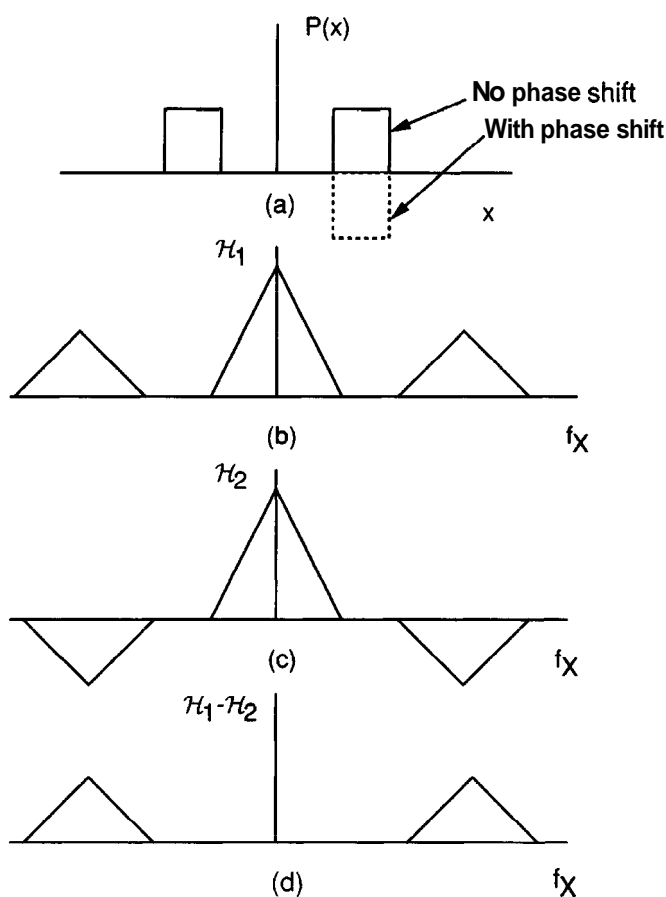


FIGURE 8.9
Two-pupil OTF synthesis.

Many variations on this theme are possible. The phase plate can be replaced by a phase modulator that shifts the temporal frequency of the light passing through one of the two apertures, and the detector can select only the difference frequency in order to eliminate the unmodulated low-frequency portions of the spectrum. Alternatively the two incoherent images can be translated by spatial light modulators into coherent images, and the amplitudes of those can be added with a 180° phase difference in an interferometer to achieve subtraction (the result is an intensity distribution representing the squared-magnitude of the difference image). In addition, synthesis of transfer functions more general than a simple **bandpass** filter are possible. For a more complete discussion, together with more references, see Ref. [182], Chapter 3.

By considering incoherent filtering systems that include the effects of diffraction, we have focused on the remaining serious problem that arises in such filtering, namely the nonnegativity of optical intensity and the lack of a convenient method for subtracting intensities. Even with electronic subtraction, it is often found that the low-frequency components being subtracted are very strong compared with the high-frequency information of interest, and imperfections in the subtraction operation may leave image artifacts or noise of serious proportions.

In summary, incoherent optical information processing is often simpler than coherent optical processing (particularly in the forms that use image casting), but in general is much less flexible in terms of the operations that can be achieved. We therefore turn to a consideration of information processing using coherent light.

8.3

COHERENT OPTICAL INFORMATION PROCESSING SYSTEMS

When coherent illumination is used, filtering operations can be synthesized by direct manipulation of the complex amplitude appearing in the back focal plane of a Fourier transforming lens. Examples of this type of processing have already been seen in the discussion of the phase-contrast microscope (Zernike) and the filtering of photographs (Maréchal). In this section we outline the system architectures used for coherent optical information processing, and point out some of the difficulties encountered in attempting to synthesize general complex filters.

8.3.1 Coherent System Architectures

Coherent systems, being linear in complex amplitude, are capable of realizing operations of the form

$$I(x, y) = K \left| \iint_{-\infty}^{\infty} g(\xi, \eta) h(x - \xi, y - \eta) d\xi d\eta \right|^2. \quad (8-10)$$

There are many different system configurations that can be used to realize this operation, three of which are shown in Fig. 8.10.

The system shown in part (a) of the figure is conceptually the most straightforward and is often referred to as a “4f” filtering architecture, due to the fact that there are four separate distances of length f separating the input plane from the output plane. Light from the point source S is collimated by lens L_1 . In order to minimize the length of the system, the input transparency, having amplitude transmittance $g(x_1, y_1)$, is placed against the collimating lens in plane P_1 . One focal length beyond the input is a Fourier transforming lens L_2 , in the rear focal plane (P_2) of which is placed a transparency to control the amplitude transmittance through that plane. An amplitude $k_1 G(x_2/\lambda f, y_2/\lambda f)$ is incident on this plane, where G is the Fourier transform of g and k_1 is a constant. A filter is inserted in plane P_2 to manipulate the spectrum of g . If H represents the desired transfer function, then the amplitude transmittance of the frequency-plane filter should be

$$t_A(x_2, y_2) = k_2 H\left(\frac{x_2}{\lambda f}, \frac{y_2}{\lambda f}\right). \quad (8-11)$$

The field behind the filter is thus GH . After one additional focal length, lens L_3 is placed, the purpose of which is to again Fourier transform the modified spectrum of the input, producing a final output in its rear focal plane, P_3 . Note that the output appears inverted in plane P_3 due to the imaging operation, or equivalently due to the fact that a sequence of two Fourier transforms has been used, rather than one transform followed by its inverse. This awkwardness can be remedied by reversing the final coordinate system (x_3, y_3) , as shown in the figure, in which case the output in plane P_3 is as described by Eq. (8-10). For simplicity, the focal lengths of all three lenses have been assumed

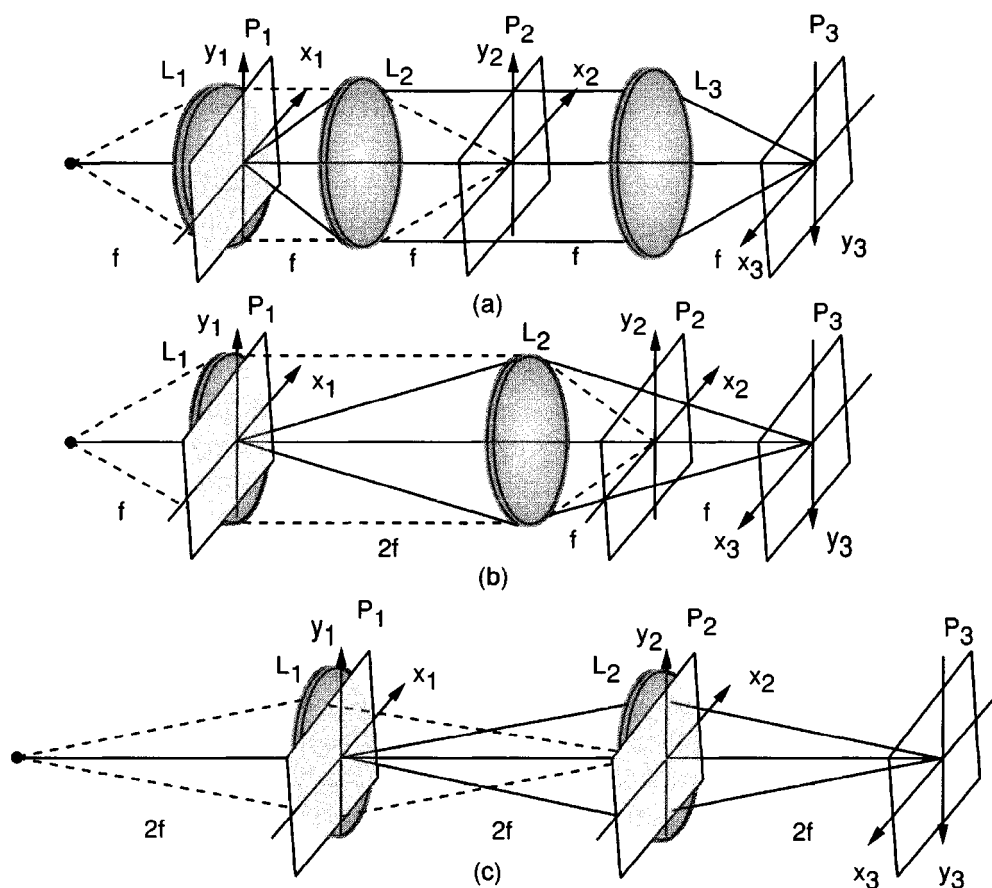


FIGURE 8.10 Architectures for coherent optical information processing.

to bf , and the total length of the system is seen to be $5f$. This architecture has the disadvantage that vignetting can occur during the first Fourier transform operation.

The system shown in part (b) of the figure has the same length as the previous system, but uses one fewer lens. Again lens L_1 collimates the light from the point source S , and again the input transparency is placed against L_1 to minimize the length of the system. Placed at distance $2f$ from the input, lens L_2 now performs both the Fourier transforming and the imaging operations, with the spectrum of the input appearing in the rear focal plane P_2 (where the Fourier filter transparency is placed) and the filtered image appearing one additional focal length beyond the focal plane, in plane P_3 . Since the object and image distances are both $2f$, the magnification of the system is unity. Note that in this geometry, the spectrum of the input has associated with it a quadratic phase factor of the form $\exp[-j\frac{k}{2f}(x_2^2 + y_2^2)]$, since the input is not in the front focal plane of the lens (cf. Eq. (5-19)). This phase factor is not of concern and is indeed needed in order to produce an image at distance $2f$ behind lens L_2 . The length of this system remains $5f$, as before.

There are two practical disadvantages of this second geometry. First, as compared with system (a), the input is now twice the distance from lens L_2 , and therefore the vignetting will be even worse than that encountered with system (a). A second

disadvantage arises from the approximations that led to Eq. (5-30) in the analysis of the coherent imaging properties of a thin lens. In that formulation, we found it necessary to assume that the amplitude of the image at any particular point consisted of contributions from only a small region surrounding the geometrical object point. If the filtering operation represented by the transfer function H is of high space-bandwidth product, then the impulse response h will extend over a sizable area, and the output of this system must be regarded as a filtered version of the function $g(x_1, y_1) \exp[j\frac{k}{4f}(x_1^2 + y_1^2)]$ rather than simply of $g(x_1, y_1)$. This problem is not encountered with system (a), which casts an image of *plane* P_1 onto a *plane* P_3 , rather than of a sphere onto a sphere. This difficulty can be corrected by adding an additional positive lens with focal length $2f$ in contact with the object, thus canceling the troubling quadratic phase factor. This additional lens also results in movement of the frequency plane from f behind lens L_2 to coincidence with that lens, but the location of image plane P_3 is not affected.

As a final example which has merit (but by no means the only other system geometry possible), consider the system shown in part (c) of the figure. Again only two lenses are used. Lens L_1 now serves as both a lens for collecting the light from the point source S and as a Fourier transforming lens. The input is placed in plane P_1 in contact with lens L_1 . This lens images the source onto the frequency plane P_2 , where the filter transparency is placed. The magnification of this imaging operation as shown is unity. The second lens, L_2 , is also placed in this plane, and images the input onto the output plane P_3 with unity magnification. Note that this system has no vignetting problems, and the quadratic phase factor across the input plane (mentioned above) is canceled by the converging illumination. The disadvantage is that the system is now of length $6f$ rather than $5f$.

Finally we mention that it is also possible to arrange a coherent system to process a stacked array of one-dimensional inputs, rather than a single two-dimensional input. An example of these so-called *anamorphic processors*⁴ is shown in Fig. 8.11. The collimating lens L_1 is followed by the input data in plane P_1 . The input data consists of an array of one-dimensional transmittance functions each running horizontally. A cylindrical lens L_2 follows, placed one focal length f from P_1 and having power only in the vertical dimension. At distance $2f$ beyond L_2 is placed a spherical lens L_3 which again has focal length f . The "frequency plane" now appears at P_2 , where an array of one-dimensional spectra is found. The lens combination L_2, L_3 has performed a double Fourier transformation in the y direction, thus imaging in the vertical direction. Since L_2 exerts no power in the x direction, the spherical lens L_3 Fourier transforms in the horizontal dimension, up to a phase factor $\exp(-j\frac{k}{f}x_2^2)$ across P_2 . This phase factor can be removed by placing a negative cylindrical lens of focal length $f/2$ immediately in front of P_2 , thus canceling the phase curvature. If the input array is the set of transmittance functions $g_k(x_1)$, $k = 1, 2, \dots, K$, then across P_2 we find displayed the corresponding set of transforms $G_k(x_2)$, $k = 1, 2, \dots, K$, with the vertical order inverted by the imaging operation.

⁴An optical system is called *anamorphic* if the focusing powers of the system in two orthogonal directions are unequal.

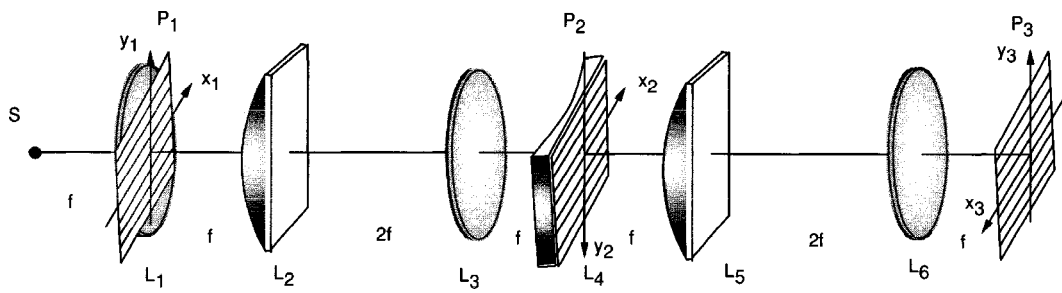


FIGURE 8.11
Example of an anamorphic processor.

A linear array of one-dimensional filters may now be introduced in plane P_2 . The lens pair L_5, L_6 again images in the y direction and Fourier transforms in the x direction, thus retaining the array structure but returning the original functions to the "space domain". The phase factor associated with the final Fourier transform is generally of no concern.

8.3.2 Constraints on Filter Realization

While coherent systems are in general more flexible and have greater data-handling capacity than most incoherent systems, nonetheless there are limitations to the types of operations that can be realized with simple frequency plane filters of the kind used earlier by Maréchal. More sophisticated techniques for realizing frequency-plane masks, based on interferometric recording, are free from some of these limitations, as will be discussed in the section to follow.

Before 1963, the conventional means for realizing a given transfer function had been the insertion of independent amplitude and phase masks in the frequency plane. The amplitude transmittance was controlled by a photographic plate, presumably immersed in a liquid gate. The phase transmittance was controlled by insertion of a transparent plate with an appropriately varying thickness. Such plates could be ruled on a substrate, much as diffraction gratings are ruled, or deposited on a flat plate using thin-film coating techniques. All such methods are rather cumbersome, and could be successfully employed only when the desired pattern of phase control was rather simple, e.g. binary and of simple geometric structure.

Figure 8.12 shows the regions of the complex plane that can be reached by the transfer functions of coherent optical systems under different constraints on the frequency-plane transparency. As shown in (a), when only an absorbing transparency is used, the reachable region is limited to the positive real axis between 0 and 1. If binary phase control is added to this absorbing transparency, then the reachable region is extended to the region -1 to 1 on the real axis, as shown in (b). If a pure phase filter is used, with arbitrary achievable values of phase, then the values of the transfer function would be restricted to the unit circle, as shown in (c). Finally part (d) of the figure shows the region of the complex plane that one would generally desire to reach if there were no constraints, namely the entire unit circle.

It should be noted that, for even a very simple impulse response (such as one in the shape of the character "P", for example), the corresponding transfer function was (1) difficult to calculate (prior to the development of the fast Fourier transform algorithm for digital computation of spectra) and (2) far too complicated to be synthesized by these rather simple techniques.

In summary, the most severe limitation to the traditional coherent processor (prior to the invention of the methods to be discussed in the next section) arose from the difficulty of simultaneously controlling the amplitude and phase transmittances in any but very simple patterns. Thus coherent optical filters were limited to those that had very simple transfer functions. It was not until 1963, with the invention of the interferometrically recorded filter, that this serious limitation was largely overcome, extending the domain of complex filters that could be realized to those with simple impulse responses.

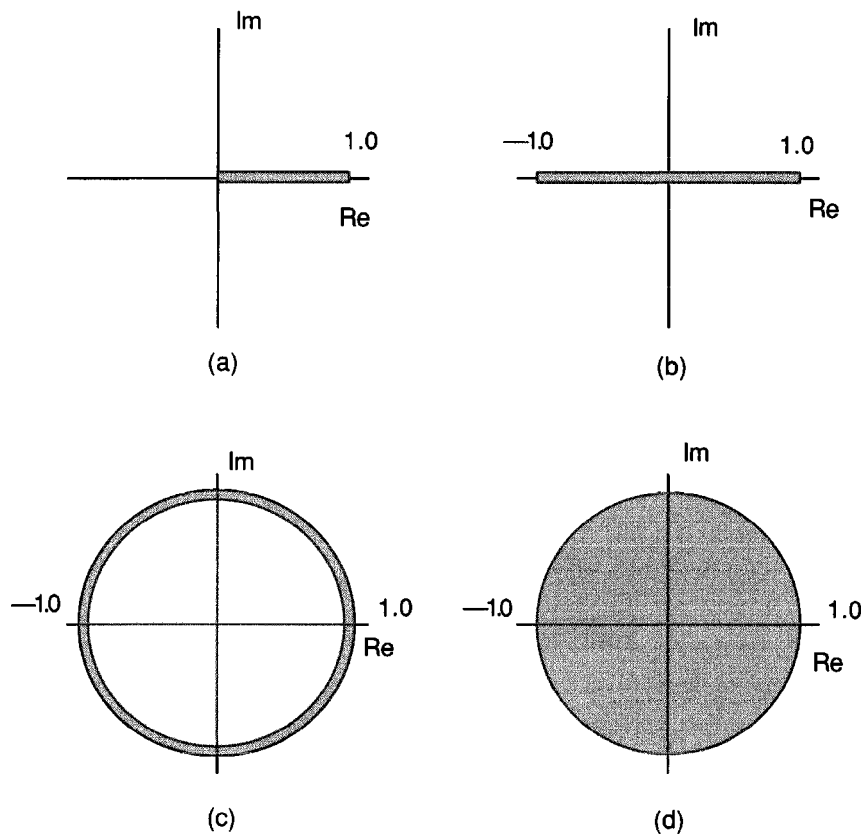


FIGURE 8.12 Reachable regions of the frequency plane for (a) a purely absorbing filter, (b) an absorbing filter and binary phase control, (c) a pure phase filter, and (d) a filter that achieves arbitrary distributions of absorption and phase control.

8.4 THE VANDERLUGT FILTER

In 1963, A.B. VanderLugt of the University of Michigan's Radar Laboratory proposed and demonstrated a new technique for synthesizing frequency-plane masks for coherent optical processors [290], [291].⁵ The frequency-plane masks generated by this technique have the remarkable property that they can effectively control both the amplitude and phase of a transfer function, in spite of the fact that they consist only of patterns of absorption. By means of this technique, it is possible to largely overcome the two limitations to coherent processing systems mentioned above.

⁵Historically, this type of filter had been preceded by a related but less general technique, known as the *hard-clipped filter*, which was a filter generated by computer and is the first example of what now might be called a *phase-only* filter. While the hard-clipped filter was used in radar signal processing as early as 1961, due to classification it did not appear in the open literature until 1965 [179]. The fundamental idea that an interferometric recording of the Fourier transform of the impulse response could realize a complex filter with a desired transfer function or its conjugate is attributable to C. Palermo (private communication, E.N. Leith).

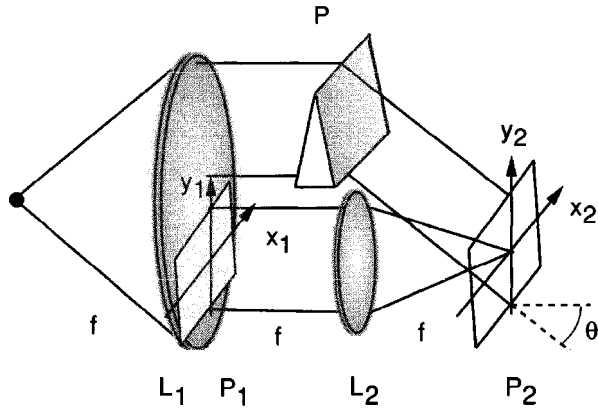


FIGURE 8.13
Recording the frequency-plane mask for a VanderLugt filter.

8.4.1 Synthesis of the Frequency-Plane Mask

The frequency-plane mask of the VanderLugt filter is synthesized with the help of an interferometric (or holographic — see Chapter 9) system, such as that shown in Fig. 8.13. The lens L_1 collimates the light from the point source S . A portion of this light strikes the mask P_1 , which has an amplitude transmittance that is proportional to the desired *impulse response* h . The lens L_2 Fourier transforms the amplitude distribution h , yielding an amplitude distribution $\frac{1}{\lambda f} H(\$, \$)$ incident on the recording medium,⁶ usually film. In addition, a second portion of the collimated light passes above the mask P_1 , strikes a prism P , and is finally incident on the recording plane at angle θ , as shown.

The total intensity incident at each point on the recording medium is determined by the interference of the two mutually coherent amplitude distributions present. The tilted plane wave incident from the prism produces a field distribution

$$U_r(x_2, y_2) = r_o \exp(-j2\pi\alpha y_2), \quad (8-12)$$

where the spatial frequency α is given by

$$\alpha = \frac{\sin \theta}{\lambda}. \quad (8-13)$$

The total intensity distribution may therefore be written

$$\begin{aligned} I(x_2, y_2) &= \left| r_o \exp(-j2\pi\alpha y_2) + \frac{1}{\lambda f} H\left(\frac{x_2}{\lambda f}, \frac{y_2}{\lambda f}\right) \right|^2 \\ &= r_o^2 + \frac{1}{\lambda^2 f^2} \left| H\left(\frac{x_2}{\lambda f}, \frac{y_2}{\lambda f}\right) \right|^2 + \frac{r_o}{\lambda f} H\left(\frac{x_2}{\lambda f}, \frac{y_2}{\lambda f}\right) \exp(j2\pi\alpha y_2) \\ &\quad + \frac{r_o}{\lambda f} H^*\left(\frac{x_2}{\lambda f}, \frac{y_2}{\lambda f}\right) \exp(-j2\pi\alpha y_2). \end{aligned} \quad (8-14)$$

⁶Here and frequently in what follows, we drop a multiplicative factor $1/\lambda f$ associated with the optical Fourier transform, with the justification that we can always change the phase reference for convenience.

Note that if the complex function H has an amplitude distribution A and a phase distribution ψ , that is, if

$$H\left(\frac{x_2}{\lambda f}, \frac{y_2}{\lambda f}\right) = A\left(\frac{x_2}{\lambda f}, \frac{y_2}{\lambda f}\right) \exp\left[j\psi\left(\frac{x_2}{\lambda f}, \frac{y_2}{\lambda f}\right)\right],$$

then the expression for \mathcal{I} can be rewritten in the form

$$\begin{aligned} \mathcal{I}(x_2, y_2) = & r_o^2 + \frac{1}{\lambda^2 f^2} A^2\left(\frac{x_2}{\lambda f}, \frac{y_2}{\lambda f}\right) \\ & + \frac{2r_o}{\lambda f} A\left(\frac{x_2}{\lambda f}, \frac{y_2}{\lambda f}\right) \cos\left[2\pi\alpha y_2 + \psi\left(\frac{x_2}{\lambda f}, \frac{y_2}{\lambda f}\right)\right]. \end{aligned} \quad (8-15)$$

This form illustrates the means by which the interferometric process allows the recording of a complex function H on an intensity-sensitive detector: amplitude and phase information are recorded, respectively, as amplitude and phase modulations of a *high-frequency carrier* that is introduced by the relative angular tilt of the "reference" wave from the prism.

There are, of course, other optical systems that will produce the same intensity distribution as that of Eq. (8-15). Figure 8.14 illustrates two additional possibilities. System (a) consists of a modified Mach-Zehnder interferometer. By tilting the mirror M_1 , a tilted plane wave is produced at the film plane. In the lower arm of the interferometer, the lens L_2 again Fourier transforms the desired impulse response. The final beam splitter allows the addition of these two waves at the recording plane.

System (b), which is a modified Rayleigh interferometer, provides a third means for producing the same intensity distribution. The collimating lens L_1 is followed by a smaller lens L_2 , which focuses a portion of the collimated light to a bright spot in the front focal plane of lens L_3 . When the spherical wave generated by this "reference point" passes through L_3 , it is collimated to produce a tilted plane wave at the recording plane. The amplitude transmitted by the impulse response mask is Fourier transformed in the usual fashion. Thus an intensity distribution similar to Eq. (8-15) is again produced at the recording plane.

As a final step in the synthesis of the frequency-plane mask, the exposed film is developed to produce a transparency which has an amplitude transmittance that is proportional to the intensity distribution that was incident during exposure. Thus the amplitude transmittance of the filter is of the form

$$\begin{aligned} t_A(x_2, y_2) \propto & r_o^2 + \frac{1}{\lambda^2 f^2} |H|^2 + \frac{r_o}{\lambda f} H \exp(j2\pi\alpha y_2) \\ & + \frac{r_o}{\lambda f} H^* \exp(-j2\pi\alpha y_2). \end{aligned} \quad (8-16)$$

Note that, aside from the simple complex-exponential factor, the third term of the amplitude transmittance is proportional to H and therefore exactly the form required to synthesize a filter with impulse response h . It remains to be demonstrated how that particular term of the transmittance can be utilized and the other terms excluded.

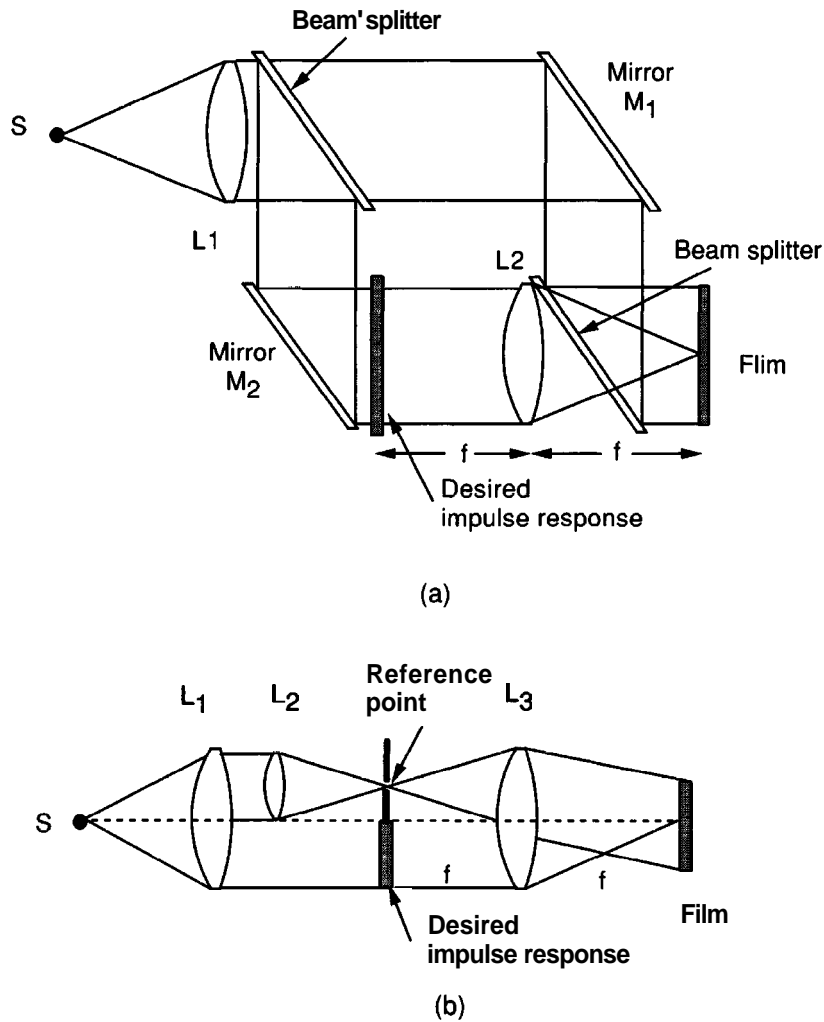


FIGURE 8.14 Two alternative systems for producing the frequency-plane transparency (a) Modified Mach-Zehnder interferometer; (b) modified Rayleigh interferometer.

8.4.2 Processing the Input Data

Once the frequency-plane mask has been synthesized, it may be inserted in any of the processing systems shown previously in Fig. 8.10. To be specific, we focus on the system shown in part (a) of that figure. If the input to be filtered is $g(x_1, y_1)$, then incident on the frequency-plane mask is a complex amplitude distribution given by $\frac{1}{\lambda f} G\left(\frac{x_2}{\lambda f}, \frac{y_2}{\lambda f}\right)$. The field strength transmitted by the mask then obeys the proportionality

$$U_2 \propto \frac{r_o^2 G}{\lambda f} + \frac{1}{\lambda^3 f^3} |H|^2 G + \frac{r_o}{\lambda^2 f^2} H G \exp(j2\pi\alpha y_2) + \frac{r_o}{\lambda^2 f^2} H^* G \exp(-j2\pi\alpha y_2).$$

The final lens L_3 of Fig. 8.10(a) optically Fourier transforms U_2 . Taking note of the reflected coordinate system in plane P_3 as well as the scaling constants present in the Fourier transform operation, the field strength in that plane is found to obey the proportionality

$$\begin{aligned}
 U_3(x_3, y_3) \propto & r_o^2 g(x_3, y_3) + \frac{1}{\lambda^2 f^2} \left[h(x_3, y_3) \otimes h^*(-x_3, -y_3) \otimes g(x_3, y_3) \right] \\
 & + \frac{r_o}{\lambda f} \left[h(x_3, y_3) \otimes g(x_3, y_3) \otimes \delta(x_3, y_3 + \alpha \lambda f) \right] \\
 & + \frac{r_o}{\lambda f} \left[h^*(-x_3, -y_3) \otimes g(x_3, y_3) \otimes \delta(x_3, y_3 - \alpha \lambda f) \right]. \quad (8-17)
 \end{aligned}$$

The third and fourth terms of this expression are of particular interest. Noting that

$$\begin{aligned}
 h(x_3, y_3) \otimes g(x_3, y_3) \otimes \delta(x_3, y_3 + \alpha \lambda f) \\
 = \iint_{-\infty}^{\infty} h(x_3 - \xi, y_3 + \alpha \lambda f - \eta) g(\xi, \eta) d\xi d\eta, \quad (8-18)
 \end{aligned}$$

we see that the third output term yields a *convolution* of h and g , centered at coordinates $(0, -\alpha hf)$ in the (x_3, y_3) plane. Similarly, the fourth term may be rewritten as

$$\begin{aligned}
 h^*(-x_3, -y_3) \otimes g(x_3, y_3) \otimes \delta(x_3, y_3 - \alpha \lambda f) \\
 = \iint_{-\infty}^{\infty} g(\xi, \eta) h^*(\xi - x_3, \eta - y_3 + \alpha \lambda f) d\xi d\eta, \quad (8-19)
 \end{aligned}$$

which is the *crosscorrelation* of g and h , centered at coordinates $(0, \alpha hf)$ in the (x_3, y_3) plane.

Note that the first and second terms of Eq. (8-17), which are of no particular utility in the usual filtering operations, are centered at the origin of the (x_3, y_3) plane. Thus it is clear that if the "carrier frequency" α is chosen sufficiently high, or equivalently if the reference wave is introduced at a sufficiently steep angle, the convolution and crosscorrelation terms will be deflected (in opposite directions) sufficiently far off-axis to be viewed independently. To find the convolution of h and g , the observer simply examines the distribution of light centered about the coordinates $(0, -\alpha hf)$. To find the crosscorrelation of h and g , the observation is centered at coordinates $(0, \alpha hf)$.

To illustrate the requirements placed on α more precisely, consider the widths of the various output terms illustrated in Fig. 8.15. If the maximum width of h in the y direction is W_h and that of g is W_g , then the widths of the various output terms are as follows:

1. $r_o^2 g(x_3, y_3) \rightarrow W_g$
2. $\frac{1}{\lambda^2 f^2} [h(x_3, y_3) \otimes h^*(-x_3, -y_3) \otimes g(x_3, y_3)] \rightarrow 2W_h + W_g$
3. $\frac{r_o}{\lambda f} [h(x_3, y_3) \otimes g(x_3, y_3) \otimes \delta(x_3, y_3 + \alpha \lambda f)] \rightarrow W_h + W_g$
4. $\frac{r_o}{\lambda f} [h^*(-x_3, -y_3) \otimes g(x_3, y_3) \otimes \delta(x_3, y_3 - \alpha \lambda f)] \rightarrow W_h + W_g.$

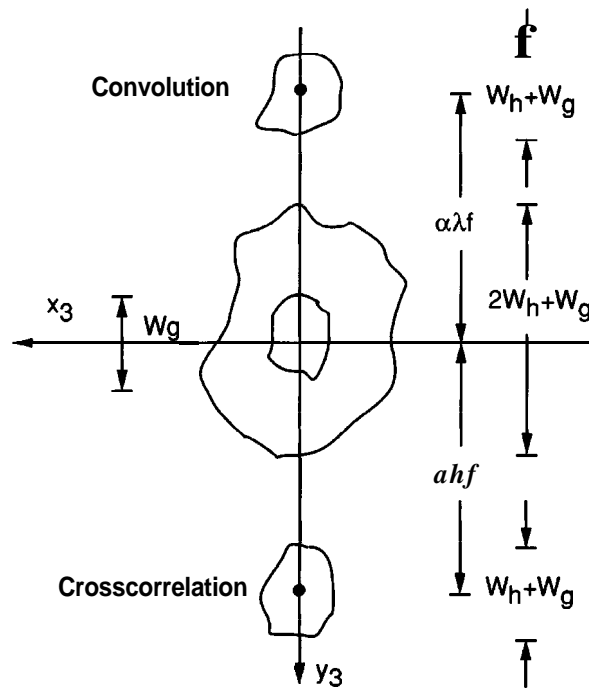


FIGURE 8.15 Locations of the various terms of the processor output.

From the figure it is clear that complete separation will be achieved if

$$\alpha > \frac{1}{\lambda f} \left(\frac{3W_h}{2} + W_g \right),$$

or equivalently, if

$$\theta > \frac{3}{2} \frac{W_h}{f} + \frac{W_g}{f}, \tag{8-20}$$

where the small-angle approximation $\sin \theta \approx \theta$ has been used.

8.4.3 Advantages of the VanderLugt Filter

The use of a VanderLugt filter removes the two most serious limitations to conventional coherent optical processors. First, when a specified impulse response is desired, the task of finding the associated transfer function is eliminated; the impulse response is Fourier transformed *optically* by the system that synthesizes the frequency-plane mask. Second, the generally complicated complex-valued transfer function is synthesized with a single *absorbing* mask; the phase transmittance through the frequency plane need no longer be controlled in a complicated manner. The absorbing mask is simply immersed in a liquid gate to eliminate all relative phase shifts.

The VanderLugt filter remains very sensitive to the exact position of the frequency-plane mask, but no more sensitive than the conventional coherent processor. The recording of the modulated high-frequency carrier requires a higher-resolution film than might

otherwise be used to synthesize the mask, but films with adequate resolution are readily available (e.g. Kodak Spectroscopic Plates) and this requirement poses no particular problem.

Note that the VanderLugt technique offers an important new flexibility to coherent processing. Whereas previously the realization of the frequency-plane mask was the major practical problem, the difficulties are now transferred back to the *space domain*. The difficulties are in general much less severe in the space domain, for the impulse responses required are often simple, and the necessary masks can be constructed by conventional photographic techniques. Thus the VanderLugt filter extends the use of coherent processors to an otherwise unattainable realm of operations. Many of the most promising applications fall in this realm.

8.5 THE JOINT TRANSFORM CORRELATOR

Before considering applications of coherent optical processing, an alternative method for performing complex filtering using a spatial carrier for encoding amplitude and phase information is considered. This method is due to Weaver and Goodman [295], and has become known as the *joint transform correlator*, although like the VanderLugt filter, it is equally capable of performing convolutions and correlations.

This type of filter differs from the VanderLugt filter in that *both* the desired impulse response *and* the data to be filtered are presented simultaneously during the recording process, rather than just presenting the desired impulse response. The transparency so constructed is then illuminated with a simple plane wave or spherical wave to obtain the filter outputs.

Consider the recording in Fig. 8.16(a). Lens L_1 collimates the light from the point source S. This collimated light then illuminates a pair of transparencies residing in the same plane, designated in the figure by their amplitude transmittances, h for the desired impulse response and g for the data to be filtered. For simplicity this input plane is taken to be the front focal plane of the Fourier transforming lens L_2 , **but** in fact this distance is arbitrary (vignetting will be eliminated if the inputs are placed in contact with lens, rather than in front of it). The Fourier transform of the composite input appears in the rear focal plane of L_2 , where the incident intensity is detected by either a photographic medium or a photosensitive spatial light modulator.

The field transmitted through the front focal plane is given by

$$U_1(x_1, y_1) = h(x_1, y_1 - Y/2) + g(x_1, y_1 + Y/2)$$

where the separation between the centers of the two inputs is Y . In the rear focal plane of the lens we find the Fourier transform of this field,

$$U_2(x_2, y_2) = \frac{1}{\lambda f} H\left(\frac{x_2}{\lambda f}, \frac{y_2}{\lambda f}\right) e^{-j\pi y_2 Y/\lambda f} + \frac{1}{\lambda f} G\left(\frac{x_2}{\lambda f}, \frac{y_2}{\lambda f}\right) e^{+j\pi y_2 Y/\lambda f}.$$

Taking the *squared magnitude* of this field, the *intensity incident on the recording plane* is found to be

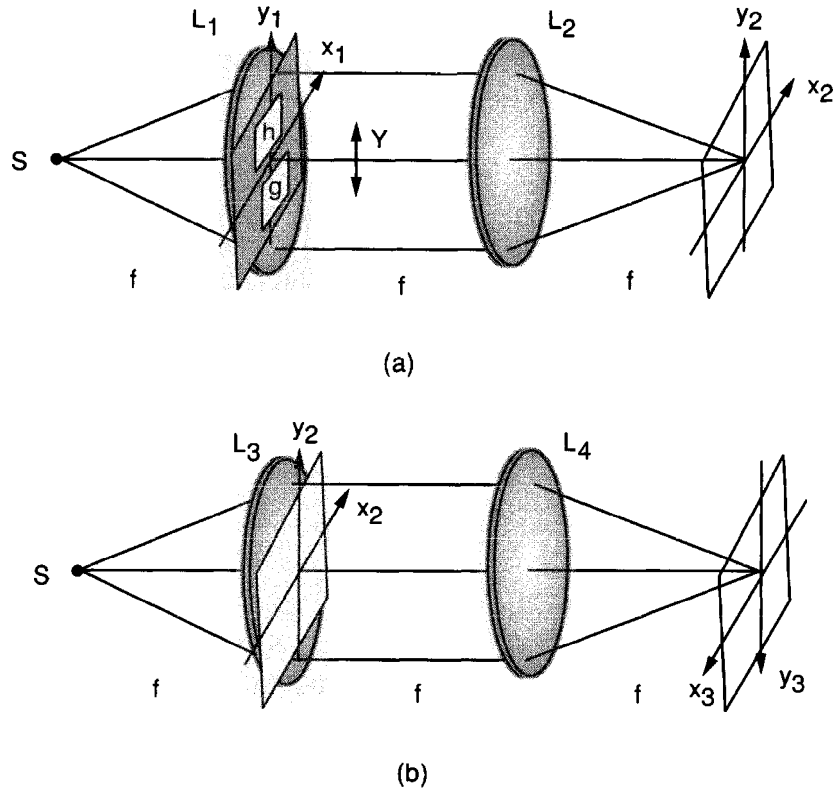


FIGURE 8.16
 The joint transform correlator. (a) Recording the filter, (b) obtaining the filtered output.

$$\begin{aligned}
 I(x_2, y_2) = & \frac{1}{\lambda^2 f^2} \left[\left| H\left(\frac{x_2}{\lambda f}, \frac{y_2}{\lambda f}\right) \right|^2 + \left| G\left(\frac{x_2}{\lambda f}, \frac{y_2}{\lambda f}\right) \right|^2 \right. \\
 & + H\left(\frac{x_2}{\lambda f}, \frac{y_2}{\lambda f}\right) G^*\left(\frac{x_2}{\lambda f}, \frac{y_2}{\lambda f}\right) e^{-j2\pi y_2 Y / \lambda f} \\
 & \left. + H^*\left(\frac{x_2}{\lambda f}, \frac{y_2}{\lambda f}\right) G\left(\frac{x_2}{\lambda f}, \frac{y_2}{\lambda f}\right) e^{+j2\pi y_2 Y / \lambda f} \right] \quad (8-21)
 \end{aligned}$$

The transparency that results from this recording is assumed to have an amplitude transmittance that is proportional to the intensity that exposed it. After processing, this transparency is illuminated by collimated light and the transmitted field is Fourier transformed by a positive lens L_4 , assumed to have the same focal length f as the lens used in the recording process (see Fig. 8.16(b)). The field in the front focal plane of this final Fourier transforming lens L_4 consists of four terms, each of which is proportional to one of the terms in Eq. (8-21). Taking account of scaling factors and coordinate inversions, the field in the rear focal plane of L_4 is

$$\begin{aligned}
U_3(x_3, y_3) = \frac{1}{\lambda f} & \left[h(x_3, y_3) \otimes h^*(-x_3, -y_3) + g(x_3, y_3) \otimes g^*(-x_3, -y_3) \right. \\
& + h(x_3, y_3) \otimes g^*(-x_3, -y_3) \otimes \delta(x_3, y_3 - Y) \\
& \left. + h^*(-x_3, -y_3) \otimes g(x_3, y_3) \otimes \delta(x_3, y_3 + Y) \right]. \quad (8-22)
\end{aligned}$$

Again it is the third and fourth terms of the expression for the output that are of most interest. We can rewrite them as

$$\begin{aligned}
h(x_3, y_3) \otimes g^*(-x_3, -y_3) \otimes \delta(x_3, y_3 - Y) \\
= \iint_{-\infty}^{\infty} h(\xi, \eta) g^*(\xi - x_3, \eta - y_3 + Y) d\xi d\eta \quad (8-23)
\end{aligned}$$

and

$$\begin{aligned}
h^*(-x_3, -y_3) \otimes g(x_3, y_3) \otimes \delta(x_3, y_3 + Y) \\
= \iint_{-\infty}^{\infty} g(\xi, \eta) h^*(\xi - x_3, \eta - y_3 - Y) d\xi d\eta. \quad (8-24)
\end{aligned}$$

Both of these expressions are crosscorrelations of the functions g and h . One output is centered at coordinates $(0, -Y)$ and the other at coordinates $(0, Y)$. The second output is a mirror reflection of the first about the optical axis.

To obtain a **convolution** of the functions h and g , it is necessary that one of them (and only one) be introduced in the processor of Fig. 8.16(a) with a mirror reflection about its own **origin**.⁷ For example, if originally we introduced the function $h(x_1, y_1 - Y/2)$, this input should be changed to $h(-x_1, -y_1 + Y/2)$, which is again centered at $Y/2$ but now is reflected about its own origin. The result will be two output terms, centered at $(0, Y)$ and $(0, -Y)$ in the output plane, each of which is a convolution of g and h . One term is identical with the other, but reflected about the optical axis.

Separation of the correlation (or convolution) terms from the uninteresting on-axis terms requires adequate separation of the two inputs at the start. If W_h represents the width of h and W_g is the width of g , both measured in the y direction, then separation of the desired terms can be shown to occur if

$$Y > \max \{W_h, W_g\} + \frac{W_g}{2} + \frac{W_h}{2}, \quad (8-25)$$

as is to be shown in Prob. 8-13.

The joint transform correlator is in some cases more convenient than the VanderLugt geometry, although both are widely used. Precise alignment of the filter transparency is required for the VanderLugt geometry, while no such alignment is necessary for the

⁷Strictly speaking, the function should also be conjugated, but in practice the functions g and h are usually real.

joint transform correlator. In addition, the joint transform approach has been found advantageous for real-time systems, i.e. systems that are required to rapidly change the filter impulse response. The price paid for the joint transform geometry is generally a reduction of the space-bandwidth product of the input transducer that can be devoted to the data to be filtered, since a portion of that space-bandwidth product must be assigned to the filter impulse response. See Ref. [201] for further comparison of the two approaches.

8.6 APPLICATION TO CHARACTER RECOGNITION

A particular application of optical information processing that has been of interest for many years is found in the field of character recognition. As we shall see, this application affords an excellent example of desired processing operations with simple impulse responses but not necessarily simple transfer functions. The carrier-frequency filter synthesis methods are therefore particularly well suited for this application.

8.6.1 The Matched Filter

The concept of the *matched filter* plays an important role in pattern recognition problems. By way of definition, a linear space-invariant filter is said to be matched to a particular signal $s(x, y)$ if its impulse response $h(x, y)$ is given by

$$h(x, y) = s^*(-x, -y). \quad (8-26)$$

If an input $g(x, y)$ is applied to a filter matched to $s(x, y)$, then the output $v(x, y)$ is found to be

$$\begin{aligned} v(x, y) &= \iint_{-\infty}^{\infty} h(x - \xi, y - \eta) g(\xi, \eta) d\xi d\eta \\ &= \iint_{-\infty}^{\infty} g(\xi, \eta) s^*(\xi - x, \eta - y) d\xi d\eta \end{aligned} \quad (8-27)$$

which is recognized to be the crosscorrelation function of g and s .

Historically the concept of the matched filter first arose in the field of signal detection; if a signal of known form, buried in "white" noise, is to be detected, then a matched filter provides the linear operation which maximizes the ratio of instantaneous signal power (at a particular time) to average noise power [286]. However, in the present application, the input patterns or characters will be assumed noiseless, and the use of a particular filtering operation must be justified on other grounds.

Considerable insight into the matched filtering operation is provided by an optical interpretation, as illustrated in Fig. 8.17. Suppose that a filter, matched to the input

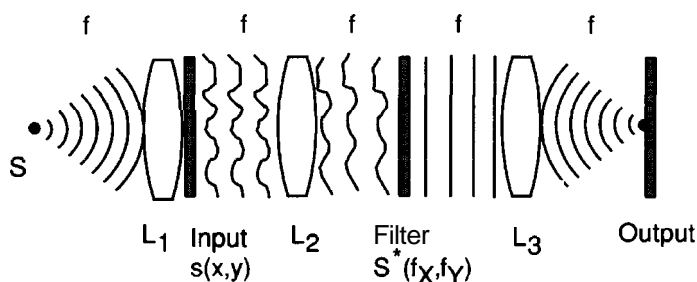


FIGURE 8.17
Optical interpretation of the
matched-filtering operation.

signal $s(x, y)$, is to be synthesized by means of a frequency-plane mask in the usual coherent processing geometry. Fourier transformation of the impulse response (8-26) shows that the required transfer function is

$$H(f_x, f_y) = S^*(f_x, f_y), \quad (8-28)$$

where $H = \mathcal{F}\{h\}$ and $S = \mathcal{F}\{s\}$. Thus the frequency plane filter should have an amplitude transmittance proportional to S^* .

Consider now the particular nature of the field distribution transmitted by the mask when the signal s (to which the filter is matched) is present at the input. Incident on the filter is a field distribution proportional to S , and transmitted by the filter is a field distribution proportional to SS^* . This latter quantity is entirely real, which implies that the frequency-plane filter exactly cancels all the curvature of the incident wavefront S . Thus the transmitted field consists of a plane wave (generally of nonuniform intensity), which is brought to a bright focus by the final transforming lens. When an input signal other than $s(x, y)$ is present, the wavefront curvature will in general not be canceled by the frequency-plane filter, and the transmitted light will not be brought to a bright focus by the final lens. Thus the presence of the signal s can conceivably be detected by measuring the intensity of the light at the focal point of the final transforming lens.

If the inputs is not centered on the origin, the bright point in the output plane simply shifts by a distance equal to the misregistration distance, a consequence of the space invariance of the matched filter (cf. Prob. 8-12).

8.6.2 A Character-Recognition Problem

Consider the following character-recognition problem: The input g to a processing system may consist of any one of N possible alphanumeric characters, represented by s_1, s_2, \dots, s_N , and the particular character present is to be determined by the processor. As will now be demonstrated, the identification process can be realized by applying the input to a bank of N filters, each matched to one of the possible input characters.

A block diagram of the recognition machine is shown in Fig. 8.18. The input is simultaneously (or sequentially) applied to the N matched filters with transfer functions $S_1^*, S_2^*, \dots, S_N^*$. The response of each filter is normalized by the square root of the total energy in the character to which it is matched. This normalization, which can be accomplished electronically after detection of the filter outputs, takes account of the fact that the various input characters will generally not be of equal energy. Finally, the

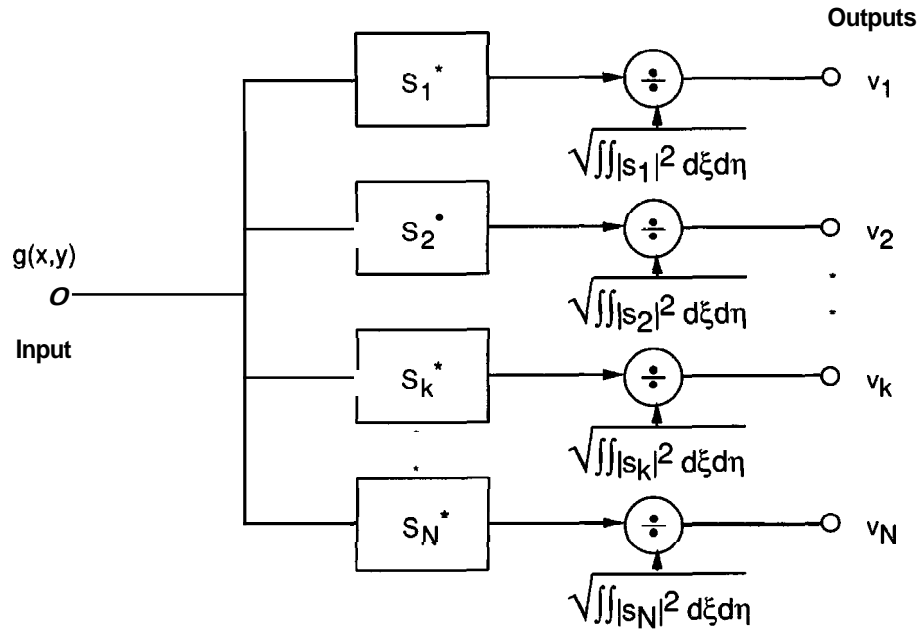


FIGURE 8.18
Block diagram of a character-recognition system.

squared moduli of the outputs $|v_1|^2, |v_2|^2, \dots, |v_N|^2$ are compared at the particular points where their maximum outputs would be anticipated (assuming that the character to which they are matched is present in each case). As will now be demonstrated, if the particular character

$$g(x, y) = s_k(x, y)$$

is actually present at the input, then the particular output $|v_k|^2$ will be the largest of the N responses.

To prove this assertion, first note that, from Eq. (8-27), the peak output $|v_k|^2$ of the correct matched filter is given by

$$|v_k|^2 = \frac{\left[\iint_{-\infty}^{\infty} |s_k|^2 d\xi d\eta \right]^2}{\iint_{-\infty}^{\infty} |s_k|^2 d\xi d\eta} = \iint_{-\infty}^{\infty} |s_k|^2 d\xi d\eta. \quad (8-29)$$

On the other hand, the response $|v_n|^2$ ($n \neq k$) of an incorrect matched filter is given by

$$|v_n|^2 = \frac{\left[\iint_{-\infty}^{\infty} s_k s_n^* d\xi d\eta \right]^2}{\iint_{-\infty}^{\infty} |s_n|^2 d\xi d\eta} \quad (8-30)$$

However, from Schwarz's inequality, we have

$$\left| \iint_{-\infty}^{\infty} s_k s_n^* d\xi d\eta \right|^2 \leq \iint_{-\infty}^{\infty} |s_k|^2 d\xi d\eta \iint_{-\infty}^{\infty} |s_n|^2 d\xi d\eta.$$

It follows directly that

$$|v_n|^2 \leq \iint_{-\infty}^{\infty} |s_k|^2 d\xi d\eta = |v_k|^2, \quad (8-31)$$

with equality if and only if

$$s_n(x, y) = \kappa s_k(x, y).$$

From this result it is evident that the matched filter does provide *one* means of recognizing which character, of a set of possible characters, is actually being presented to the system. It should be emphasized that this capability is not unique to the matched filter. In fact it is often possible to modify (mismatch) all the filters in such a way that the discrimination between characters is improved. Examples of such modifications include: (1) overexposing the low-frequency portion of a VanderLugt filter transparency so as to suppress the influence of those frequencies in the decision process (see, for example, Ref. [284], pp. 130–133); (2) eliminating the amplitude portion of the transfer functions of the matched filters and retaining only phase information [149]; and (3) modifying the nonlinearity of the normally square-law detection process in the joint-transform correlator to enhance discrimination between patterns [154], [155].

Not all pattern-recognition problems are of the type described above. For example, rather than trying to distinguish between several possible known patterns, we may wish simply to detect the presence or absence of a single known object in a larger image. Such a problem is closer to what the matched filter is known to do well, namely detect a known pattern in the presence of background noise, but has the added difficulty that the orientation and possibly the scale size of the target may not be under the same level of control that is present in the character recognition problem. We return in a later subsection to discussing some of the difficulties of the matched filter approach to such problems.

8.6.3 Optical Synthesis of a Character-Recognition Machine

The matched filter operation can readily be synthesized by means of either the VanderLugt technique or the joint transform technique discussed earlier. Our discussion here is directed at the VanderLugt-type system, but the reader may wish to contemplate how the equivalent system could be realized with the joint transform geometry.

Recall that one of the outputs of the VanderLugt filtering operation is itself the crosscorrelation of the input pattern with the original pattern from which the filter was synthesized. By restricting attention to the proper region of the output space, the matched filter output is readily observed.

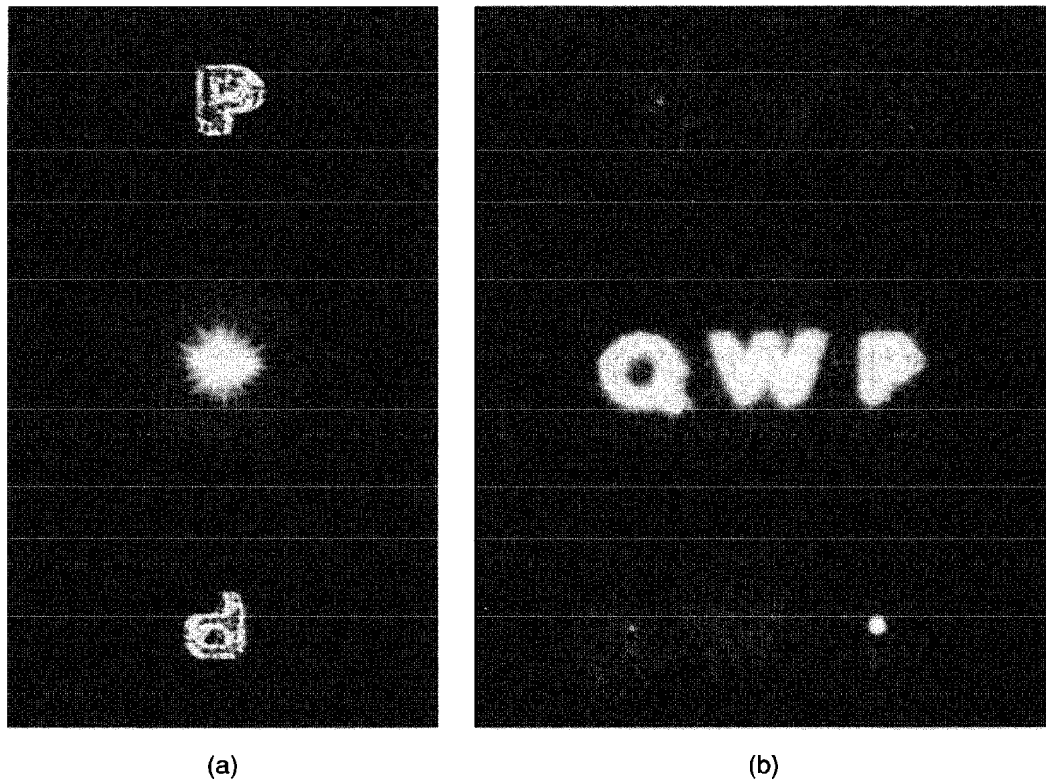


FIGURE 8.19 Photographs of (a) the impulse response of a VanderLugt filter, and (b) the response of the matched filter portion of the output to the letters Q, W, and P.

Figure 8.19(a) shows a photograph of the impulse response of a VanderLugt filter which has been synthesized for the character P. The upper portion of response will generate the convolution of the input data with the symbol P, while the lower response will generate the crosscorrelation of the input with the letter P. The central portion of the response is undesired and not of interest.

Figure 8.19(b) shows the response of the matched filter portion of the output to the letters Q, W, and P. Note the presence of the bright point of light in the response to P, indicating the high correlation between the input letter and the letter to which it is matched.

To realize the entire *bank* of matched filters illustrated in Fig. 8.18, it would be possible to synthesize N separate VanderLugt filters, applying the input to each filter sequentially. Alternatively, if N is not too large, it is possible to synthesize the entire bank of filters on a single frequency-plane filter. This can be done by frequency-multiplexing, or recording the various frequency-plane filters with different carrier frequencies on a single transparency. Figure 8.20(a) illustrates one way of recording the multiplexed filter. The letters Q, W, and P are at different angles with respect to the reference point, and as a consequence, the crosscorrelations of Q, W, and P with the input character appear at different distances from the origin, as illustrated in Fig. 8.20(b).

The number of different filters that can be realized by this technique is limited by the dynamic range that can be achieved in the frequency-plane filter. Synthesis of nine separate impulse responses in a single mask was demonstrated by VanderLugt at an early date (see Ref. [284], pp. 133–139).

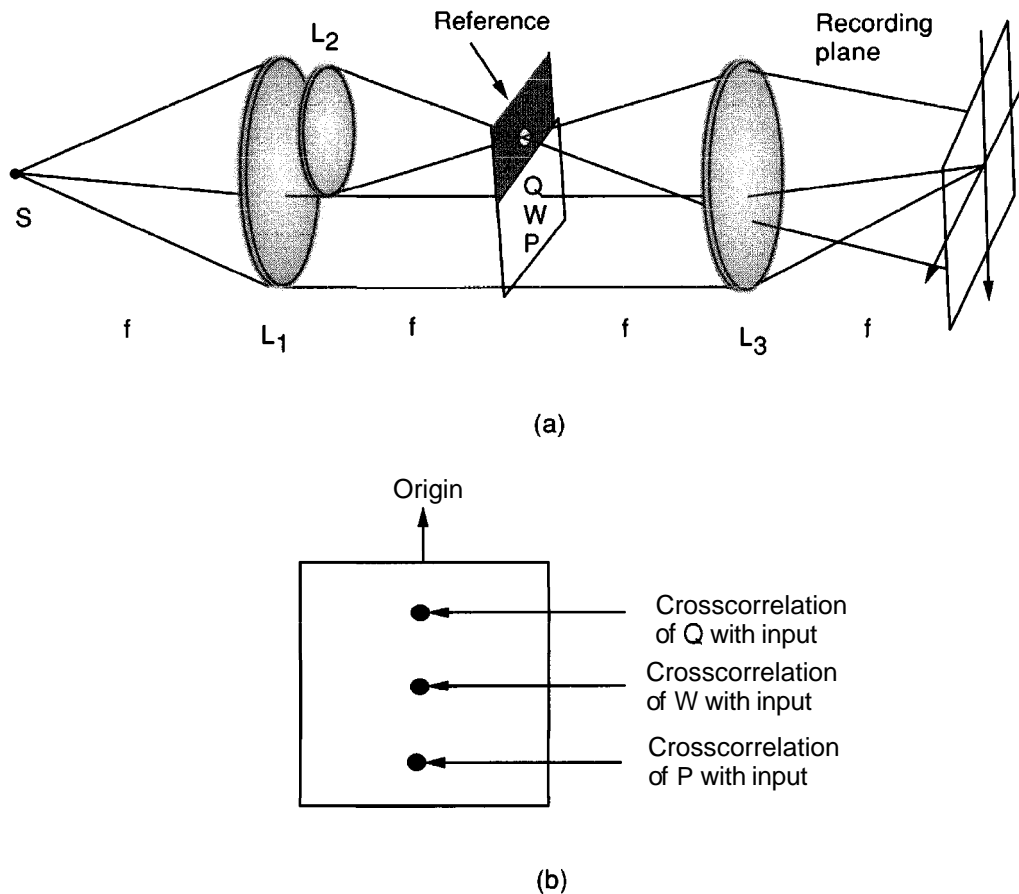


FIGURE 8.20
Synthesis of a bank of matched filters with a single frequency-plane filter. (a) Recording the frequency-plane filter; (b) format of the matched filter portion of the output.

8.6.4 Sensitivity to Scale Size and Rotation

The coherent optical pattern-recognition technique described above suffers from certain deficiencies that are shared by all matched-filter approaches to the pattern recognition problem. Specifically, such filters are too sensitive to scale size changes and rotations of input patterns. When an input pattern is presented with an angular orientation or a scale size that is different from those of the pattern to which the filter is matched, the response of the correct matched filter is reduced, and errors arise in the pattern recognition process. The degree of sensitivity of a matched filter to rotation and scale-size depends to a large extent on the structure of the pattern to which it is matched. For example, a matched filter for the character **L** is obviously much more rotation-sensitive than that for the letter **O**. One solution that has been used is to make a bank of matched filters, each of which is matched to the pattern of interest with a different rotation **and/or** scale size. If any of these matched filters have a large output, then the pattern of interest is known to have been presented to the input.

We turn attention in the next section to a few of the many other techniques that have been explored as possible solutions to scale-size and rotation variations of the object.

8.7

OPTICAL APPROACHES TO INVARIANT PATTERN RECOGNITION

There exists a vast number of different pattern-recognition approaches that are aimed at reducing or eliminating sensitivity to extraneous parameters, such as scale size, and rotation. Note that the classical matched filtering approach is insensitive to one extraneous parameter, translation, in the sense that wherever the object may be in the input field, a bright spot will appear at its corresponding location in the output plane. The strength of that spot is not affected by pure translation. This is a property that one would like to preserve in any approach to reducing sensitivity to other types of object variation.

One approach to handling patterns with different scale sizes and rotations is to synthesize a matched filter for an object of fixed size and rotation, and to perform a mechanical search, rotating and **magnifying/demagnifying** the input to the system. Mechanical searches are awkward and time-consuming, and therefore are not considered further here.

Note that many methods for achieving invariance have been developed for digital processing of images, and only some of these are naturally well suited for optical implementation. In some cases there must be a heavy intrusion of electronic digital processing, a process that can slow down the natural speed (obtained from parallelism) of a purely optical solution. In the end, the complexity of digital and optical solutions to a given approach must be carefully and critically assessed to determine whether there is really a practical motivation to pursue the optical solution.

In what follows we will only briefly touch on three different approaches to invariant pattern recognition that have received considerable attention in the literature. Space limitations do not allow a complete discussion of these methods, so our goal is to introduce the reader to the basic underlying ideas.

8.7.1 Mellin Correlators

While Fourier-based correlators such as discussed above are extremely sensitive to both magnification and rotation of the object, there exists a different transform, closely related to the Fourier transform, that exhibits a certain invariance to object magnification. We refer to the Mellin transform (see [32], p. 254). For simplicity we introduce the Mellin transform in one-dimensional form, although it is easily generalized to two dimensions.

The Mellin transform of a function $g(\xi)$ is defined by

$$M(s) = \int_0^{\infty} g(\xi) \xi^{s-1} d\xi, \quad (8-32)$$

where in the most general case, s is a complex variable. A simple relation between the Fourier transform and the Mellin transform can be discovered if the complex variable s is restricted to the imaginary axis, i.e. $s = j2\pi f$. A substitution of variables $\xi = e^{-x}$ yields the following expression for the Mellin transform of g ,

$$M(j2\pi f) = \int_{-\infty}^{\infty} g(e^{-x}) e^{-j2\pi f x} dx. \quad (8-33)$$

which we see is nothing but the Fourier transform of the function $g(e^{-x})$. We conclude that it is possible to perform a Mellin transform with an optical Fourier transforming system provided the input is introduced in a "stretched" coordinate system, in which the natural space variable is logarithmically stretched ($x = -\ln \xi$). Such a stretch can be introduced, for example, by driving the deflection voltage of a cathode ray tube through a logarithmic amplifier and writing onto an SLM with the resulting stretched signal. There also exist optical methods for coordinate distortion that can be used for this task [42].

The particular interest in the Mellin transform arises because its magnitude is independent of scale-size changes in the input. To prove this fact, let M_1 represent the Mellin transform of $g(\xi)$ and let M_a represent the Mellin transform of $g(a\xi)$, where $0 < a < \infty$. A value of a greater than unity implies a demagnification of g and a value of a between zero and unity implies a magnification. The Mellin transform of $g(a\xi)$ can now be found as follows:

$$\begin{aligned} M_a(j2\pi f) &= \int_0^{\infty} g(a\xi) \xi^{j2\pi f-1} d\xi = \int_0^{\infty} g(\xi') \left(\frac{\xi'}{a}\right)^{j2\pi f-1} \frac{d\xi'}{a} \\ &= a^{-j2\pi f} \int_0^{\infty} g(\xi') \xi'^{j2\pi f-1} d\xi', \end{aligned} \quad (8-34)$$

where the change of variables $\xi' = a\xi$ was made. Taking the magnitude of M_a and noting that the term $|a^{-j2\pi f}| = 1$ proves that $|M_a|$ is independent of scale size a . The independence of the Mellin magnitude with respect to object scale size, coupled with the fact that the Mellin transform can be performed as a Fourier transform of a stretched input, will be shown to suggest one way to achieve independence from scale size.

The second parameter we wish to eliminate is rotation of the object. As the basis for this elimination, we note that rotation of an object by a certain angle is equivalent to translation in one dimension if the object is presented in polar coordinates, provided that the center chosen for the polar coordinate system coincides with the center of rotation of the object. There is a subtlety that arises from the fact that on a scale that varies from 0 to 2π radians, rotation by angle θ may result in portions of the object shifting by θ , while other portions of the object may "wrap around the angular coordinate and appear at a position corresponding to $2\pi - \theta$ ". This problem can be removed if the angle coordinate is allowed to cover two or more periods of length 2π , in which case the "wrap around" problem can be minimized or eliminated.

The approach to achieving simultaneous scale and rotation invariance, which was pioneered by Casasent and Psaltis (see [50] and references therein), can now be described. A two dimensional object $g(\xi, \eta)$ is entered into the optical system in a distorted polar coordinate system, the distortion arising from the fact that the radial coordinate is stretched by a logarithmic transformation. The optical system that follows is a matched filtering system, for which the filters have been made under the same coordinate transformations to which the input was subjected. The output intensity will translate with

rotation, but will not drop in strength under either scale-size or rotational changes of the input.

Unfortunately the achievement of scale-size and rotation invariance by the method described above leads to a loss of the original translation invariance that characterized the conventional matched filter. To achieve simultaneous invariance to all three parameters, it is possible to replace the input function g by the magnitude of its Fourier transform $|G|$, which is invariant to translation of g [50]. This magnitude function is subjected to the same coordinate changes discussed above, and the filter should be matched to the magnitude of the Fourier transform of the pattern of interest, again subject to the coordinate transformations outlined. Some loss of correlator performance can be expected due to the fact that the phase of the Fourier transform of the input and the matched filter transfer function have both been discarded.

8.7.2 Circular Harmonic Correlation

An approach that focuses on the problem of invariance to object rotation is based on a circular harmonic decomposition of the object [151], [150]. For an excellent overview, see Ref. [9].

The circular harmonic expansion rests on the fact that a general two-dimensional function $g(r, \theta)$, expressed in polar coordinates, is periodic in the variable θ , with period 2π . As a consequence it is possible to express g in a Fourier series in the angular variable,

$$g(r, \theta) = \sum_{m=-\infty}^m g_m(r) e^{jm\theta}, \quad (8-35)$$

where the Fourier coefficients are functions of radius,

$$g_m(r) = \frac{1}{2\pi} \int_0^{2\pi} g(r, \theta) e^{-jm\theta} d\theta. \quad (8-36)$$

Each term in Eq. (8-35) is referred to as a "circular harmonic" of the function g . Note that if the function $g(r, \theta)$ undergoes a rotation by angle α to produce $g(r, \theta - \alpha)$, the circular harmonic expansion becomes

$$g(r, \theta - \alpha) = \sum_{m=-\infty}^{\infty} g_m(r) e^{-jm\alpha} e^{jm\theta} \quad (8-37)$$

and thus the m th circular harmonic is subjected to a phase change of $-m\alpha$ radians.

Consider now the crosscorrelation of the functions g and h , which in rectangular coordinates is written

$$R(x, y) = \iint_{-\infty}^{\infty} g(\xi, \eta) h^*(\xi - x, \eta - y) d\xi d\eta. \quad (8-38)$$

Of particular interest is the value of the crosscorrelation at the origin, which in rectangular and polar coordinates can be written

$$R_0 = R(0, 0) = \iint_{-\infty}^{\infty} g(\xi, \eta) h^*(\xi, \eta) d\xi d\eta = \int_0^{\infty} r dr \int_0^{2\pi} g(r, \theta) h^*(r, \theta) d\theta. \quad (8-39)$$

The particular case of the crosscorrelation between the function $g(r, \delta)$ and an angularly rotated version of the same function, $g(r, \delta - \alpha)$, yields

$$R_\alpha = \int_0^{\infty} r dr \int_0^{2\pi} g^*(r, \theta) g(r, \theta - \alpha) d\theta \quad (8-40)$$

which, when the function $g^*(r, \delta)$ is expanded in a circular harmonic expansion, is equivalently expressed as

$$R_\alpha = \int_0^{\infty} r \left[\sum_{m=-\infty}^{\infty} g_m^*(r) \int_0^{2\pi} g(r, \theta - \alpha) e^{-jm\theta} d\theta \right] dr. \quad (8-41)$$

But

$$\frac{1}{2\pi} \int_0^{2\pi} g(r, \theta - \alpha) e^{-jm\theta} d\theta = g_m(r) e^{-jm\alpha},$$

and therefore

$$R_\alpha = 2\pi \sum_{m=-\infty}^{\infty} e^{-jm\alpha} \int_0^{\infty} r |g_m(r)|^2 dr. \quad (8-42)$$

From this result we see that each of the circular harmonic components of the crosscorrelation undergoes a *different* phase shift $-m\alpha$.

If a particular circular harmonic component of R_α , say the M th, is extracted digitally, then from the phase associated with that component it is possible to determine the angular shift that one version of the object has undergone. Of more relevance to us here, if an optical filter that is matched to the M th circular harmonic component of a particular object is constructed, perhaps using digital techniques, and placed in an optical correlation system, then if that same object is entered as an input to the system with any angular rotation, a correlation peak of strength proportional to $\int_0^{\infty} r |g_M(r)|^2 dr$ will be produced, independent of rotation. Hence an optical correlator can be constructed that will recognize that object independent of rotation.

The price paid for rotation invariance is that the strength of the correlation peak is smaller than what would be obtained for the crosscorrelation with an unrotated version of the object when all of the circular harmonics are used simultaneously. The reduction in peak correlation intensity incurred by use of only the M th circular harmonic component is easily shown to be given by

$$\kappa_M = \frac{\int_0^{\infty} r |g_M(r)|^2 dr}{\sum_{m=-\infty}^{\infty} \int_0^{\infty} r |g_m(r)|^2 dr}. \quad (8-43)$$

It should be mentioned that the circular harmonic expansion of a function depends on the particular point chosen for the center of that expansion, and the quality of the correlation peaks obtained depends on making a "good choice of center. This problem

has been addressed and procedures for determining an appropriate center have been found (see Ref. [9]).

8.7.3 Synthetic Discriminant Functions

The final method we will discuss for achieving invariant pattern recognition is through the use of what are known as "synthetic discriminant functions" (SDF). This method has its roots in a number of earlier ideas, particularly in the work of Braunecker and Lohmann [36], [35] and in the ideas of Caulfield and Haimes [54]. However, it has been carried to its present state by D. Casasent and his students (e.g. see [48] and [49]).

The SDF approach is a method for constructing a single pattern-recognition filter that has its correlation properties tailored in advance by means of a certain "training set" of images, whose desired correlations with the reference filter are known in advance. The members of the training set may be distorted versions of a single object, where the distortions correspond to scale change and rotation, they may be more generally distorted versions of that object, or they may be examples of other objects for which we desire to have zero filter output. Let the training set of N images be represented by $\{g_n(x, y)\}$ where $n = 1, 2, \dots, N$. For a particular training image, we may want the correlation with our filter's impulse response $h(x, y)$ to be unity (i.e. that particular training image is a distorted version of the ideal image), and in some cases we may wish it be zero (i.e. that particular training image represents a distorted version of an entirely different ideal image). We will divide the set $\{g_n\}$ into two subsets, $\{g_n^+\}$ for which we wish the correlation to be unity, and $\{g_n^-\}$ for which we wish the correlation to be zero. Thus we have the constraints

$$\iint_{-\infty}^{\infty} g_n^+(x, y) h(x, y) dx dy = 1 \quad \iint_{-\infty}^{\infty} g_n^-(x, y) h(x, y) dx dy = 0. \quad (8-44)$$

To obtain a filter impulse $h(x, y)$ that will have the desired correlations with the training set, we first expand (symbolically) that impulse response in a series using the training images as basis functions,

$$h(x, y) = \sum_{n=1}^N a_n g_n(x, y), \quad (8-45)$$

where the a_n are for the moment unknown. Now consider the correlation of any one member of the training set, say $g_k(x, y)$ with the filter function $h(x, y)$,

$$c_k = \iint_{-\infty}^{\infty} g_k^*(x, y) h(x, y) dx dy = \sum_{n=1}^N a_n \iint_{-\infty}^{\infty} g_k^*(x, y) g_n(x, y) dx dy, \quad (8-46)$$

where we have substituted the previous expansion for $h(x, y)$, and c_k is known to be either zero or unity, depending on which class of inputs g_k is drawn from. Letting p_{kn} represent the correlation between g_k and g_n , we see that

$$c_k = \sum_{n=1}^N a_n p_{kn}. \quad (8-47)$$

Now by considering all N members of the training set, we establish a total set of N linear equations in the N unknowns a_n , each similar to Eq. (8-47), but for a different value of k . The entire collection of these equations can be expressed in a single matrix equation

$$\mathbf{P} \vec{a} = \vec{c}, \quad (8-48)$$

where \vec{a} and \vec{c} are column vectors of length N , and \mathbf{P} is an $N \times N$ matrix of correlations between the training images,

$$\vec{a} = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_N \end{bmatrix} \quad \vec{c} = \begin{bmatrix} c_1 \\ c_2 \\ \vdots \\ c_N \end{bmatrix} \quad \mathbf{P} = \begin{bmatrix} p_{11} & \cdots & p_{1N} \\ p_{21} & \cdots & p_{2N} \\ \vdots & \vdots & \vdots \\ p_{N1} & \cdots & p_{NN} \end{bmatrix}. \quad (8-49)$$

Note that the vector \vec{c} is a column vector of known values (each element is zero or one in the case we are considering, but clearly they can also have more general values), the matrix \mathbf{P} contains known elements (calculated in advance), and we seek knowledge of the vector \vec{a} , for this will allow us to specify the desired impulse response of our filter according to Eq. (8-45). This unknown vector can be found by inverting the matrix \mathbf{P} and multiplying the inverse by the vector \vec{c} (using a digital computer for these calculations)

$$\vec{a} = \mathbf{P}^{-1} \vec{c}. \quad (8-50)$$

Thus we have described a method for constructing a filter which will produce prescribed correlations between a group of images in a training set. The theory does not directly indicate what the response might be to an image that is not a member of the training set, but the method still provides a useful design procedure to obtain filters with a significant degree of invariance to various image parameters. The theory presented here can clearly be generalized in many ways, but space constraints prevent us from delving further into this subject.

8.8 IMAGE RESTORATION

A common problem in image processing, and one that has been studied extensively in the context of optical information processing, is image restoration, by which we mean the restoration of an image that has been blurred by a known linear, invariant point-spread function. In this section we summarize some of the past work on this problem. The reason for doing so is only partly because of the extensive past work. Equally important, there are lessons that have been learned in this application, particularly about clever use of the properties of wavefront modulation devices, that can be applied to other unrelated problems.

8.8.1 The Inverse Filter

Let $o(x, y)$ represent the intensity distribution associated with an incoherent object, and let $i(x, y)$ represent the intensity distribution associated with a blurred image of that object. For simplicity we assume that the magnification of the imaging system is unity and we define the image coordinates in such a way as to remove any effects of image inversion.

We assume that the blur the image has been subjected to is a linear, space-invariant **transformation**, describable by a known space-invariant point-spread function $s(x, y)$. Thus, in the simplest description of the problem, the object and image are related by

$$i(x, y) = \iint_{-\infty}^{\infty} o(\xi, \eta) s(x - \xi, y - \eta) d\xi d\eta. \quad (8-51)$$

We seek to obtain an estimate $\hat{o}(x, y)$ of $o(x, y)$, based on the measured image intensity $i(x, y)$ and the known point-spread function $s(x, y)$. In other words, we wish to invert the blurring operation and recover the original object.

An unsophisticated solution to this problem is quite straightforward. Given the relationship between object and image in the frequency domain,

$$\mathcal{F}\{i(x, y)\} = \mathcal{F}\{s(x, y) \otimes o(x, y)\} = S(f_x, f_y) O(f_x, f_y), \quad (8-52)$$

it seems obvious that the spectrum of the original object can be obtained by simply dividing the image spectrum by the known OTF of the imaging system,

$$\hat{O}(f_x, f_y) = \frac{I(f_x, f_y)}{S(f_x, f_y)}. \quad (8-53)$$

An equivalent statement of this solution is that we should pass the detected image $i(x, y)$ through a linear space-invariant filter with transfer function

$$H(f_x, f_y) = \frac{1}{S(f_x, f_y)}. \quad (8-54)$$

Such a filter is commonly referred to as an "inverse filter", for obvious reasons.

This straightforward solution has several serious defects:

1. Diffraction limits the set of frequencies over which the transfer function $S(f_x, f_y)$ is nonzero to a finite range. Outside this range, $S = 0$ and its inverse is ill defined. For this reason, it is necessary to limit the application of the inverse filter to those frequencies lying within the diffraction-limited passband.
2. Within the range of frequencies for which the diffraction-limited transfer function is nonzero, it is possible (indeed likely) that transfer function S will have isolated zeros. Such is the case for both a serious defocusing error and for many kinds of motion blur (see Prob. 8-14). The value of the restoration filter is undefined at the frequencies where these isolated zeros occur. Another way of stating this problem is that the restoration filter would need a transfer function with infinite dynamic range in order to properly compensate the spectrum of the image.

3. The inverse filter takes no account of the fact that there is inevitably noise present in the detected image, along with the desired signal. The inverse filter boosts the most those frequency components that have the worst signal-to-noise ratios, with the result that the recovered image is usually dominated by noise.

The only solution to the last of the problems raised above is to adopt a new approach to determining the desired restoration filter, an approach that includes the effects of noise. One such approach is described in the following, and it will be seen to solve the first two problems as well.

8.8.2 The Wiener Filter, or the Least-Mean-Square-Error Filter

A new model for the imaging process is now adopted, one that takes into account explicitly the presence of noise. The detected image is now represented by

$$i(x, y) = o(x, y) \otimes s(x, y) + n(x, y), \quad (8-55)$$

where $n(x, y)$ is the noise associated with the detection process. In addition to the presence of the noise term, which must be regarded as a random process, we also treat the object $o(x, y)$ as a random process in this formulation (if we knew what the object is, we would have no need to form an image of it, so the object that is present is regarded as one realization of a random process). We assume that the power spectral densities⁸ (i.e. the distributions of average power over frequency) of the object and the noise are known, and are represented by $\Phi_o(f_x, f_y)$ and $\Phi_n(f_x, f_y)$. Finally, the goal is to produce a linear restoration filter that minimizes the mean-square difference between the true object $o(x, y)$ and the estimate of the object $\hat{o}(x, y)$, i.e. to minimize

$$\epsilon^2 = \text{Average} \left[|o - \hat{o}|^2 \right]. \quad (8-56)$$

The derivation of the optimum filter would take us too far afield, so we content ourselves with presenting the result and referring the reader to another source [119]. The transfer function of the optimum restoration filter is given by

$$H(f_x, f_y) = \frac{S^*(f_x, f_y)}{|S(f_x, f_y)|^2 + \frac{\Phi_n(f_x, f_y)}{\Phi_o(f_x, f_y)}}. \quad (8-57)$$

This type of filter is often referred to as a *Wiener filter*, after its inventor, Norbert Wiener.

Note that at frequencies where the signal-to-noise ratio is high ($\Phi_n/\Phi_o \ll 1$), the optimum filter reduces to an inverse filter,

$$H \approx \frac{S^*}{|S|^2} = \frac{1}{S},$$

while at frequencies where the signal-to-noise ratio is low ($\Phi_n/\Phi_o \gg 1$), it reduces to a strongly attenuating matched filter,

⁸For a detailed discussion of the concept of power spectral density, see [123], Section 3.3.

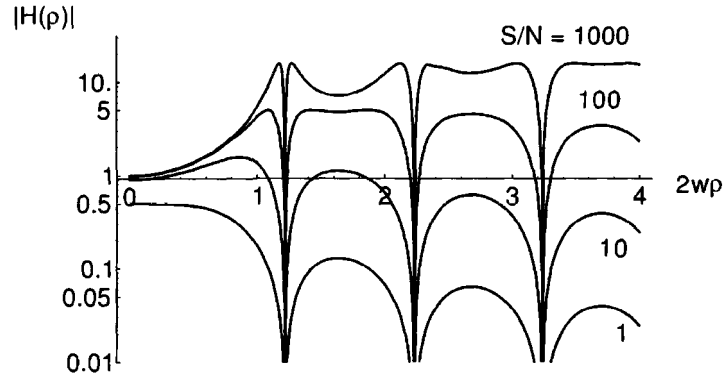


FIGURE 8.21

Magnitudes of the transfer function of a Wiener filter. The image is assumed to have been blurred by a point-spread function consisting of a circular disk of radius w . The signal-to-noise ratio is varied from 1000 to 1. The phase of the filter changes between 0 and π radians between alternate zeros of this transfer function.

$$H \approx \frac{\Phi_o}{\Phi_n} S^*$$

Figure 8.21 shows plots of the magnitude of the transfer function of the restoration filter under the assumption of a severe focusing error and white (i.e. flat) power spectra for the signal and the noise. Several different signal-to-noise ratios are represented. Note that at high signal-to-noise ratio, the Wiener filter reduces the relative strength of the low frequencies and boosts the relative strength of the high frequencies. At low signal-to-noise ratio, all frequencies are reduced.

Note that at frequencies outside the diffraction-limited **passband** of the imaging system, no object information is present, and therefore the noise-to-signal ratio is infinite. Hence the Wiener filter makes no attempt to restore object frequency components that are simply not present in the image, a very sensible strategy.

8.8.3 Filter Realization

Many methods exist for optically realizing inverse and Wiener restoration filters. We discuss only two such methods, one relatively obvious, the other not at all obvious. Both depend on the use of VanderLugt-type filters. In both cases we suppose that there is available a transparency that has recorded the known impulse response of the blurred system. This transparency could have been obtained by imaging a point source through the blurred system, or could have been generated by computer. We also assume that this transparency has been made in such a way that its amplitude transmittance, t_A , is proportional to $s(x, y)$.

Inverse filter

The first method is one that attempts to realize an inverse filter [275]. Using the recording of the blur, we record two transparencies which will be sandwiched (i.e. placed in close contact) to form the frequency plane filter. Referring back to Fig. 8.20(a), one component of the filter is of the VanderLugt type, recorded interferometrically as shown, but with an input that consists only of the known blur function s . This filter captures both the amplitude and phase associated with the transfer function of the blur, S . A second transparency is recorded in the same geometry, but with the reference point source blocked, thus capturing information only about the intensity $|S|^2$.

The transmittance of the VanderLugt filter consists of four terms, as before, and only one of these is of interest in this problem. We again focus on the term proportional to S^* , the same term that was of interest in the case of the matched filter. With exposure in the linear region of the t_A vs. E curve and with proper processing, this component of amplitude transmittance can be written

$$t_{A1} \propto S^*(f_X, f_Y).$$

The second transparency is exposed in the linear region of the H&D curve and processed with a photographic γ equal to 2. The result is an amplitude transmittance

$$t_{A2} \propto \frac{1}{|S(f_X, f_Y)|^2}.$$

When these two transparencies are placed in close contact, the amplitude transmittance of the pair is

$$t_A = t_{A1} t_{A2} = \frac{S^*(f_X, f_Y)}{|S(f_X, f_Y)|^2} = \frac{1}{S(f_X, f_Y)},$$

which is the transfer function of an inverse filter.

In addition to all the difficulties associated with an inverse filter that were mentioned earlier, this method suffers from other problems related to the photographic medium. The dynamic range of amplitude transmittance over which this filter can function properly is quite limited. The problem is evident if we consider only the second filter, which was recorded in the linear region of the H&D curve. If we wish this filter to behave as desired over a 10 : 1 dynamic range of $|S|$, this requires proper behavior over a 100 : 1 range of $1/|S|^2$. But since the amplitude transmittance of this filter is proportional to $1/|S|^2$, the intensity transmittance is proportional to $1/|S|^4$, and a 10 : 1 change of S implies a 10,000 : 1 change of intensity transmittance. To properly control this filter over the range of interest would require controlling the density accurately over a range of 0 to 4. Densities as high as 4 can seldom be achieved in practice, and even a density of 3 requires some special effort. For this reason, the dynamic range of $|S|$ over which the filter functions properly is severely limited in practice.

Wiener filter

A superior approach to realizing an image restoration filter is one that generates a Wiener filter, and does so with considerably more dynamic range than the previous

method afforded. Such a method was introduced by Ragnarsson [239]. There are several novel aspects to this approach to filter realization:

1. Diffraction, rather than absorption, is used to attenuate frequency components.
2. Only a single interferometrically generated filter is required, albeit one with an unusual set of recording parameters.
3. The filter is bleached and therefore introduces only phase shifts in the transmitted light.

Certain postulates underlie this method of recording a filter. First, it is assumed that the maximum phase shift introduced by the filter is much smaller than 2π radians, and therefore

$$t_A = e^{j\phi} \approx 1 + j\phi.$$

In addition, it is assumed that the phase shift of the transparency after bleaching is linearly proportional to the silver density present before bleaching,

$$\phi \propto D.$$

This assumption is true to a very good approximation if a nontanning bleach is used, for such a bleach returns metallic silver to a transparent silver salt, and the density of that transparent material determines the phase shift introduced by the bleached transparency. Finally, it is assumed that the filter is exposed and processed such that operation is in the linear part of the H&D curve, where density is linearly proportional to the logarithm of exposure, i.e. where

$$D = \gamma \log E - D_o.$$

Note that this is not the usual region of operation used for other interferometrically generated filters, which are typically recorded in the linear portion of the t_A vs. E curve.

The three postulates above lead to certain conclusions regarding the mathematical relationship between changes of exposure and resulting changes of amplitude transmittance. To discover this relationship, first note that a change of logarithmic exposure induces a proportional change of amplitude transmittance, as evidenced by the chain

$$\Delta t_A \propto \Delta \phi \propto \Delta D \propto \Delta(\log E),$$

which is implied by the above hypotheses. In addition, if the exposure pattern consists of a strong average exposure \bar{E} and a weaker varying exposure ΔE , then

$$\Delta(\log E) \approx \frac{\Delta E}{\bar{E}},$$

making

$$\Delta t_A \propto \frac{\Delta E}{\bar{E}}. \quad (8-58)$$

With the above information as background, attention is turned to the process of recording the **deblurring** filter. The recording geometry is that of a **VanderLugt** filter, exactly as illustrated previously in Fig. 8.20(a), but with only the function $s(x, y)$ present in the input transparency. The exposure produced by this interferometric recording is

$$E(x, y) = T \left\{ A^2 + a^2 \left| S \left(\frac{x}{\lambda f}, \frac{y}{\lambda f} \right) \right|^2 + 2Aa \left| S \left(\frac{x}{\lambda f}, \frac{y}{\lambda f} \right) \right| \cos \left[2\pi\alpha x + \phi \left(\frac{x}{\lambda f}, \frac{y}{\lambda f} \right) \right] \right\}, \quad (8-59)$$

where A is the square root of the intensity of the reference wave at the film plane, a is the square root of the intensity of the object wave at the origin of the film plane,⁹ α is again the carrier frequency introduced by the off-axis reference wave, ϕ is the phase distribution associated with the blur transfer function S , and T is the exposure time.

An additional unusual attribute of the Ragnarsson filter is the fact that it is recorded with the object wave much stronger at the origin of the film plane than the reference wave, i.e.

$$A^2 \ll a^2.$$

Because of this condition, we make the following associations with the average exposure \bar{E} and the varying component of exposure ΔE ,

$$\begin{aligned} \bar{E} &= \left[A^2 + a^2 \left| S \left(\frac{x}{\lambda f}, \frac{y}{\lambda f} \right) \right|^2 \right] T \\ \Delta E &= 2AaT \left| S \left(\frac{x}{\lambda f}, \frac{y}{\lambda f} \right) \right| \cos \left[2\pi\alpha x + \phi \left(\frac{x}{\lambda f}, \frac{y}{\lambda f} \right) \right] \end{aligned} \quad (8-60)$$

Choosing the term of transmittance of the processed transparency that is proportional to S^* , we have

$$\Delta t_A \propto \frac{\Delta E}{\bar{E}} \propto \frac{S^*}{K + |S|^2}, \quad (8-61)$$

where

$$K = \frac{A^2}{a^2},$$

(often called the beam ratio), which is precisely the amplitude transmittance required for a Wiener filter when the signal and noise have flat power spectra with a ratio of noise power to signal power of K at all frequencies.

Both Ragnarsson [239] and Tichenor and Goodman [283] have demonstrated restorations with dynamic ranges of 100 : 1 in $|S|$ using this technique. Figure 8.22 shows photographs of the blur impulse response, the magnitude of the deblur impulse response, and the impulse response of the cascaded blur and deblur filters, illustrating the restoration of a blurred point source. The **deblurring** operation becomes highly

⁹For simplicity, we have assumed that the transfer function S has been normalized to unity at the origin.

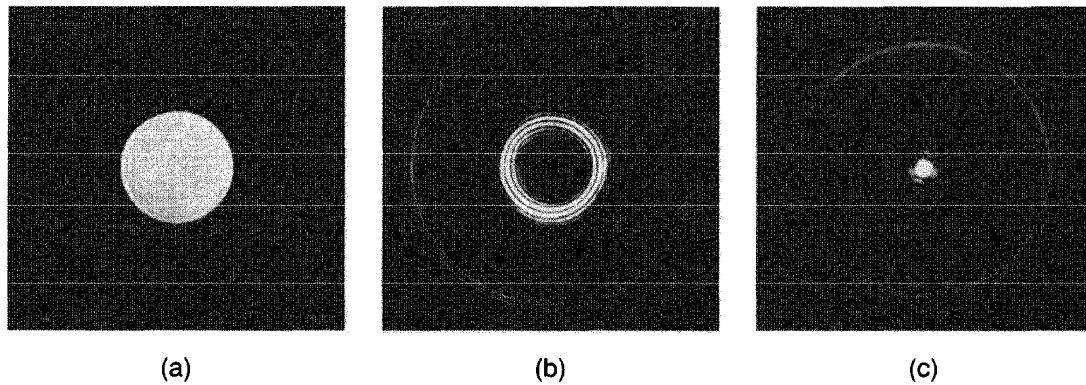


FIGURE 8.22

Deblurring of the blur point-spread function. (a) The original blur, (b) the magnitude of the deblur point-spread function, and (c) the point-spread function of the blur-deblur sequence. [Courtesy of D.A. Tichenor]

sensitive to optical noise at the input of the processor as the dynamic range of the deblurring operation increases. For example, dust specks and small phase perturbations on the input transparency generate deblur impulse responses in the output image which eventually mask the desired image detail [122].

8.9

PROCESSING SYNTHETIC-APERTURE RADAR (SAR) DATA

One of the most successful applications of optical information processing during the 1960s and 1970s was to processing of data gathered by *synthetic-aperture radars*. While optical processing techniques have been largely replaced by digital processing for these problems since the 1970s, nonetheless the ideas developed for optical processing of such data form an important intellectual heritage in the field of optical information processing. Many excellent discussions of SAR can be found in the literature, including several books (see, for example, [117] and [102]). See also [74], on which the early parts of this discussion are based.

8.9.1 Formation of the Synthetic Aperture

With reference to Fig. 8.23, consider a side-looking radar system carried by an aircraft flying with constant speed v_a along a linear flight path in the x direction. Suppose that the function of the radar is to obtain a high-resolution map of the microwave reflectivity of the terrain across an area adjacent to the flight path. Resolution in slant range from the flight path is obtained by transmitting pulsed radar signals and recording the returned signals as a function of time, i.e. by pulse-echo timing. Resolution in azimuth, or equivalently, along the direction of the flight path, could in principle be obtained by using a radar beam of extremely narrow azimuthal extent. However, the azimuthal resolution obtainable at range R from an antenna of linear extent D is roughly $\lambda_r R/D$.

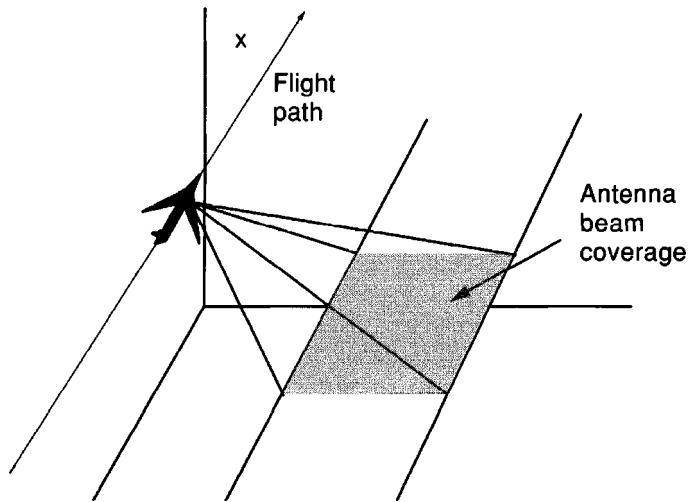


FIGURE 8.23
Synthetic-aperture radar
geometry.

Since the microwave wavelength λ , is typically four or five orders of magnitude larger than an optical wavelength, antennas so large as to be impractical would be required to obtain resolutions comparable with those of optical photo-reconnaissance systems.

A solution to this problem is offered by the synthetic-aperture technique. Let the aircraft carry a small, *broadbeam* antenna which points in a fixed side-looking direction with respect to the aircraft. Radar pulses are transmitted from a uniformly spaced sequence of positions along the flight path, and the time records of both the amplitude and the phase of the radar returns received at these positions are recorded. Each such signal may be regarded as the signal that would be obtained from a single element of a large antenna array, and the various recorded waveforms need only be properly combined to synthesize an effective aperture that may be hundreds or even thousands of meters long.

In order to maintain coherence across the various elements of the synthetic array, it is necessary that there be a common phase reference for the signals measured at all positions along the flight path. This reference is provided by a highly stable local oscillator carried in the aircraft which is used in the detection of all received signals.

Note that to realize the longest possible synthetic array, the radar antenna must illuminate a given point on the terrain for the longest possible portion of the flight path. Thus the broader the beamwidth of the radar antenna, the higher the resolution that can potentially be obtained from the received data. It is possible to show that the best resolution obtainable in the final map of the terrain reflectivity is approximately equal to one-half the dimension of the antenna carried by the aircraft (see Prob. 8-17).

8.9.2 The Collected Data and the Recording Format

To examine the signal-collecting process in more detail, consider the geometry illustrated in Fig. 8.24. The distance along the flight path is represented by the coordinate x .

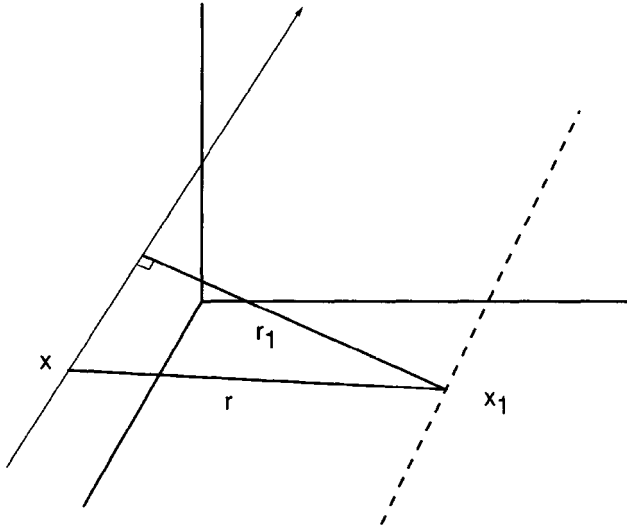


FIGURE 8.24
Flight-path geometry.

For simplicity, we assume that a simple point scatterer exists at coordinate x_1 , which lies at a perpendicular distance r_1 from the flight path. r is the distance of the aircraft from this point scatterer, a function of time. For additional simplicity, we assume that the waveform transmitted by the radar is a steady sinusoid of frequency f_r . The pulsed nature of the actual signal transmitted results simply in a periodic sampling of the signals predicted under the sinusoidal assumption. While the pulsed nature of the transmitted signal must be taken into account when considering imaging in range, it is not important when considering only imaging in azimuth, and for the moment azimuthal imaging (the direction in which the aperture is synthesized) is of primary concern.

The signal returned to the aircraft from the point scatterer under consideration can be represented by the time-varying phasor

$$s_1(t) = \sigma_1 \exp\left[-j2\pi f_r \left(t - \frac{2r}{c}\right)\right] \quad (8-62)$$

where f_r is the RF frequency of the radar, c is the velocity of light, and σ_1 is a complex amplitude factor which depends on such parameters as transmitted power, target reflectivity and phase shift, and inverse fourth-power propagation attenuation. The distance r may be expressed in terms of r_1 , x_1 , and x (the flight-path coordinate) by

$$r = \sqrt{r_1^2 + (x - x_1)^2} \approx r_1 + \frac{(x - x_1)^2}{2r_1}, \quad (8-63)$$

yielding a signal

$$s_1(t) = \sigma_1(x_1, r_1) \exp\left\{-j\left[2\pi f_r t - \frac{4\pi r_1}{\lambda_r} - \frac{2\pi(x - x_1)^2}{\lambda_r r_1}\right]\right\}. \quad (8-64)$$

The motion of the aircraft links the variable x to the time variable t through the prescription

$$x = v_a t.$$

If the terrain at slant range r_1 along the flight path is regarded as consisting of a collection of many scatterers, the total returned signal can be written

$$\begin{aligned} s(t) &= \sum_n s_n(t) \\ &= \sum_n \sigma_n(x_n, r_1) \exp \left\{ -j \left[2\pi f_r t - \frac{4\pi r_1}{\lambda_r} - \frac{2\pi(v_a t - x_n)^2}{\lambda_r r_1} \right] \right\}, \end{aligned} \quad (8-65)$$

where each $s_n(t)$ is the signal received from a different scatterer. The returned signal is synchronously demodulated, using the stable internal local oscillator referred to earlier. This operation translates the center frequency of the return from the microwave frequency f_r to a new lower frequency f'_r , yielding

$$s'(t) = \sum_n |\sigma_n(x_n, r_1)| \cos \left[2\pi f'_r t - \frac{4\pi r_1}{\lambda_r} - \frac{2\pi}{\lambda_r r_1} (v_a t - x_n)^2 + \phi_n \right], \quad (8-66)$$

where ϕ_n is the phase associated with the complex quantity σ_n , and we have abandoned the phasor notation in favor of real notation.

Note that the signal received from a given point scatterer is in fact a sinusoidal one with a phase that varies quadratically with time. Equivalently, the instantaneous frequency of this signal is "chirping", with a chirp rate that depends on the aircraft speed as well as the distance of the scatterer from the flight path. The chirping frequency is caused by doppler shifts experienced as the aircraft approaches, passes, and recedes from the scatterer. The received signal starts at a high frequency when the aircraft velocity is nearly directly towards the scatterer, the frequency drops as the component of the velocity vector pointing towards the scatterer grows smaller, and the frequency shift vanishes when the velocity of the aircraft has no component towards the scatterer. The aircraft then begins to recede from the scatterer, with a component of velocity away from the scatterer that increases with distance. This, then, is the origin of the doppler shifts that cause the chirping received signal. It is the entire received doppler history that must be operated on to form an azimuthal image of the scatterer.

In the early days of optical processing of SAR data, it was common to record the received signals on photographic film in a format suitable for optical processing. Film was thus used both as a medium for high-density storage of the received radar data and later as a means for inserting the data into a coherent optical processor. The demodulated signal was used to intensity-modulate a cathode-ray tube (CRT), with the electron beam swept vertically during the time interval of interest after each transmitted pulse. If film is drawn past the CRT face with horizontal velocity v_f in a direction normal to the sweeping electron beam, the recording format shown in Fig. 8.25 is obtained (the electron beam is blanked as it returns for the next range sweep). The vertical lines represent successive range sweeps, while the azimuthal position of a given scatterer with respect to the radar varies along the horizontal direction.

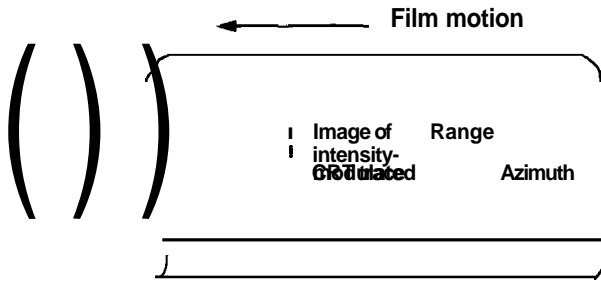


FIGURE 8.25
Recording format.

8.9.3 Focal Properties of the Film Transparency

The photographic recording serves as the input to a coherent optical processor with a special architecture to be discussed. For the moment we limit attention to the focusing properties of the film when illuminated with coherent light, turning later to consideration of the proper processing architectures for obtaining an image of the terrain.

Again limit attention to a single slant range r_1 , thus considering only the data recorded along a line $\eta = \eta_1$ on the film, and again neglect the pulsed nature of the transmitted signal. With proper care in exposure and chemical processing, the azimuthal history of the received signal can be made to generate a photographic record with amplitude transmittance given by

$$t_A(\xi, \eta_1) = t_b + \chi \sum_n |\sigma_n(x_n, r_1)| \times \cos \left[2\pi f_X \xi - \frac{4\pi r_1}{\lambda_r} - \frac{2\pi}{\lambda_r r_1} \left(\frac{v_a}{v_f} \xi - x_n \right)^2 + \phi_n \right], \quad (8-67)$$

where t_b is a bias transmittance introduced to allow the recording of the bipolar radar signals, ξ is the horizontal coordinate on the film, and χ is a constant that is proportional to the slope of the t_A vs. E curve of the photographic film. In writing Eq. (8-67), use has been made of the relation

$$\xi = v_f t,$$

from which it follows that the carrier frequency f_X on the film is given by

$$f_X = \frac{f_r'}{v_f}.$$

It is also worth noting that the vertical coordinate η_1 where the azimuthal signal from a particular point scatterer is recorded, and slant range r_1 of that scatterer from the flight path are related through

$$\eta_1 = 2 \frac{v_e}{c} r_1, \quad (8-68)$$

where v_e is the vertical speed of the CRT spot during the recording process. This relation is easily proven by equating the position of the vertically scanning spot and the range from which a signal is being received at any particular chosen time t .

By decomposing the cosine of Eq. (8-67) into two complex-exponential factors, the transmittance may be expressed as the sum of the bias and two additional terms of the form

$$t_\alpha(\xi, \eta_1) = \frac{\chi}{2} \sum_n \sigma'_n(x_n, r_1) \exp \left\{ j \left[2\pi f_X \xi - \frac{2\pi}{\lambda_r r_1} \left(\frac{v_a}{v_f} \right)^2 \left(\xi - \frac{v_f}{v_a} x_n \right)^2 \right] \right\} \quad (8-69)$$

and

$$t_\beta(\xi, \eta_1) = \frac{\chi}{2} \sum_n \sigma_n'^*(x_n, r_1) \exp \left\{ -j \left[2\pi f_X \xi - \frac{2\pi}{\lambda_r r_1} \left(\frac{v_a}{v_f} \right)^2 \left(\xi - \frac{v_f}{v_a} x_n \right)^2 \right] \right\}, \quad (8-70)$$

where the constant phase $4\pi r_1/\lambda_r$, as well as the phases ϕ_n , have been absorbed into the definition of the σ_n' .

Restricting attention to only one of the point scatterers, say the one with index $n = N$, the appropriate component of t_α is

$$t_\alpha^{(N)}(\xi, \eta_1) = \frac{\chi}{2} \sigma_N'(x_N, y_1) \exp(j2\pi f_X \xi) \\ \times \exp \left[-j \frac{2\pi}{\lambda_r r_1} \left(\frac{v_a}{v_f} \right)^2 \left(\xi - \frac{v_f}{v_a} x_N \right)^2 \right]. \quad (8-71)$$

The first exponential term, having a linear phase dependence, introduces a simple tilt of the phase front of this component of light. The angle θ of the tilt from the transparency plane may be determined from the relation

$$\sin \theta = \lambda_o f_X \quad (8-72)$$

where λ_o is the wavelength of the light.

Turning to the second exponential factor, we note its close resemblance to the amplitude transmittance function of a positive cylindrical lens, centered at coordinate $\xi = \xi_o$,

$$t_l(\xi) = \exp \left[-j \frac{\pi}{\lambda_o f_1} (\xi - \xi_o)^2 \right], \quad (8-73)$$

where f_1 is the focal length. Equating the last term of (8-71) with (8-73), we find that this component of t_α behaves like a positive cylindrical lens with focal length

$$f_1 = \frac{1}{2} \frac{\lambda_r}{\lambda_o} \left(\frac{v_f}{v_a} \right)^2 r_1, \quad (8-74)$$

and with lens center (axis) located at coordinate

$$\xi = \frac{v_f}{v_a} x_N. \quad (8-75)$$

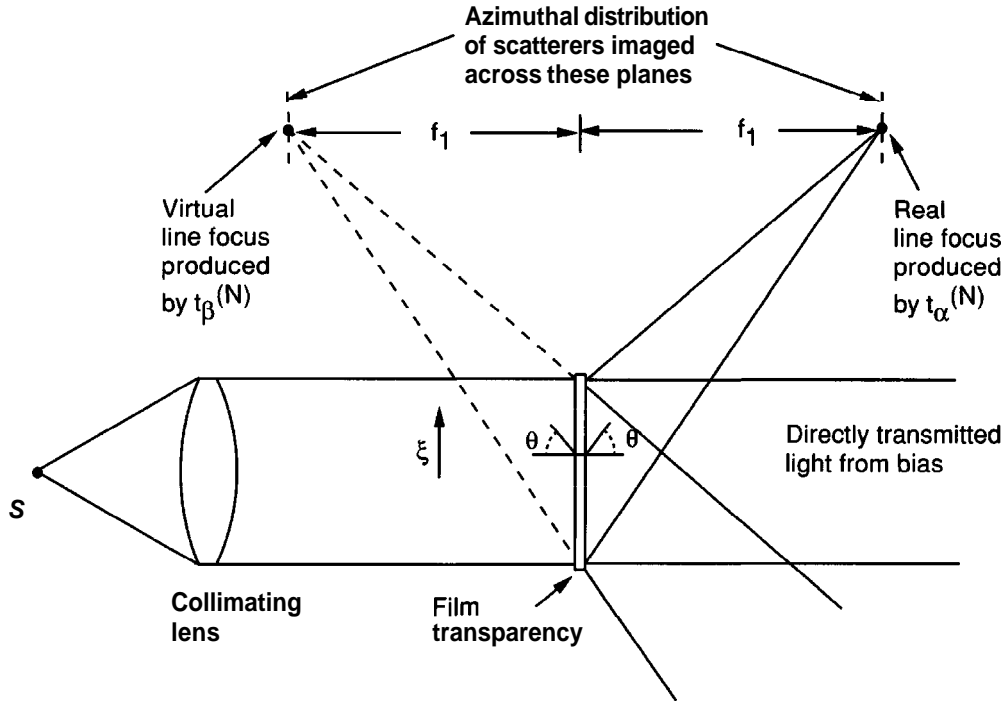


FIGURE 8.26
Light transmitted by the line $\eta = \eta_1$ of the film transparency.

In a similar fashion, the Nth component of t_β ,

$$t_\beta^{(N)}(\xi, \eta_1) = \frac{\chi}{2} \sigma_N'^*(x_N, y_1) \exp(-j2\pi f_X \xi) \times \exp \left[j \frac{2\pi}{\lambda r_1} \left(\frac{v_a}{v_f} \right)^2 \left(\xi - \frac{v_f}{v_a} x_N \right)^2 \right], \quad (8-76)$$

has an exponential factor which introduces a wavefront tilt in the opposite direction, i.e. at angle $-\theta$, and a second exponential factor which is identical with the amplitude transmittance of a *negative* cylindrical lens, again centered at $\xi = (v_f/v_a)x_N$ and with a focal length given by the negative of Eq. (8-74).

Figure 8.26 illustrates the three components of light transmitted by this slit-segment of film for the case of a single point scatterer. The bias transmittance t_b allows the incident optical wave to pass through the transparency, uniformly attenuated but otherwise unchanged in the ξ direction.¹⁰ The components $t_\alpha^{(N)}$ and $t_\beta^{(N)}$ of the transmittance may be regarded as generating a pair of "images" of the point scatterer in the following sense: the component $t_\alpha^{(N)}$ focuses light to a bright line focus (rising out of the page in Fig. 8.26) to the right of the transparency, while the component $t_\beta^{(N)}$ produces a wave

¹⁰Of course the transmitted optical wave is expanding in the η direction (i.e. out of the paper in Fig. 8.26) due to diffraction by the narrow slit of film being considered. However, we ignore the η behavior for the moment, returning to it later.

that appears to originate from a line source to the left of the transparency (see Fig. 8.26). If a multitude of point scatterers is present at various locations along the flight path at range r_1 , each generates its own pair of real and virtual line foci when illuminated with coherent light. The relative azimuthal positions of the point scatterers determine the relative positions of the centers of the lens-like structures on the film, and therefore are preserved in the relative positions of the corresponding line foci. Thus an entire image of the azimuthal distribution of scatterers at range r_1 is recreated across appropriate planes in front of and behind the transparency. We emphasize again that this image is spread in the η direction, since the film exerts no focal power in that direction.

8.9.4 Forming a Two-Dimensional Image

We ultimately wish to form an image, not only of the azimuthal distribution of scatterers, but also of their distribution in range. The azimuth history we have been discussing above (corresponding to the particular range r_1) resides at a particular η_1 coordinate on film corresponding to a scaled version of its actual range from the flight path (see Eq. (8-68)). Thus it is necessary to *image* the η variations of film transmittance directly onto the plane of focus of the azimuthal signals. This task is complicated by the fact that the focal length of the azimuthal variations on film is a function of the particular range r_1 under consideration. To construct the final radar image, it is evidently necessary to image the η variations onto a *tilted* plane in which the azimuthal foci occur.

This task can be accomplished with the optical system of Fig. 8.27. A positive *conical* lens (called an “axicon” – cf. Prob. 5-3) is inserted immediately behind the transparency. The transmittance function of this lens is

$$t_l(\xi, \eta) = \exp\left[-j \frac{\pi}{\lambda_o f(\eta)} \xi^2\right] \quad (8-77)$$

and its focal length depends linearly on the η coordinate (or equivalently on the range coordinate through Eq. (8-68)) according to

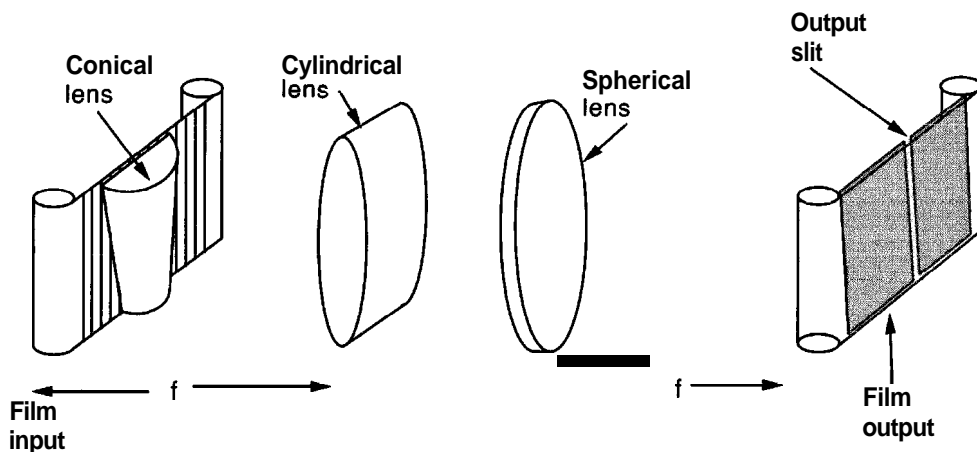


FIGURE 8.27

Optical system for recording an image. A slit is used in the output plane.

$$f(\eta) = \frac{1}{2} \frac{\lambda_r}{\lambda_o} \left(\frac{v_f}{v_a} \right)^2 r_1 = \frac{1}{4} \left(\frac{\lambda_r}{\lambda_o} \right) \left(\frac{c}{v_e} \right) \left(\frac{v_f}{v_a} \right)^2 \eta. \quad (8-78)$$

Because the focal length of the conical lens varies linearly with the range coordinate, it removes all the *virtual* line sources (for all ranges) to infinity, i.e. the quadratic-phase factor of the lens cancels the quadratic-phase factors of all of the azimuthal histories of signals received from all ranges. Thus the entire tilted plane of the azimuthal image is moved to infinity. Azimuthal information is retained through the *angles* at which the infinitely distant virtual line sources lie.

Next a cylindrical lens with power only in the vertical or range dimension is placed one focal length from the film. This lens creates an image of the vertical (η) structure on the film, i.e. the range information, at infinity. The azimuthal and range images now coincide at infinity, and must be brought back from infinity to form a real image in a single plane. This is accomplished by a positive spherical lens, following the cylindrical lens, placed one focal distance from the final observation plane. A single real image, with both range information and azimuthal information in focus, now lies in the back focal plane of the spherical lens. However, this image has one serious defect. The focal powers exerted by the optical system on azimuthal records corresponding to different ranges are all different. As a consequence, the azimuthal magnification of the final image varies with range, with greater magnification at ranges where the focal length of the *axicon* was shorter. The result is a seriously distorted image in the output plane. To overcome this distortion, a vertical slit is inserted in the output plane, and the output film strip is moved linearly past this slit in synchronism with the motion of the input film strip (but in general with a different speed). Since only a vertical range strip is being recorded at one time, the varying azimuthal magnifications have no effect, and an undistorted image of the terrain is recorded on the film strip.

Thus through the use of a reasonably sophisticated optical system, a full image of the microwave reflectivity of the ground strip has been recorded by the optical processing system. The primary deficiency of the system is that, at any one time, it provides only a linear strip image, rather than a full two-dimensional image, so the full two-dimensional processing power of the optical system has not been utilized. This defect is remedied by a second, even more sophisticated optical processing system, to be discussed next.

8.9.5 The Tilted Plane Processor

The tilted plane processor [180] overcomes the chief deficiency of the simpler processor described above, namely its inability to use the full two-dimensional output of the optical processing system. The tilted plane processor is arguably one of the most sophisticated optical information processing systems yet constructed.

The nature of the problem to be solved is clarified with the help of Fig. 8.28, showing the tilted azimuthal image planes and the untilted range plane. The goal is to bring one of the tilted azimuthal planes into coincidence with the range plane, and to do so in such a way that both range and azimuthal magnification are constant across the output. To understand how this is done with the tilted plane processor, it is necessary to first

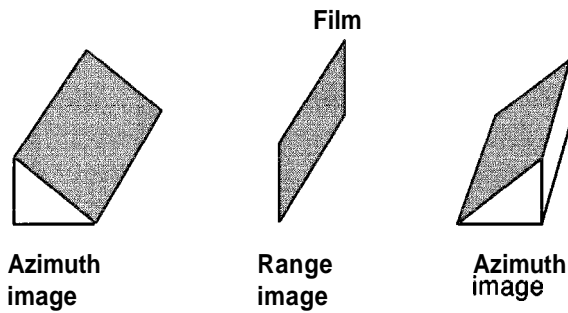


FIGURE 8.28 Tilted azimuth planes and untilted range plane.

digress into a short discussion of the three-dimensional imaging properties of telescopic systems.

Three-dimensional imaging properties of telescopic systems

Consider a simple telescopic system such as shown in Fig. 8.29. Lens L_1 is spherical with focal length f_1 and lens L_2 is spherical with focal length f_2 . An object is placed in the front focal plane of L_1 and an image is observed at the rear focal plane of L_2 .

The transverse magnification m_t from object to image plane for this system is easily shown to be

$$m_t = -\frac{f_2}{f_1}, \tag{8-79}$$

where the minus sign accounts for image inversion. Of equal interest to us here is the axial magnification, m_a , which applies for displacements along the optical axis. It is straightforward to show that the axial magnification of a telescopic system is given by

$$m_a = m_t^2 = \left(\frac{f_2}{f_1}\right)^2. \tag{8-80}$$

The transverse magnification is independent of the transverse coordinates of the object point in question, and the axial magnification is independent of the axial position of the object point. As a consequence, as illustrated in Fig. 8.30 any rectangular parallelepiped in object space is transformed into a parallelepiped having a different shape in the image space. As illustrated in the figure, with a magnification less than unity, the

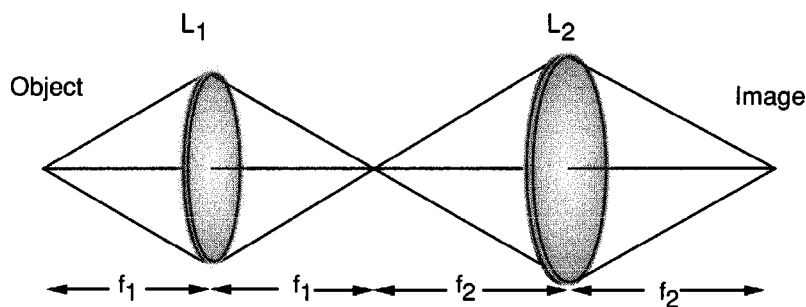


FIGURE 8.29 Two-lens telescopic system.

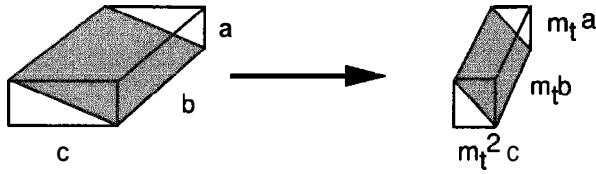


FIGURE 8.30
Demagnification of a tilted plane for $f_1 > f_2$.

tilt of an object plane is reduced, and it is made more upright in the image plane, while preserving a constant transverse magnification. This property can be used to bring one of the azimuthal image planes into a nearly vertical position, in which case a slightly tilted range plane can be made coincident with it.

An anamorphic telescope

In order to bring the azimuthal and range image planes into a common plane of focus, a telescope with different focal properties in the two transverse dimensions is required. Such a telescope is called anamorphic, and can be constructed with combinations of spherical and cylindrical lenses. Figure 8.31 shows one such system from the side and from the top. For both views, the lens L_0 is simply a collimating lens that provides plane-wave illumination for the optical system that follows.

Considering first the top view, the only lenses with power in the range direction are the spherical lenses L_1 and L_3 , which form a telescope with $f_1 = f_3$. The input film

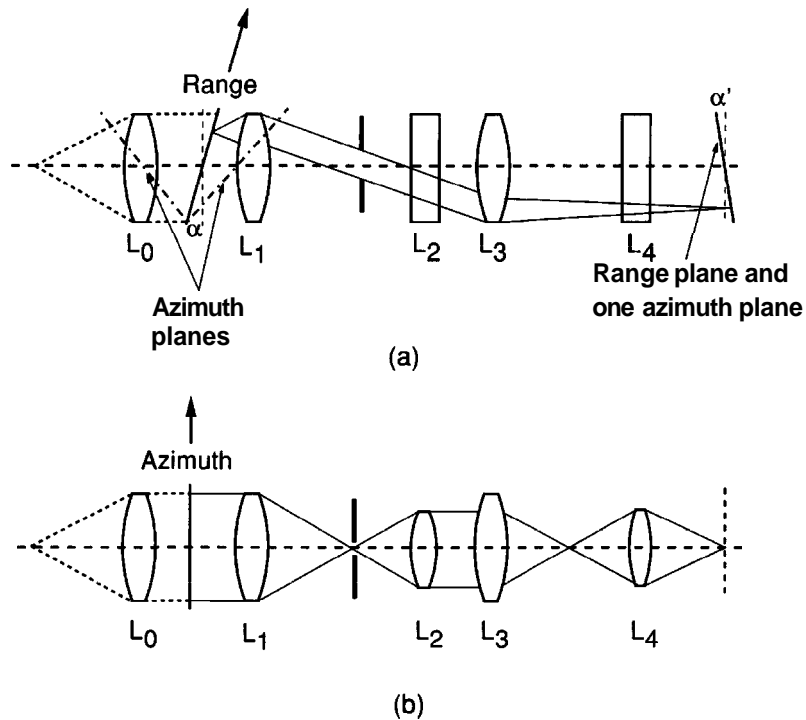


FIGURE 8.31
Anamorphic telescope comprising a tilted plane processor. (a) Side view (range direction is vertical), and (b) top view (azimuth direction is vertical).

is seen to be tilted through angle α in the input plane, and due to the demagnification of the two spherical lenses, tilted by a smaller angle α' in the output plane. Thus the system images in the range dimension from one tilted plane onto a second plane with reduced tilt. The output film is tilted to coincide with this tilted plane. Considering next the side view, all four lenses have power in the azimuth dimension and they form a system that images (for fixed range) any azimuthal focal point (remember such points do not lie in the film plane) onto the recording film. The input and output are not tilted when viewed from this perspective. By properly choosing the parameters of the system (i.e. f_1 , f_2 , f_3 , and f_4 and the tilts of the input and output films) [180], it is possible to achieve equal magnifications in the range and azimuth dimensions at all points in the output aperture. Thus the entire two-dimensional processing ability of the system is utilized; much more light is brought to the image plane than for the system with the slit, which in turn means that film can be moved through the system faster. Finally, the use of the full two-dimensional aperture results in the reduction of coherent artifacts that often arise from dust specks on lenses and other imperfections in a coherent optical system.

Other forms of the tilted plane processor are also possible. For details the reader should consult Ref. [180]. Figure 8.32 shows a processed synthetic-aperture radar image obtained from a tilted plane processor. The image shows the Huron river near Ann Arbor, Michigan, with a ground resolution of 5 ft by 7 ft.

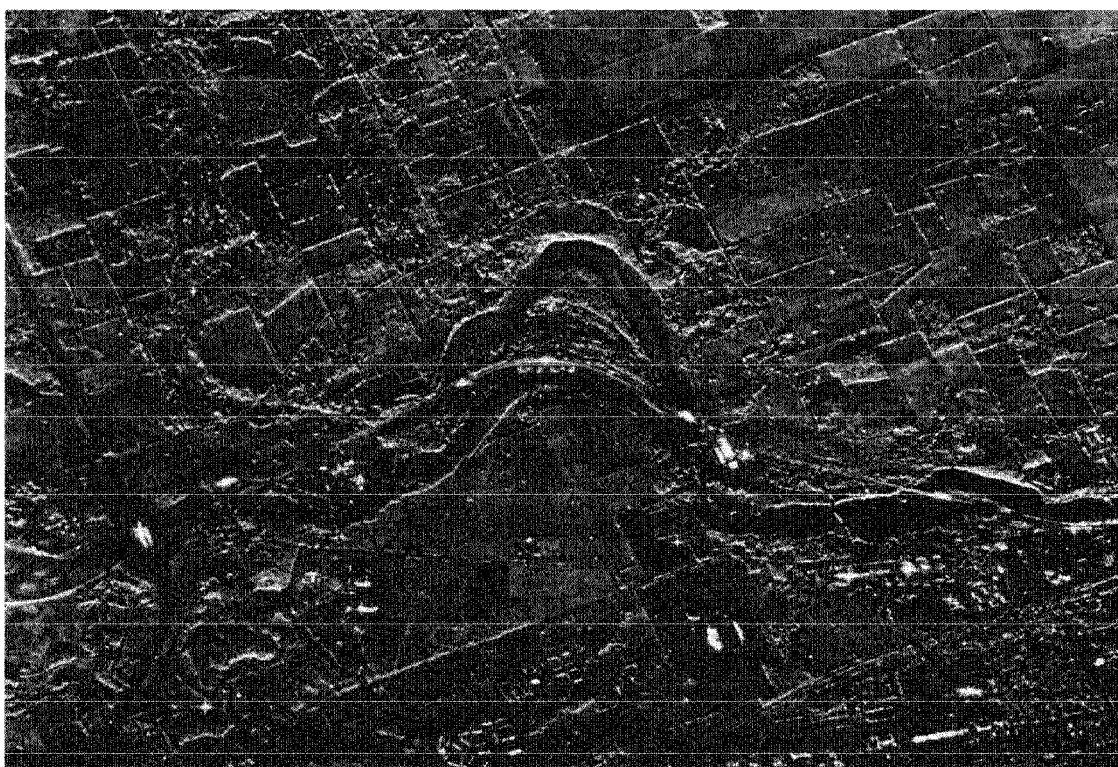


FIGURE 8.32

Synthetic-aperture radar image obtained from a tilted plane processor. The image shows a portion of the Huron river and surrounding farmland near Ann Arbor, Michigan. [Courtesy of the Environmental Research Institute of Michigan.]

8.10 ACOUSTO-OPTIC SIGNAL PROCESSING SYSTEMS

The means by which a temporal electrical signal can be converted into a moving spatial optical signal with the help of an acousto-optic cell was discussed in Section 7.2.6. Attention is turned here to the use of such cells as input transducers for various types of signal processing systems. Since virtually all modern work in this area utilizes microwave signals in crystals, we focus attention exclusively on **Bragg** cells as the input transducers. While systems based on **Raman-Nath** diffraction were important in the early days of acousto-optic signal processing [264], [6], they are virtually non-existent today.

Our discussion is of necessity brief, but we will describe three different system architectures. One, the Bragg cell spectrum analyzer, can be used to analyze the spectrum of broadband microwave signals. Attention is then turned to two types of acousto-optic correlators, the space-integrating correlator and the time-integrating correlator.

8.10.1 Bragg Cell Spectrum Analyzer

The ease with which Fourier transforms can be performed in coherent light suggests that a system that combines an acousto-optic input transducer with a coherent optical Fourier transform system can function as a spectrum analyzer for **wideband** and **high-frequency** electrical signals. Figure 8.33 shows the basic structure of such a spectrum analyzer.

Consider a high-frequency signal represented by the electrical voltage

$$v(t) = A(t) \cos[2\pi f_c t - \psi(t)] = \text{Re} \left\{ A(t) e^{j\psi(t)} e^{-j2\pi f_c t} \right\} = \text{Re} \left\{ s(t) e^{-j2\pi f_c t} \right\}, \quad (8-81)$$

where $s(t) = A(t) e^{j\psi(t)}$ is the complex representation of signal.

With reference to Eq. (7-34) and Fig. 8.33, let the coordinate y_1 refer to the plane where the transmitted field exits the Bragg cell, and let the coordinate y_2 refer to the plane in the rear focal plane of the Fourier transforming lens. When the above signal is applied to an acousto-optic cell, and the cell is illuminated at the Bragg angle by a collimated, monochromatic wave, there results a transmitted wavefront into the +1 diffracted order given by

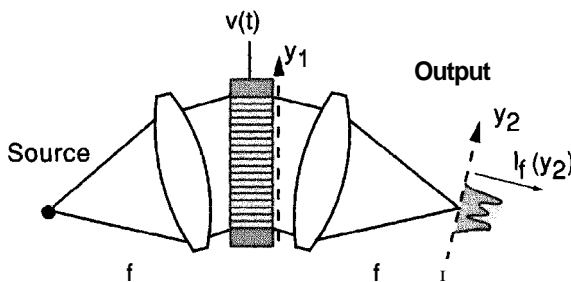


FIGURE 8.33
Bragg cell spectrum analyzer.

$$U(y_1; t) = C s\left(\frac{y_1}{V} + t - \tau_o\right) e^{-j2\pi y_1/\Lambda} \text{rect}\frac{y_1}{L}$$

where C is a constant and we have neglected the temporal frequency shift by f_c , since it has no impact on our calculations.

This optical signal now passes through a positive lens tilted at the Bragg angle, as shown in Fig. 8.33. Noting that the linear phase factor in y_1 is canceled by the tilt of the lens, the spatial distribution of fields appearing in the back focal plane of the lens will be (aside from quadratic-phase factors in y_2 which we can neglect)

$$U_f(y_2; t) = C' \int_{-\infty}^{\infty} s\left(\frac{y_1 + Vt - V\tau_o}{V}\right) \text{rect}\frac{y_1}{L} \exp\left(-j\frac{2\pi y_1 y_2}{\lambda f}\right) dy_1. \quad (8-82)$$

This is a Fourier transform of a product of two functions, so the convolution theorem will apply, and we consider each of the two spectra individually. Consider the Fourier transform of the scaled signal first; we have

$$\mathcal{F}\left\{s\left(\frac{y_1 + Vt - V\tau_o}{V}\right)\right\} = V S(V f_Y) \exp[j2\pi f_Y V(t - \tau_o)] \quad (8-83)$$

where $S = \mathcal{F}\{s\}$. The presence of the term depending on time t in this result is an indication that *every spatial frequency component is oscillating with a different optical frequency*. Considering the rect function next, we have

$$\mathcal{F}\left\{\text{rect}\frac{y_1}{L}\right\} = L \text{sinc } L f_Y.$$

For the moment, neglect the finite length of the Bragg cell, allowing L to become arbitrarily large, in which case the sinc function approaches a δ function. The optical intensity incident in the focal plane will then be (neglecting multiplicative constants)

$$I_f(y_2) = \left| S\left(\frac{V y_2}{\lambda f}\right) \exp\left[j\frac{2\pi}{\lambda f} V y_2(t - \tau_o)\right] \right|^2 = \left| S\left(\frac{V y_2}{\lambda f}\right) \right|^2. \quad (8-84)$$

The intensity distribution measured by an array of time-integrating detectors will therefore be proportional to the *power spectrum* of the input signal, and the acousto-optic system acts as a *spectrum analyzer*.

The relationship between position y_2 in the focal plane and temporal frequency f_t of the input signal can be found by first noting that the center frequency f_c of the electrical signal corresponds to the origin of the y_2 plane (we choose the origin to make this true). As we move in the positive y_2 direction, we are moving to lower temporal frequencies (zero temporal frequency corresponds to the direction of the zero order of the acoustic grating). From the scaling factors present in the equation above we find that the temporal input frequency corresponding to coordinate y_2 is

$$f_t = f_c - \frac{V y_2}{\lambda f}$$

However, when only the time integrated intensity of the light is detected, the temporal frequency of the light is of no consequence.

When the length of the Bragg cell is not infinite, convolution with the sinc function in the spectral domain cannot be neglected. This convolution in effect smooths the measured spectrum, establishing the frequency resolution obtainable. The minimum resolvable difference in temporal frequency is readily shown to be approximately the reciprocal of the total time delay stored in the Bragg cell window, i.e.

$$\Delta f_t = V/L.$$

The technology of high-performance Bragg cell spectrum analyzers is well developed. Center frequencies lie in the hundreds of MHz to the 1- to 3-GHz range, and time-bandwidth products (equal to the number of resolvable spectral elements) from several hundred to more than 1,000 have been reported.

8.10.2 Space-Integrating Correlator

Bragg cells can also be used as real-time inputs to convolvers and correlators. Historically the first such systems were based on what is now called a space-integrating architecture. Consider the acousto-optic system shown in Fig. 8.34. This system contains one Bragg cell, which is used for converting a temporal voltage $v_1(t)$ into a complex distribution of field $s_1(\frac{y_1}{V} + t - \tau_0)$. Here s_1 is the complex representation of an amplitude and phase modulated voltage, analogous to the representation of Eq. (8-81). Due to the Bragg effect, the cell is assumed to transmit only the zero order and the +1 order.

The second input is provided by a fixed transparency which contains an amplitude and phase modulated grating. If $s_2 = B \exp(j\chi)$ represents the second signal, with which s_1 is to be correlated, then the amplitude transmittance of the transparency should ideally be chosen to be

$$\begin{aligned} t_A(y_2) &= \frac{1}{2} \{ 1 + B(y_2) \cos[2\pi f_0 y_2 - \chi(y_2)] \} \\ &= \frac{1}{2} + \frac{B(y_2)}{4} e^{-j\chi(y_2)} e^{j2\pi f_0 y_2} + \frac{B(y_2)}{4} e^{j\chi(y_2)} e^{-j2\pi f_0 y_2}. \end{aligned} \quad (8-85)$$

Such a grating could be computer generated. Alternatively a transparency with the same two first-order grating components could be recorded interferometrically, in a manner

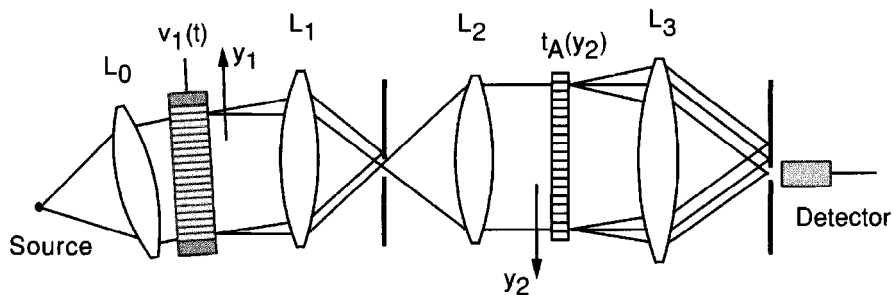


FIGURE 8.34
Acousto-optic space-integrating correlator.

analogous to that used for the **VanderLugt** filter. It is assumed to be a thin grating, so a zero order and two first orders are generated.

The optical system following the Bragg cell contains a stop in the spatial frequency domain that blocks the zero-order transmitted component and passes the first-order diffracted component. Lenses L_1 and L_2 together image the amplitude distribution corresponding to the first diffraction order onto the fixed transparency, with inversion. The y_2 coordinate system is inverted to account for the inversion associated with the imaging operation. Lens L_3 is used to bring the -1 order component diffracted by the fixed grating to focus on a pinhole, which is followed by a nonintegrating photodetector.

The operation performed by lens L_3 , the pinhole, and the detector can be expressed as a spatial integration of the product of the two complex functions of interest. The current generated by the detector is therefore (up to multiplicative constants)

$$i_d(t) = \left| \int_{-\infty}^{\infty} s_1 \left(\frac{y_2 + Vt - V\tau_o}{V} \right) s_2^*(y_2) \text{rect} \frac{y_2}{L} dy_2 \right|^2 \quad (8-86)$$

where L is again the length of the Bragg cell, and the complex conjugate of s_2 occurs because we have chosen the -1 diffraction order of the fixed grating. As time progresses, the scaled signal s_1 slides through the Bragg cell and the relative delay between s_1 and s_2 changes, thus providing the values of the correlation between the two signals for different delays. The correlation operation takes place only within the window provided by the finite length of the Bragg cell.

The distinguishing characteristics of the space-integrating correlator are that the correlation integration is over space and the various values of relative delay occur sequentially in time.

8.10.3 Time-Integrating Correlator

An entirely different approach to realization of an acousto-optic correlator is provided by a system that interchanges the roles of time and space vis-à-vis the space-integrating correlator. Such an approach was first conceived of by Montgomery [216]. A different architecture that accomplishes a similar operation was demonstrated by Sprague and Koliopoulos [272]. This general approach to correlation is known as "time-integrating correlation".

Figure 8.35 shows one architecture of such a correlator. Two RF voltages $v_1(t)$ and $v_2(t)$ are applied to different Bragg cells in close proximity, arranged so that the resulting acoustic signals propagate in opposite directions. The lenses L_1 and L_2 form a standard double Fourier transform system. A light ray entering the first Bragg cell exits as a zero-order ray and a -1 order ray.¹¹ Both the zero order and the -1 order transmitted by the first Bragg cell are split by the second cell, which itself applies zero-order and

¹¹When the voltage is applied at the top of the cell, downwards deflection corresponds to the $+1$ order and upwards deflection to the -1 order. When the voltage is applied to the bottom of the cell, the opposite is true.

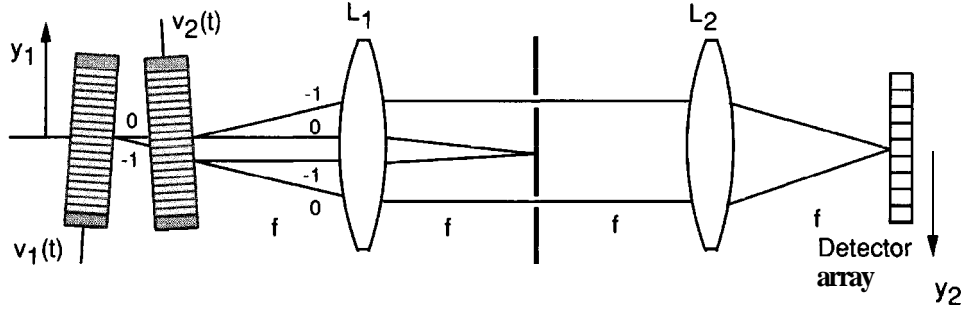


FIGURE 8.35
Time-integrating correlator.

-1 order diffraction to each of those incident rays. A stop in the rear focal plane of L_1 passes only rays that have undergone the sequence $0 \rightarrow -1$ order diffractions or $-1 \rightarrow 0$ order diffractions, blocking the rays that have undergone two zero-order or two -1 order diffractions. Note that the optical frequencies of the two beams passed are identical because they have undergone the same diffractions, although in opposite order. The two optical signals passed by the aperture-stop are then brought back together on an array of time-integrating detectors, which is situated in a plane where the product of the amplitude transmittances of the two Bragg cells is imaged.

Note that each element of the detector array measures the intensity associated with a different vertical location on the two Bragg cells, and as a consequence, for each detector element there is a different relative time delay between the two signals driving the cells, due to the opposite directions of acoustic wave propagation. If $s_1(y_1)$ represents the complex representation of the signal in the first cell, and $s_2(y_1)$ is the complex representation of the signal in the second cell, a detector at location y_2 measures the finite time integral of the squared magnitude of the sum of the two fields that have traveled the two different paths to the detector. Neglecting multiplicative constants, the integral in question is given by

$$E(y_2) = \int_{\Delta T} \left| s_1^* \left(-\frac{y_2}{V} + t + \tau_o \right) e^{-j2\pi\alpha_c y_2} + s_2 \left(\frac{y_2}{V} + t + \tau_o \right) e^{j2\pi\alpha_c y_2} \right|^2 dt, \quad (8-87)$$

where the linear exponential terms of opposite sign account for the opposite angles with which the two components arrive at the detector plane, and ΔT is the finite integration time.

Considering the various parts of this integral, the terms

$$E_1 = \int_{\Delta T} \left| s_1 \left(-\frac{y_2}{V} + t + \tau_o \right) \right|^2 dt$$

$$E_2 = \int_{\Delta T} \left| s_2 \left(\frac{y_2}{V} + t + \tau_o \right) \right|^2 dt$$

will approach constants as the integration time ΔT grows large. The remaining term is

$$\begin{aligned}
E_3 &= \int_{\Delta T} \left[s_1 \left(t + \tau_o - \frac{y_2}{V} \right) e^{j2\pi\alpha_c y_2} s_2^* \left(t + \tau_o + \frac{y_2}{V} \right) e^{j2\pi\alpha_c y_2} + cc \right] dt \\
&= 2\text{Re} \left\{ e^{j4\pi\alpha_c y_2} \int_{\Delta T} s_1 \left(t + \tau_o - \frac{y_2}{V} \right) s_2^* \left(t + \tau_o + \frac{y_2}{V} \right) dt \right\} \\
&= 2\text{Re} \left\{ e^{j4\pi\alpha_c y_2} \int_{\Delta T'} s_1(t') s_2^* \left(t' + \frac{2y_2}{V} \right) dt' \right\}, \tag{8-88}
\end{aligned}$$

where a simple change of variables has been made in the last line, $\Delta T'$ is of the same duration as ΔT , but shifted in accord with the variable change, and cc stands for the complex conjugate of the previous term.

Let the complex function $c(\tau)$ represent the complex finite-time crosscorrelation of s_1 and s_2 ,

$$c(\tau) = \int_{\Delta T'} s_1(t') s_2^*(t' + \tau) dt' = |c(\tau)| e^{j\phi(\tau)}.$$

Then the last line of Eq. (8-88) becomes

$$E_3 = 2\text{Re} \left\{ e^{j4\pi\alpha_c y_2} c(2y_2/V) \right\} = 2|c(2y_2/V)| \cos[4\pi\alpha_c y_2 - \phi(2y_2/V)]. \tag{8-89}$$

If the detectors are small compared with the period $\frac{1}{2\alpha_c}$ of the spatial carrier frequency, the fringe pattern incident on the detector will be sampled at or above the Nyquist rate, and the complex correlation information will be captured by the detector array. Since the detector array is of the charge-coupled-device (**CCD**) type, the measured intensities are read out serially from the array. As a result there is an AC output from the CCD array, which can be isolated from the DC output components with a **highpass** or **bandpass** filter. The amplitude of the AC component represents the magnitude of the complex correlation coefficient, and the phase is the phase of that coefficient. By using an envelope detector, the magnitude of the complex correlation can be measured. To measure the phase, synchronous detection must be used. Generally it is the magnitude information that is of most interest.

Note that for this architecture, each detector element measures the correlation for a different relative delay $2y_2/V$ between the two signals. The time-bandwidth product of the correlation measurement is determined by the integration time of the detector array, and is no longer limited to the delay time of the acoustic cell; rather, the delay time determines the range of relative delays that can be explored. In practice the integration time is limited by accumulation of dark current and by detector saturation that is ultimately introduced by the constant terms E_1 and E_2 .

The architecture described above is that of Montgomery [216]. The architecture of Sprague and Koliopoulos [272] differs in that only a single Bragg cell is used and the second signal is introduced by temporal modulation of the optical source. The reader should consult the reference for details.

8.10.4 Other Acousto-Optic Signal Processing Architectures

A multitude of other acousto-optic signal processing architectures exist, but they will not be covered here. We mention in particular various extensions of acousto-optic systems to two-dimensional processing (see [293], Chapter 15, for examples). Applications of acousto-optic processing to numerical or digital computation are omitted here. An application of Bragg cells to discrete processing is described in the section to follow.

8.11 DISCRETE ANALOG OPTICAL PROCESSORS

Until now, we have considered only optical systems that process continuous analog optical signals. Attention is now turned to another class of systems, namely those that process *discrete* analog optical signals. Discrete signals arise in many different applications. For example, an array of sensors collects a discrete set of measurements. Those measurements may be changing continuously with time, but because there is a discrete array of sensors, only a discrete array of data is available at any one time. In addition, it is often necessary to discretize continuous data in order to subject it to processing. Thus discrete data can arise in a multitude of different ways.

The discreteness does not imply that the data is digital. Quite the contrary, the data of interest here has analog values, which have not been quantized, but there is a finite set of such data to be processed. All of the optical processing systems we shall describe are analog processing systems, in keeping with our earlier restrictions.

8.11.1 Discrete Representation of Signals and Systems

Any continuous signal s dependent on a time coordinate t and/or space coordinates (x, y) can be sampled in a discrete array of data representable by a *vector* of values

$$\vec{s} = \begin{bmatrix} s_1 \\ s_2 \\ \vdots \\ s_N \end{bmatrix}. \quad (8-90)$$

If the sample values are taken sufficiently close together and if the function s is band-limited or nearly bandlimited, we know that it would be possible to reconstruct s either exactly (in the bandlimited case) or with high accuracy (in the almost bandlimited case). Therefore the vector \vec{s} is a suitable representation of the original data. Note that if the signal s arose from a discrete array of sensors, then each component of \vec{s} may be a function of time.

For discrete signals, the superposition integral becomes a matrix-vector product (see [182], Chapter 6). Thus the output \vec{g} (M samples) of a linear system having \vec{s} at the input (N samples) is represented by

$$\vec{g} = \mathbf{H} \vec{s}, \quad (8-91)$$

where \mathbf{H} is a matrix with M rows and N columns,

$$\mathbf{H} = \begin{bmatrix} h_{11} & h_{12} & \cdots & h_{1N} \\ h_{21} & h_{22} & \cdots & h_{2N} \\ \vdots & \vdots & \vdots & \vdots \\ h_{M1} & h_{M2} & \cdots & h_{MN} \end{bmatrix} \quad (8-92)$$

and

$$\vec{g} = \begin{bmatrix} g_1 \\ g_2 \\ \vdots \\ g_M \end{bmatrix}. \quad (8-93)$$

Note that there are $M \times N$ analog multiplications and $M \times N$ analog additions¹² needed to perform the above computation.

Thus in discrete signal processing, the matrix-vector product is as fundamental as the superposition integral, or as its special case, the convolution integral. It is therefore of great interest to devise methods for performing such operations optically, preferably using the parallelism of optical systems to full advantage.

8.11.2 A Serial Matrix-Vector Multiplier

The first important optical processor directed at the problem of processing discrete data was the *serial* incoherent matrix-vector processor of Bocker [23], [38]. As the name implies, this processor was aimed at serially processing samples of data, and used the parallelism of the optical system to perform the M analog multiplications in parallel. A description of its operation now follows.

Figure 8.36 illustrates the operation of this system. For simplicity we assume initially that all elements of the input vector \vec{s} and all elements of the system matrix \mathbf{H} are nonnegative real numbers. Methods of generalization will be discussed later. Discrete analog data is entered into the system as a sequential set of current pulses applied to a light-emitting diode (LED). Each such current pulse has an amplitude that is proportional to the amplitude of one of the elements of the vector \vec{s} . In response the LED emits a series of light pulses, each with an intensity proportional to an element of the signal vector. These light pulses are allowed to diverge and flood a two-dimensional matrix mask, each element of which has an intensity transmittance proportional to one of the elements of the system matrix. Transmitted by that mask is an intensity proportional to the product of the applied signal pulse s_k and all the elements of the system matrix, i.e. a matrix of light intensities proportional to $s_k \mathbf{H}$.

¹²Strictly speaking, only $(M - 1) \times (N - 1)$ additions are needed, but we count generation of the first component of a sum as addition with zero.

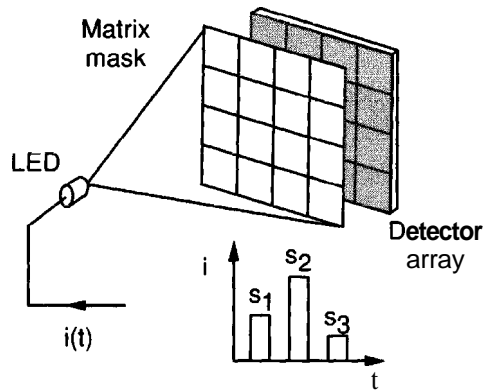


FIGURE 8.36
Serial incoherent matrix-vector multiplier.

The light transmitted by the matrix mask falls upon a two-dimensional **charge-coupled-device (CCD) detector**, which is operated in an unusual mode. The light pulses transmitted by the matrix mask are converted to charge packets residing in the wells associated with each of the discrete detectors. Before the arrival of the next light pulse, all charges in every row are shifted (or clocked) in parallel to the right by one well. The next signal pulse then illuminates the matrix mask, and the CCD detector array detects a new set of signals, which in turn generate charges that are added to the charge packets already residing in the wells. The charges are clocked to the right again and the process repeats until the last signal pulse has been processed.

Concentrate attention on the column of charges which starts initially on the far left of the detector array and ends up after the last signal pulse on the far right of the array. After detection of the first pulse, the cell in the j th row of this first column has accumulated charge proportional to $g_{j1} = h_{j1} s_1$. This charge is clocked to the right one cell, and after detection of the signal pulse s_2 the total charge accumulated is the sum of the first and second charge packets, $g_{j2} = h_{j1} s_1 + h_{j2} s_2$. The process continues, until after N cycles the charge column on the right contains a set of M charges proportional to the elements of the vector \vec{g} .

Thus the parallelism of the system arises through the generation of M useful analog products simultaneously on the CCD array. Many more charge packets are generated per cycle, but only M of them will ultimately be used. In addition, M useful electrical analog additions take place per cycle. After N cycles the entire vector \vec{g} has been accumulated, and it can be read out of the CCD serially.

The distinguishing characteristics of this system are that it accepts discrete data serially, it produces discrete results serially, and it performs M optical multiplications and M electrical additions in parallel.

8.11.3 A Parallel Incoherent Matrix-Vector Multiplier

A fully parallel incoherent matrix-vector processor is illustrated in Fig. 8.37 [124] which temporarily omits the details of the optical elements used. This architecture has come to be known as the "Stanford matrix-vector multiplier" and its use is now nearly as wide-spread as the **VanderLugt** filter in optical signal processing. This system is fun-

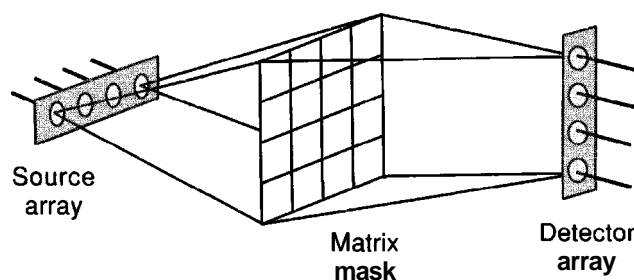


FIGURE 8.37
A fully parallel incoherent matrix-vector multiplier.

damentally faster than the previous serial one, due to the entry of all elements of the signal vector \vec{s} simultaneously, in one clock cycle.

The optics before the mask are arranged so that the light from any one input source, which may be an LED or a laser diode, is spread vertically and imaged horizontally, so that it fills a single vertical column of the mask. Each source thus illuminates a different column. The optics following the mask are arranged so that light from each row of the matrix mask is focused horizontally and imaged vertically, so that it falls upon a single detector element in the output detector array. Thus the light transmitted by each row of the mask is summed optically on a unique detector element. The detectors used here do not integrate charge, but rather respond as fast as possible, generating output signals that vary in unison with the variations of the light intensities falling upon them.

In effect, the input vector \vec{s} is spread vertically so that each output detector can measure an inner product of that vector with a different row vector stored in the matrix mask. For this reason, such a processor is sometimes called an "inner product processor".

There are several different ways to construct an optical system that will achieve the operations indicated diagrammatically in Fig. 8.37. Figure 8.38 shows one such arrangement of elements. Note that because different operations are desired in the horizontal and vertical dimensions, both before and after the matrix mask, the optical system must be anamorphic, with combinations of spherical and cylindrical lenses. The operation of this optical system is as follows. Each of the lenses, whether spherical or cylindrical, has a focal length f . The combination of a spherical and cylindrical lens in close contact has a focal length that is $f/2$ in the direction for which the cylinder has power, and f in the direction for which the cylinder has no power. Thus such a pair will collimate light diverging in the direction with weaker power, and image light diverging in the direction of stronger power. As a result, the lens combination L_1, L_2 collimates the light diverging vertically from an input source, but images in the horizontal direction, thereby illuminating a column of the matrix mask. Similarly, the lens combination L_3, L_4 images a row of the mask onto the vertical position of a single detector element, but collimates or spreads the light from a single column of the matrix mask. Ideally the detector elements should be long in the horizontal direction, allowing detection of most of the light across a row of the mask, but a long detector has high capacitance, and such capacitance limits the bandwidth of the electronic channel that follows.

The parallel matrix-vector multiplier performs all $N \times M$ multiplications and additions in a single clock cycle. A clock cycle can be very short, for example 10 nsec, depending on the amount of light available from each source. Lasers can be used as

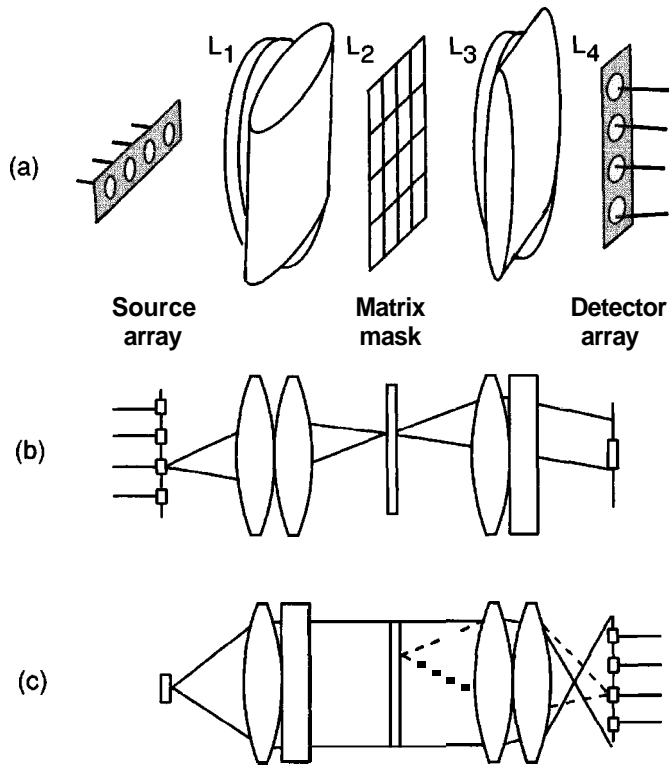


FIGURE 8.38
Optical elements comprising the parallel matrix-vector multiplier:
(a) perspective view, (b) top view, and (c) side view.

the sources, in spite of the fact that incoherent addition is used by this system, due to the fact that all additions are of light from different lasers, which for most types of semiconductor lasers will be mutually incoherent on the time scale of a clock cycle.

Many applications of the parallel matrix-vector multiplier architecture have been proposed and demonstrated. These include the construction of an optical crossbar switch [84], the iterative inversion of matrices [236], the construction of **Hopfield** neural networks [97], and others. The architecture is a useful workhorse of the optical information processing field.

8.11.4 An Outer Product Processor

A fundamentally different architecture for discrete operations is the *outer product processor* of Athale and Collins [11], which is a method for performing a matrix-matrix multiplication.

Suppose we wish to multiply two 3×3 matrices **A** and **B** to produce a product 3×3 matrix **C**, where

$$\mathbf{A} = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix} \quad \mathbf{B} = \begin{bmatrix} b_{11} & b_{12} & b_{13} \\ b_{21} & b_{22} & b_{23} \\ b_{31} & b_{32} & b_{33} \end{bmatrix} \quad \mathbf{C} = \begin{bmatrix} c_{11} & c_{12} & c_{13} \\ c_{21} & c_{22} & c_{23} \\ c_{31} & c_{32} & c_{33} \end{bmatrix}. \tag{8-94}$$

Straightforward manipulations show that it is possible to express **C** as a sum of outer products of the column vectors of **A** and the row vectors of **B** as follows:

$$\mathbf{C} = \begin{bmatrix} a_{11} \\ a_{21} \\ a_{31} \end{bmatrix} \begin{bmatrix} b_{11} & b_{12} & b_{13} \end{bmatrix} + \begin{bmatrix} a_{12} \\ a_{22} \\ a_{32} \end{bmatrix} \begin{bmatrix} b_{21} & b_{22} & b_{23} \end{bmatrix} + \begin{bmatrix} a_{13} \\ a_{23} \\ a_{33} \end{bmatrix} \begin{bmatrix} b_{31} & b_{32} & b_{33} \end{bmatrix}. \quad (8-95)$$

The sum of outer products can be achieved optically with the system shown in Fig. 8.39. Two two-dimensional SLMs are used, each operating as an array of one-dimensional SLMs. Lens L_0 collimates the light from the source S . That light is incident on a set of independently addressable SLM rows, the outputs of which are imaged by spherical lens L_1 onto a second SLM, consisting of a set of independently addressable SLM columns. Finally, lens L_2 images the light transmitted by the second SLM onto a two-dimensional array of time-integrating detectors.

The operation of the system in this simple example is described as follows. The first column vector of \mathbf{A} is entered into the first SLM and the first row vector of \mathbf{B} is entered into the second SLM. The detector array then stores charge proportional to the first of the individual outer products found in Eq. (8-95). The first SLM is now filled with the second column vector of \mathbf{A} and the second SLM with the second row vector of \mathbf{B} . The light incident on the detector array now adds charge proportional to the second outer product in Eq. (8-95). The process repeats one more time with the third column vector of \mathbf{A} and the third row vector of \mathbf{B} . The total stored charge array, now proportional to the elements of the product matrix \mathbf{C} , is read out of the detector.

Neglecting the time required to dump the detected charges, the speed of operation of this approach for a general product of an $N \times M$ matrix (i.e. with N rows and M columns) \mathbf{A} and an $M \times N$ matrix \mathbf{B} would be one cycle for each $N \times N$ outer product component, and M such cycles to accumulate the entire output matrix. During each cycle, N^2 multiplies and additions take place. For $N = M$, the degree of parallelism is similar to that of the parallel matrix-vector multiplier discussed earlier.

8.11.5 Other Discrete Processing Architectures

A multitude of other discrete processing architectures have been proposed and in some cases demonstrated in the past. We mention in particular the systolic approach of Caulfield et al. [55]. Also of interest is the systolic processor demonstrated by Guilfoyle [135], although this processor was aimed at "numerical" processing, rather than analog processing, and therefore is not within the realm of our coverage here.

8.11.6 Methods for Handling Bipolar and Complex Data

Until now we have assumed that the elements of both the input vector and the system matrix are nonnegative real numbers, an assumption that assures compatibility with the nonnegative real character of incoherent optical signals. To handle bipolar data or complex data, two different methods can be utilized, either individually or together. For the purposes of this discussion, we focus on the parallel matrix-vector multiplier, although the methods are more widely applicable.

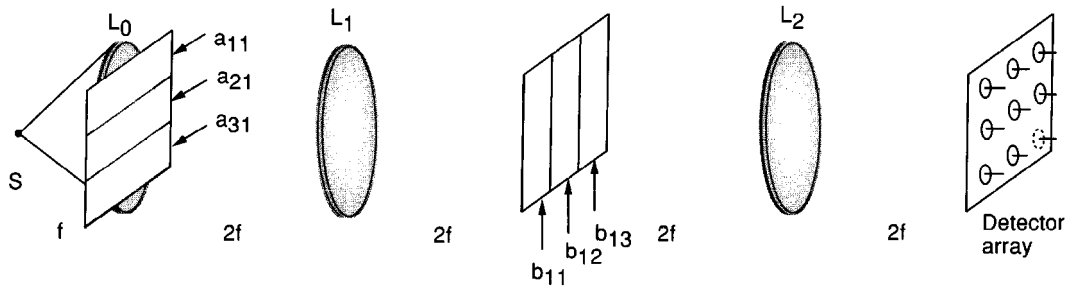


FIGURE 8.39
Outer product processor.

The first method places all bipolar signals on a bias, with the bias chosen large enough so that all elements of the input vector and all elements of the system matrix remain nonnegative. The biasing operation can be represented mathematically by noting that the input vector is now the sum of the signal vector \vec{s} and a bias vector \vec{b} (all elements assumed identical), and the system matrix is likewise the sum of two matrices, \mathbf{H} and \mathbf{B} , where the elements of the bias matrix are also assumed to be identical. The output of the system now becomes

$$(\mathbf{H} + \mathbf{B})(\vec{s} + \vec{b}) = \mathbf{H}\vec{s} + \mathbf{H}\vec{b} + \mathbf{B}\vec{s} + \mathbf{B}\vec{b}. \quad (8-96)$$

If the bias matrix and the bias vector are known and constant over time, then the last term can be subtracted from the output electronically. In addition, the matrix \mathbf{H} is known a priori, so the product $\mathbf{H}\vec{b}$ can be calculated in advance and subtracted from any result. However, the vector \vec{s} is not known in advance, and therefore it is generally necessary to measure its inner product with a row vector of the bias matrix, perhaps by adding a simple extra bias row to the matrix \mathbf{H} and one extra element to the detector array.

An alternative approach to handling bipolar elements is to represent the input vector and the system matrix as the *difference* of two nonnegative vectors or two nonnegative matrices, respectively. Thus $\mathbf{H} = \mathbf{H}_+ - \mathbf{H}_-$ and $\vec{s} = \vec{s}_+ - \vec{s}_-$, where the matrix \mathbf{H}_+ contains positive elements only in those locations where \mathbf{H} contains positive elements, and zero elsewhere, and \mathbf{H}_- contains positive elements equal to the magnitude of any negative elements of \mathbf{H} and zero for all other elements, with a similar construction procedure for \vec{s}_+ and \vec{s}_- . In addition, the output vector \vec{g} can be similarly decomposed. It is now easily shown that the nonnegative components of the output vector are related to the similar components of the input vector and the system matrix by

$$\begin{aligned} \vec{g}_+ &= \mathbf{H}_+ \vec{s}_+ + \mathbf{H}_- \vec{s}_- \\ \vec{g}_- &= \mathbf{H}_- \vec{s}_+ + \mathbf{H}_+ \vec{s}_-. \end{aligned} \quad (8-97)$$

A simpler way of stating this relation is to stack \vec{s}_+ and \vec{s}_- in a longer column vector, and to do the same for the two parts of \vec{g} , yielding

$$\begin{bmatrix} \vec{g}_+ \\ \vec{g}_- \end{bmatrix} = \begin{bmatrix} \mathbf{H}_+ & \mathbf{H}_- \\ \mathbf{H}_- & \mathbf{H}_+ \end{bmatrix} \begin{bmatrix} \vec{s}_+ \\ \vec{s}_- \end{bmatrix} \quad (8-98)$$

From this result we can see that a doubling of the two dimensions of the matrix mask to accommodate the larger matrix above, and a doubling of the length of the input vector, will allow the two components of \vec{g} to be computed without the use of biases. Those two output vectors must then be subtracted electronically, element by element.

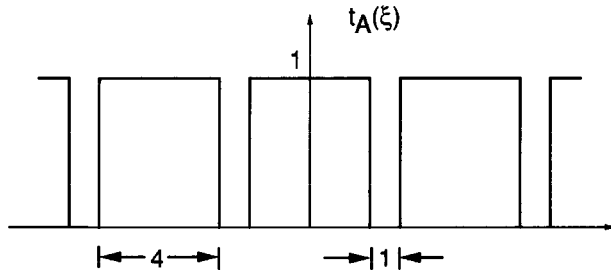
When complex elements of the input vector and matrix are important, then the most straightforward approach is to quadruple the dimensions of the input vector, the output vector, and the matrix, thus allowing positive and negative real parts and positive and negative imaginary parts to be handled properly. More efficient decompositions can also be found [130].

PROBLEMS-CHAPTER 8

8-1. An object has a periodic amplitude transmittance described by

$$t_A(\xi, \eta) = t_A(\xi) \cdot 1$$

where $t_A(\xi)$ is shown in Fig. P8.1. The object is placed in the object plane of the optical system shown in Fig. 8.1, and a tiny completely opaque stop is introduced on the optical axis in the focal plane, blocking only the spot on the optical axis. Sketch the intensity distribution observed in the image plane.



ξ **FIGURE P8.1**

8-2. The central dark ground method for observing phase objects is achieved by placing a tiny opaque stop on the optical axis in the focal plane to block the undiffracted light. Assuming that the variable component of phase shift through the object is always small compared with 2π radians, find the observed image intensity in terms of the object phase delay.

8-3. The schlieren method for observing phase objects is achieved by introduction of a knife edge in the focal plane to block half of the diffracted light. The amplitude transmittance through the focal plane may be written

$$t_f(x, y) = \frac{1}{2}(1 + \text{sgn}x).$$

(a) Assuming a magnification of unity and neglecting image inversion, show that the image amplitude U_i is related to the object amplitude U_o by

$$U_i(u, v) = \frac{1}{2} \left[U_o(u, v) + \frac{j}{\pi} \int_{-\infty}^{\infty} \frac{U_o(\xi, v)}{u - \xi} d\xi \right].$$

(b) Let the field transmitted by the object be of the form

$$U_o(\xi, \eta) = e^{j\phi_o} \exp[j\Delta\phi(\xi, \eta)]$$

where $\Delta\phi(\xi, \eta) \ll 2\pi$. Show that the image intensity can be approximated as

$$I_i(u, v) \approx \frac{1}{4} \left[1 - \frac{2}{\pi} \int_{-\infty}^{\infty} \frac{\Delta\phi(\xi, v)}{u - \xi} d\xi \right]$$

(c) Find and sketch the image intensity distribution when

$$\Delta\phi = \Phi \text{rect}\left(\frac{\xi}{U}\right)$$

with the constant $\Phi \ll 2\pi$.

- 8-4.** Find an expression for the image intensity observed when the phase-shifting dot of the Zernike phase-contrast microscope is also partially absorbing, with intensity transmittance equal to a ($0 < a < 1$).
- 8-5.** A certain coherent processing system has an input aperture that is 3 cm wide. The focal length of the initial transforming lens is 10 cm, and the wavelength of the light is $0.6328 \mu\text{m}$. With what accuracy must a frequency-plane mask be positioned in the focal plane, assuming that the mask has a structure comparable in scale size with the smallest structure in the spectrum of the input?
- 8-6.** It is desired to remove an additive periodic intensity interference of the form $I_N(\xi, \eta) = \frac{1}{2} [1 + \cos 2\pi f_o \xi]$ from a photograph taken by an imaging system. A coherent “4f” optical processing system will be used for that removal. The wavelength of the coherent light is λ . The image was recorded on photographic film (size $L \times L$) using the linear region of the **H&D** curve. A purely absorbing positive transparency with a photographic gamma of -2 was made, and that transparency is to be inserted in the input plane of the optical processing system. Specify the absorbing mask you would place in the frequency plane of the coherent optical processor in order to remove the interference. Consider especially:
- Where should the absorbing spots be placed?
 - What size would you make the absorbing spots?
 - What would you do at frequency ($f_X = 0, f_Y = 0$)?

Note: Neglect any effect the mask might have on the nonperiodic signal that is also present at the input.

- 8-7.** A grating with amplitude transmittance $t_A(x, y) = \frac{1}{2} [1 + \cos(2\pi f_o x)]$ is placed at the input to a standard “4f” coherent optical processing system of the kind illustrated in Fig. 8.10(a). Specify the transfer function (as a function of f_X) of a pure phase spatial filter that will completely suppress the spatial frequency component of output intensity having spatial frequency f_o . Assume normally incident plane wave illumination, monochromatic light, and neglect the effects of finite lens apertures.
- 8-8.** A transparent object with complex amplitude transmittance $t_A(x, y)$ is placed immediately in front of a positive spherical lens. The object is normally illuminated with a monochromatic plane wave, and a photographic transparency records the intensity distribution across the back focal plane. A positive transparency with a gamma of -2 is produced. The developed transparency is then illuminated by a plane wave, and the same positive lens is inserted directly behind the transparency. What is the relationship between the amplitude transmittance of the original object and the intensity distribution observed across the back focal plane of the lens in the second step of the process?
- 8-9.** A phase object with amplitude transmittance $t_A(x_1, y_1) = \exp[j\phi(x_1, y_1)]$ is present in the object plane of a coherent imaging system. In the back focal plane of the system, an attenuating plate (of uniform thickness) with intensity transmittance

$$\tau(x_2, y_2) = \alpha(x_2^4 + 2x_2^2 y_2^2 + y_2^4)$$

is introduced. How is the resulting image intensity related to the object phase?

- 8-10.** Consider the optical system shown in Fig. P8.10. A transparency with a real and non-negative *amplitude* transmittance $s_1(\xi, \eta)$ is placed in plane P_1 and coherently illuminated by a monochromatic, unit-intensity, normally incident plane wave. Lenses L_1 and L_2 are spherical with common focal length f . In plane P_2 , which is the rear focal plane of L_1 , a *moving* diffuser is placed. The effect of the moving diffuser can be considered to be the conversion of spatially coherent incident light into spatially incoherent transmitted light, without changing the intensity distribution of the light in plane P_2 and without appreciably broadening the spectrum of the light. In plane P_3 , in contact with L_2 , is placed a second transparency, this one with *amplitude* transmittance $s_2(x, y)$. Find an expression for the intensity distribution incident on plane P_4 .

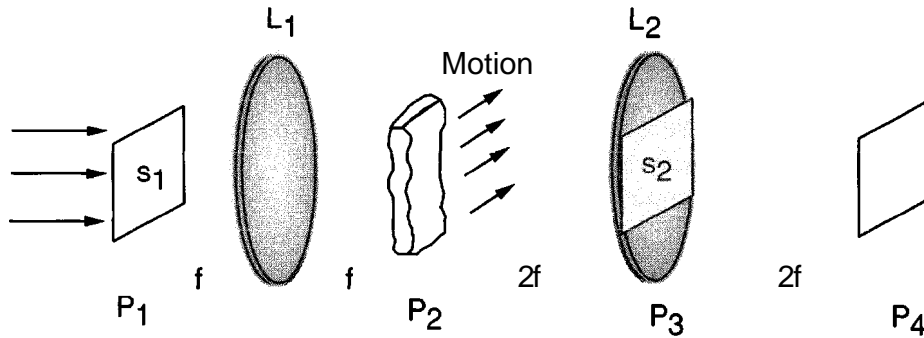


FIGURE P8.10

- 8-11.** The VanderLugt method is used to synthesize a frequency-plane filter. As shown in Fig. P8.11(a), a "signal" transparency with amplitude transmittance $s(x, y)$ is placed immediately against a positive lens (rather than in the front focal plane) and a photographic plate records the intensity in the back focal plane. The amplitude transmittance of the developed plate is made proportional to exposure, and the resulting transparency is placed in the system of part (b) of the figure. Assuming that the appropriate portions of the output

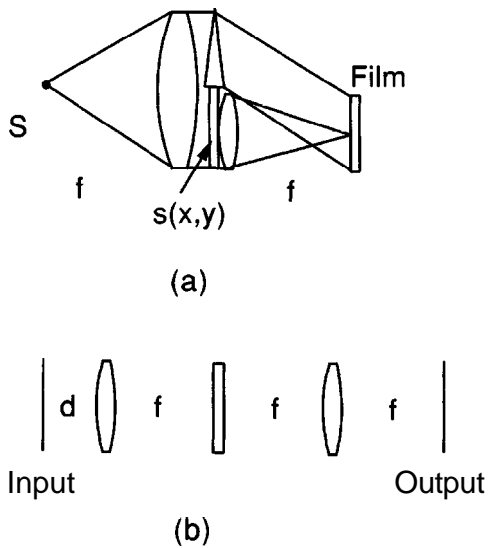


FIGURE P8.11

plane are examined in each case, what should the distance d between the object plane and the first lens of the filtering system be in order to synthesize:

- (a) A filter with impulse response $s(x, y)$?
 - (b) A filter with impulse response $s^*(-x, -y)$?
- 8-12.** Given a standard **VanderLugt** matched filtering system, prove that the output correlation spot shifts with any shift of the input signal against which we are correlating the reference. Assume that the magnification of the system is unity.
- 8-13.** Prove that the inequality of Eq. (8-25) must be satisfied if the various output terms of the joint transform correlator are to be separated.
- 8-14.** A certain image is blurred by camera motion such that the image incident on the recording film has moved linearly with velocity V on the film during a T -second exposure time.
- (a) Specify the point-spread function and the optical transfer function of the blur.
 - (b) Specify and plot the magnitude of the transfer function of an inverse filter that will in principle remove the blur.
 - (c) Assuming a constant ratio of signal power spectrum to noise power spectrum of 10, specify and plot the transfer function of a Wiener filter that will serve as a better deblurring filter than the simple inverse filter.
 - (d) If you have access to a computer, calculate and plot the impulse response of the Wiener filter of part (c).
- 8-15.** Consider an ideal "perfect" periodic transmitting object with amplitude transmittance $p(x, y)$ having period L in both the x and y directions. On this object there exists an opaque defect, with a size much smaller than the period of the object, but nonetheless much larger than the smallest structure contained in $p(x, y)$. We wish to create an optical filter that will enhance the brightness of the defect with respect to the periodic object, thus enhancing our ability to detect the defect. Ideally we would like to completely suppress the periodic portion of the image and pass only a bright image of the defect.
- (a) Describe how you might make a spatial filter that would accomplish the task described above. Be as specific as possible.
 - (b) Suppose that your filter were able to completely eliminate the discrete frequency components associated with the ideal periodic object, but also pass essentially all the light caused to leave these locations by the defect. Find an approximate expression for the image intensity that would be obtained at the output. As an aid to your analysis, let the amplitude transmittance of the defect be described by $1 - d(x, y)$, where the function $d(x, y)$ is unity within the defect and zero outside it. Remember that the defect and the periodic object should be treated as two separate diffracting structures in close contact with one another. You may neglect the finite sizes of the lenses and any vignetting effects.

- 8-16.** You are to construct a coherent optical "grade change" filter that will change the letter F into the letter A. The filtering system is of the standard " $4f$ " type. Describe in detail how you would construct such a filter. Be specific. How would you expose the photographic plate in making the filter? What behavior of the photographic transparency would you try to achieve? Where would you look in the output plane of the processor? Give as many details as you can.
- 8-17.** A synthetic-aperture radar system carries an antenna of dimension D along the direction of the flight path. Therefore on any particular transmitted pulse, that antenna illuminates a patch on the terrain that is of approximate width $\lambda_r r/D$ where r is the slant range.
- Approximately how long (in the direction of the flight path) is the synthetic aperture that such a system can generate?
 - If the radar is transmitting a microwave frequency f_r and moving with linear velocity v_a , find the frequency of the radiation returned from a fixed target and received in the aircraft, first when the target initially enters the radiation pattern of the transmitting/receiving antenna, and second when it is just leaving that radiation pattern. From these results deduce the total bandwidth of received radiation seen by the radar from one target.
 - Taking into account the fact that both the transmitter and the receiver are moving, show that the approximate size of a minimum resolvable spot achieved in the image of the terrain (in the direction of the flight path) can be as small as approximately $D/2$, and therefore the smaller the antenna carried, the better the resolution of the system (neglecting noise).
- 8-18.** With reference to Eq. (7-34), show that the optical frequency of the light at coordinate y_2 in the spatial frequency domain of the Bragg cell spectrum analyzer is offset from the optical frequency of the source by an amount that is exactly equal to the temporal frequency of the RF spectral component represented at that coordinate.

Holography

In 1948, Dennis **Gabor** [106] proposed a novel two-step, **lensless** imaging process which he called *wavefront reconstruction* and which we now know as *holography*. **Gabor** recognized that when a suitable coherent reference wave is present simultaneously with the light diffracted by or scattered from an object, then information about both the amplitude and phase of the diffracted or scattered waves can be recorded, in spite of the fact that recording media respond only to light intensity. He demonstrated that, from such a recorded interference pattern (which he called a *hologram*, meaning a "total recording"), an image of the original object can ultimately be obtained.

While **Gabor's** imaging technique received only mild interest in its early days, the 1960s saw dramatic improvements in both the concept and the technology, improvements that vastly extended its applicability and practicality. In 1971 **Gabor** received the Nobel prize in physics for his invention.

In this chapter we examine the basic principles behind holography, explore the many modern variations upon **Gabor's** original theme, and survey some of the important applications that have been found for this novel imaging technique. Several excellent books devoted to holography exist. The classic text is that of Collier, Burckhardt, and Lin [70]. For another excellent and authoritative treatment see the book by **Hariharan** [139]. Other broad books include those by Smith [265], Develis and Reynolds [83], Caulfield [53], and Saxby [254].

9.1 HISTORICAL INTRODUCTION

Gabor was influenced in his early studies of holography by previous work of W.L. Bragg in X-ray crystallography (see, for example, [34]), but was primarily motivated by possible applications of his newfound technique to electron holography. **Gabor** followed

his original proposal with two more lengthy papers ([107], [108]) published in 1949 and 1951, considering the possible application of holography to microscopy. While for practical reasons he was unable to realize his envisioned application, the improvements developed in the 1960s led to many applications that **Gabor** could not possibly have foreseen.

In the 1950s, a number of authors, including G.L. Rogers [245], H.M.A. El-Sum [93], and A.W. Lohmann [197], significantly extended the theory and understanding of holography. It was not, however, until the early 1960s that a revolution in holography began. Again it was workers at the University of Michigan's Radar Laboratory, in particular E.N. Leith and J. Upatnieks [188], who recognized the similarity of **Gabor's** lensless imaging process to the synthetic-aperture-radar problem and suggested a modification of his original technique that greatly improved the process. At virtually the same time, Y.N. Denisyuk [82], working in what was then the Soviet Union, created a remarkable synthesis of the ideas of both **Gabor** and French physicist G. Lippmann to invent the thick reflection hologram, which he perfected to an advanced state.

The Michigan workers soon coupled their new developments with the emerging technology of lasers in order to perform lensless three-dimensional photography [190]. The quality and realism of the three-dimensional images obtained by holography were largely responsible for the development of a great popular interest in the field. Today it is common to find museums or galleries specializing in holography in many of the great cities of the world. However, contrary to popular impression, many of the most interesting and useful properties of holography are quite independent and separate from the three-dimensional imaging capability, as we shall see in some detail in later sections.

9.2 THE WAVEFRONT RECONSTRUCTION PROBLEM

The fundamental problem addressed by holography is that of recording, and later reconstructing, both the amplitude and the phase of an optical wave arriving from a coherently illuminated object. This problem is sufficiently general to be of interest for electromagnetic waves in all regions of the spectrum, as well as for acoustic and seismic waves. Our considerations here, however, will be largely restricted to the optical problem.

9.2.1 Recording Amplitude and Phase

As indicated above, the wavefront-reconstruction problem must consist of two distinct operations: a recording or detection step, and a reconstruction step. For the moment we focus on the first of these two operations.

Since the wavefronts of concern are coherent, it is necessary to detect information about both the amplitude and phase of the waves. However, all recording media respond only to light intensity. It is therefore required that the phase information somehow be converted to intensity variations for recording purposes. A standard technique for

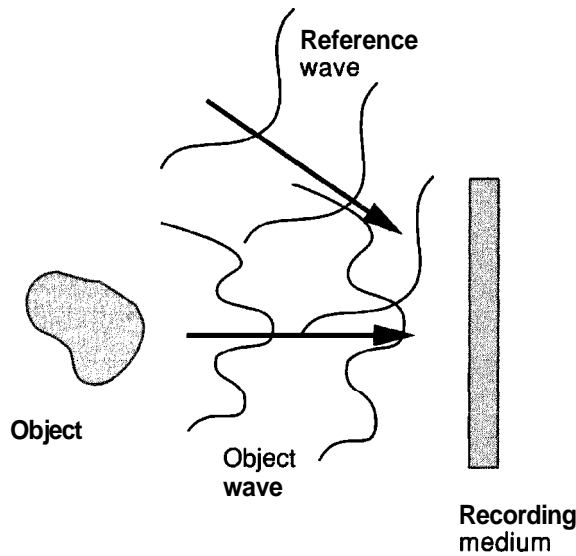


FIGURE 9.1
Interferometric recording.

accomplishing this task is *intetferometry*; that is, a second wavefront, mutually coherent with the first and of known amplitude and phase, is added to the unknown wavefront, as shown in Fig. 9.1. The intensity of the sum of two complex fields then depends on both the amplitude and phase of the unknown field. Thus if

$$a(x, y) = |a(x, y)| \exp[-j\phi(x, y)] \quad (9-1)$$

represents the wavefront to be detected and reconstructed, and if

$$A(x, y) = |A(x, y)| \exp[-j\psi(x, y)] \quad (9-2)$$

represents the "reference" wave with which $a(x, y)$ interferes, the intensity of the sum is given by

$$\mathcal{I}(x, y) = |A(x, y)|^2 + |a(x, y)|^2 + 2|A(x, y)||a(x, y)| \cos[\psi(x, y) - \phi(x, y)]. \quad (9-3)$$

While the first two terms of this expression depend only on the intensities of the individual waves, the third depends on their relative phases. Thus information about both the amplitude and phase of $a(x, y)$ has been recorded. The issue as to whether it is sufficient information to reconstruct the original wavefront remains to be dealt with. At this point we have not specified any detailed character of the reference wave $A(x, y)$. Properties that the reference wave must satisfy in order to enable reconstruction of $a(x, y)$ will become evident as the discussion progresses. The recording of the pattern of interference between an "object" wave and a "reference" wave may be regarded as a hologram.

9.2.2 The Recording Medium

The material used to record the pattern of interference described above will be assumed to provide a linear mapping of intensity incident during the detection process into amplitude transmitted by or reflected from the material during the reconstruction process. Usually both light detection and wavefront modulation are performed by photographic

film or plate. The linear relation required is then provided by operation in the linear portion of the t_A vs. E curve of the emulsion. However, many other materials suitable for holography exist, including photopolymers, dichromated gelatin, photorefractive materials, and others (see Section 9.8). It is even possible to detect the interference pattern electronically and reconstruct the wavefront with a digital computer. However, photographic materials remain the most important and widely used recording medium.

Thus we assume that the variations of exposure in the interference pattern remain within a linear region of the t_A vs. E curve. In addition, it is assumed that the MTF of the recording material extends to sufficiently high spatial frequencies to record all the incident spatial structure (effects of removing some of these ideal assumptions are examined in Section 9.10). Finally we assume that the intensity $|A|^2$ of the reference wave is uniform across the recording material, in which case the amplitude transmittance of the developed film or plate can be written

$$t_A(x, y) = t_b + \beta' (|a|^2 + A^*a + Aa^*), \quad (9-4)$$

where t_b is a uniform "bias" transmittance established by the constant reference exposure, and β' is the product of the slope β of the t_A vs. E curve at the bias point and the exposure time. Note that, as in Section 7.1, β' is a negative number for a negative transparency, and a positive number for a positive transparency.

9.2.3 Reconstruction of the Original Wavefront

Once the amplitude and phase information about the object wave $a(x, y)$ have been recorded, it remains to reconstruct that wave. Suppose that the developed transparency is illuminated by a coherent *reconstruction* wave $B(x, y)$. The light transmitted by the transparency is evidently

$$\begin{aligned} B(x, y)t_A(x, y) &= t_bB + \beta'aa^*B + \beta'A^*Ba + \beta'ABa^* \\ &= U_1 + U_2 + U_3 + U_4. \end{aligned} \quad (9-5)$$

Note that if B is simply an exact duplication of the original uniform reference wavefront A , the third term of this equation becomes

$$U_3(x, y) = \beta'|A|^2a(x, y). \quad (9-6)$$

Since the intensity of the reference wave is uniform, it is clear that reconstructed wave component U_3 is, up to a multiplicative constant, an exact duplication of the original wavefront $a(x, y)$, as shown in Fig. 9.2(a).

In a similar fashion, if $B(x, y)$ happens to be chosen as the *conjugate* of the original reference wave, i.e. as $A^*(x, y)$, the fourth term of the reconstructed field becomes

$$U_4(x, y) = \beta'|A|^2a^*(x, y), \quad (9-7)$$

which is proportional to the *conjugate* of the original wavefront. This case is illustrated in Fig. 9.2(b).

Note that in either case, the particular field component of interest (that is, U_3 when $B = A$ and U_4 when $B = A^*$) is accompanied by three additional field components,

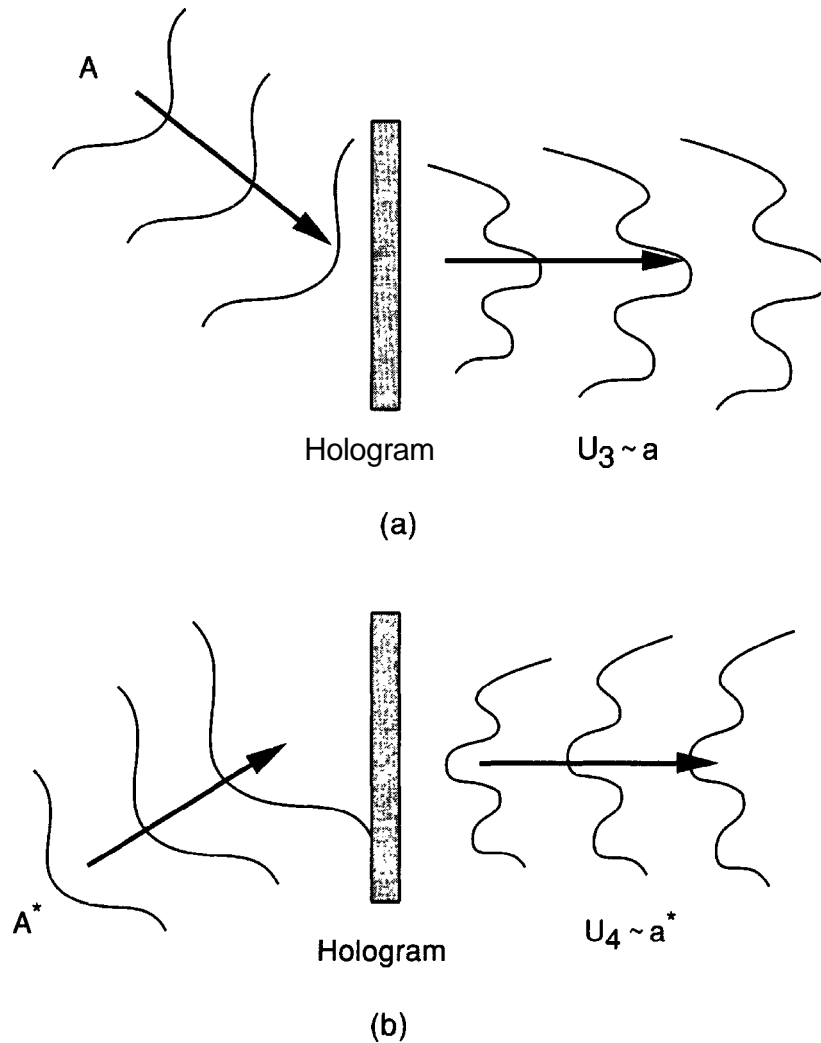


FIGURE 9.2
Wavefront reconstruction with (a) the original reference wave A as illumination, and (b) the conjugate reference wave A^* as illumination.

each of which may be regarded as extraneous interference. Evidently, if a usable duplication of the object wave $a(x, y)$ (or of $a^*(x, y)$) is to be obtained, some method for separating the various wave components of transmitted light is required.

9.2.4 Linearity of the Holographic Process

The characteristic behavior hypothesized for the recording material in Eq. (9-4) corresponds to a highly nonlinear mapping of fields incident during exposure into fields transmitted after development. It would therefore appear, at first glance, that linear systems concepts can play no role in the theory of holography. While the overall mapping introduced by the film is nonlinear, nonetheless the mapping of object field $a(x, y)$ into

the transmitted field component $U_3(x, y)$ is entirely linear, as evidenced by the proportionality of Eq. (9-6). Similarly, the mapping of $a(x, y)$ into the transmitted field component $U_4(x, y)$, as represented by Eq. (9-7), is a linear one. Thus if the object field $a(x, y)$ is regarded as an input, and the transmitted field component $U_3(x, y)$ (or $U_4(x, y)$) is regarded as an output, the system so defined is a linear one. The nonlinearity of the detection process manifests itself in the generation of several output terms, but there is no nonlinear distortion of the one term of interest, assuming that the exposure variations remain in the linear region of the t_A vs. E curve.

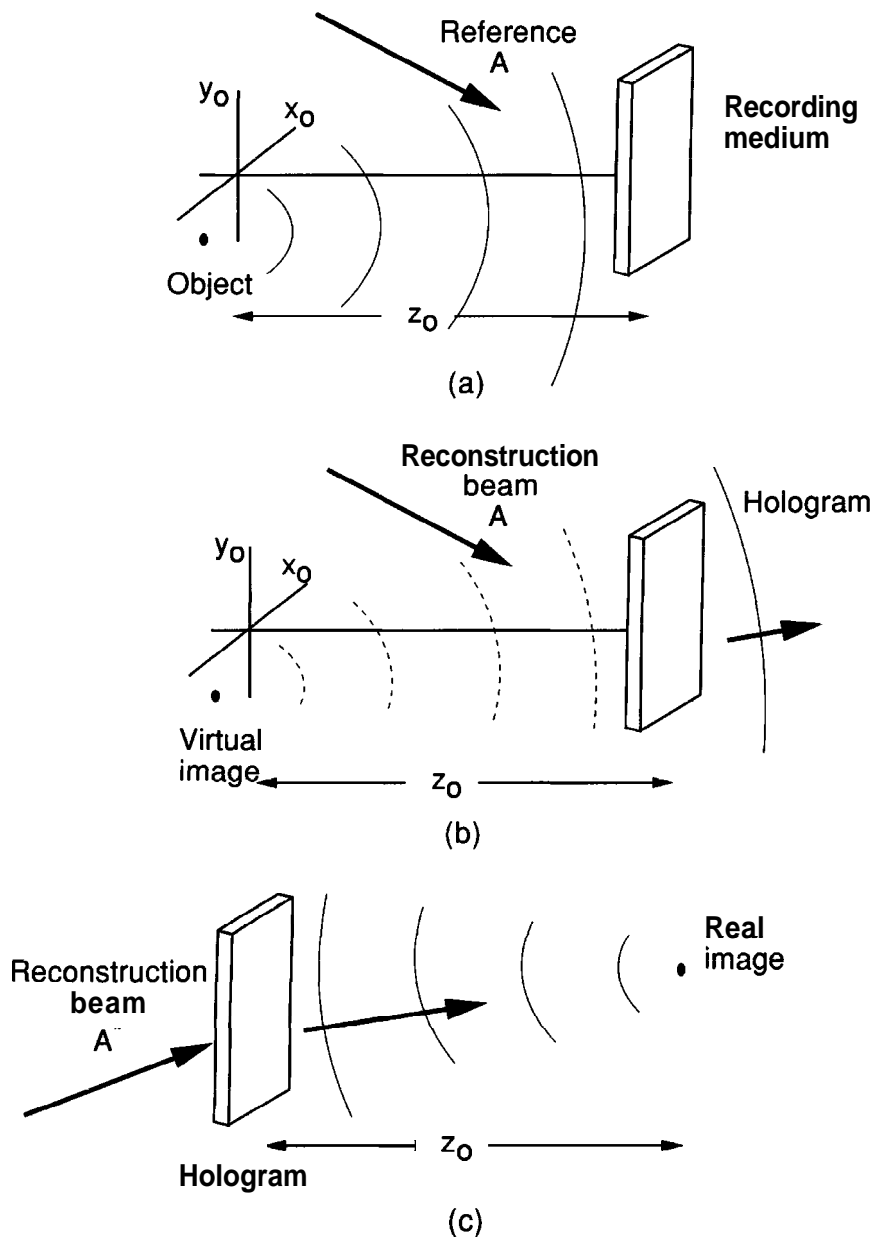


FIGURE 9.3

Imaging by wavefront reconstruction. (a) Recording the hologram of a point-source object; (b) generation of the virtual image; (c) generation of the real image.

9.2.5 Image Formation by Holography

To this point we have considered only the problem of reconstructing a wavefront which arrived at a recording medium from a coherently illuminated object. It requires but a small change in point of view to regard the wavefront reconstruction process as a means of imageformation.

To adopt this point of view, note that the wave component $U_3(x, y)$ of Eq. (9-6), being simply a duplication of the original object wavefront $a(x, y)$, must appear to the observer to be diverging from the original object, in spite of the fact that the object has long since been removed. Thus when the reference wave $A(x, y)$ is used as the illumination during reconstruction, the transmitted wave component $U_3(x, y)$ may be regarded as generating a virtual image of the original object. This case is illustrated in Fig. 9.3(a),(b) for the particular case of a simple point-source object.

In a similar fashion, when the conjugate of the reference wave, $A^*(x, y)$, is used as the illumination during reconstruction, the wave component $U_4(x, y)$ of Eq. (9-7) also generates an image, but this time it is a real image which corresponds to an actual focusing of light in space. To prove this assertion, we invoke the linearity property discussed above, considering an object which consists of a single point source. The corresponding result for a more complicated object may then be found by linear superposition of point-source solutions.

Incident on the recording medium we have the sum of the reference wave $A(x, y)$ and a simple spherical object wave,

$$a(x, y) = a_o \exp \left[jk \sqrt{z_o^2 + (x - \hat{x}_o)^2 + (y - \hat{y}_o)^2} \right] \quad (9-8)$$

where (\hat{x}_o, \hat{y}_o) are the (x, y) coordinates of the object point, and z_o is its normal distance from the recording plane. Illuminating the developed hologram with a reconstruction wave $A^*(x, y)$, we obtain the transmitted wave component

$$\begin{aligned} U_4(x, y) &= \beta' |A|^2 a^*(x, y) \\ &= \beta' |A|^2 a_o^* \exp \left[-jk \sqrt{z_o^2 + (x - \hat{x}_o)^2 + (y - \hat{y}_o)^2} \right], \end{aligned} \quad (9-9)$$

which is a spherical wave that converges towards a real focus at distance z_o to the right of the hologram, as shown in Fig. 9.3(c). A more complicated object may be considered to be a multitude of point sources of various amplitudes and phases; and by the linearity property, each such point source generates its own real image as above. Thus a real image of the entire object is formed in this fashion.

Note that the amplitude of the wave described by Eq. (9-9) is proportional to a_o^* , the conjugate of the original object point-source amplitude. Similarly, for a more complicated object, the real image generated by the hologram is always the complex conjugate of the original object amplitude. Such a change of phase does not affect image intensity, but it can be important in certain applications that utilize both the amplitude and phase of the image.

It should again be emphasized that we have considered, in each case, only one of the four wave components transmitted by the hologram. This approach is acceptable if,

by proper choice of reference wave, the undesired components are suppressed or are separated from the image of interest. When this is not the case, the interference of the various components of transmitted light must be taken into account.

9.3 THE GABOR HOLOGRAM

Keeping in mind the preceding general discussion, we now consider the **wavefront-reconstruction** process in the form originally proposed and demonstrated by Gabor. In Section 9.4, we turn to modifications of the process which improve its imaging capabilities.

9.3.1 Origin of the Reference Wave

The geometry required for recording a *Gabor hologram* is illustrated in Fig. 9.4. The object is assumed to be highly transmissive, with an amplitude transmittance

$$t(x_o, y_o) = t_o + \Delta t(x_o, y_o), \quad (9-10)$$

where t_o is a high average level of transmittance, Δt represents the variations about this average, and

$$|\Delta t| \ll |t_o|. \quad (9-11)$$

When such an object is coherently illuminated by the collimated wave shown in Fig. 9.4, the transmitted light consists of two components: (1) a strong uniform plane wave passed by the term t_o , and (2) a weak scattered wave generated by the transmittance variations $\Delta t(x_o, y_o)$. The intensity of the light incident on the recording medium at distance z_o from the object may be written

$$\begin{aligned} \mathcal{I}(x, y) &= |A + a(x, y)|^2 \\ &= |A|^2 + |a(x, y)|^2 + A^* a(x, y) + A a^*(x, y), \end{aligned} \quad (9-12)$$

where A is the amplitude of the plane wave, and $a(x, y)$ is the amplitude of the scattered light at the recording plane.

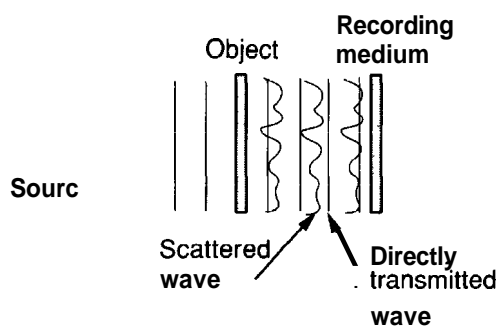


FIGURE 9.4
Recording a Gabor hologram.

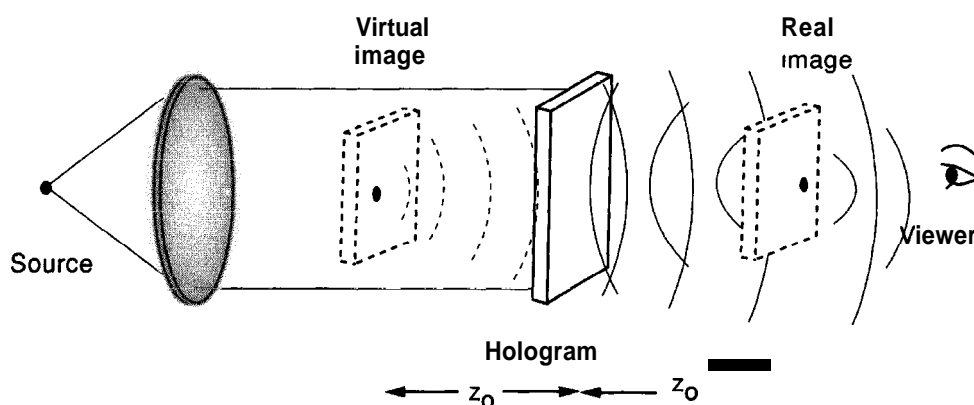


FIGURE 9.5
Formation of twin images from a Gabor hologram.

Thus the object has, in a sense, supplied the required reference wave itself through the high average transmittance t_o . The interference of the directly transmitted light with the scattered light results in a pattern of intensity that depends on both the amplitude and the phase of the scattered wave $a(x, y)$.

9.3.2 The Twin Images

The developed hologram is assumed to have an amplitude transmittance that is proportional to exposure. Thus

$$t_A(x, y) = t_b + \beta' (|a|^2 + A^* a + A a^*). \quad (9-13)$$

If the transparency is now illuminated by a normally incident plane wave with uniform amplitude B , the resulting transmitted field amplitude consists of a sum of four terms:

$$B t_A = B t_b + \beta' B |a(x, y)|^2 + \beta' A^* B a(x, y) + \beta' A B a^*(x, y). \quad (9-14)$$

The first term is a plane wave which passes directly through the transparency, suffering uniform attenuation but without scattering. The second term may be dropped as negligible by virtue of our assumption (9-11), which implies that

$$|a(x, y)| \ll A. \quad (9-15)$$

The third term represents a field component that is proportional to the original scattered wave $a(x, y)$. This wave appears to originate from a virtual image of the original object located at distance z_o from the transparency, as shown in Fig. 9.5. Similarly, the fourth term is proportional to $a^*(x, y)$ and, in accord with our earlier discussions, leads to the formation of a real image at distance z_o on the opposite side of the transparency from the virtual image (again, see Fig. 9.5).

Thus the Gabor hologram generates simultaneous real and virtual images of the object transmittance variations $A t$, both images being centered on the hologram axis. These so-called twin images are separated by the axial distance $2z_o$, and are accompanied by a coherent background $B t_b$.

Note from Eq. (9-14) that positive and negative transparencies yield different signs for the image-forming waves with respect to the background (β' is positive for a positive transparency and negative for a negative transparency). In addition, for any one of these two cases, the real image wave is the conjugate of the virtual image wave, and depending on the phase structure of the object, further contrast reversals are possible when one of these waves interferes with the constant background. For an object with constant phase, a positive hologram transparency is found to produce a positive image, and a negative hologram transparency is found to produce a negative image.

9.3.3 Limitations of the Gabor Hologram

The **Gabor** hologram is found to suffer from certain limitations which restrict the extent of its applicability. Perhaps the most important limitation is inherent in the assumption of a highly transparent object and the consequent conclusion (9-15) that followed. If this assumption is not adopted, there exists an additional wave component

$$U_2(x, y) = \beta' B |a(x, y)|^2 \quad (9-16)$$

transmitted by the hologram which can no longer be dropped as negligible. In fact, if the object is of low average transmittance, this particular wave component may be the largest transmitted term, and as a consequence may entirely obliterate the weaker images. Thus with a **Gabor** hologram it is possible to image an object consisting of, for example, opaque letters on a transparent background, but not transparent letters on an opaque background. This restriction seriously hampers the use of **Gabor** holograms in many potential applications.

A second serious limitation lies in the generation of overlapping twin images, rather than a single image. The problem lies not with the presence of twin images per se, but rather with their inseparability. When the real image is brought to focus, it is always accompanied by an out-of-focus virtual image. Likewise an observer viewing the virtual image sees simultaneously a defocused image arising from the real-image term. Thus, even for highly transparent objects, the quality of the images is reduced by the twin image problem. A number of methods have been proposed for eliminating or reducing the twin-image problem, (e.g. see [197]), including one technique originated by **Gabor** himself [111]. The most successful of these methods has been that of Leith and Upatnieks [188], which we discuss in detail in the next section.

9.4 THE LEITH-UPATNIEKS HOLOGRAM

Leith and Upatnieks suggested and demonstrated a modification of **Gabor's** original recording geometry that solved the twin image problem and vastly extended the applicability of holography. This type of hologram will be called the *Leith-Upatnieks* hologram, and is also known as an *onset-reference* hologram. The major change between this type of hologram and the **Gabor** hologram is that, rather than depending on the light directly transmitted by the object to serve as a reference wave, a separate and

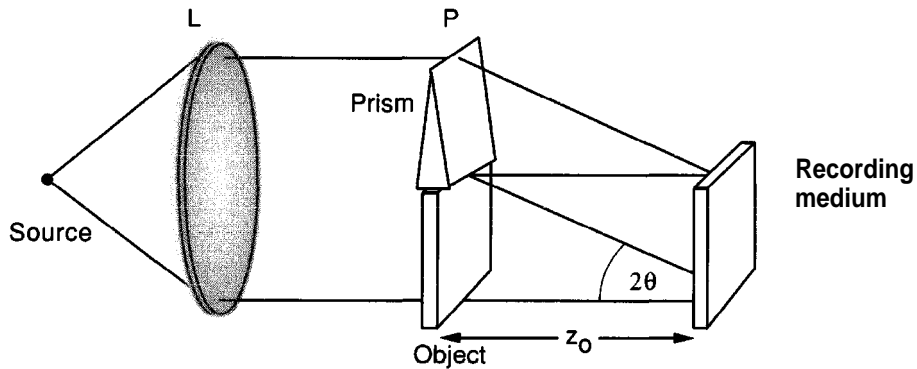


FIGURE 9.6
Recording a Leith-Upatnieks hologram.

distinct reference wave is introduced. Furthermore the reference is introduced at an offset angle, rather than being collinear with the object-film axis.

The first successful demonstration of this type of hologram, reported in [188], was carried out without a laser source. However, it was not until the technique was combined with highly coherent laser illumination that its full potential became evident [190], [189].

9.4.1 Recording the Hologram

One possible geometry for recording a Leith-Upatnieks hologram is illustrated in Fig. 9.6. The light from a point source of illumination is collimated by the lens L. A portion of the resulting plane wave strikes the object, which is taken to be a transparency with a general amplitude transmittance $t(x_o, y_o)$. A second portion of the plane wave strikes a prism P located above the object and is deflected downwards at angle 2θ with respect to the normal to the recording plane.¹ Thus at the recording surface we find the sum of two mutually coherent waves, one consisting of light transmitted by the object, and the other consisting of a tilted plane wave. The amplitude distribution incident on the recording plane may be written

$$U(x, y) = A \exp(-j2\pi\alpha y) + a(x, y), \quad (9-17)$$

where the spatial frequency α of the reference wave is given by

$$\alpha = \frac{\sin 2\theta}{\lambda}. \quad (9-18)$$

The intensity distribution across the recording plane is evidently

$$\begin{aligned} \mathcal{I}(x, y) = & |A|^2 + |a(x, y)|^2 \\ & + A^* a(x, y) \exp(j2\pi\alpha y) + A a^*(x, y) \exp(-j2\pi\alpha y). \end{aligned} \quad (9-19)$$

¹The reason for calling this angle 2θ rather than θ will become evident when we consider fringe orientation through the depth of a thick emulsion.

An alternative more revealing form may be obtained by writing $a(x, y)$ explicitly as an amplitude and phase distribution,

$$a(x, y) = |a(x, y)| \exp[-j\phi(x, y)] \quad (9-20)$$

and combining the last two terms of (9-19) to yield

$$\mathcal{I}(x, y) = |A|^2 + |a(x, y)|^2 + 2|A||a(x, y)| \cos[2\pi\alpha y - \phi(x, y)]. \quad (9-21)$$

This expression demonstrates that the amplitude and phase of the light arriving from the object have been recorded, respectively, as amplitude and phase modulations of a spatial carrier of frequency α . If the carrier frequency is sufficiently high (we shall see shortly just how high it must be), the amplitude and phase distributions can be unambiguously recovered from this pattern of interference.

9.4.2 Obtaining the Reconstructed Images

In the usual fashion, the photographic plate is developed to yield a transparency with an amplitude transmittance proportional to exposure. Thus the film transmittance may be written

$$t_A(x, y) = t_b + \beta' [|a(x, y)|^2 + A^* a(x, y) \exp(j2\pi\alpha y) + A a^*(x, y) \exp(-j2\pi\alpha y)]. \quad (9-22)$$

For convenience we represent the four terms of transmittance by

$$\begin{aligned} t_1 &= t_b & t_3 &= \beta' A^* a(x, y) \exp(j2\pi\alpha y) \\ t_2 &= \beta' |a(x, y)|^2 & t_4 &= \beta' A a^*(x, y) \exp(-j2\pi\alpha y). \end{aligned} \quad (9-23)$$

For the present we assume that the hologram is illuminated by a normally incident, uniform plane wave of amplitude B , as illustrated in Fig. 9.7. The field transmitted by the hologram has four distinct components, each generated by one of the transmittance terms of Eq. (9-23):

$$\begin{aligned} U_1 &= t_b B & U_3 &= \beta' B A^* a(x, y) \exp(j2\pi\alpha y) \\ U_2 &= \beta' B |a(x, y)|^2 & U_4 &= \beta' B A a^*(x, y) \exp(-j2\pi\alpha y). \end{aligned} \quad (9-24)$$

The field component U_1 is simply an attenuated version of the incident reconstruction illumination, and therefore represents a plane wave traveling down the optical axis. The second term U_2 is spatially varying and therefore has plane wave components traveling at various angles with respect to the optical axis. However, as we shall see in more detail shortly, if the bandwidth of $a(x, y)$ is sufficiently small compared with the carrier frequency α , the energy in this wave component remains sufficiently close to the optical axis to be spatially separated from the images of interest.

The wave component U_3 is proportional to the original object wavefront multiplied by a linear exponential factor. Proportionality to a implies that this term generates a virtual image of the object at distance z_o to the left of the transparency, while the linear exponential factor $\exp(j2\pi\alpha y)$ indicates that this image is deflected away from the optical axis at angle 2θ , as shown in Fig. 9.7. Similarly, wave component U_4 is

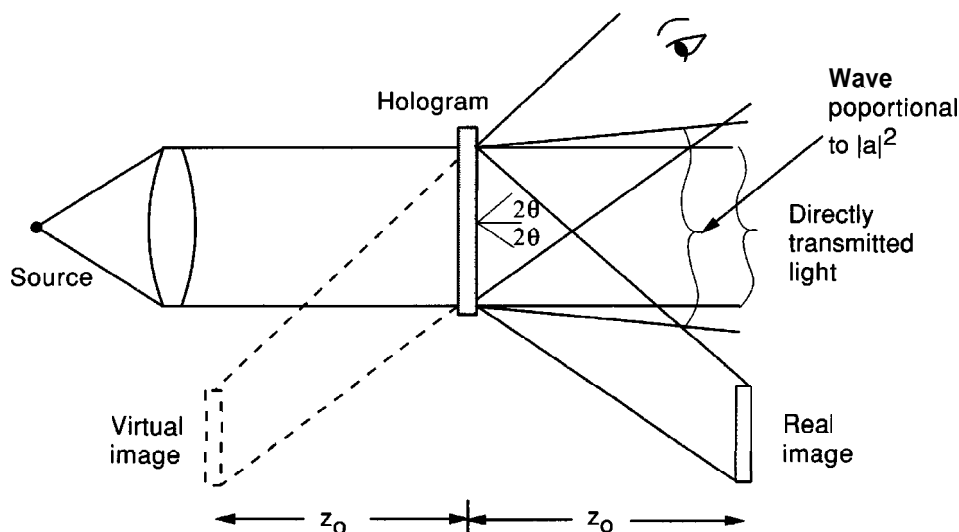


FIGURE 9.7
Reconstruction of images from a Leith-Upatnieks hologram.

proportional to the conjugate wavefront a^* , which indicates that a real image forms at distance z_o to the right of the transparency. The presence of the linear exponential factor $\exp(-j2\pi\alpha y)$ indicates that the real image is deflected at angle -2θ from the optical axis, as again can be seen in Fig. 9.7.

The most important observation to be derived from these results is that, while twin images are again generated by the holographic process, they have been angularly separated from each other and from the wave components U_1 and U_2 . This separation comes about due to the use of a reference wave with an angular offset; indeed, successful isolation of each of the twin images requires the use of an angle between object and reference which is chosen larger than some lower limit (the minimum reference angle will be discussed in more detail shortly). When this angle exceeds the minimum allowable angle, the twin images are not contaminated by each other nor by other wave components.

Note in addition that since the images may be viewed without the presence of a coherent background generated by the object transparency, the particular sign associated with the wave components U_3 and U_4 of Eq. (9-24) is immaterial. The transparency may be either a positive or a negative; in each case a positive image is obtained. For practical reasons it is generally preferable to use negatives directly, thus avoiding the two-step process usually required for making a positive transparency.

Finally we should point out that we have chosen to illuminate the hologram with a normally incident plane wave, which is neither a duplication of the original reference wave nor its complex conjugate, yet we have obtained a real and a virtual image simultaneously. Evidently our conditions concerning the required nature of the reconstruction illumination were overly restrictive. However, when we consider the effects of the thickness of the emulsion on the reconstructed wavefronts, the exact nature of the reconstruction illumination will become more important. As will be discussed in Section 9.7, it then becomes critical that the hologram be illuminated with a duplicate of the original reference wave to obtain one image, and the complex conjugate of the reference wave to obtain the other image.

9.4.3 The Minimum Reference Angle

Returning to the reconstruction geometry of Fig. 9.7, if the twin images are to be separated from each other and from the light transmitted with directions close to the optical axis, the offset angle 2θ of the reference beam with respect to the object beam must be greater than some minimum angle $2\theta_{\min}$. To find this minimum, it suffices to determine the minimum carrier frequency a for which the spatial frequency spectra of t_3 and t_4 (that is the virtual-image and real-image terms of hologram transmittance) do not overlap each other and do not overlap the spectra of t_1 and t_2 . If there is no overlap, then in principle the hologram amplitude transmittance can be Fourier transformed with the help of a positive lens, the unwanted spectral components can be removed with appropriate stops in the focal plane, and a second Fourier transformation can be performed to yield just that portion of the transmitted light that leads to the twin images.²

Consider the spatial frequency spectra of the various terms of transmittance listed in Eq. (9-23). Neglecting the finite extent of the hologram aperture, we have directly that

$$G_1(f_X, f_Y) = \mathcal{F}\{t_1(x, y)\} = t_b \delta(f_X, f_Y). \quad (9-25)$$

Using the autocorrelation theorem, we also have

$$G_2(f_X, f_Y) = \mathcal{F}\{t_2(x, y)\} = \beta' G_a(f_X, f_Y) \star G_a(f_X, f_Y) \quad (9-26)$$

where $G_a(f_X, f_Y) = \mathcal{F}\{a(x, y)\}$ and the \star indicates the autocorrelation operation. Finally we have

$$\begin{aligned} G_3(f_X, f_Y) &= \mathcal{F}\{t_3(x, y)\} = \beta' A^* G_a(f_X, f_Y - a) \\ G_4(f_X, f_Y) &= \mathcal{F}\{t_4(x, y)\} = \beta' A G_a^*(-f_X, -f_Y - a). \end{aligned} \quad (9-27)$$

Now note that the bandwidth of G_a is identical with the bandwidth of the object, for the two spectra differ only by the transfer function of the propagation phenomenon, which (neglecting the evanescent wave cutoff) is the pure phase function of Eq. (3-70). Suppose that the object has no spatial frequency components higher than B cycles/mm. Thus the spectrum $|G_a|$ might be as shown in Fig. 9.8(a). The corresponding spectrum of the hologram transmittance is illustrated in Fig. 9.8(b). The term $|G_1|$ is simply a δ function at the origin in the (f_X, f_Y) plane. The term $|G_2|$, being proportional to the autocorrelation function of $|G_a|$, extends to frequencies as high as $2B$. Finally, $|G_3|$ is simply proportional to $|G_a|$, displaced to a center frequency $(0, a)$, while $|G_4|$ is proportional to a reflected version of $|G_a|$ centered at frequency $(0, -a)$.

Examination of Fig. 9.8(b) shows that $|G_3|$ and $|G_4|$ can be isolated from $|G_2|$ if

$$\alpha \geq 3B \quad (9-28)$$

or equivalently if

$$\sin 2\theta \geq 3Bh. \quad (9-29)$$

²Spatial filtering operations are seldom used in practice to separate the twin images. If the reference angle satisfies the requirements to be derived here, the images will separate of their own accord due to the different directions of propagation of the respective wave components (cf. Fig. 9.7). However, spatial-filtering arguments do provide a conceptually simple way of finding sufficient conditions for separation.

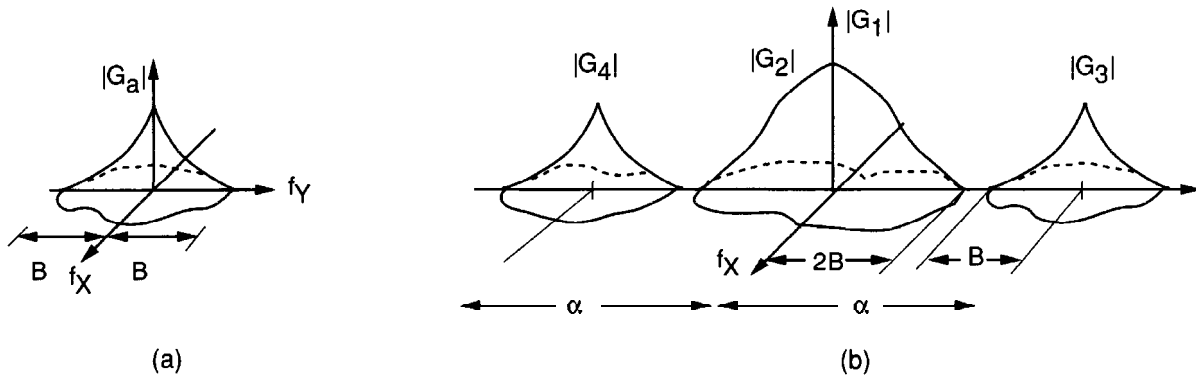


FIGURE 9.8
Spectra of (a) the object and (b) the hologram.

Evidently the minimum allowable reference angle is given by

$$2\theta_{\min} = \sin^{-1} 3B\lambda. \quad (9-30)$$

When the reference wave is much stronger than the object wave, this requirement can be relaxed somewhat. The term G_2 is generated physically by interference of light from each object point with light from all other object points, while G_3 and G_4 arise from interference between the object and reference waves. When the object wave is much weaker than the reference wave (i.e. when $|a| \ll |A|$), the term G_2 is of much smaller magnitude than G_1 , G_3 , or G_4 , and can be dropped as negligible. In this case the minimum reference angle is that which barely separates G_3 and G_4 from each other, or

$$2\theta_{\min} = \sin^{-1} (B\lambda). \quad (9-31)$$

9.4.4 Holography of Three-Dimensional Scenes

In 1964, Leith and Upatnieks reported the first successful extension of holography to three-dimensional imagery [190]. Success in this endeavor rested to a large degree on the availability of the HeNe laser, with its excellent temporal and spatial coherence.

Figure 9.9(a) illustrates the general geometry used for recording holograms of three-dimensional scenes. Coherent light illuminates the scene of interest. In addition, a portion of the illumination strikes a "reference" mirror placed next to the scene. Light is reflected from the mirror directly to the photographic plate, where it serves as a reference wave, interfering with light reflected from the scene itself. Thus the photographic plate records a hologram of the three-dimensional scene.

To reconstruct a three-dimensional image of the scene, two different geometries are recommended, one for viewing the virtual image and the other for viewing the real image. As indicated in Fig. 9.9(b), to view the virtual image we illuminate the hologram with an exact duplicate of the original reference wave, in which case the virtual image appears fixed in three-dimensional space behind the photographic plate at exactly the same location where the object was originally located. Since the wavefronts originally incident on the plate have been duplicated during the reconstruction process, the

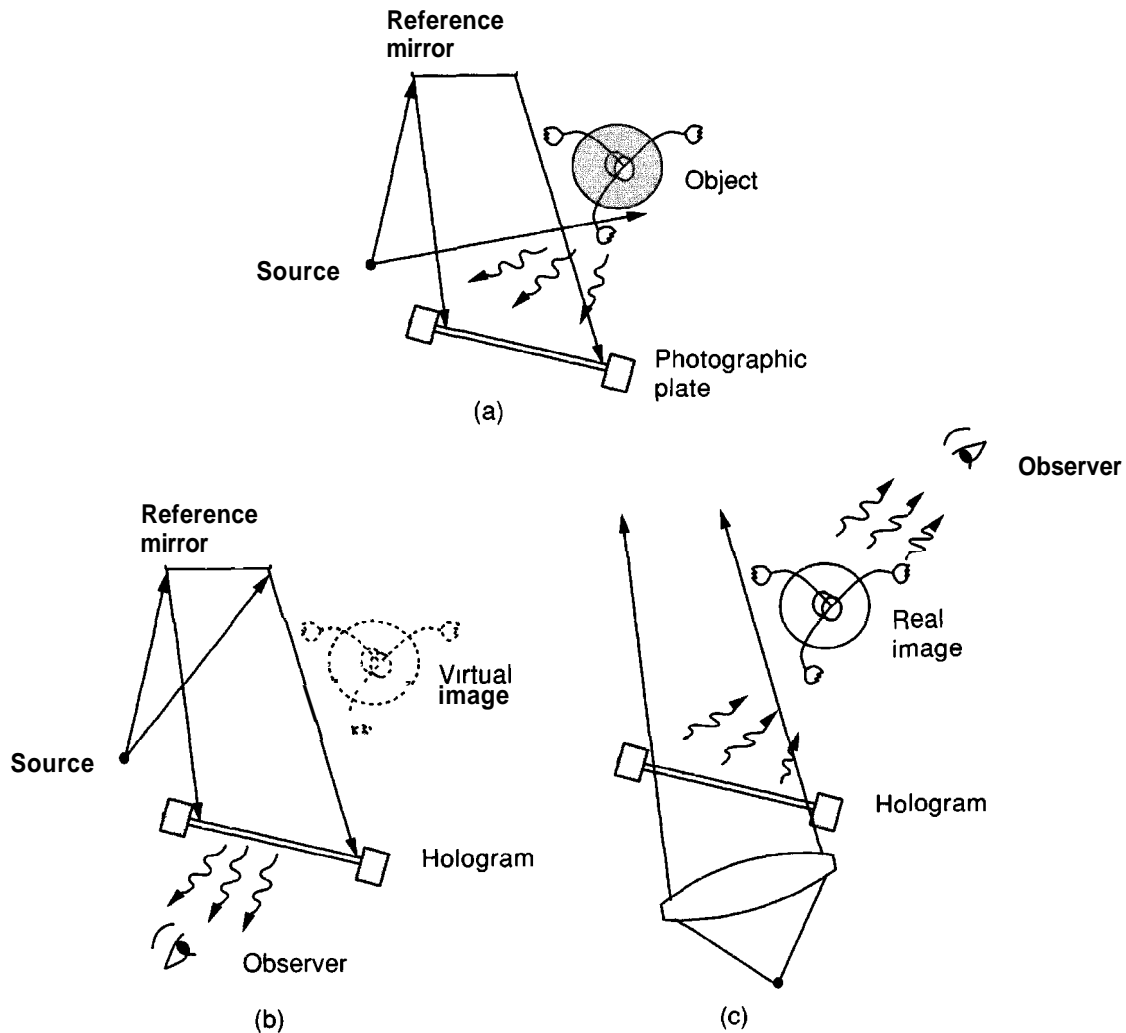


FIGURE 9.9

Holographic imaging of a three-dimensional scene. (a) Recording the hologram; (b) reconstructing the virtual image; (c) reconstructing the real image.

image retains all three-dimensional properties of the object. In particular, it is possible to "look behind" objects in the foreground simply by changing one's viewing position or perspective.

The real image is best viewed when we illuminate the hologram in a different manner. Let the reconstruction wave be a wave that duplicates the reference wave in all respects except one, namely it is traveling backwards towards the original location of the reference source as if time had been reversed during the recording process. This wave can be referred to as an "anti-reference" wave, and can be thought of as being obtained by reversing the direction of the local \vec{k} vector of the reference wave at each point on the hologram. The result is a reconstruction wave with a complex distribution of field that is the complex conjugate of the original reference wave, i.e. $A^*(x, y)$. Under such illumination, the real image forms in space between the photographic plate and the observer, as shown in Fig. 9.9(c). For three-dimensional objects, the real image has

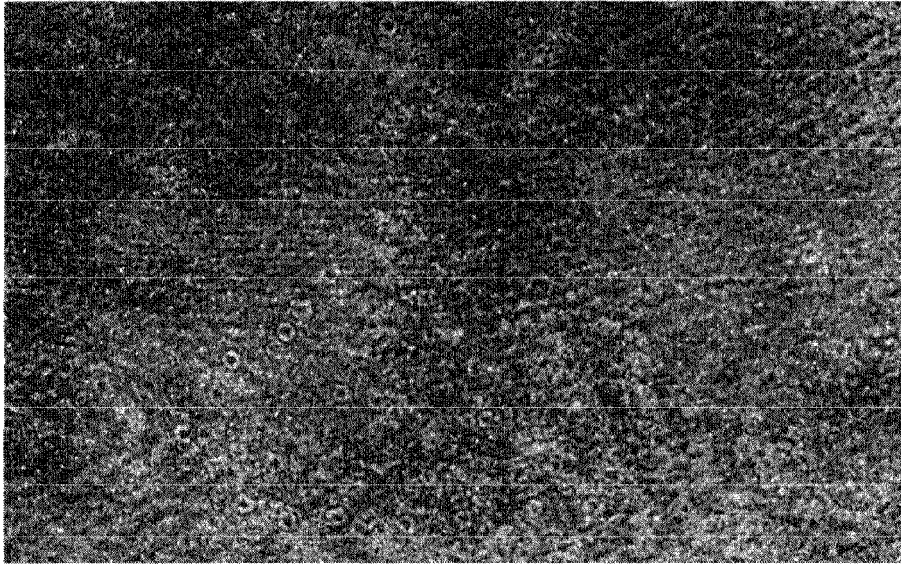


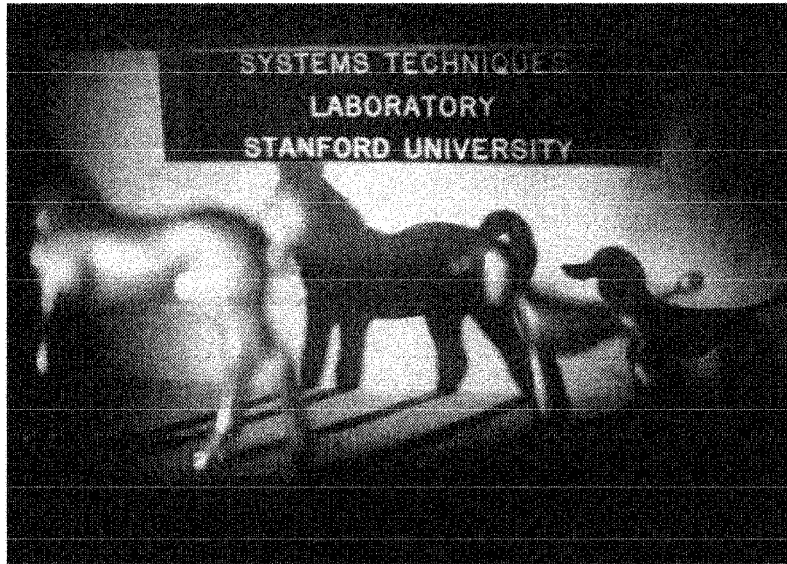
FIGURE 9.10
Photograph of a portion of a hologram of a diffuse three-dimensional scene.

certain properties that make it less useful than the virtual image in many applications. First, points on the object that were closest to the photographic plate (and therefore closest to an observer of the original scene) appear in the real image closest to the photographic plate again, which in this case is farthest from the *observer* (cf. Fig. 9.9(c)). Thus to an observer of the real image, the parallax relations are not those associated with the original object, and the image appears (in a certain peculiar sense that must be personally observed to be fully appreciated) to be "inside out". Images of this type are said to be pseudoscopic, while images with normal parallax relations (like the virtual image) are said to be orthoscopic.

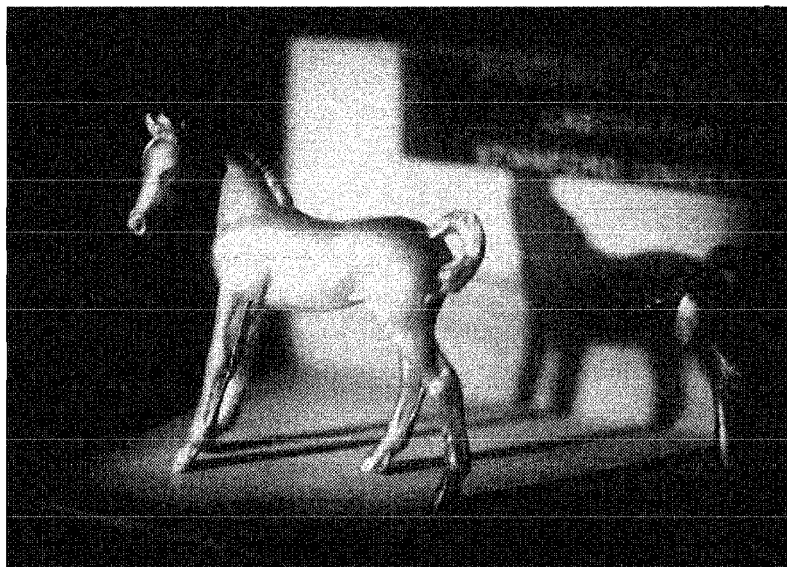
As a second disadvantage of the real image, if photographic film is inserted directly into that image in an attempt to record it directly, the experimenter soon discovers that (for holograms of reasonable size) the depth of focus is generally so small that a recognizable recording can not be obtained. This problem can be alleviated by illuminating only a small portion of the hologram, in which case the depth of focus is increased and a usable two-dimensional image can be recorded. If the illuminating spot on the hologram is moved, then the apparent perspective of the two-dimensional image changes. Thus every small region of a large hologram is capable of producing a real image of the original object with a different perspective!

Figure 9.10 shows a photograph of a portion of a hologram of a diffusely reflecting three-dimensional scene. Note that there is nothing recognizable in the structure recorded on the hologram. In fact, most of the observable structure is irrelevant to the reconstruction in the sense that it arises from imperfections in the optical apparatus (e.g. from dust specks on mirrors and lenses). The structure that generates the reconstructed images is far too fine to be resolved in this photograph.

To illustrate the truly three-dimensional nature of the reconstructed images, we refer to Fig. 9.11, which shows two photographs of the virtual image. In 9.11(a), the camera is focused on the background of the virtual image; the sign in the background



(a)



(b)

FIGURE 9.11
Photographs showing the three-dimensional character of the virtual image reconstructed from a hologram.

is sharply in focus, while the figurines in the foreground are out of focus. Note also that the tail of the horse obscures the head of the shadow of the horse. The camera is now refocused on the foreground and moved to change perspective, with the result shown in Fig. 9.11(b). The foreground is now in focus and the background out of focus. The tail of the horse no longer obscures the head of the shadow of the horse, a consequence of the change of perspective. Thus the camera has succeeded in "looking behind" the tail by means of a simple lateral movement.

9.4.5 Practical Problems in Holography

There are several problems that any practitioner of holography faces and must overcome in order to successfully make a hologram. To become better acquainted with the practice of holography, the reader may wish to consult Ref. [254].

Historically, the extremely short coherence lengths of optical sources available before the advent of the laser seriously constrained the types of holograms that could be recorded. Today the availability of high-quality laser sources has vastly alleviated this problem. However, the experimenter must still take some precautions, for the coherence of lasers is not perfect. For example, it is good practice to measure the distances the reference beam and the object beam travel from source to photographic plate and to equalize the lengths of these paths as closely as possible.

The process of recording a hologram is an exercise in interferometry. As with any interferometric experiment, if clean and sharp interference fringes are to be recorded, it is essential that all path-length differences for interfering light be kept stable to within a fraction of an optical wavelength **during** the duration of the exposure period. The higher the power available from the laser source, the shorter the required exposure time and the less severe the stability requirements become. The exposure time required depends on a multitude of factors, including the transmissivity or reflectivity of the object, the distances and geometry involved, and the particular film or plate used to record the hologram. Pulsed lasers with pulse durations as short as a few nanoseconds have been used in some instances, and CW exposures as long as several hours have been used in some cases.

Some of the most stringent experimental requirements are associated with the recording of holograms of three-dimensional scenes. Photographic emulsions with extremely high resolution are required in such cases (see Section 9.8 for a more complete discussion of this point). One of the most commonly used emulsions for this use, Kodak Spectroscopic Plate Type 649F, has a resolution better than 2000 lines-pairs (cycles)/mm and an equivalent ASA speed of about 0.03.³ It is invariably true that high-resolution emulsions are extremely insensitive.

An additional problem of some significance is the limited dynamic range of photographic recording materials. The amplitude transmittance vs. exposure curve is linear over only a limited range of exposure. It is desirable to choose an average exposure that falls at the midpoint of this linear region. However, when the object is, for example, a transparency with a rather coarse structure, there may exist significant areas on the hologram with exposures falling well outside the linear region. As a consequence of this nonlinearity, degradation of the reconstructed images can be expected (see Section 9.10.2 for further discussion). The dynamic range problem can be largely overcome by a technique first demonstrated by Leith and Upatnieks [190]. The object is illuminated through a diffuser, which spreads the light passed by any one point on the object to cover the entire hologram. Thus a bright spot on the object will no longer generate a

³The **ASA** speed is a standard way of rating the photographic sensitivity of an emulsion. The lower the **ASA** number the less sensitive the emulsion. The **ASA** speeds of emulsions used in conventional photography are often of the order of 100 or higher.

strong Fresnel diffraction pattern on part of the hologram, but rather contributes a more uniform distribution of light. Attendant with the advantageous reduction of dynamic range of the exposing light pattern is another advantage: since each object point contributes to every point on the hologram, an observer looking at a reconstructed image through only a portion of the hologram will always see the entire image. As might be expected, the virtual image appears to be backlit with diffuse illumination.

9.5 IMAGE LOCATIONS AND MAGNIFICATION

To this point we have considered primarily collimated reference and reconstruction waves. In practice these waves are more commonly spherical waves diverging from or converging toward particular points in space. Therefore this section is devoted to an analysis of the holographic process in this more general case. We begin by determining image locations, and then utilize these results to find the axial and transverse magnifications characteristic of the imaging process. The section then concludes with an example.

9.5.1 Image Locations

Referring to Fig. 9.12(a), we suppose that the reference wave is generated by a point source located at coordinates (x_r, y_r, z_r) . Since the mapping of object amplitudes into image amplitudes is linear, provided the reference offset angle is sufficiently large to separate the twin images from each other and from other unwanted terms of trans-

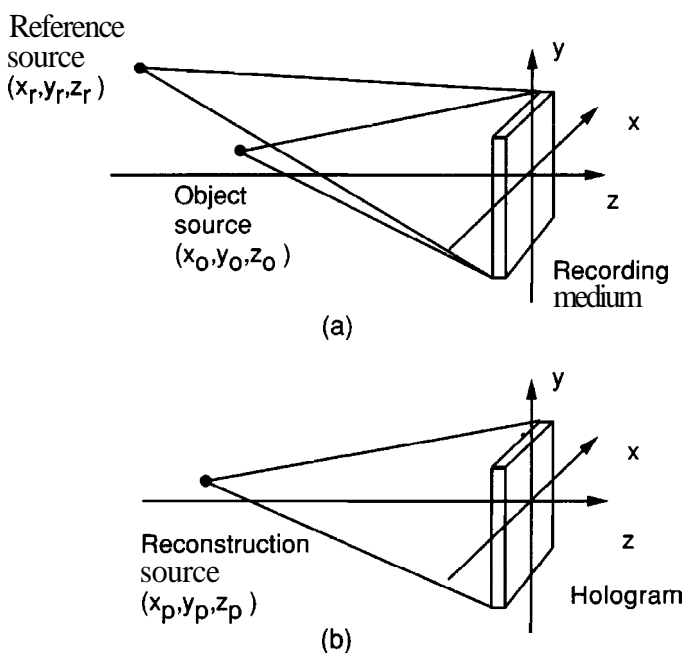


FIGURE 9.12
Generalized (a) recording and (b)
reconstruction geometries.

mitted light, it suffices to consider a single object point source located at coordinates (x_o, y_o, z_o) . Note from the figure that, for our choice of the location of the center of the coordinate system, both z_r and z_o are negative numbers for point sources lying to the left of the hologram recording plane (i.e. for diverging spherical waves), and positive numbers for points lying to the right of that plane (i.e. for converging spherical waves).

During the reconstruction step, illustrated in Fig. 9.12(b), the hologram is assumed to be illuminated by a spherical wave originating from a point source at coordinates (x_p, y_p, z_p) . Again z_p is negative for a diverging wave and positive for a converging wave.

To achieve the maximum possible generality, we allow for the possibility that the recording and reconstruction processes may involve radiation with different wavelengths. Such is the case, for example, in microwave holography, in which the hologram is recorded with microwaves and reconstructed using visible light. The recording wavelength will be represented by λ_1 , and the reconstruction wavelength by λ_2 .

Our analysis will determine the **paraxial** approximations to the (twin) image locations for an object point source at the given coordinates. An extended coherent object may then be regarded as a collection of many mutually coherent point sources.

Using quadratic-phase approximations to the spherical waves involved: the total field incident on the recording plane may be written

$$U(x, y) = A \exp \left\{ -j \frac{\pi}{\lambda_1 z_r} [(x - x_r)^2 + (y - y_r)^2] \right\} + a \exp \left\{ -j \frac{\pi}{\lambda_1 z_o} [(x - x_o)^2 + (y - y_o)^2] \right\}, \quad (9-32)$$

where \mathbf{A} and \mathbf{a} are complex constants representing the amplitudes and relative phases of the two spherical waves. The corresponding intensity distribution in the pattern of interference between the two waves is

$$\begin{aligned} \mathcal{I}(x, y) &= |A|^2 + |a|^2 \\ &+ \mathbf{A} \mathbf{a}^* \exp \left\{ -j \frac{\pi}{\lambda_1 z_r} [(x - x_r)^2 + (y - y_r)^2] + j \frac{\pi}{\lambda_1 z_o} [(x - x_o)^2 + (y - y_o)^2] \right\} \\ &+ \mathbf{A}^* \mathbf{a} \exp \left\{ j \frac{\pi}{\lambda_1 z_r} [(x - x_r)^2 + (y - y_r)^2] - j \frac{\pi}{\lambda_1 z_o} [(x - x_o)^2 + (y - y_o)^2] \right\}. \end{aligned} \quad (9-33)$$

If the amplitude transmittance of the developed transparency is proportional to exposure, then the two important terms of transmittance are

⁴Note that diverging spherical waves have a negative sign in the exponent, whereas in the past they have had a positive sign. This is because the values of z being used here are negative, whereas in previous cases they were positive. It remains true that, if the overall sign in the exponent is positive, the spherical wave is diverging, and if it is negative, the wave is converging.

$$\begin{aligned}
t_3 &= \beta' \mathbf{A} \mathbf{a}^* \exp \left\{ -j \frac{\pi}{\lambda_1 z_r} [(x - x_r)^2 + (y - y_r)^2] + j \frac{\pi}{\lambda_1 z_o} [(x - x_o)^2 + (y - y_o)^2] \right\} \\
t_4 &= \beta' \mathbf{A}^* \mathbf{a} \exp \left\{ j \frac{\pi}{\lambda_1 z_r} [(x - x_r)^2 + (y - y_r)^2] - j \frac{\pi}{\lambda_1 z_o} [(x - x_o)^2 + (y - y_o)^2] \right\}.
\end{aligned} \tag{9-34}$$

The hologram is illuminated with a spherical wave, which in the paraxial approximation is described by

$$U_p(x, y) = B \exp \left\{ -j \frac{\pi}{\lambda_2 z_p} [(x - x_p)^2 + (y - y_p)^2] \right\}. \tag{9-35}$$

The two wavefronts of interest behind the transparency are found by multiplying (9-34) and (9-35), yielding

$$\begin{aligned}
U_3(x, y) &= t_3 B \exp \left\{ -j \frac{\pi}{\lambda_2 z_p} [(x - x_p)^2 + (y - y_p)^2] \right\} \\
U_4(x, y) &= t_4 B \exp \left\{ -j \frac{\pi}{\lambda_2 z_p} [(x - x_p)^2 + (y - y_p)^2] \right\}.
\end{aligned} \tag{9-36}$$

To identify the nature of these transmitted waves, we must examine their (x, y) dependence. Since only linear and quadratic terms in x and y are present, the two expressions U_3 and U_4 may be regarded as quadratic-phase approximations to spherical waves leaving the hologram. The presence of linear terms simply indicates that the waves are converging towards or diverging from points that do not lie on the z axis. It remains to determine the exact locations of these real or virtual points of convergence.

Since the waves emerging from the hologram are given by a product of quadratic-phase exponentials, they must be representable as quadratic-phase exponentials themselves. Thus we can identify the coordinates (x_i, y_i, z_i) of the images if we compare the expanded equations (9-36) with a quadratic-phase exponential of the form

$$U_i(x, y) = K \exp \left\{ -j \frac{\pi}{\lambda_2 z_i} [(x - x_i)^2 + (y - y_i)^2] \right\}. \tag{9-37}$$

From the coefficients of the quadratic terms in x and y we conclude that the axial distance z_i of the image points is

$$z_i = \left(\frac{1}{z_p} \pm \frac{\lambda_2}{\lambda_1 z_r} \mp \frac{\lambda_2}{\lambda_1 z_o} \right)^{-1} \tag{9-38}$$

where the upper set of signs applies for one image wave and the lower set of signs for the other. When z_i is negative, the image is virtual and lies to the left of the hologram, while when z_i is positive, the image is real and lies to the right of the hologram.

The x and y coordinates of the image points are found by equating the linear terms in x and y in Eqs. (9-36) and (9-37), with the result

$$\begin{aligned}
 x_i &= \mp \frac{\lambda_2 z_i}{\lambda_1 z_o} x_o \pm \frac{\lambda_2 z_i}{\lambda_1 z_o} x_r + \frac{z_i}{z_p} x_p \\
 y_i &= \mp \frac{\lambda_2 z_i}{\lambda_1 z_o} y_o \pm \frac{\lambda_2 z_i}{\lambda_1 z_o} y_r + \frac{z_i}{z_p} y_p.
 \end{aligned}
 \tag{9-39}$$

Equations (9-38) and (9-39) provide the fundamental relations that allow us to predict the locations of images of point sources created by the holographic process. Depending on the geometry, it is possible for one image to be real and the other virtual, or for both to be real or both virtual (see Prob. 9-2).

9.5.2 Axial and Transverse Magnifications

The axial and transverse magnifications of the holographic process can now be found from the equations derived above for image locations. The transverse magnification is easily seen to be given by

$$M_t = \left| \frac{\partial x_i}{\partial x_o} \right| = \left| \frac{\partial y_i}{\partial y_o} \right| = \left| \frac{\lambda_2 z_i}{\lambda_1 z_o} \right| = \left| 1 - \frac{z_o}{z_r} \mp \frac{\lambda_1 z_o}{\lambda_2 z_p} \right|^{-1}
 \tag{9-40}$$

Similarly, the axial magnification is found to be

$$M_a = \left| \frac{\partial z_i}{\partial z_o} \right| = \left| \frac{\partial}{\partial z_o} \left(\frac{1}{z_p} \pm \frac{\lambda_2}{\lambda_1 z_r} \mp \frac{\lambda_2}{\lambda_1 z_o} \right)^{-1} \right| = \frac{\lambda_1}{\lambda_2} M_t^2.
 \tag{9-41}$$

Note that in general the axial and transverse magnifications will not be identical. This can be very important when we consider the imaging of three-dimensional objects, as we shall do shortly, for the difference between these magnifications will create a three-dimensional distortion of the image. There does exist one additional parameter that can be used to combat such distortions, namely, it is possible to *scale* the hologram itself between the recording and reconstruction process. For example, if the hologram were formed with microwaves or acoustic waves, it would be possible to plot out the hologram at any scale size we choose, and record a transparency of the hologram with magnification or demagnification. If m is the magnification ($m > 1$) or demagnification ($m < 1$) to which the hologram is subjected, then we can show that the transverse and axial magnifications take the form

$$\begin{aligned}
 M_t &= m \left| 1 - \frac{z_o}{z_r} \mp m^2 \frac{\lambda_1 z_o}{\lambda_2 z_p} \right|^{-1} \\
 M_a &= \frac{\lambda_1}{\lambda_2} M_t^2.
 \end{aligned}
 \tag{9-42}$$

9.5.3 An Example

Consider an example in which we record a hologram at wavelength $\lambda_1 = 10$ cm in the microwave region of the spectrum, and reconstruct images in the visible region of the

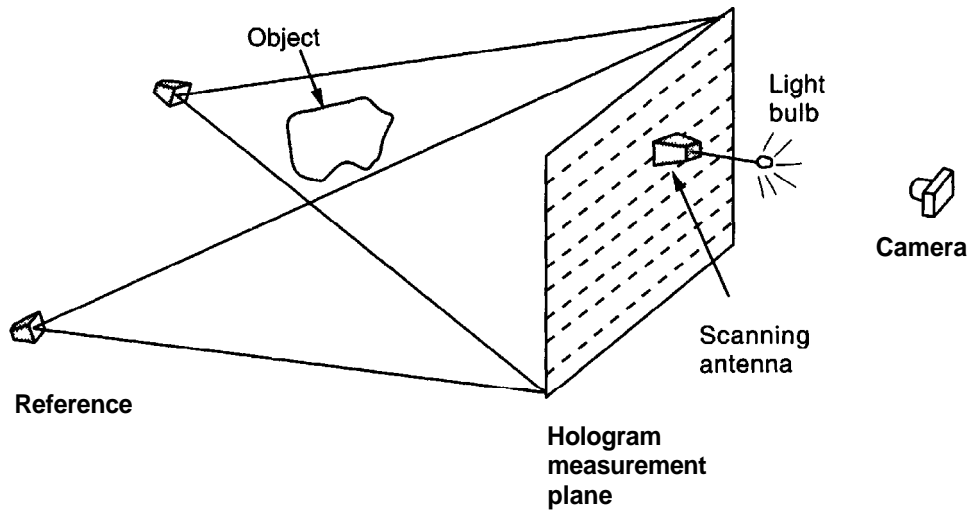


FIGURE 9.13

Recording a microwave hologram. The sources that provide the object and reference illuminations are derived from the same microwave signal generator to assure coherence.

spectrum with $\lambda_2 = 5 \times 10^{-5}$ cm. Figure 9.13 illustrates how the experiment might be performed. A microwave source illuminates an object with variable microwave transmittance, perhaps a three-dimensional structure which partially absorbs microwaves. A mutually coherent microwave source provides a reference wave which interferes with the radiation diffracted by the object. Some distance away a scanning rig with a microwave horn antenna measures the microwave intensity impinging on a larger aperture. To be specific, suppose that the size of the scanned aperture is 10 m \times 10 m. Attached to the scanning microwave horn antenna is a light bulb, which is driven with a **current** that is proportional to the incident microwave power at each point in the scanned aperture. A camera records a time exposure of the brightness pattern of the light bulb as it scans across the microwave aperture, and this recorded photograph generates an optical transparency that is inserted into an optical system and illuminated at the visible wavelength quoted above.

If we imagined a physically impossible case of a photographic transparency that is as large as the total scanned microwave aperture, and we suppose that the microwave reference wave supplied in the recording process is a plane wave ($z_r = \infty$) and the optical reconstruction wave is a plane wave ($z_p = \infty$), then application of Eqs. (9-42) yields the following transverse and axial magnifications:

$$M_t = 1$$

$$M_a = 2 \times 10^5.$$

As can be seen from these numbers, there is an enormous amount of distortion of the image, with the transverse magnification being more than five orders of magnitude smaller than the axial magnification.

Now suppose that we modify this experiment such that the photograph is optically reduced to be only 5 μm on a side, which corresponds to a demagnification of $m = \frac{\lambda_2}{\lambda_1} = 5 \times 10^{-6}$. Again the reference wave and the reconstruction wave are assumed to

be plane waves. In this case we find that the transverse and axial magnifications are given by

$$M_t = 5 \times 10^{-6}$$

$$M_a = 5 \times 10^{-6}.$$

Thus we see that the two magnifications have been made equal to each other by means of the scaling of the hologram by the wavelength ratio, thereby removing the three-dimensional distortion. Unfortunately in the process the image has been made so small (5×10^{-6} times smaller than the original object) that we may need a microscope to examine it, in which case the microscope will introduce distortions similar to what we have removed from the holographic process.

The above example is somewhat contrived, but it does illustrate some of the problems that can be encountered when the recording and reconstruction wavelengths are significantly different. Such is often the case for acoustic holography and microwave holography. For holography at very short wavelengths such as the ultraviolet and X-ray regions of the spectrum, the problem is reversed, and the hologram must be scaled upwards in size in order to avoid distortions.

9.6 SOME DIFFERENT TYPES OF HOLOGRAMS

Attention is now turned to a brief guided tour through several different kinds of holograms. There are many different aspects with respect to which holograms may differ, and this has led to a rather confused classification system, in which a given hologram may in fact be properly classified in two or more different classes at the same time. There is nothing fundamentally wrong with this, as long as we understand what the different classes mean. In what follows we do not include the categorization "thin" vs. "thick" as a classification, only because these differences will be discussed in detail in later sections.

9.6.1 Fresnel, Fraunhofer, Image, and Fourier Holograms

Our first dimension of classification is one that distinguishes between the diffraction or imaging conditions that exist between the object and the photographic plate where the hologram is recorded. Thus we say that a hologram is of the *Fresnel* type if the recording plane lies within the region of Fresnel diffraction of the illuminated object, whereas it is of the *Fraunhofer* type if the transformation from object to hologram plane is best described by the Fraunhofer diffraction equation.

In some cases a hologram is recorded in what must be called an image plane, and such a hologram would be referred to as an *image* hologram. This geometry is most frequently used when the object is three-dimensional but perhaps not extremely deep in the third dimension. The middle of the object can then be brought to focus in the plane of the photographic plate, and the resulting image obtained from the hologram will appear to float in space at the hologram, with parts of the object extending forwards and backwards from the hologram.

A category that applies primarily to transparency objects is the *Fourier* hologram, for which the recording plane resides in a plane that will yield the Fourier transform of the object amplitude **transmittance**. Thus with a normally illuminated object transparency in the front focal plane of a lens and the recording plane in the rear focal plane of the lens, the relation between fields in the two planes will be that of a Fourier transform. For such a hologram, the light from each point on the object interferes with the reference beam (assumed planar) to create a sinusoidal fringe with a vector spatial frequency that is unique to that object point. The transformation from object points into sinusoidal fringes of unique spatial frequencies is thus characteristic of the Fourier transform hologram. To view the reconstructed images, we can place the hologram in front of a positive lens, illuminate it with a normally incident plane wave, and look for images in the rear focal plane. Note that both of the twin **images** come to focus in the same plane for such a hologram, as can be verified by applying Eq. (9-38).

Finally, for transparency objects one sometimes sees mention of a hologram that is called a *lensless Fourier transform* hologram. The name is a misnomer, for the geometry usually requires a lens, but not a Fourier transforming lens as is used in the **ordinary Fourier** transform geometry. Rather, as shown in Fig. 9.14, the reference wave is brought to focus in the plane of the object transparency, and then diverges to the recording plane without passing through any optical elements. Likewise, the wave transmitted by the object propagates to the recording plane without the intervention of any optical elements. The interference pattern is then recorded. The distance between the object and the hologram recording plane is immaterial. The reason for associating the words *Fourier transform* with such a hologram, when in fact no Fourier transform actually takes place, can be understood by considering the interference pattern generated by light from a single point on the object. Both the object wave and the reference wave **are diverging** spherical waves with the same curvature, and as a consequence when they interfere, the pattern of intensity is (within the **paraxial approximation**) a sinusoidal fringe of a vector spatial frequency that is unique to that object point. This is the same property that holds for a true Fourier transform hologram, and hence the mention of *Fourier* in the name for this type of hologram. The difference between this type of hologram and the true Fourier transform hologram lies in the spatial phases that are associated with the various sinusoidal fringes, which in this case are not the phases of the Fourier transform of the object. The twin images can be observed if the fields transmitted by the hologram are Fourier transformed by a positive lens. Again, if the

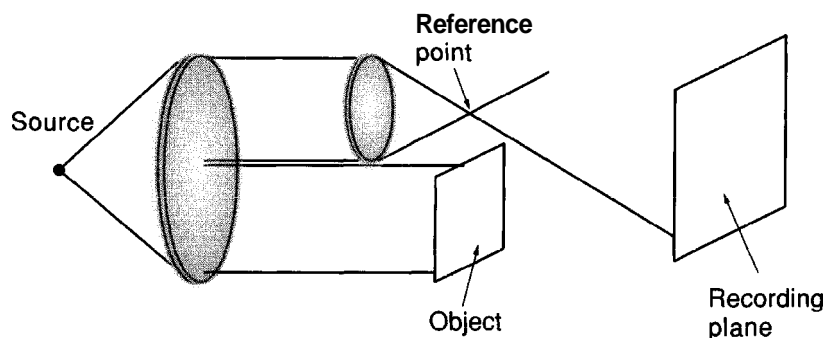


FIGURE 9.14
Recording a lensless Fourier transform hologram.

hologram is illuminated with a plane reconstruction wave, both images appear in the focal plane of the transforming lens.

The encoding of object points into sinusoidal fringes of constant frequency should be contrasted with the *Fresnel* hologram, in which each object point is encoded into a portion of a frequency-chirped sinusoidal fringe (a sinusoidal zone plate) with an entire range of spatial frequency components present. The Fourier transform and **lensless** Fourier transform holograms make the most efficient use of the space bandwidth product of the hologram.

9.6.2 Transmission and Reflection Holograms

The majority of the holograms discussed so far have been of the *transmission* type. That is, we view the images in light that has been transmitted through the hologram. Such holograms are comparatively tolerant to the wavelength used during the reconstruction process (although the amount of tolerance depends on the thickness of the emulsion), in the sense that a bright image can be obtained without exactly duplicating the wavelength used during exposure. However this also leads to chromatic blur when a transmission hologram is illuminated with white light, so some filtering of the source is generally required.

Another important class of holograms is that of *rejection* holograms, for which we view the images in light that is reflected from the hologram. The most widely used type of reflection hologram is that invented by Y. Denisyuk in 1962 [82]. The method for recording such a hologram is illustrated in Fig. 9.15(a). In this case there is only one illumination beam, which supplies both the object illumination and the reference beam simultaneously. As shown in the figure, the object is illuminated *through the holographic plate*. The incident beam first falls upon the holographic emulsion, where it serves as a reference wave. It then passes through the photographic plate and illuminates the object, which usually is three dimensional. Light is scattered backwards from the object, towards the recording plane, and it passes through the emulsion traveling in a direction that is approximately opposite to that of the original incident beam. Within the emulsion the two beams interfere to produce a standing interference pattern with extremely fine fringes. As will be seen in Section 9.7, the period of the sinusoidal fringe formed when two plane waves traveling at angle 2θ with respect to each other interfere is given by

$$\Lambda = \frac{2\pi}{|K|} = \frac{\lambda}{2 \sin\left(\frac{2\theta}{2}\right)}. \quad (9-43)$$

When $2\theta = 180^\circ$, as is the case for the reflection hologram, the fringe period is half of an optical wavelength in the the **emulsion**.⁵ As will be seen in Section 9.7, the fringes are oriented such that they bisect the angle between directions of travel of the reference

⁵Note also that the optical wavelength in the emulsion is smaller than the vacuum wavelength by a factor $1/n$, which for $n \approx 1.5$ is a factor of $2/3$.

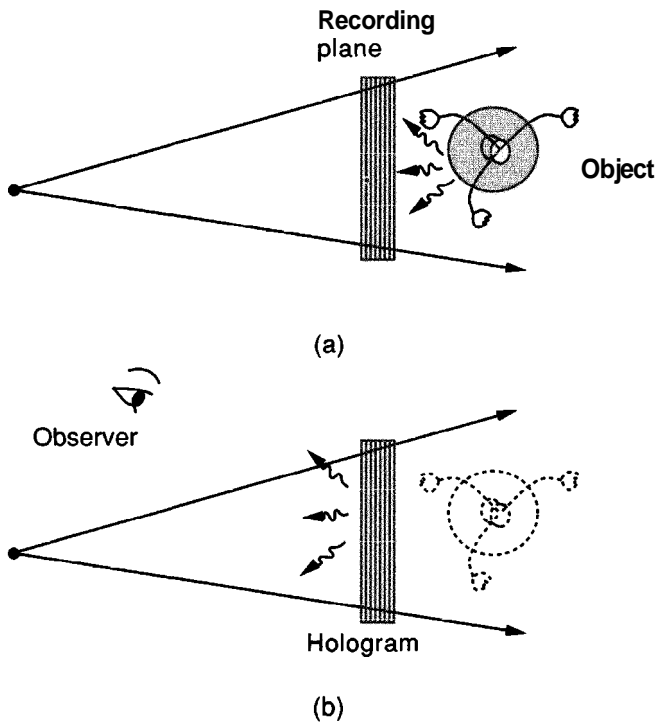


FIGURE 9.15
 (a) Recording a reflection hologram,
 and (b) reconstructing an image in
 reflected light.

and the object waves, and are therefore approximately *parallel to the surface of the emulsion* for a reflection hologram.

Figure 9.15(b) shows how the virtual image would be viewed for a reflection hologram. The hologram is illuminated by a duplication of the original reference wave, and a duplicate of the object wave is generated, which in this case is a reflected wave. The observer looks into the reflected wave and sees the virtual image in the original location of the object, behind the hologram. Figure 9.16 shows a photograph of the virtual image reconstructed from a reflection hologram that is being illuminated by white light.

This type of hologram can be illuminated with white light, for the hologram is highly wavelength selective, and the wavelength that satisfies the Bragg condition will automatically be reflected, while others will not. In this regard it should be noted that photographic emulsions usually suffer some shrinkage during the chemical processing and drying steps, and therefore the color of the light reflected from this type of hologram will usually be different than that used during recording. For example, a hologram recorded with red light may reflect green light. Such effects can be compensated by intentionally swelling the emulsion by means of proper chemical treatment.

9.6.3 Holographic Stereograms

At a relatively early stage in the development of holography, several ideas emerged for using the holographic process to capture a multitude of images that were recorded by conventional photography and to create the illusion of three dimensions through the stereo effect. The function of holography in these schemes is to allow the observer to see different images, taken from different perspectives, in each eye, thereby

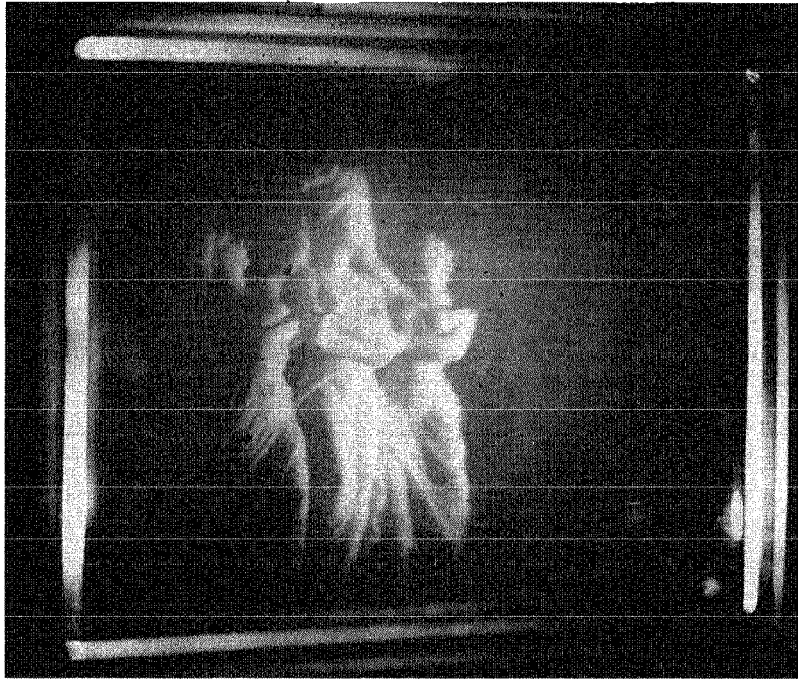


FIGURE 9.16
 Photograph of a virtual image reconstructed from a reflection
 hologram.

creating the stereo effect. The fact that the process begins with ordinary photography, and does not require that the original scene be illuminated by a laser, is a distinct advantage. A laser is required in the hologram-recording process. References include [205], [81], and [242].

One method for recording such a hologram is illustrated in Fig. 9.17 [81]. A series of black and white photographs are taken of the subject from a sequence of horizontal

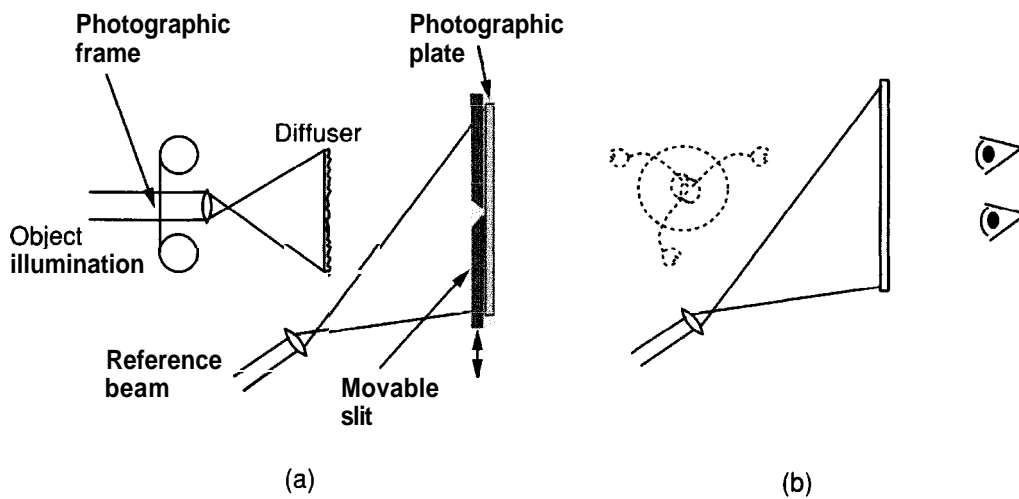


FIGURE 9.17
 Recording a holographic stereogram (top view). (a) Recording the holograms, and (b)
 viewing the image.

positions, each with its own unique perspective. Each frame of the sequence is then projected with light from a laser onto a translucent screen. A reference beam is introduced and a hologram is recorded through a movable slit. As the photographic frame is advanced, the slit is moved, with the result that a multitude of holograms are recorded side-by-side, each hologram capable of reconstructing an image of the original object taken from a different horizontal perspective. If the resulting hologram is illuminated in its entirety by a duplicate of the reference wave, the observer will look through a different holographic stripe with each eye, and therefore each eye will see the subject from a different perspective, creating a three-dimensional image through the stereo effect.

An alternative approach [242] uses angular multiplexing in thick holograms to superimpose a multitude of overlapping holograms on a photographic plate. Each eye views the holographic plate from a slightly different angle, and as a consequence the Bragg effect leads to the reconstruction of a different image seen by each eye, and a resulting three-dimensional image is created.

9.6.4 Rainbow Holograms

An important advance in the field of display holography was the invention of the *rainbow hologram* by *S. Benton* [19]. This invention provides a method for utilizing white light as the illumination when viewing the hologram, and does so by minimizing the blur introduced by color dispersion in a transmission hologram, at the price of giving up parallax information in one dimension. The ability to view holographic images in white light was a vital step on the road to making holography suitable for display applications.

The method entails a two-step process, in which an initial hologram is made, and then a second hologram is made using the first hologram as part of the process. The first step in the process is to record a hologram of a three-dimensional scene in the usual way, in particular using monochromatic or nearly monochromatic light, as is illustrated in Fig. 9.18(a). The light from the reference source R_1 and light scattered by the object O interfere to form holographic recording H_1 . This recording is processed in the usual way, and a hologram results. We now illuminate this first hologram with an monochromatic "anti-reference" wave, i.e. a wave that duplicates the original reference wave, except that the direction of travel is reversed, as illustrated in Fig. 9.18(b). A real image of the original object is produced by hologram H_1 , and the location of that real image coincides with the original location of the object.

Now a new element is introduced in the reconstruction geometry of Fig. 9.18(b), namely a narrow horizontal slit immediately following hologram H_1 . The light passing through this slit again reconstructs a real image of the original object, but this time vertical parallax is eliminated; the image that is formed is the one that would have been seen from the particular vertical location of the slit. Having created this real image, a second hologram H_2 is recorded, this time a hologram of the image produced by the first hologram, with a new reference wave being used, in particular a reference wave that is a converging spherical wave. Again the light used in the recording process is monochromatic, and the pattern of interference is between the reference wave from R_2 and the light that has passed through the focus of the real image and traveled on to the recording plane, as shown in Fig. 9.18(b). H_2 is the final hologram created by this process.

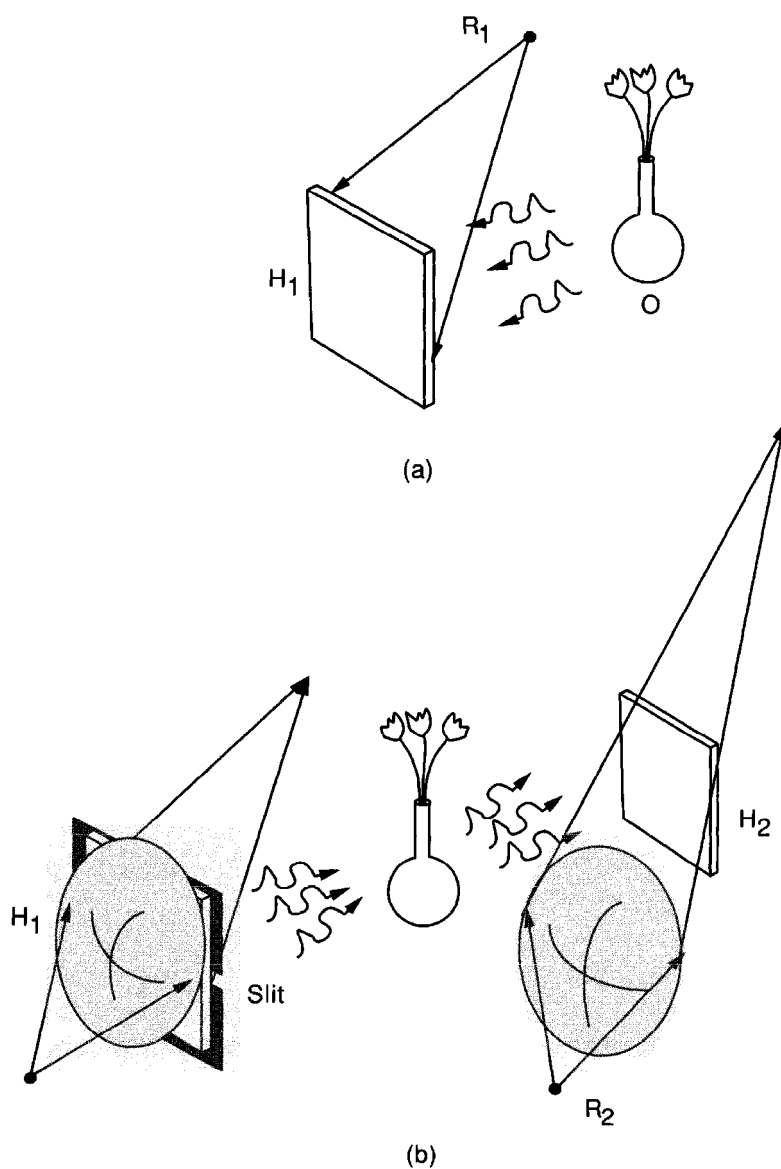


FIGURE 9.18
The rainbow hologram. (a) The first recording step, and (b) the second recording step.

The hologram obtained by the method described above is now illuminated with a diverging spherical wave, which is in fact the "anti-reference" wave for the converging reference wave from R_2 , as shown in Fig. 9.19(a).

The hologram forms a real image⁶ of the original object, but beyond that image, closer to the viewer, there is also formed an image of the slit that was introduced when

⁶Because this is a pseudoscopic image of a pseudoscopic image, it is orthoscopic from the perspective of the viewer.

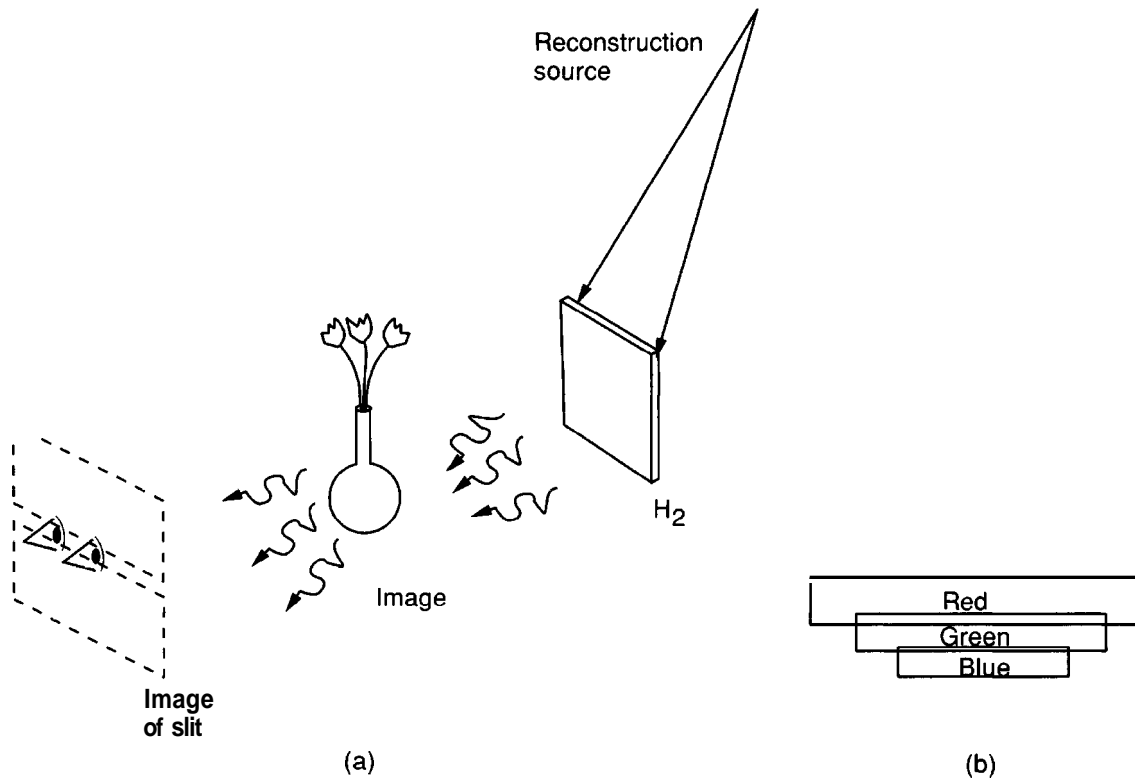


FIGURE 9.19

Reconstruction of the image from a rainbow hologram; (a) Reconstruction geometry, (b) slit sizes at different wavelengths.

hologram H_2 was recorded. Now if the reconstruction source in this last step emits white light, then the dispersion of the hologram will in fact generate a blur of images of both the object and the slit. In particular, each different narrow band of colors of the reconstruction source will create an image of the slit at a different vertical location (and with a different magnification), with red light having been diffracted vertically the most and blue light the least. An observer located in the plane of the images of the slit will in effect look through a slit appropriate for only a narrow color band, and will intercept no light outside of this color band. Thus the image will appear free from color blur, and will have a color that depends on exactly where the observer's eyes are located in the vertical dimension. Tall observers will see an image with a different color (and a somewhat different magnification) than short observers. Thus the dispersive properties built into hologram H_2 have been used to advantage to eliminate color blur and to allow the image to be viewed with a white light source. As shown in Fig. 9.19(b), the slit position varies with color in both vertical position and magnification.

9.6.5 Multiplex Holograms

Another major type of hologram that has been widely used for display purposes is the *multiplex hologram* invented by Lloyd Cross [72]. Good descriptions of the multiplex hologram process can be found in Ref. [254].

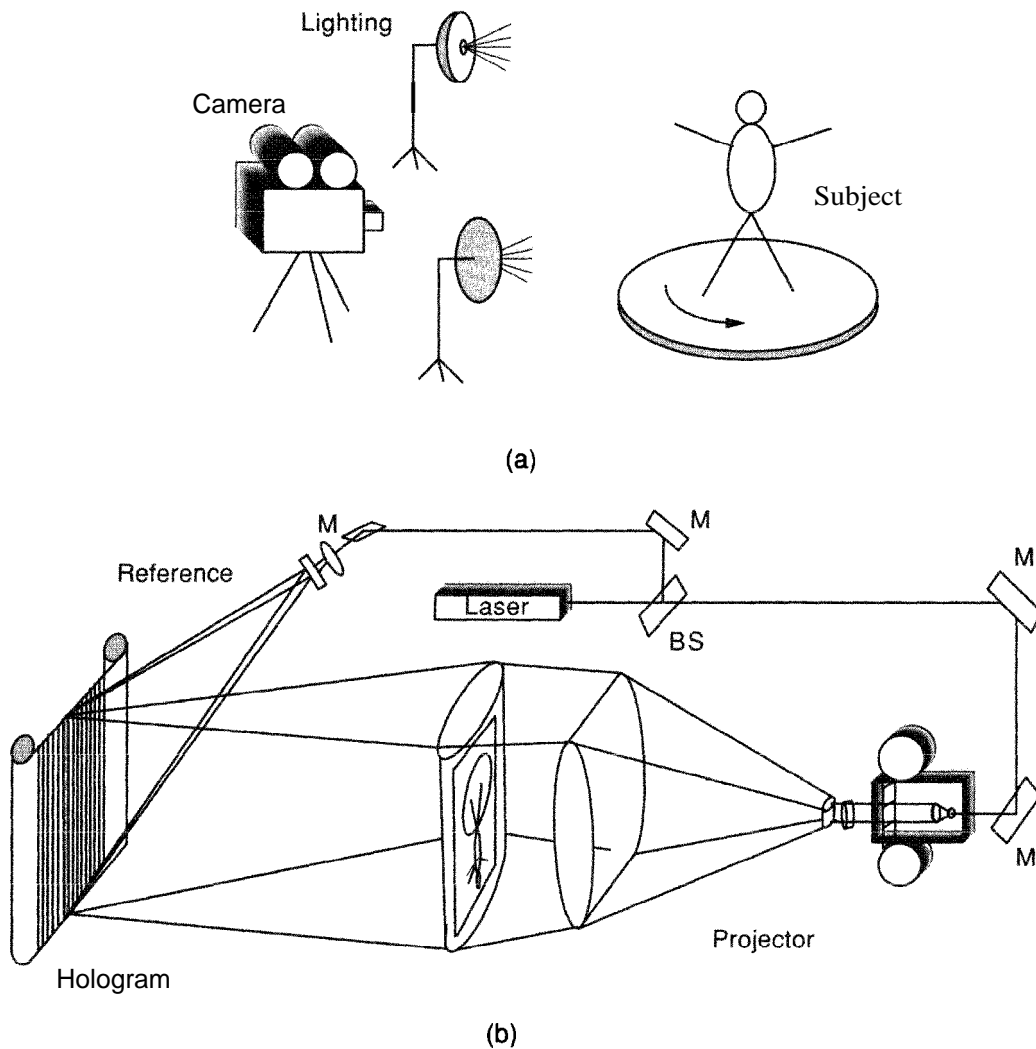


FIGURE 9.20
Constructing a multiplex hologram. (a) Recording the still-frame sequence, and (b) recording the multiplex hologram. M indicates a mirror, BS a beam splitter. The reference wave arrives at the recording plane from above.

The process begins with a series of still-frame photographs, typically made with a motion picture camera operated a single frame at a time. Figure 9.20(a) shows the process. A subject is placed on a slowly rotating platform and still-frame photographs are taken, typically at a rate of 3 frames for every degree of rotation of the subject. Thus for a hologram that is to offer a 120° viewing angle, a total of 360 images are recorded. During the rotation process, the subject may undergo some motion or movement of its own, a motion that will eventually be evident in the image viewing process.

The sequence of photographs obtained as above is now fed through a special projector, as shown in Fig. 9.20(b). Using light from a laser, the images are projected onto a large cylindrical lens, which focuses the light down to a narrow vertical stripe on the large film strip that is to record the hologram. At the same time a portion of the laser light is brought above the projection system and focused to a vertical stripe that coincides with the projected image and provides a reference beam that is offset from the object

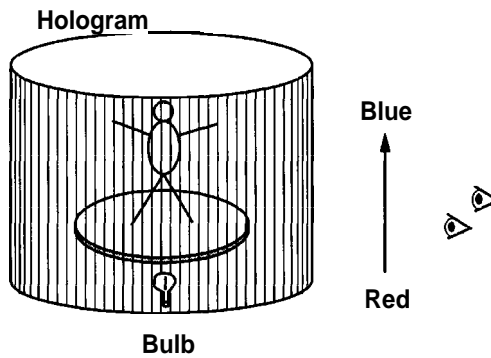


FIGURE 9.21
Viewing the image with a multiplex hologram.

beam by an angle in the vertical dimension. Thus a narrow vertical stripe hologram is recorded for a given still-frame photo, with a carrier frequency running vertically. The film is now advanced to the next still-frame photo, and the holographic film is moved so that a new vertical stripe on the film will be exposed, usually partially overlapping the adjacent stripe. Through a sequence of such exposures, the 360 holograms are recorded. Note that each still-frame photo, and therefore each holographic stripe, contains image information taken from a different perspective in the original photographic process.

To view a three-dimensional image with the hologram after processing, the film is bent into a cylindrical shape and illuminated from within the cylinder with a **white-light** source, properly positioned in the vertical direction to account for the reference angle used during recording (see Fig. 9.21). An illumination source with a clear bulb and a vertical filament is required in order to avoid image blur. The observer looks into the hologram and sees an apparently three-dimensional image within the cylinder. The white light is dispersed in the vertical dimension by the holographic gratings, with red light suffering more downward deflection than blue light. An observer looking into the hologram will automatically perform two types of selection. First the vertical position of the observer's head will place him or her in a certain narrow region of the color spectrum, so color filtering is performed simply by geometry. Second, the two eyes of the observer will look through different regions of the multiplex hologram, and therefore will be looking predominantly through two different holographic stripes, each of which yields an image of the original object taken from a different perspective. As a consequence the stereo effect creates the illusion that the object is three-dimensional. As the observer moves horizontally, the image appears to be stationary in space and the perspective changes accordingly. If the subject goes through some form of motion while the original platform is being rotated, then corresponding motion of the three-dimensional image will be seen as the viewer moves around the hologram, or alternatively as the hologram is rotated. Note that, as with the rainbow hologram, vertical parallax is not present in the image.

9.6.6 Embossed Holograms

Embossing has become a highly refined and advanced technique for replicating compact disks and video disks, which have structures of the same order of size as an optical wavelength. The same techniques can be applied to the replication of holograms, with substantial cost savings as compared with optical methods of duplication. The ability

to produce holograms inexpensively has led to their use, for example, in security cards, credit cards, magazines, books, and in some cases on monetary bills. We outline here the various steps that are involved in creating an embossed hologram.

The first step in the process is to record a hologram of the subject of interest, on photoresist. With a proper choice of photoresist, the resolution is fully adequate to the task at hand. Usually a rather powerful argon-ion laser is used in the recording step. The exposed photoresist is then developed, leading to a relief pattern that constitutes the photoresist master hologram.

A metal master hologram is now made from the photoresist hologram by means of an electroforming process. A silver spray is applied to the photoresist surface, making it conducting. The master is then immersed in a plating tank, together with a bar of pure nickel, and current is passed through the tank with the result that a thin layer of nickel is plated on top of the photoresist master. The layer of nickel, which forms the metal master, is then separated from the photoresist. It is now possible to use the metal master in a second electroforming process, in which a second-generation metal submaster can be made from the original. The process can be repeated to make many metal submasters, which will serve as stampers in the reproduction process.

With the metal submasters in hand it is now possible to initiate the embossing process. There are several different methods for embossing, including flat-bed embossing, roll embossing, and hot stamping. In all cases the metal submaster is heated to an elevated temperature, and used to stamp the hologram pattern, usually into a polyester material. Often the embossed pattern is metallized to create a reflection hologram.

Without doubt, of all the holograms in existence today, the largest number are of the embossed type, for only with embossing can the cost of reproducing holograms be brought down to the levels needed for extremely high-volume applications.

9.7 THICK HOLOGRAMS

Just as for acousto-optic spatial light modulators (see Section 7.2.6), holograms behave differently depending on the relation between the period of the finest fringe they contain and the thickness of the recording medium. It is therefore common to categorize holograms as *thick* or *thin*, depending on this relation. Like the acoustic waves in an acousto-optic SLM, a hologram is a grating. Unlike the acousto-optic case, the grating is stationary rather than moving, and it may also be highly absorbing, partially absorbing, or nonabsorbing, depending on the conditions of exposure and the photographic processing to which it has been subjected.

If we consider a hologram consisting of a single sinusoidal grating with grating planes normal to the surface of the emulsion, it behaves as a thick or thin grating depending on the value of the Q parameter of Eq. (7-29), which is repeated here,

$$Q = \frac{2\pi\lambda_0 d}{n\Lambda^2}, \quad (9-44)$$

where λ_0 is the vacuum wavelength of the light used during reconstruction, n is the refractive index of the emulsion after processing, Λ is the period of the sinusoidal grating,

and d is the emulsion thickness. Again, for $Q > 2\pi$ the grating is considered "thick", while for $Q < 2\pi$ the grating is "thin".

The most common photographic plates used in holography have thicknesses of the order of $15\ \mu\text{m}$, while the fringes formed in holograms may be only a few wavelengths, or in some cases as small as half a wavelength, depending on the angle between the reference wave and the object wave. Therefore a hologram of an object with any significant angular subtense at the hologram plane will contain at least some fringes that exhibit the properties of a thick grating. Hence Bragg diffraction effects must be considered in most cases.

In this section we consider in more detail the properties of the gratings recorded by the holographic process, and the requirements for achieving high diffraction efficiency from such gratings. Finally we determine the diffraction efficiencies of thick holograms and compare them with those of thin holograms. An excellent and far more detailed treatment of this subject will be found in [267].

9.7.1 Recording a Volume Holographic Grating

Consider the very simple case of a plane reference wave and a plane object wave incident on an emulsion of nonnegligible thickness. These two simple waves may be regarded as generating a simple holographic grating.

With reference to Fig. 9.22, it is assumed for simplicity that the two wave normals (represented by arrows and pointing in the directions of the two \vec{k} vectors), are each inclined at angle θ to the surface normal. Wavefronts, or successive lines of zero phase, are shown dotted; the wavefronts of any one wave are spaced by a normal distance of one wavelength. Along the lines (points in this two-dimensional figure) within the emulsion where the wavefronts of the two waves intersect, the two amplitudes add in phase, yielding high exposure. As time progresses, the wavefronts move in the direction of their respective wave normals, and the lines of constructive interference move through the emulsion, tracing out *planes* of high exposure. Simple geometry shows that these planes bisect the angle 2θ between the two wave normals and occur periodically throughout the emulsion.

Describing the three-dimensional interference process mathematically, the complex amplitudes of the two waves can be represented by

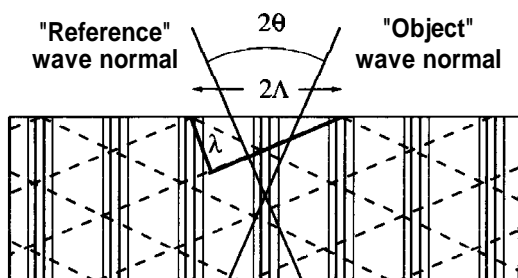
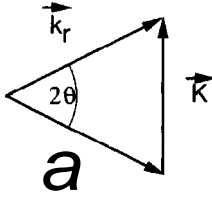


FIGURE 9.22
Recording an elementary hologram with a thick emulsion.


FIGURE 9.23

Wave vector diagram illustrating the length and direction of the grating vector.

$$\begin{aligned} U_r(\vec{r}) &= A e^{j\vec{k}_r \cdot \vec{r}} \\ U_o(\vec{r}) &= a e^{j\vec{k}_o \cdot \vec{r}}, \end{aligned} \quad (9-45)$$

where \vec{k}_r and \vec{k}_o are the wave vectors of the reference and object waves, respectively, and \vec{r} is a position vector with components (x, y, z) . The intensity distribution that results from superimposing these waves is given by

$$I(\vec{r}) = |A|^2 + |a|^2 + 2|A||a| \cos [(\vec{k}_r - \vec{k}_o) \cdot \vec{r} + \phi], \quad (9-46)$$

where ϕ is the phase difference between the phasors A and a .

At this point it is convenient to define a grating vector \vec{K} as the difference of the two wave vectors,

$$\vec{K} = \vec{k}_r - \vec{k}_o. \quad (9-47)$$

The vector \vec{K} has a magnitude that is $\frac{2\pi}{\Lambda}$, where Λ is the fringe period, and points in the direction of the difference between \vec{k}_r and \vec{k}_o . A pictorial representation of \vec{K} is given by the wave vector diagram shown in Fig. 9.23. From this figure we can deduce that the period Λ of the grating is given by

$$\Lambda = \frac{2\pi}{|\vec{K}|} = \frac{\lambda}{2 \sin \theta}, \quad (9-48)$$

as asserted in an earlier section.

If the photographic plate is developed, silver atoms will appear concentrated along the quasi-planes of high exposure, which we will call silver "platelets". The distance between these platelets is the period Λ specified above.

9.7.2 Reconstructing Wavefronts from a Volume Grating

Suppose that we attempt to reconstruct the original object plane wave by illuminating the volume grating with a reconstruction plane wave. The question naturally arises as to what angle of illumination should be used to obtain a reconstructed object wave of maximum intensity. To answer this question, we may regard each platelet of high silver concentration as a partially reflecting mirror, which diverts part of the incident wave according to the usual laws of reflection, and transmits part of the wave. If the plane-wave illumination is incident on the silver platelets at angle α , as shown in Fig. 9.24, then the reflected wave will travel in the direction satisfying the law of reflection. However, such reflections occur at all the platelets, and if the various reflected plane waves are to add in phase, then it is essential that the various path lengths traveled by waves

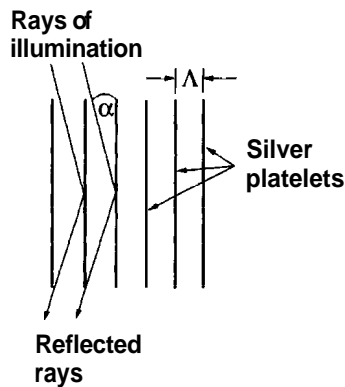


FIGURE 9.24
Reconstruction geometry.

reflected from adjacent platelets differ by precisely one optical wavelength.⁷ With reference to the figure, simple geometry shows that this requirement will be satisfied only if the angle of incidence satisfies the **Bragg condition**,

$$\sin \alpha = \pm \frac{\lambda}{2\Lambda}. \quad (9-49)$$

Comparison of Eqs. (9-48) and (9-49) shows that maximum intensity of the diffracted wave will be obtained only if

$$\alpha = \begin{cases} \pm \theta \\ \pm(\pi - \theta). \end{cases} \quad (9-50)$$

This result is a very important one, for it indicates the condition necessary to obtain a reconstructed plane wave of maximum intensity. Actually, this equation defines a cone of reconstruction angles that will yield the desired results. It is only necessary that the wave vector \vec{k}_p of the reconstruction wave be inclined at angle θ to the planes of the silver platelets. Figure 9.25 shows the allowable cones of incident and diffracted wave vectors. As the reconstruction (or "playback") wave vector \vec{k}_p moves around circle shown, the wave vector of the diffracted light \vec{k}_i moves with it such that the k-vector diagram always closes. Note that it is possible to interchange the roles of \vec{k}_p and \vec{k}_i in this figure and still satisfy the Bragg condition.⁸ The fact that an entire cone of incident k-vectors will diffract strongly from a given volume grating is referred to as "Bragg degeneracy".

9.7.3 Fringe Orientations for More Complex Recording Geometries

The discussion above has focused on the simplest case of a hologram formed by the interference of two plane waves that have equal but opposite angles with respect to the normal to the surface of the recording medium. This case is less restrictive than it

⁷The path-length difference could be any integer number of wavelengths. We consider only a single wavelength, which corresponds to the *first* order diffracted wave.

⁸We use the subscript i on the k-vector of diffracted wave because in most circumstances it is an "image" wave. The i does not stand for "incident".

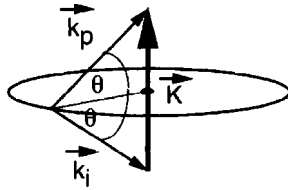


FIGURE 9.25
Cone of incident wave vectors that satisfies the Bragg condition.

might appear, for it is possible to consider two arbitrary wavefronts to be locally planar and their interference to be similar to the interference of two plane waves in any local region, albeit with a different tilt angle with respect to the recording medium than has been assumed here. In all such cases the general principle governing the orientation of the fringe pattern is the same as for the simple case examined: *the fringes formed in the recording medium are always oriented locally to bisect the angle between the two interfering waves within the medium.*⁹

Application of the principle stated above allows one to accurately predict the fringe structures expected in any given case. Figure 9.26 shows several cases of interest, including plane waves interfering to produce slant fringes, plane waves and spherical waves interfering, and waves interfering from opposite sides of the recording medium, a case that leads to a *reflection* hologram.

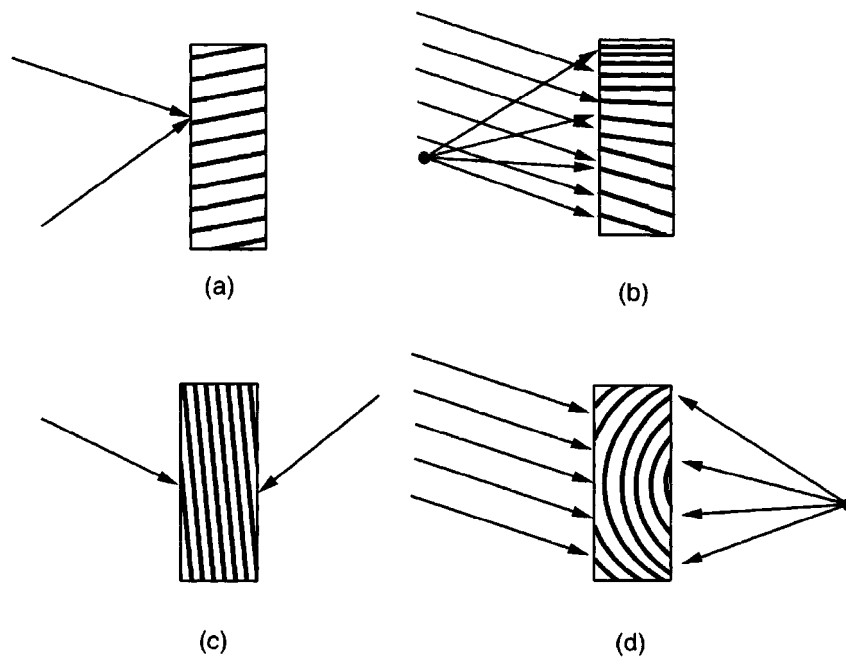
Another general case worth considering is that of two equiphase point sources, perhaps at different distances from the recording medium, generating interference fringes. The fringe peaks form along surfaces for which the difference of distances from the two point sources is an integer multiple of an optical wavelength. Such surfaces are hyperboloids, and any slice through the surface shows hyperboloidal lines of fringe peaks, as shown in Fig. 9.27. Note that if our distance from the two sources is much greater than their separation, and if we examine the fringes over a region that is small compared with the distance from the sources, the fringes will appear to be of an approximately constant spatial frequency determined by the angular separation of the sources, viewed from the recording plane. Notice also that the fringe spacing is smallest when the spherical waves are approaching one another from opposite directions. When the angle between reference and object reaches 180° , Eq. (9-48) implies that the fringe spacing is $\lambda_o/2n$, where n is the refractive index of the recording medium.

9.7.4 Gratings of Finite Size

The theoretical treatments of volume gratings are, for simplicity, often based on the assumption that the grating is infinite in extent. Such is never the case in practice, of course, so it is important to understand the consequences of finite grating size. Such gratings are confined spatially to the finite volume occupied by the recording medium, and usually that volume has a highly asymmetric *shape*.¹⁰ For example, photographic emulsions are usually very much thinner than their lateral extent.

⁹Remember that the angle between two waves within the recording medium is different than the angle between them external to that medium, due to the generally higher refractive index of the recording medium.

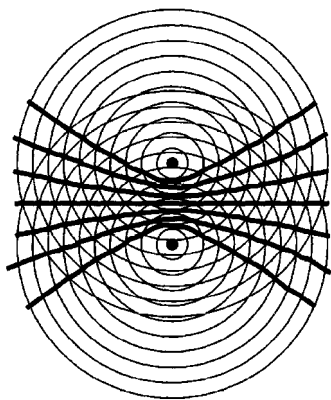
¹⁰An exception is found for nonlinear crystals, which may have sizes that are comparable in all three dimensions.

**FIGURE 9.26**

Orientation of interference fringes within a recording medium. (a) Two plane waves forming slant fringes, (b) a plane wave and a spherical wave, (c) two plane waves impinging from opposite sides of the emulsion, and (d) a plane wave and a spherical wave impinging from opposite sides of the recording medium.

We now present an analysis which is at best a rough approximation to the full description of the effects of finite grating size. The approach is an approximation primarily because it neglects the effects of absorption on the readout beam, but it does provide some physical intuition regarding some of the properties of thick holograms.

For this purpose we use three-dimensional Fourier analysis to express a finite-size grating as a superposition of a multitude of infinite-size gratings, each having a different

**FIGURE 9.27**

Slice through the hyperboloids of fringe maxima for the case of two point sources. The dark lines represent interference fringes, while the lighter lines are the wavefronts.

\vec{K} vector. Suppose that $g(\vec{r})$ represents the local refractive index of a volume phase grating, or the local absorption coefficient of a volume amplitude grating. It is convenient to represent g with a three-dimensional Fourier integral,

$$g(\vec{r}) = \iiint_{-\infty}^{\infty} G(\vec{k}) e^{j\vec{k}\cdot\vec{r}} d^3\vec{k} \quad (9-51)$$

where $G(\vec{k})$ describes the amplitude and phase of k -vector components contained in g , and $d^3\vec{k} = dk_x dk_y dk_z$.

In the special case of a sinusoidal fringe of constant grating vector \vec{K}_g , the form of g is

$$g(\vec{r}) = [1 + m \cos(\vec{K}_g \cdot \vec{r} + \phi_0)] \operatorname{rect} \frac{x}{X} \operatorname{rect} \frac{Y}{Y} \operatorname{rect} \frac{z}{Z}, \quad (9-52)$$

where ϕ_0 is an unimportant spatial phase of the grating, m is the modulation of the grating, and the recording medium has been assumed to have dimensions X, Y, Z in the three rectangular coordinate directions.

The grating-vector spectrum of the above spatially bounded fringe is easily found to be

$$G(\vec{K}) = \left[\delta(\vec{K}) + \frac{1}{2} \delta(\vec{K} - \vec{K}_g) + \frac{1}{2} \delta(\vec{K} + \vec{K}_g) \right] \otimes XYZ \operatorname{sinc} \frac{Xk_x}{2\pi} \operatorname{sinc} \frac{Yk_y}{2\pi} \operatorname{sinc} \frac{Zk_z}{2\pi}. \quad (9-53)$$

The result of this convolution is a blurring of the grating-vector tip into a continuum of grating vectors surrounding the ideal location of the tip of the grating vector for an infinite grating. This blurring operation then leads to the possibility that the k -vector triangle required by the Bragg effect can be closed in many different ways, perhaps at some cost in terms of the strength of the diffracted wave. If the k -vector triangle closes within the central portion of the **primary** lobe of the three-dimensional sinc function above, then the diffraction efficiency should still be near its maximum possible value.

Figure 9.28 shows the effects of the grating-vector cloud on k -vector closure in two different cases. In all cases, the angle of illumination of the grating is assumed to be identical with the angle of the reference wave used during the recording process. In 9.28(a), the grating has been recorded by plane waves incident from the same side of the recording medium, which is assumed much thinner in the z direction than in the other directions. Since the grating-vector cloud is extended in the direction normal to the recording surface, this geometry is quite tolerant to changes of the wavelength of the reconstruction beam (i.e. the length of k_p) relative to that used during recording, but less tolerant to changes of the direction of the reconstruction beam. In 9.28(b), the object and reference waves have come from opposite sides of the emulsion, producing a grating that operates by reflection rather than transmission. In this case the grating-vector blur extends predominantly in a direction along the grating-vector direction. This orientation leads to tolerance with respect to the angle of illumination, and less wavelength tolerance than in the previous case. The degree of tolerance to angular or wavelength changes depends, in both cases, on the thickness of the grating as well as

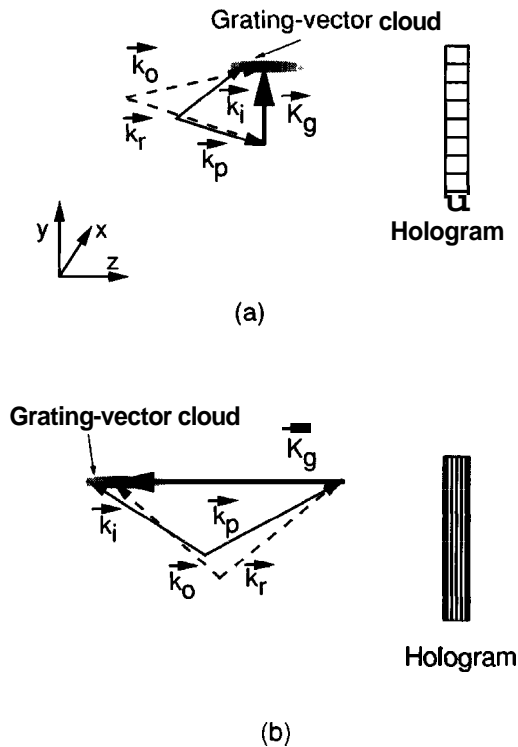


FIGURE 9.28

Grating-vector clouds and their effect on closing the k-vector triangle. The dotted vectors correspond to k vectors when the grating is recorded, and the solid vectors correspond to the k vectors when reconstruction takes place. Changes of the lengths of the k vectors correspond to reconstruction at a different wavelength than was used for recording. In part (a), a change of the length of k_p does not prevent closure of the k-vector diagram. In part (b), a change of the angle of k_p does not prevent closure.

on the period of the fringes, but it is generally true that transmission gratings are more tolerant to wavelength changes than are reflection gratings, and reflection gratings are more tolerant to angle changes than are transmission gratings.

A more exact understanding of the tolerance of volume gratings to the angles and wavelengths of illumination requires a more sophisticated analysis. An example of such an analysis is the coupled mode theory that follows.

9.7.5 Diffraction Efficiency-Coupled Mode Theory

It is extremely important to know the diffraction efficiencies that can be theoretically expected from thick holograms of various types. To find these efficiencies, as well as the tolerance of various types of gratings to the particular angle and wavelength used for reconstruction, it is necessary to embark on a significant analysis. Many methods for calculating these quantities have been found. Most entail some form of approximation, and some are more accurate than others. For an in-depth discussion of a variety of different methods, see Ref. [267]. However, the most widely used method is the coupled mode theory, pioneered by Kogelnik [169], [170] in holography. This is the approach that we shall use here, although we shall follow [139], Chapter 4, most closely. See also Ref. [114] for another useful reference.

The general geometry is illustrated in Fig. 9.29. In this general case, the grating within the emulsion is tilted at angle ψ with respect to the normal to the surface of the recording medium and has grating period $\Lambda = 2\pi/K$. The reconstruction wave is incident at angle θ to that same normal.

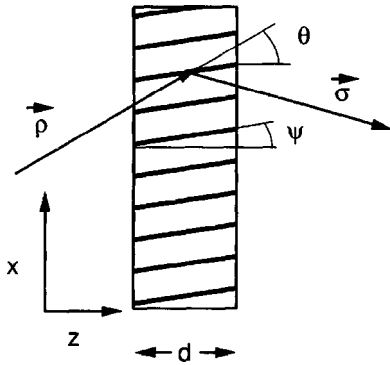


FIGURE 9.29
Geometry for analysis of a thick hologram.

The analysis

The analysis begins with the scalar wave equation,

$$\nabla^2 U + k^2 U = 0, \quad (9-54)$$

valid in a source free region for monochromatic light. The wave number in the most general case is complex-valued, $k = (2\pi n/\lambda_0) + ja$, where a is the absorption constant and λ_0 is the vacuum wavelength. The refractive index n and the absorption constant a within the grating are assumed to vary in sinusoidal fashions according to

$$\begin{aligned} n &= n_0 + n_1 \cos \vec{K} \cdot \vec{r} \\ \alpha &= \alpha_0 + \alpha_1 \cos \vec{K} \cdot \vec{r}, \end{aligned} \quad (9-55)$$

where $\vec{r} \sim (x, y, z)$ and \vec{K} is the grating vector. The hologram is assumed to lie with its faces parallel to the (\mathbf{x}, \mathbf{y}) plane and to be of thickness d in the z dimension.

A number of assumptions are needed for simplification of the problem of solving the wave equation. First it is assumed that the hologram is thick enough that only two waves need be considered within the grating. One is the reconstruction or playback wave $U_p(\vec{r})$, which is gradually depleted by diffraction and absorption, and the other is the first-order Bragg-matched grating order $U_i(\vec{r})$. We assume that the total field within the grating is composed of a sum of these two waves, and we accordingly write that field as

$$\begin{aligned} U(\vec{r}) &= U_p(\vec{r}) + U_i(\vec{r}) \\ &= R(z) e^{j\vec{\rho} \cdot \vec{r}} + S(z) e^{j\vec{\sigma} \cdot \vec{r}}, \end{aligned} \quad (9-56)$$

where the symbols $\vec{\rho}$ and $\vec{\sigma}$ are conventionally used in place of what would be \vec{k}_p and \vec{k}_i , respectively, in our previous notation. We assume that the wave vector $\vec{\rho}$ of R is that of the playback wave in the absence of coupling, and that the wave vector $\vec{\sigma}$ of the diffracted wave is given by

$$\vec{\sigma} = \vec{\rho} - \vec{K}. \quad (9-57)$$

In addition, it is assumed that absorption in a distance of one wavelength is small and that the variation of the refractive index is small compared to its mean,

$$\begin{aligned} n_0 k_o &\gg \alpha_0 \\ n_0 k_o &\gg \alpha_1 \\ n_0 &\gg n_1, \end{aligned} \quad (9-58)$$

where k , is the vacuum wave number, $k_o = 2\pi/\lambda_o$.

It is now possible to expand and simplify k^2 for use in the wave equation as follows:¹¹

$$\begin{aligned} k^2 &= [k_o(n_0 + n_1 \cos \vec{K} \cdot \vec{r}) + j(\alpha_0 + \alpha_1 \cos \vec{K} \cdot \vec{r})]^2 \\ &\approx B^2 + 2jB\alpha_0 + 4\kappa B \cos \vec{K} \cdot \vec{r}, \end{aligned} \quad (9-59)$$

where liberal use of the approximations (9-58) has been made, $B = k_o n_0$, and κ is the coupling constant, given by

$$\kappa = \frac{1}{2}(k_o n_1 + j\alpha_1). \quad (9-60)$$

The next step is to substitute the assumed solution (9-56) and the expression for k^2 above into the wave equation (9-54). During the substitution, $R(z)$ and $S(z)$ are assumed to be slowly varying functions of z so that their second derivatives can be dropped, the term $\cos \vec{K} \cdot \vec{r}$ is expanded into its two complex-exponential components, and $\vec{\sigma}$ is replaced according to (9-57). Terms with wave vectors $\vec{\sigma} - \vec{K} = \vec{\rho} - 2\vec{K}$ and $\vec{\rho} + \vec{K} = \vec{\sigma} + 2\vec{K}$ are dropped, since they correspond to propagation directions that are far from satisfying the Bragg condition. Finally, equating the sum of all terms multiplying $\exp[j\vec{\rho} \cdot \vec{r}]$ to zero and similarly for the sum of all terms multiplying $\exp[j\vec{\sigma} \cdot \vec{r}]$, we find that $R(z)$ and $S(z)$ must individually satisfy the following equations in order for the wave equation to be satisfied:

$$\begin{aligned} c_R \frac{dR}{dz} + \alpha_0 R &= j\kappa S \\ c_S \frac{dS}{dz} + (\alpha_0 - j\zeta) S &= j\kappa R, \end{aligned} \quad (9-61)$$

where ζ is called the "detuning parameter", given by

$$\zeta = \frac{B^2 - |\vec{\sigma}|^2}{2B}, \quad (9-62)$$

and the quantities c_R and c_S are given by

$$\begin{aligned} c_R &= \frac{\rho z}{B} = \cos \theta \\ c_S &= \frac{\sigma z}{B} = \cos(\theta - 2\psi), \end{aligned} \quad (9-63)$$

where θ and ψ are defined in Fig. 9.29.

¹¹Note that the definition of α used here is the reciprocal of the propagation distance within which the field drops to $1/e$ of its original value. The intensity drops to $1/e^2$ in this same distance.

The quantity ζ is a measure of the "Bragg mismatch" of the reconstructed wave, and deserves further discussion. Equation (9-57) is a statement of the Bragg matching condition. Using this equation, we see

$$\begin{aligned} B^2 - |\vec{\sigma}|^2 &= B^2 - (\vec{\rho} - \vec{K}) \cdot (\vec{\rho} - \vec{K}) \\ &= B^2 - |\vec{\rho}|^2 + 2\vec{\rho} \cdot \vec{K} - K^2 \\ &= 2\rho K \cos(\psi + \pi/2 - \theta) - K^2 \\ &= 2\rho K \sin(\theta - \psi) - K^2, \end{aligned} \quad (9-64)$$

where $K = |\vec{K}|$ and $p = |\vec{\rho}| = B = k_o n_o$. Thus

$$\zeta = \frac{B^2 - |\vec{\sigma}|^2}{2B} = K \left[\sin(\theta - \psi) - \frac{K}{2k} \right]. \quad (9-65)$$

Note that the quantity in brackets will be zero when the Bragg condition is satisfied. Consider now a departure from the Bragg matched conditions caused by a combination of a small mismatch in the illumination angle $\theta' = \theta_B - \Delta\theta$ and a small mismatch in the wavelength $\lambda' = \lambda - \Delta\lambda$. Substitution into (9-65) yields the following expression for the detuning parameter in terms of the angular and wavelength mismatches:

$$\zeta = K \left[\Delta\theta \cos(\theta_B - \psi) - \frac{\Delta\lambda}{2\lambda} \right]. \quad (9-66)$$

It can now be clearly seen that mismatch due to wavelength error grows as the grating period Λ shrinks, and therefore wavelength selectivity will be maximum for counterpropagating object and reference beams, which produce a reflection hologram. Selectivity to angular mismatch can be shown (see Prob. 9-10) to be maximum when the reference and object beams are separated by an angle of 90° . With the help of Eq. (9-66), we can estimate the value of the detuning parameter for any combination of angular or wavelength mismatch.

Returning to the coupled wave equations, note that the equation for S contains a driving or forcing term on the right that depends on the incident wave R . It is this term that leads to a transfer of energy from the incident wave to the diffracted wave. If the coupling constant κ is zero, no such coupling will occur. The detuning parameter ζ , if sufficiently large, will swamp the driving term in R , leading to a spoiling of the coupling phenomena due to phase mismatch through the coupling region. In addition, the equation for the amplitude of the incident wave contains a driving term that depends on the diffracted wave, leading to coupling from that wave back into the incident wave.

For all specific solutions discussed in the following, we assume that the grating is unslanted. For a transmission grating, this implies that $\psi = 0$ while for a reflection grating, $\psi = 90^\circ$.

Solution for a thick phase transmission grating

For a pure phase grating we have $\alpha_0 = \alpha_1 = 0$. For a transmission geometry, the boundary conditions to be applied to the differential equations (9-61) are $R(0) = 1$ and $S(0) = 0$. The solution for the diffracted wave S at the exit of the grating ($z = d$) then takes the form

$$S(d) = je^{j\chi} \frac{\sin\left(\Phi \sqrt{1 + \frac{\chi^2}{\Phi^2}}\right)}{\sqrt{1 + \frac{\chi^2}{\Phi^2}}}, \quad (9-67)$$

where¹²

$$\Phi = \frac{\pi n_1 d}{\Lambda \cos \theta}$$

$$\chi = \frac{\zeta d}{2 \cos \theta} = \frac{Kd}{2 \cos \theta} \left[\Delta \theta \cos(\theta - \psi) - \frac{\Delta \lambda}{2\Lambda} \right]. \quad (9-68)$$

The diffraction efficiency of the grating is given by

$$\eta = \frac{|S(d)|^2}{|R(0)|^2} = \frac{\sin^2\left(\Phi \sqrt{1 + \chi^2/\Phi^2}\right)}{1 + \chi^2/\Phi^2}. \quad (9-69)$$

When the grating is illuminated at the Bragg angle with the same wavelength used during recording, the parameter χ is identically zero, yielding for the diffraction efficiency

$$\eta_B = \sin^2 \Phi. \quad (9-70)$$

The diffraction efficiency is seen to increase initially with increasing thickness of the grating, reach a maximum of 100%, fall to zero, rise to 100%, etc., periodically. Since the grating is lossless, the power in the undiffracted wave oscillates similarly but with minima and maxima interchanged. The first maximum of 100% for the diffraction efficiency is reached when $\Phi = \pi/2$, or when

$$\frac{d}{\cos \theta} = \frac{\lambda}{2n_1}. \quad (9-71)$$

Figure 9.30 shows the oscillations of the diffracted power and the undiffracted power as a function of the parameter Φ .

When the grating is illuminated off of the Bragg angle or with a different wavelength than was used during recording, the parameter χ is nonzero. Figure 9.31 shows a three-dimensional plot illustrating efficiency as a function of the both Φ and χ . It can be seen that for any fixed value of Φ , an increase in χ leads to a loss of diffraction efficiency, although oscillations of a diminishing magnitude occur for some values of Φ . This figure is useful when either Φ or χ is fixed and we wish to understand the effect of changing the other parameter. Note, however, that both parameters are proportional to the grating thickness d , so if the behavior as a function of thickness is of interest, a slice through the surface at some angle with respect to the Φ axis is required.

¹²In expressions involving both λ and θ , it is possible to take them to have their values outside the emulsion or inside the emulsion, as long as the same condition holds for both (cf. Prob. 9-7(a)).

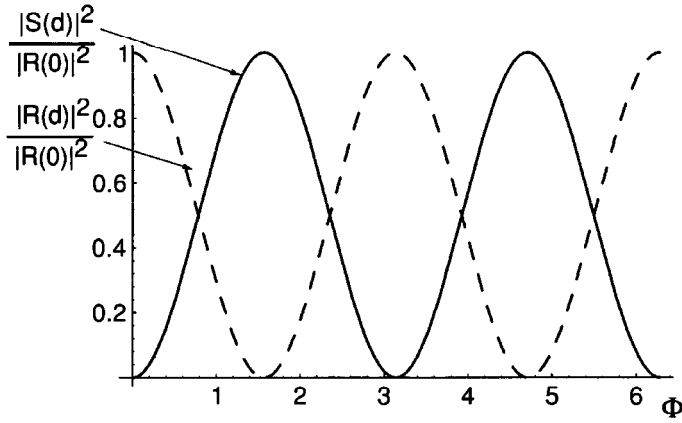


FIGURE 9.30
Normalized intensities of the diffracted and undiffracted waves as a function of Φ for the Bragg matched case.

Solution for a thick amplitude transmission grating

For an unslanted amplitude grating, the index modulation n_1 is zero and α_1 is nonzero. The appropriate boundary conditions are the same as for the phase transmission grating, $R(0) = 1$ and $S(0) = 0$. The solution for the diffracted amplitude at the grating output is now given by

$$S(d) = -\exp\left(-\frac{\alpha_0 d}{\cos \theta}\right) e^{j\chi} \frac{\sinh\left(\Phi_a \sqrt{1 - \chi^2/\Phi_a^2}\right)}{\sqrt{1 - \chi^2/\Phi_a^2}}, \tag{9-72}$$

where \sinh is a hyperbolic sine function,

$$\Phi_a = \frac{\alpha_1 d}{2 \cos \theta}, \tag{9-73}$$

and again

$$\chi = \frac{\zeta d}{2 \cos \theta}. \tag{9-74}$$

For Bragg matched conditions, $\chi = 0$ and the diffraction efficiency is given by

$$\eta_B = \exp\left(-\frac{2\alpha_0 d}{\cos \theta}\right) \sinh^2\left(\frac{\alpha_1 d}{2 \cos \theta}\right). \tag{9-75}$$

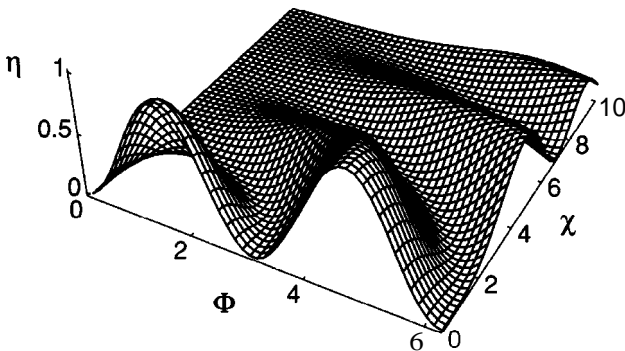


FIGURE 9.31
Diffraction efficiency of a thick phase transmission grating when Bragg mismatch is present.

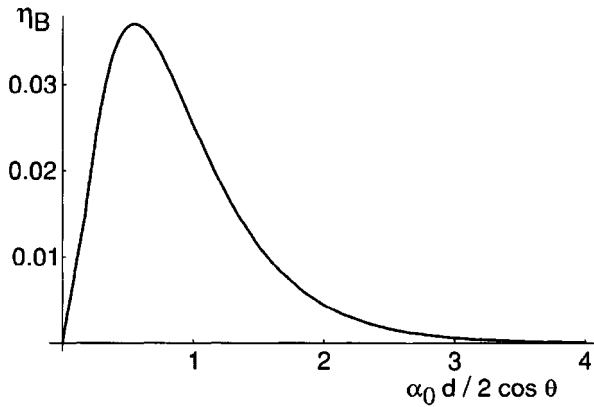


FIGURE 9.32
Maximum possible Bragg matched diffraction efficiency vs. Φ'_a for a thick amplitude transmission grating.

This solution is a product of two functions, the first of which simply represents the attenuation of light due to the average absorption coefficient α_0 as it propagates through the distance $d \cos \theta$ in the hologram. The second term represents the rising effect of diffraction as the wave propagates through increasing thickness. The absorption can never be negative in a passive medium, and therefore the modulation of absorption can never exceed the average absorption, $\alpha_1 \leq \alpha_0$. Because of this constraint, the two terms balance one another in such a way that there is an optimum thickness where diffraction efficiency is maximized.

Diffraction efficiency will be maximized if the attenuation modulation is taken to have its largest possible value, $\alpha_1 = \alpha_0$. Defining $\Phi'_a = \frac{\alpha_0 d}{2 \cos \theta}$, this maximum diffraction efficiency can be expressed as

$$\eta_B = \exp(-4\Phi'_a) \sinh^2(\Phi'_a) \tag{9-76}$$

under Bragg matched conditions, a plot of which is shown in Fig. 9.32. This expression takes on a maximum value of 0.037 or 3.7% for $\Phi'_a = 0.55$.

Figure 9.33 shows a three-dimensional plot of the maximum possible diffraction efficiency (again, $\alpha_1 = \alpha_0$) with the quantity Φ'_a running from the left and the quantity χ running into the right, thus illustrating the effects of Bragg mismatch on the diffraction efficiency. Note that when Φ'_a is in the vicinity of the value that yields maximum diffraction efficiency, values of χ of the order of 2.5 will drive the diffraction efficiency to near zero.

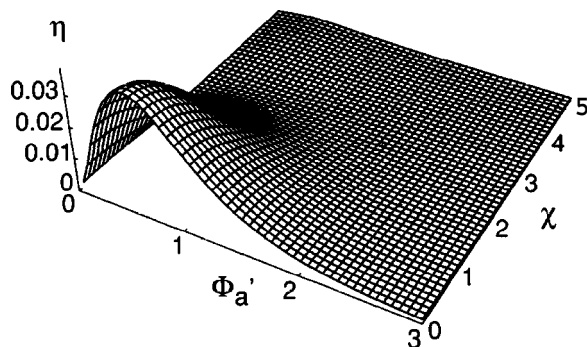


FIGURE 9.33
Diffraction efficiency of a thick amplitude transmission grating with Bragg mismatch.

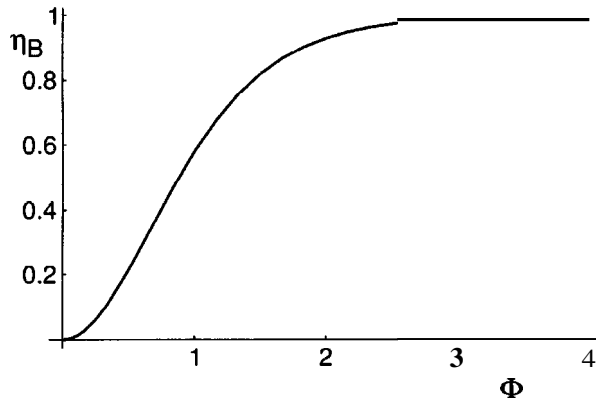


FIGURE 9.34
Diffraction efficiency of a thick Bragg matched phase reflection grating.

Solution for a thick phase reflection grating

For a reflection grating, the grating planes run nearly parallel with the face of the recording medium. In what follows we assume for simplicity that the grating is unslanted, i.e. that the grating planes are exactly parallel to the surface ($\psi = 90^\circ$). The boundary conditions change, now being $R(0) = 1$ and $S(d) = 0$ (i.e. the diffracted wave is now growing in the "backwards" direction). Again for a pure phase grating, $\alpha_0 = \alpha_1 = 0$. The solution of the coupled mode equations for the amplitude of the diffracted wave is now

$$S(0) = -j \left[-j \frac{\chi}{\Phi} + \sqrt{1 - \frac{\chi^2}{\Phi^2}} \coth \left(\Phi \sqrt{1 - \frac{\chi^2}{\Phi^2}} \right) \right]^{-1}, \quad (9-77)$$

where Φ and χ are again given by Eqs. (9-68) and \coth is a hyperbolic cotangent. The diffraction efficiency then becomes¹³

$$\eta = \left[1 + \frac{1 - \frac{\chi^2}{\Phi^2}}{\sinh^2 \left(\Phi \sqrt{1 - \frac{\chi^2}{\Phi^2}} \right)} \right]^{-1}. \quad (9-78)$$

Under Bragg matched conditions, $\chi = 0$, and the diffraction efficiency can be expressed as

$$\eta_B = \tanh^2 \Phi, \quad (9-79)$$

where \tanh is a hyperbolic tangent. Figure 9.34 shows a plot of this diffraction efficiency vs. the parameter Φ . As can be seen, the diffraction efficiency asymptotically approaches 100% as the parameter Φ increases.

The behavior of the diffraction efficiency with Bragg mismatch is illustrated in the three-dimensional plot of Fig. 9.35. In this figure we have interchanged the directions

¹³When evaluating this equation under the condition $\chi > \Phi$, use must be made of the fact that $\sinh iu = i \sin u$.

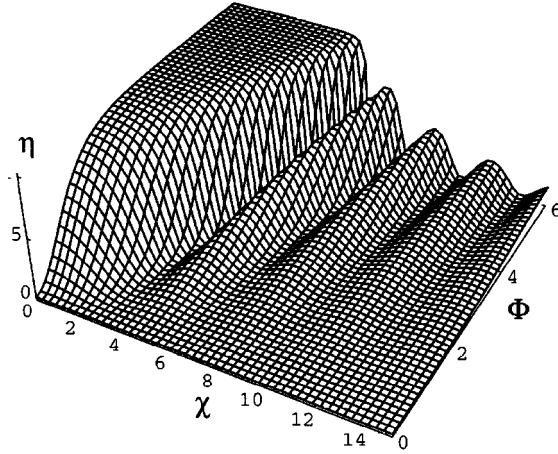


FIGURE 9.35
Diffraction efficiency of a thick phase reflection grating when Bragg mismatch is present.

of Φ and χ from those of the other cases shown previously in order to better reveal the structure of the surface.

Solution for a thick amplitude reflection grating

The final case of interest is that of a thick amplitude reflection grating. The boundary conditions are the same as for the previous case, but now $n_1 = 0$ and diffraction is caused by variations a_1 of the absorption coefficient.

Solution of the coupled-wave equations now yields the following expression for the diffracted amplitude:

$$S(0) = -j \left[-j \frac{\chi_a}{\Phi_a} + \sqrt{1 - \frac{\chi_a^2}{\Phi_a^2}} \coth \left(\Phi_a \sqrt{1 - \frac{\chi_a^2}{\Phi_a^2}} \right) \right]^{-1} \quad (9-80)$$

where Φ_a is again given by Eq. (9-73) and

$$\chi_a = \frac{\alpha_0 d}{\cos \theta} + \frac{j \zeta d}{2 \cos \theta}. \quad (9-81)$$

Under Bragg matched conditions, ζ goes to zero. Again maximum diffraction efficiency will be achieved if the variations of absorption have their largest possible value, $\alpha_1 = \alpha_0$. Under these conditions, $\chi_a/\Phi_a = 2$, and

$$\eta_B = \left[2 + \sqrt{3} \coth(\sqrt{3} \Phi_a) \right]^{-2} \quad (9-82)$$

which is shown plotted vs. Φ_a in Fig. 9.36. The diffraction efficiency is seen to asymptotically approach its maximum value of 0.072, or 7.2%.

Under Bragg mismatched conditions, again with the largest possible modulation of absorption, the following expression holds for χ_a :

$$\chi_a = 2\Phi_a + j\chi, \quad (9-83)$$

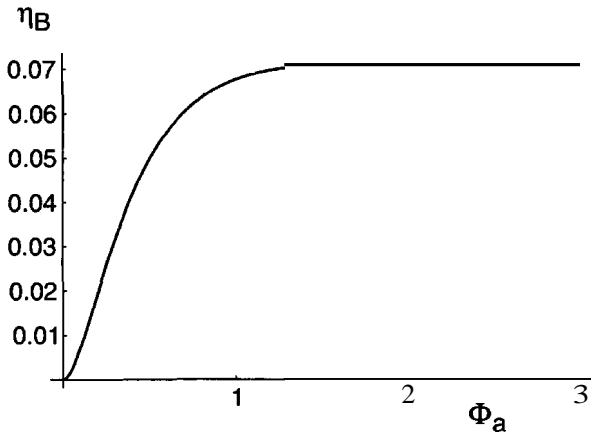


FIGURE 9.36 Bragg matched diffraction efficiency of a thick amplitude reflection grating.

where χ is as given in Eq. (9-68) earlier. Thus the expression for the Bragg mismatched diffraction efficiency can be written

$$\eta = \left| 2 + j \frac{\chi}{\Phi_a} + \sqrt{\left(2 + j \frac{\chi}{\Phi_a}\right)^2 - 1} \coth \left(\Phi_a \sqrt{\left(2 + j \frac{\chi}{\Phi_a}\right)^2 - 1} \right) \right|^{-2} \quad (9-84)$$

Figure 9.37 illustrates the dependence of diffraction efficiency on Φ_a and χ , again with the display rotated to make its structure most clear. The broadening tolerance to Bragg mismatch as the parameter Φ_a increases is a result of the increasing absorption of the grating, and a resulting decrease of its effective thickness.

Summary of maximum possible diffraction efficiencies

In Table 9.1 the various maximum possible diffraction efficiencies possible with thick gratings of various kinds are summarized. For comparison purposes, recall that for a thin sinusoidal amplitude grating the maximum possible diffraction efficiency is 6.25% and for a thin sinusoidal phase grating the maximum is 33.8%.

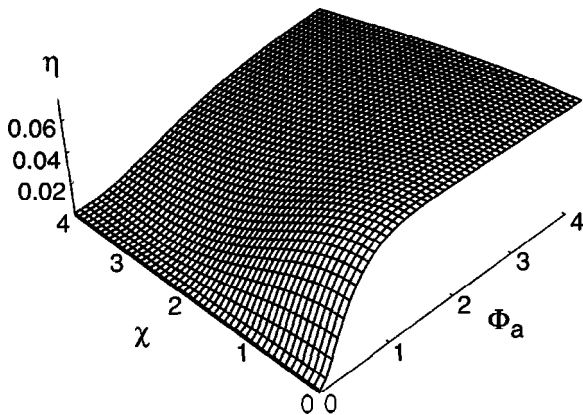


FIGURE 9.37 Diffraction efficiency of a thick amplitude reflection hologram when Bragg mismatch is present.

TABLE 9.1

Maximum possible diffraction efficiencies of volume sinusoidal gratings.

Phase transmission	Amplitude transmission	Phase reflection	Amplitude reflection
100%	3.7%	100%	7.2%

9.8 RECORDING MATERIALS

Holograms have been recorded in a vast number of different materials during the history of holography. In this section we will offer a brief review of some of the more important recording materials. Unfortunately space limitations prevent a complete coverage here. For further information, the reader may wish to consult any of the textbooks on holography already cited. A particularly relevant reference is [266].

9.8.1 Silver Halide Emulsions

The most widely used recording materials for holograms are certainly those based on silver halide photographic technology. A detailed review of such materials, with particular reference to holography, can be found in [22]. It should be noted at the start that a major distinguishing characteristic of holographic materials is that they must be capable of recording extremely fine structures, as compared with the structures encountered in ordinary photography. The spatial frequency response of holographic recording materials often exceeds 2000 cycles/mm, whereas in conventional photography, a spatial frequency response of 200 cycles/mm is considered high. A corollary to this fact is that high resolution is always accompanied by low sensitivity. High resolution is achieved by constructing emulsions with very small grain sizes, but a certain number of photons must be absorbed by each grain to make it developable. It follows that the energy densities needed to expose high-resolution materials are much greater than those required for low-resolution materials.

There are four major classes of silver-halide emulsions that have found use in holography, classes that can be distinguished by the manufacturer. The first material to be widely used was a high-resolution plate manufactured by Kodak, especially the Spectroscopic Plate 649F, which is sensitive across the entire visible spectrum, including the red (emulsions with this broad wavelength sensitivity are called panchromatic). The emulsion thickness for this plate is about 15 μm . A common measure of sensitivity is the energy density required to achieve an amplitude transmittance of 0.5. For 649F plate about 80 $\mu\text{J}/\text{cm}^2$ are required. The same emulsion can be obtained on film, but with a smaller thickness (5 μm). Other Kodak materials are also available, as will be summarized in Table 9.2.

The second manufacturer is Agfa-Gevaert, and in particular its Scientia series, including plates numbered 8E75 HD (red sensitive, < 5000 cycles/mm), 8E56 (blue-green sensitive, < 5000 cycles/mm), 10E75 (red sensitive, < 2800 cycles/mm), and

TABLE 9.2
Properties of some silver halide materials.

Material	Emulsion thickness [μm]	Spectral sensitivity [nm]	Sensitivity [$\mu\text{J}/\text{cm}^2$]	Resolving power [cycles/mm]
Kodak				
649-F (plate)	17	< 700	80	> 2,000
649-F (film)	6	< 700	80	> 2,000
649-GH	6	< 560	100	> 2,000
120	6	< 700	40	> 2,500
SO-173	6	< 700	40	> 2,500
Agfa-Gevaert				
8E75 HD	7	< 750	10	< 5,000
10E75	7	< 750	1	< 2,800
8E56 HD	7	< 560	25	< 5,000
Ilford				
FT340T (plate)	6	< 700	200	< 7,000
Hotec R (film)	7	< 700	20	< 7,000
SP695T (plate)	6	< 560	100	< 5,000
SP672 (film)	7	< 560	100	< 7,000

10E56 (blue-green sensitive, < 2000 cycles/mm). These materials have been optimized for particular laser wavelengths in each case.

The third manufacturer, Ilford Limited, has a number of new holographic silver halide plates and films, as indicated in Table 9.2.

The last class of materials are those made in the former Soviet Union and in Eastern Europe. These manufacturers took a rather different path than those in the West, actively developing emulsions especially for holography. For a discussion of these emulsions, with many references, see [267] and [22]. We do not include these materials in Table 9.2 because they are not widely available.

Table 9.2, after [22], presents a summary of the relevant properties of the silver-halide materials from the manufacturers mentioned above.

9.8.2 Photopolymer Films

Photopolymer films provide a recording medium with two major virtues: (1) the holograms obtained are primarily phase holograms, and (2) the films can be coated with considerable thickness (as thick as 8 mm). The thick phase holograms that result can have excellent efficiency.

The modulation mechanism for these holograms is a combination of thickness variation and internal refractive index change. The recording material is a photopolymerizable monomer, i.e. a monomer that experiences polymerization or cross-linking under exposure to light. After the initial polymerization of the portions of the monomer exposed to light, diffusion of the remaining monomer takes place away from areas of

high concentration (low exposure). A final uniform exposure polymerizes the remaining monomer, but due to the previous diffusion, the distribution of polymer is now nonuniform, leading to the phase modulation properties of the hologram. Changes of refractive index of 0.2% to 0.5% are possible.

Work on recording volume holograms in photopolymers began in the late 1960s at the Hughes Research Laboratories [67]. Further important work included that of Booth [26], [27], Colburn and Haines [69], and many others. See [267], pp. 293–298, for a more detailed consideration of the history of this material in volume holography.

Photopolymer materials are either self-developing or dry-developing, for example by exposure to UV light. Resolutions are excellent but sensitivities are low, typically a few mJ/cm^2 . E.I. Dupont de Nemours & Co. markets one such material under the name *OmniDex*, and Polaroid offers another material named DMP-128.

9.8.3 Dichromated Gelatin

Dichromated gelatin films are widely used to record extremely efficient volume phase holograms, particularly of the reflection type. Diffraction efficiencies in excess of 90% are readily achieved repeatably.

A gelatin film containing a small amount of a dichromate, such as $(\text{NH}_4)_2\text{Cr}_2\text{O}_7$, is found to harden under exposure to light. The process is a form of molecular cross-linking, similar to that observed in polymer films. Since dichromated gelatin plates are not available commercially, users must make their own photosensitive plates from gelatin films, typically coated on a glass plate. The methods used for preparing such plates and developing them are quite complex and must be performed with great care. A description of these methods can be found, for example, in [139], [254], and [267].

Particularly important publications on this material from an historical point-of-view include [258], [195], [56], [211], and others. Again a more detailed discussion of the history can be found in [267], pp. 278–286.

A number of theories have been proposed to explain the physical mechanism that takes place in the dichromated gelatin plates. Currently the best-accepted theory [51] is that a large number of very tiny vacuoles, with sub-wavelength dimensions, form in unhardened areas of the film. The density of vacuoles changes the local refractive index, allowing smooth and continuous variations of phase shift.

Recording using dichromated gelatin films is carried out typically at 488 nm or 514.5 nm wavelengths in the blue and green, respectively. Emulsion thickness may be of the order of $15\ \mu\text{m}$, and exposures required are of the order of 50 to $100\ \text{mJ}/\text{cm}^2$, a very high exposure indeed.

9.8.4 Photorefractive Materials

A number of crystals, including lithium niobate (LiNbO_3), barium titanate (BaTiO_3), bismuth silicon oxide (BSO), bismuth germanium oxide (BGO), potassium tantalum niobate (KTN), and strontium barium nitrate (SBN), exhibit a combination of sensitivity to light and an electro-optic effect. This combined effect has come to be known as

the photorefractive *effect*, and the materials that exhibit it are known as photorefractives or photorefractive materials. For an excellent background on these materials and their applications in optics, see [137] and [138].

Early work on photorefractive holograms took place at the Bell Laboratories [57] and the RCA Laboratories [2], [3]. Considerable advances in theoretical understanding were developed in this early work. The reader should consult the more general references cited above for a more complete historical picture.

The mechanisms by which incident optical **interference** patterns are stored as changes of the local refractive index of these materials are extremely complex and in some cases not completely understood. Index change is known to occur as a result of charge transport and the electro-optic effect. The charge transport results from **photoexcitation** of trapped carriers, transport of these carriers, and re-trapping at new locations. In some materials (e.g. SBN) the transport mechanism is diffusion, while in others (e.g. LiNbO₃), it may be the photovoltaic effect under some circumstances and diffusion under others. After charge transport and re-trapping, internal electric fields will result from the charge redistribution. These internal electric fields cause, through the electro-optic effect, local changes of the refractive index experienced by polarized light.

Figure 9.38 illustrates an incident sinusoidal intensity pattern and the resulting distributions of charge, electric field, and refractive index. The charge carriers, which in this case carry positive charge, migrate to the nulls of the intensity pattern, establishing a charge distribution that is 180° out-of-phase with the incident intensity distribution. Electric field is proportional to the spatial derivative of charge, and hence the electric field distribution is 90° out-of-phase with the charge distribution and (in the opposite sense) with the intensity distribution. Assuming the linear electro-optic effect, the refractive index change is proportional to the electric field, and a volume index grating, spatially phase-shifted by 90° from the exposure pattern, results.

The 90° phase shift between the exposure pattern and the pattern of refractive index change plays an important role in the transfer of energy between the two interfering beams during the exposure process. The two interfering beams create, in any incremental distance Δz normal to the grating fringes, a weak phase grating with amplitude transmittance of the form

$$t_A(x, y) = \exp\left[j2\pi \frac{n(x, y) \Delta z}{\lambda_o}\right]. \quad (9-85)$$

Since the grating is weak, the argument of the exponential is small, and

$$t_A(x, y) = \exp\left[j2\pi \frac{\Delta n \Delta z}{\lambda_o} \sin 2\pi x/\Lambda\right] \approx 1 + j2\pi \frac{\Delta n \Delta z}{\lambda_o} \sin 2\pi x/\Lambda, \quad (9-86)$$

where Δn is the peak refractive index change in the grating, and Λ is the grating period. Note in particular the 90° phase difference between the zero order (represented by unity) and the combination of the two first orders (represented by the sinusoid) in the last expression. For one of the first-order diffracted components, the spatial shift of the index grating with respect to the incident intensity pattern compensates for the similar phase shift in the above equation, with the result that strong coupling and energy transfer can occur between the two incident beams. In this fashion a strong incident beam can couple to a weak incident beam such that the component diffracted in the direction of the weak beam is amplified by energy transfer from the strong beam.

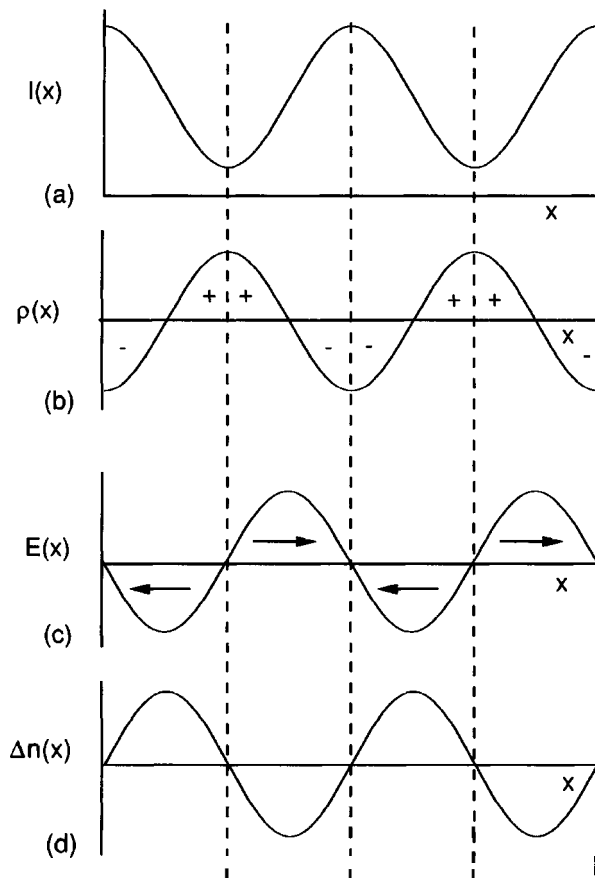


Figure 9.38 Between (a) an incident sinusoidal intensity pattern and the resulting distributions of (b) charge, (c) electric field, and (d) refractive index change in a photorefractive material.

It is often found desirable to apply an external voltage across the photorefractive crystal in a direction orthogonal to the grating planes to induce a drift component of charge transfer. Such a voltage is found to strengthen the diffraction efficiency of the crystal for low spatial frequency components of the interference pattern, whereas without the applied field the diffraction efficiency may be poorer for low frequencies than for high frequencies.

Many photorefractive crystals are extremely slow when compared with photographic emulsions, at least for exposures with typical CW lasers. In fact their response time depends on the rate at which energy is delivered to them, and therefore a recording can be made in a very short time (e.g. a few nsec) with a powerful pulsed laser.

The chief difficulty found with the use of photorefractive materials as a medium for holography is the fact that the reconstruction beam will partially or totally erase the stored hologram as the image is read out. While in some simple cases it is possible to read out the image with a different wavelength than was used for recording, in particular a wavelength to which the crystal is not sensitive, this is not possible in general due to the impossibility of properly Bragg matching the entire set of gratings that were recorded for a complex object when there is a wavelength change. Various methods for "fixing" the recorded holograms have been investigated and remain an area of active research.

Photorefractive crystals have found applications in interferometry, adaptive optics, holographic memories, and optical signal and image processing. They have formed

the basis for certain types of spatial light modulators. For a review of many of these applications, see [138].

9.9 COMPUTER-GENERATED HOLOGRAMS

The vast majority of holograms are made using interference of coherent light, as described in previous sections. However, a significant amount of study has been given to methods for creating holograms by means of calculations on a digital computer, which are then transferred to a transparency by means of a plotting or printing device. The advantage gained by such a process is that one can create images of objects that in fact never existed in the real physical world. We thus become limited in the creation of images (two-dimensional or three-dimensional) only by our ability to describe that image mathematically, our ability to compute the hologram numerically in a reasonable amount of time, and our ability to transfer the results of that computation to a suitable transparent medium, such as photographic film or plate.

The process of creating a computer-generated hologram can be broken down into three separate parts. First is the computational part, which involves calculation of the fields that the object would produce in the hologram plane if it existed. It is these fields, or an approximation to them, that we wish the hologram to generate. This portion of the problem has itself two distinct parts: (1) a decision as to how many sampling points we should use for the object and the hologram (we can calculate only a discrete set of samples of the desired field starting from a discrete representation of the object); and (2) the carrying out of the correct discrete Fresnel or Fourier transform on the object fields, which is usually accomplished with a fast Fourier transform algorithm.

The second part of the process is the choice of a suitable representation of the complex fields in the hologram plane. The result of the calculation mentioned above is usually a discrete set of samples of a complex field, each sample point having both a magnitude and a phase. In general we cannot create structures that directly control both the amplitude and the phase of the amplitude transmittance in arbitrary ways, so some form of encoding of these quantities into a form suitable for representation on a transparency must be chosen.

The third part of the problem is the transfer of the encoded representation of the fields to a transparency. This plotting or printing operation is constrained by the properties of available computer output devices, whether they be pen plotters, laser printers, or electron-beam lithography machines. In fact, the choice of an encoding step is often influenced by the properties of the plotting device that will be used, so the second and third parts of the problem are not entirely independent. Most plotting devices are capable of writing small rectangles at various locations on the output plane. In some cases those rectangles can be written with gray scale, while in others they are restricted to binary values, *i.e.* transparent or opaque.

Many different methods for creating computer-generated holograms have been discovered, but we are limited by space constraints here to discuss only a few of the most important kinds. For more complete discussions, see [185] and [304]. It should be noted that computer-generated holograms are almost invariably thin holograms, due to the

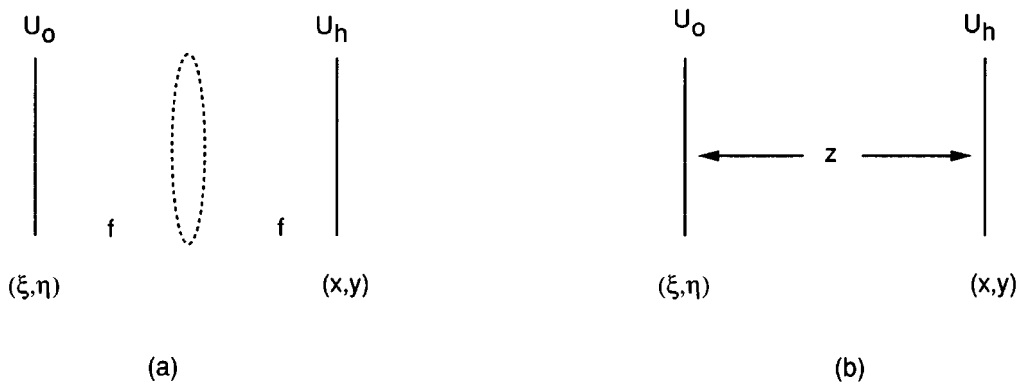


FIGURE 9.39 Geometries when (a) the desired hologram field is the Fourier transform of the object field, and (b) the desired hologram field is the Fresnel transform of the object field.

constraints imposed by the various methods that are used to write such holograms onto transparencies.

9.9.1 The Sampling Problem

The process of holography, whether analog or digital, invariably involves the creation of a complex field in the hologram plane, a field that we wish to regenerate during the wavefront reconstruction process. For computer-generated holograms, we must calculate that field using a digital computer; of necessity the field must be sampled, and complex values computed at each sample point. How many samples of the field must we compute?

To answer this question, we must consider two different situations illustrated in Fig. 9.39. For the case illustrated in part (a) of the figure, we wish to create a hologram field that is the *Fourier transform* of the desired object field. In part (b) of the figure, the goal is to create a hologram field that is the *Fresnel transform* of the object field. We consider each of the cases illustrated in this figure separately.

Fourier hologram

Consider the case of a Fourier hologram first. In this case, a fictitious lens of infinite size and focal length f exists between the object and the hologram field. Since the two fields exist in the front and back focal planes of this lens, the hologram field $U_h(x, y)$ and the object field $U_o(\xi, \eta)$ are related by a Fourier transform,

$$U_h(x, y) = \frac{1}{\lambda f} \iint_{-\infty}^{\infty} U_o(\xi, \eta) \exp\left[-j\frac{2\pi}{\lambda f}(\xi x + \eta y)\right] d\xi d\eta \quad (9-87)$$

where f is the focal length of the lens. The number of samples required in the plane of the hologram field is determined by the bandwidth of that field, as dictated by the Whittaker-Shannon sampling theorem (cf. Sect. 2.4). The size of the object determines the bandwidth of the hologram field in this case, and therefore in turn determines the

number of hologram field samples required. If the dimensions of the object are $L_\xi \times L_\eta$, then the spectrum of the field in the hologram plane is contained within a centered rectangle with dimensions $2B_X \times 2B_Y$, where

$$\begin{aligned} 2B_X &= \frac{L_\xi}{\lambda f} \\ 2B_Y &= \frac{L_\eta}{\lambda f}. \end{aligned} \quad (9-88)$$

The rectangular sampling lattice¹⁴ in the hologram plane should therefore have spacings

$$\begin{aligned} \Delta x &= \frac{1}{2B_X} = \frac{\lambda f}{L_\xi} \\ \Delta y &= \frac{1}{2B_Y} = \frac{\lambda f}{L_\eta}. \end{aligned} \quad (9-89)$$

If the extent of the field in the hologram plane is to be $L_X \times L_Y$, then the number of samples required in that plane will be

$$\begin{aligned} N_X &= \frac{L_X}{\Delta x} = \frac{L_X L_\xi}{\lambda f} \\ N_Y &= \frac{L_Y}{\Delta y} = \frac{L_Y L_\eta}{\lambda f}. \end{aligned} \quad (9-90)$$

It is a straightforward matter to show that the number of samples required in the object plane (from which the hologram field will be calculated) is identically the same as for the hologram field.

Fresnel hologram

Consider now the case of the Fresnel hologram illustrated in part (b) of the figure. The hologram field is no longer related to the object field by a simple Fourier transform, since the lens is missing.

The hologram field and the object field are now related by the Fresnel diffraction integral,

$$U_h(x, y) = e^{j\frac{\pi}{\lambda z}(x^2+y^2)} \iint_{-x}^{\infty} U_o(\xi, \eta) e^{j\frac{\pi}{\lambda z}(\xi^2+\eta^2)} \exp\left[-j\frac{2\pi}{\lambda z}(x\xi + y\eta)\right] d\xi d\eta. \quad (9-91)$$

The relation between the bandwidth of the hologram field and the size of the object is not as obvious as it was in the case of the Fourier hologram.

An excellent approximation to the bandwidth of the hologram field can be obtained from the following argument. The object in this case is viewed as being the function $U_o(\xi, \eta) \exp[j\frac{\pi}{\lambda z}(\xi^2 + \eta^2)]$. The presence of a phase distribution across the object does

¹⁴If the spectrum of the desired field in the hologram plane is not efficiently contained within a centered rectangle, or if a rectangular lattice is not the most efficient means of packing spectral islands in the frequency domain, more efficient forms of the sampling theorem can be used, but we will not dwell on this point here.

not affect its intensity distribution, which is the quantity we wish to re-create from the hologram field. Therefore Eq. (9-91) can be viewed as expressing the hologram field as the product of the quadratic-phase factor in (x, y) and the Fourier transform of the modified object field. We can now consider the bandwidth of each of these terms individually, and by invoking the convolution theorem, we can add the two bandwidths to give us the total bandwidth in the hologram plane.

The bandwidth that arises from the Fourier transform of the modified object will be identical to that obtained in the case of the Fourier hologram, for the presence of the quadratic-phase factor in the object plane has not changed the width of the object. The bandwidth of the quadratic-phase factor in (x, y) can be approximated using local spatial frequencies and taking into account the finite extent of the hologram field (cf. Eq. (2-39)). The local frequencies in the x direction are easily shown to vary between $\pm \frac{L_X}{2\lambda z}$ and those in the y direction between $\pm \frac{L_Y}{2\lambda z}$. The total bandwidth of the hologram field is now found by adding these bandwidths to those obtained in the Fourier hologram case, with the result

$$\begin{aligned} 2B_X &= \frac{L_\xi + L_X}{\lambda z} \\ 2B_Y &= \frac{L_\eta + L_Y}{\lambda z}. \end{aligned} \quad (9-92)$$

Note that the bandwidth now depends on both the extent of the object field and the extent of the hologram field. The sampling intervals in the hologram plane must now be taken to be

$$\begin{aligned} \Delta x &= \frac{\lambda z}{L_\xi + L_X} \\ \Delta y &= \frac{\lambda z}{L_\eta + L_Y}, \end{aligned} \quad (9-93)$$

and the total number of samples in each dimension becomes

$$\begin{aligned} N_X &= \frac{L_X}{\Delta x} = \frac{L_X(L_X + L_\xi)}{\lambda z} \\ N_Y &= \frac{L_Y}{\Delta y} = \frac{L_Y(L_Y + L_\eta)}{\lambda z} \end{aligned} \quad (9-94)$$

Again the number of samples required in the object domain is equal to the number required in the hologram domain. Note that the number of samples required in the Fresnel hologram case is greater than the number required in the Fourier hologram case. See Prob. 9-12 for a further illustration of this fact.

Having determined the sampling requirements for the hologram field, we now turn to a short discussion of the computational problem associated with determining the hologram field from the object field.

9.9.2 The Computational Problem

The relations between the hologram field and the object field represented by Eqs. (9-87) and (9-91) both involve Fourier transforms. After the object and hologram fields have

been sampled, the calculations required to determine the hologram field take the form of a discrete sum. Given that the sampling spacings in the hologram and object planes are $(\Delta x, \Delta y)$ and $(\Delta \xi, \Delta \eta)$, respectively, the case of the Fourier transform hologram requires that the following be calculated:

$$U_h(p\Delta x, q\Delta y) = \sum_{m=0}^{N_X-1} \sum_{n=0}^{N_Y-1} U_o(m\Delta \xi, n\Delta \eta) \exp \left[j2\pi \left(\frac{pm}{N_X} + \frac{qn}{N_Y} \right) \right]. \quad (9-95)$$

Such a transform is referred to as a discrete Fourier transform. It is most rapidly computed using the fast Fourier transform algorithm [37], which for a total of $N_X N_Y$ points requires order $N_X N_Y \log$, $N_X N_Y$ complex multiplications and additions. The computation is fastest if the numbers of samples N_X and N_Y are chosen to be powers of 2.

Due to the constraints of the representational methods to be discussed in the next section, it is often desirable to introduce a randomly and independently chosen phase at each object point, which simulates the presence of a diffuser through which the object is illuminated. Such a step does not change the sampling requirements discussed earlier. The result is a more uniform hologram field, but the price paid is that the images reconstructed by this field will contain speckle.

For a Fresnel hologram, the procedure differs only through the postmultiplication of the discrete Fourier transform obtained by a discrete quadratic-phase function.

9.9.3 The Representational Problem

Having obtained the complex field in the hologram plane, the remaining important step is to adopt a representation of that field that can be encoded in a hologram. Just as with holograms recorded by analog optical means, it is not practical in general to attempt to control both the amplitude and the phase of the amplitude transmittance of a transparency (an exception is the so-called ROACH, to be mentioned below). Therefore some method for encoding complex amplitude into either amplitude or phase is required. We discuss various such methods in what follows. The reader should keep in mind that once a suitable hologram has been plotted by any of the means discussed, it is then necessary to photo-reduce the plot and produce a transparency that can be illuminated with coherent light.

Detour-phase holograms

The oldest and perhaps the best known method for creating holograms from computed complex fields is the so-called "detour-phase" method, invented by Brown and Lohmann [41] and Lohmann and Paris [200]. This method accepts the constraints imposed by most plotting devices, namely that it is easiest to plot binary patterns (ink or no ink) and that convenient basic building blocks are black rectangles that can be centered at any of a quantized set of locations and can be controlled in size at least in certain quantized increments.

Suppose that the final hologram transparency will be illuminated by an off-axis plane wave, and the image will be obtained with a positive lens of focal length f by looking on the optical axis in the focal plane behind the lens. Let the illuminating wave be inclined with respect to the x axis for simplicity, so that its complex field distribution incident on the hologram plane is

$$U_p(x, y) = \exp[-j2\pi\alpha x], \tag{9-96}$$

where α is equal to $\sin 2\theta/\lambda$, 2θ being the angle between the k vector of the incident beam and the normal to the hologram. Then for each value of x on the hologram plane, the optical phase of the illuminating beam has a different value.

Let the hologram plane be divided into $N_x \times N_y$ separate cells, with the width of a cell in the x direction being equal to one full period of the incident phase function, i.e. the width is α^{-1} . The width in the y direction need not necessarily be the same but for simplicity might be chosen so. Each cell defined in this way will encode one of the Fourier coefficients that was calculated with the fast Fourier transform.

Suppose that one particular Fourier coefficient is given by

$$a_{pq} = U_h(p\Delta x, q\Delta y) = |a_{pq}| \exp(j\phi_{pq}). \tag{9-97}$$

Then within that cell we will plot a black rectangle with an area proportional to $|a_{pq}|$ and with a position in the x direction such that at the center of the rectangle, the incident phase from the reconstruction beam is precisely ϕ_{pq} . Remembering that a black rectangle will be changed to a transparent rectangle after the plot is photographed, we have created a transmitted wave component from this cell that has the amplitude of the desired Fourier component and a phase equal to that of the desired Fourier component. Phase shift has been achieved by moving the center of the plotted rectangle, a method known as *detourphase*, and illustrated in Fig. 9.40. Note that our goal is to synthesize an image field of the form

$$U_f(u, v) = \sum_{p=0}^{N_x-1} \sum_{q=0}^{N_y-1} |a_{pq}| e^{j\phi_{pq}} \exp\left[j \frac{2\pi}{\lambda f} (up\Delta x + vq\Delta y)\right], \tag{9-98}$$

which expresses the image field as the sum of its Fourier components, all with proper amplitudes and phases.

To understand the approximations inherent in the detour-phase approach, we undertake a short analysis. Consider first the diffraction pattern created in the rear focal plane of the transforming lens when a single transparent rectangle of widths (w_x, w_y) exists in the hologram plane, that rectangle being describable by the function

$$t_A(x, y) = \text{rect}\left(\frac{x - x_0}{w_x}\right) \text{rect}\left(\frac{y - y_0}{w_y}\right), \tag{9-99}$$

where (x_0, y_0) is the center of the rectangle. When this rectangle is illuminated by the reconstruction wave of Eq. (9-96), the transmitted field is

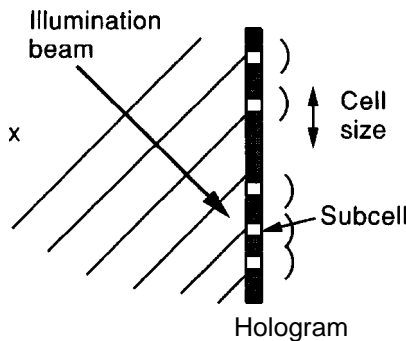


FIGURE 9.40
The detour-phase concept. The subcells are moved within a cell to control the phase of the transmitted light. Zero-phase lines of the reconstruction wavefront are shown.

$$U_i(x, y) = e^{-j2\pi\alpha x} \operatorname{rect}\left(\frac{x - x_0}{w_X}\right) \operatorname{rect}\left(\frac{y - y_0}{w_Y}\right),$$

and the optical Fourier transform of this field is given by

$$U_f(u, v) = \frac{w_X w_Y}{f} \operatorname{sinc}\left[\frac{w_X(u + Af\alpha)}{\lambda f}\right] \operatorname{sinc}\left[\frac{w_Y v}{\lambda f}\right] \exp\left\{j2\pi[(u + Af\alpha)x_0 + vy_0]\right\}, \quad (9-100)$$

where we have made use of the similarity and shift theorems.

If the width w_X of this rectangle is limited in the x direction so that it is a small fraction of the period of the reconstruction beam,

$$w_X \ll \alpha^{-1},$$

then the shift of the first sinc function can be neglected. In addition, if the region of interest in the image plane (size $L_U \times L_V$) is much smaller than the width of the sinc functions, then those functions can be replaced by unity within that region. The resulting approximation to the contribution of this single rectangle can then be written

$$U_f(u, v) = \frac{w_X w_Y}{\lambda f} e^{-j2\pi\alpha x_0} \exp\left[j\frac{2\pi}{\lambda f}(ux_0 + vy_0)\right]. \quad (9-101)$$

Now consider the result of introducing many such rectangles, one for each cell defined in the hologram plane. The cells are indexed by (p, q) , since each cell represents one Fourier coefficient of the image. For the moment we assume that all rectangles are located precisely in the center of their respective cells, but in general each may have a different set of widths (w_X, w_Y) , subject to the constraint on w_X introduced above. Thus the center of the (p, q) th cell is located at

$$\begin{aligned} (x_0)_{pq} &= p\Delta x \\ (y_0)_{pq} &= q\Delta y. \end{aligned} \quad (9-102)$$

The total reconstructed field in the image plane becomes

$$U_f(u, v) = \sum_{p=0}^{N_X-1} \sum_{q=0}^{N_Y-1} (w_X)_{pq} (w_Y)_{pq} e^{-j2\pi p} \exp\left[j\frac{2\pi}{\lambda f}(up\Delta x + vq\Delta y)\right], \quad (9-103)$$

where the period α^{-1} of the reconstruction wave must equal Δx , the width of one cell. Thus when the **subcells** are all centered in their respective cells, the phase of the first exponential term is seen to be an integer multiple of 2π , and that term can be replaced by unity. The terms represented by the second exponential are the Fourier basis functions that we are attempting to add to synthesize the final image. The amplitude of the (p, q) th Fourier component is $w_X w_Y / \lambda f$ and the phases of all components are identical. While we can properly control the amplitudes of the Fourier components by controlling their widths $(w_Y)_{pq}$ (which are not constrained by the limitation imposed on w_X in our earlier approximation), we have not yet controlled the phases of these components properly.

Phase control is introduced by moving the centers of the **subcells** in the x direction within each cell. Suppose that the center of the (p, q) th cell is now located at

$$\begin{aligned}(x_0)_{pq} &= p\Delta x + (\delta x)_{pq} \\ (y_0)_{pq} &= q\Delta y.\end{aligned}\quad (9-104)$$

With this change, the expression of Eq. (9-103) becomes

$$U_f(u, v) = \sum_{p=0}^{N_X-1} \sum_{q=0}^{N_Y-1} (w_X)_{pq}(w_Y)_{pq} e^{-j2\pi\frac{(\delta x)_{pq}}{\Delta x}} \exp\left[j\frac{2\pi}{\lambda f}(up\Delta x + u(\delta x)_{pq} + vq\Delta y)\right] \quad (9-105)$$

where an exponential term with a phase that is an integer multiple of 2π has been replaced by unity. One further approximation is needed. We assume that the width of the image region of interest, which extends in the u direction over $(L_U/2, -L_U/2)$ is sufficiently small that

$$\frac{L_U(\delta x)_{pq}}{2\lambda f} \ll 1,$$

in which case the exponential term $\exp[-j\frac{2\pi}{\lambda f}u(\delta x)_{pq}] \approx 1$, leaving the following expression for the image field:

$$U_f(u, v) = \sum_{p=0}^{N_X-1} \sum_{q=0}^{N_Y-1} (w_X)_{pq}(w_Y)_{pq} e^{-j2\pi\frac{(\delta x)_{pq}}{\Delta x}} \exp\left[j\frac{2\pi}{\lambda f}(up\Delta x + vq\Delta y)\right]. \quad (9-106)$$

This field does have the phases of the Fourier components properly controlled, provided

$$e^{-j2\pi\frac{(\delta x)_{pq}}{\Delta x}} = e^{j\phi_{pq}}.$$

Given the phase ϕ_{pq} of the (p, q) th Fourier component, the subcell in the (p, q) th cell should be centered at $(\delta x)_{pq}$ satisfying

$$-\frac{(\delta x)_{pq}}{\Delta x} = \frac{\phi_{pq}}{2\pi}. \quad (9-107)$$

In addition we choose the width $(w_Y)_{pq}$ of the (p, q) th subcell to be proportional to the desired magnitude of the (p, q) th Fourier component,

$$(w_Y)_{pq} \propto |a_{pq}|. \quad (9-108)$$

$(w_X)_{pq}$ is held constant to satisfy the previous approximation regarding the overlap of the sinc functions. Thus we have created a field in the image plane that is, to within a proportionality constant, equal to the desired field represented by Eq. (9-98). Figure 9.41 illustrates a single cell in the detour-phase hologram.

Once the desired reconstructed field is generated by the hologram, an image will appear in the rear focal plane of a positive lens placed behind the hologram. In fact, as in the case of optically recorded holograms, this type of computer-generated hologram utilizes a carrier frequency α and generates twin images. The second image can be made to appear on the optical axis of the transforming lens if the incident illumination wave is taken to be the conjugate of the previous reconstruction wave, i.e. if its angle with respect to the normal to the hologram is the negative of that in the previous case. Alternatively, the incident wave can be normal to the hologram, in which case both twin images appear with opposite displacements off axis in the rear focal plane.

Figure 9.42 shows (a) a binary detour-phase hologram and (b) an image reconstructed from that hologram.

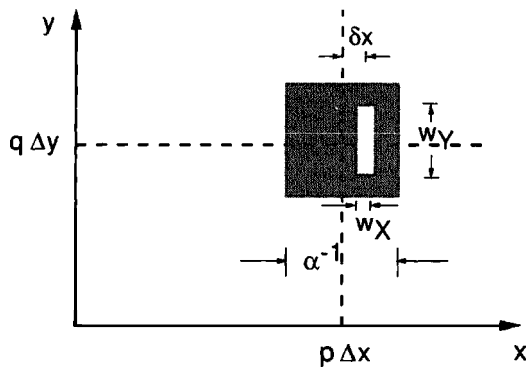
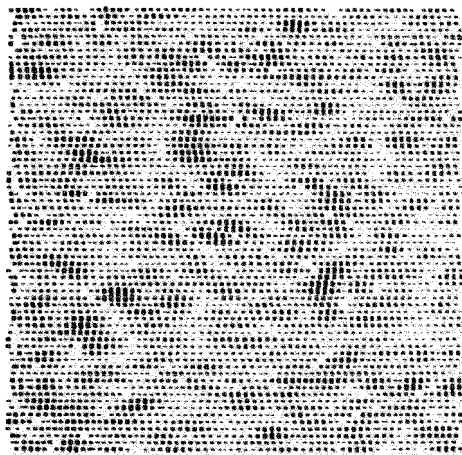
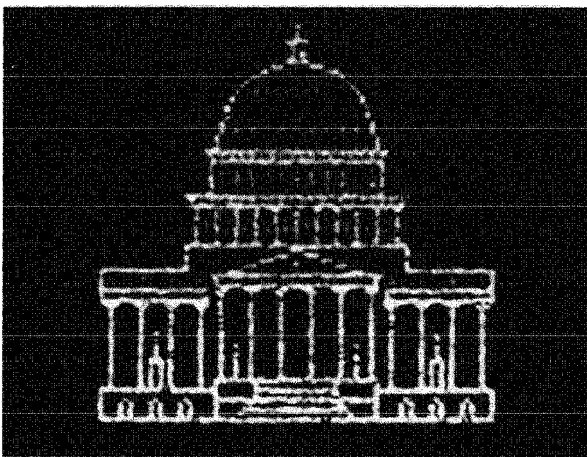


FIGURE 9.41
A single cell in a detour-phase hologram.



(a)



(b)

FIGURE 9.42
(a) Binary detour-phase hologram;
(b) image reconstructed from that
hologram. [Courtesy of A.W. Lohmann.
Copyright 1969 by International
Business Machines Corporation;
reprinted with permission.]

Note that in practice it is necessary to quantize both the amplitude and the phase of a binary detour-phase hologram, a process that leads to noise in the reconstructed image. The effects of phase quantization are particularly important and interesting [129], [77], [78].

Alternative methods of representation using the detour-phase concept exist. For example, Lee [183] utilized four fixed **subcells** per cell, each with a gray-level transmittance, the first representing the real and positive part of the Fourier coefficient (angle 0°), the second the imaginary and positive part (angle 90°), the third the real and negative part (angle 180°), and the fourth the imaginary and negative part (angle 270°). Since the real and imaginary parts are either positive or negative but not both, two of the **subcells** in every cell are normally opaque. Burckhardt [44] recognized that any point in the complex plane can be reached with only three gray-level phasors, one at 0° , one at 120° , and the third at 240° .

The Kinoform and the ROACH

An entirely different method for computer-generated hologram representation is known as the *kinoform* [193]. In this case, an assumption is made that the *phases* of the Fourier coefficients carry the majority of information about an object, and the amplitude information can be entirely eliminated. While this assumption might at first glance appear surprising, it turns out to be quite accurate if the object is a diffuse one, i.e. if the object points all are assigned random and independent phases.

Considering a Fourier geometry again, the hologram is divided up into $N_X \times N_Y$ cells, each representing one Fourier coefficient of the object. The amplitudes $|a_{pq}|$ of all Fourier coefficients are assigned value unity, and it is only the phases ϕ_{pq} that we attempt to encode in the hologram. The encoding is done by linearly mapping the phase range $(0, 2\pi)$ into a continuum of gray levels displayed by an output device such as a photographic plotter. The gray-level transparency obtained from this process is subjected to photographic bleaching. Thus each gray level is mapped into a corresponding phase shift introduced by the transparency, and if the bleaching process is well enough controlled to assure that the complete phase range $(0, 2\pi)$ is exactly and properly realized by the transparency, an excellent image can be obtained in the Fourier plane of the kinoform. In this case there is only a single image and it appears on the optical axis. The diffraction efficiency of the kinoform is very high because it is a pure phase transparency. Errors in "phase matching" the $(0, 2\pi)$ interval result in a single bright spot on the optical axis, generally in the midst of the desired image.

Figure 9.43 shows a photograph of the gray-level recording that leads to a kinoform after bleaching, and the image obtained from the same kinoform.

A related approach known as the "referenceless on-axis complex hologram" (ROACH) utilizes color film to control both the amplitude and the phase of the Fourier coefficients simultaneously [58]. Suppose we wish to create a computer-generated Fourier hologram which will reconstruct an image in red light. Let the magnitudes $|a_{pq}|$ of the Fourier coefficients first be displayed as gray levels on a black-and-white CRT display. This display is photographed through a red-transmitting filter onto a frame of reversal color film. The red-absorbing layer of the three-layer film records this exposure. Then the desired array of Fourier phases is encoded as an array of gray levels,

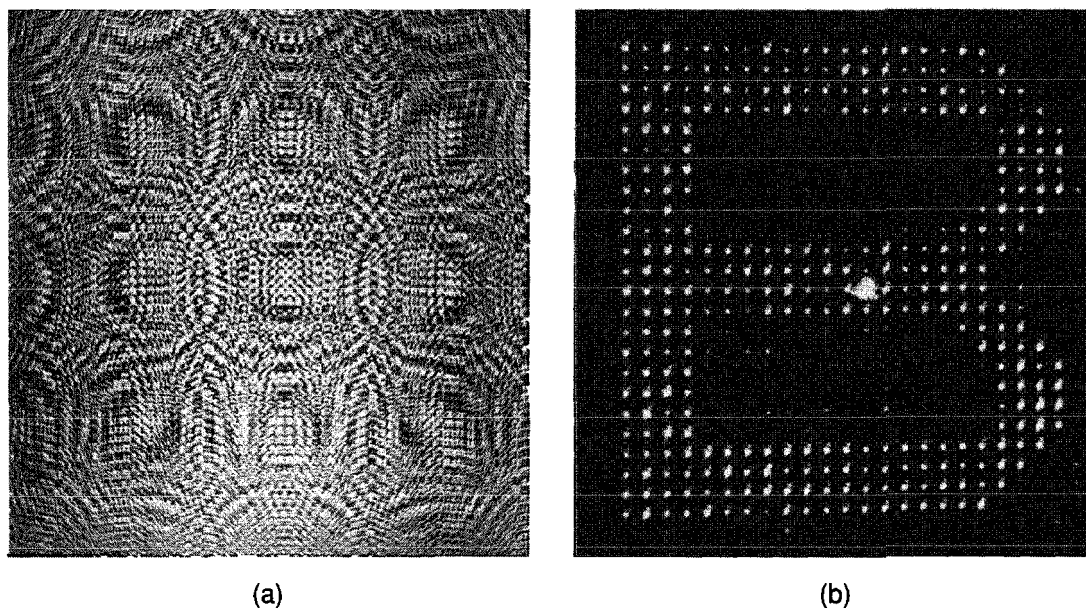


FIGURE 9.43

(a) The gray level display that leads to a kinoform, and (b) the image obtained from that kinoform. [Copyright 1969 by International Business Machines Corporation; reprinted with permission.]

as was done for the kinoform, displayed on the same CRT, and photographed through a blue-green transmitting filter, thus exposing the blue and green absorbing layers of the same frame of film used previously. After processing, the layer exposed to red light becomes absorbing in the red, but the layers exposed to blue-green light are transparent in the red. However, these layers do introduce a phase shift in the transmitted red light, due to thickness variations. Thus the color photographic transparency controls both the amplitude and the phase of the transmitted red light, and as such creates an on-axis image of the desired object. Again proper phase matching is critical, and errors in this regard result in a bright spot of light on axis.

Note that both the kinoform and the ROACH are more efficient than detour-phase holograms in their utilization of the space-bandwidth product of the plotter or display used, since only one resolution cell is required for each Fourier coefficient, whereas many cells are required for the binary hologram. However, both the kinoform and the ROACH require that the phase matching problem be solved, whereas no such problem exists for the detour-phase hologram.

Phase contour interferograms

When phase variations exceeding 2π radians are to be created by the hologram, as is often required for holographic optical elements used in optical testing, detour-phase holograms have the disadvantage that subapertures near the 2π phase boundaries may partially overlap. For such applications, other representational approaches have some advantages. We discuss here only one such approach, namely a method due to Lee [184].

We focus here on the problem of generating elements that control only the phase of the transmitted wavefront. Consider an optical element with ideal amplitude transmittance

$$t_A(x, y) = \frac{1}{2} \{ 1 + \cos[2\pi\alpha x - \phi(x, y)] \}. \quad (9-109)$$

This is a carrier frequency hologram which generates two reconstructed waves of interest, one with a pure phase distribution $\phi(x, y)$ and the other a conjugate term with the negative of this phase distribution. Since this amplitude transmittance contains a continuum of gray levels, it would not be easy to display directly on a plotter and record with high fidelity. We prefer some form of binary pattern for this purpose. Let the plotter create a contour plot of t_A , with one contour line per period, each located at a maximum of the distribution. Such contours are defined by the equation

$$2\pi\alpha x - \phi(x, y) = 2\pi n, \quad (9-110)$$

where each integer n defines a different contour line. Such a plot, when photographically reduced, has been shown by Lee to generate both the desired phase distribution and its conjugate, each in a different first diffraction order [184].

Figure 9.44 shows such a plot generated for the case of a quadratic-phase approximation to a lens, i.e. a phase distribution

$$\phi(x, y) = \frac{\pi}{\lambda f} (x^2 + y^2),$$

where the constant λf has been chosen to be unity for convenience, and $a_x = 2.5$. The photoreduction of such a pattern will yield an optical element that provides the effect of a positive lens in one first diffraction order and the effect of a negative lens in the other first order.

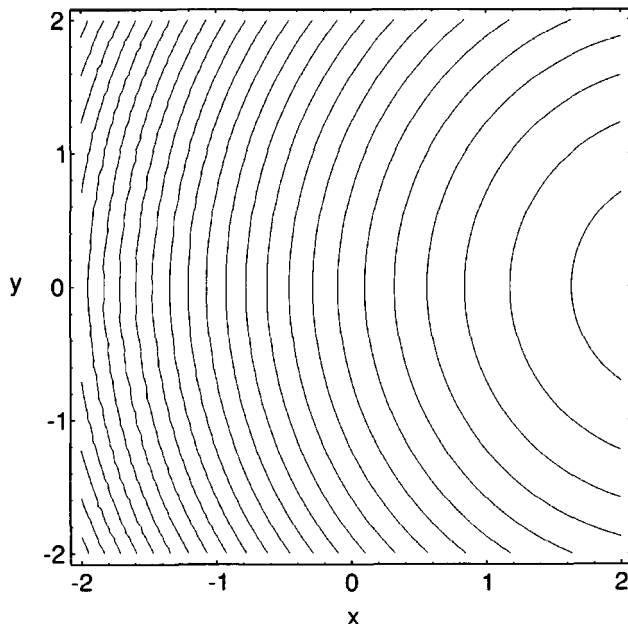


FIGURE 9.44
Plot of a phase contour interferogram for a quadratic-phase approximation to a spherical lens.

Generalizations of this procedure to allow the incorporation of amplitude information in addition to phase information have been demonstrated by Lee [186]. The reader is referred to the original reference for more details.

9.10 DEGRADATIONS OF HOLOGRAPHIC IMAGES

Holographic imaging, like other imaging approaches, suffers from certain degradations that limit the quality of the images that can be obtained. Some degradations, such as that caused by diffraction, are common to all systems. Others, while having a counterpart in conventional photography, manifest themselves in distinctly different ways in holography. In this section we review some of the common sources of image degradations and discuss the effects expected and found in holography.

Holography, like other imaging processes, can suffer from all of the classical aberrations encountered in optics. Consideration of such aberrations is beyond the scope of our treatment. The interested reader can find an excellent discussion in the work of Meier [209]. We mention only that if a hologram is reconstructed with an exact duplicate of the original reference wave at the same wavelength used for recording, and no auxiliary optical elements exist between the object and the hologram and between the hologram and the image plane, the image obtained will be aberration-free (provided there has been no swelling or shrinking of the emulsion on which the hologram was recorded).

The holographic imaging process almost always uses coherent light (for exceptions, cf. Section 9.11). Under usual circumstances, such as operation in the linear region of the t_A vs. E curve for thin holograms, the imaging process has been argued to be linear in complex amplitude, as long as attention is focused on one of the twin images. Under such circumstances it is possible to characterize the holographic process by an amplitude transfer function $H(f_X, f_Y)$. As with nonholographic systems, the amplitude transfer function is determined by the pupil function of the imaging system, and specifies the response of the system to complex-exponential components of the ideal image. Thus in the absence of effects associated with limited film MTF or film **nonlinearities**, we would characterize the holographic imaging process by an amplitude transfer function of the form

$$H(f_X, f_Y) = P(\lambda z_i f_X, \lambda z_i f_Y), \quad (9-111)$$

where P is the pupil function, usually determined by the finite size of the hologram. The amplitude transfer function fully accounts for the limitations to image quality posed by diffraction, and therefore we concentrate on other effects in what follows.

9.10.1 Effects of Film MTF

It has been seen previously that the holographic process in general places heavy requirements on the resolving power of recording materials, requirements that may not always be perfectly met in some applications. It is therefore of some interest to

understand the effects of a limited spatial frequency response (MTF) of a recording material used in holography.

We present analyses of two particularly important recording geometries. For more detailed consideration of the subject, the reader is referred to the classic work of van Ligten [288], [289].

Collimated reference wave

We examine first the effects of the MTF of the recording medium on the mapping of object amplitudes into real-image amplitudes when a plane reference wave is used.

The linearity of the holographic imaging process is not affected by the linear processes (e.g. light scattering) that reduce the spatial frequency response of the recording medium. Therefore it suffices to find the effects of the MTF on one general frequency component of the object amplitude, and to construct the more general result by linear superposition.

To this end, consider the reference and object waves to be given by

$$\begin{aligned} U_r(x, y) &= A \exp(-j2\pi\alpha y) \\ U_o(x, y) &= a \exp[-j2\pi(f_x x + f_y y)], \end{aligned} \quad (9-112)$$

respectively. Thus the object and reference waves are plane waves propagating in different directions. The distribution of exposing intensity is therefore

$$\mathcal{I}(x, y) = A^2 + |a|^2 + 2A|a| \cos\{2\pi[f_x x + (f_y - \alpha)y] + \phi\}, \quad (9-113)$$

where $\phi = \arg a$ and the phase angle of A has been taken as the phase reference. If we represent the MTF of the emulsion as $M(f_x, f_y)$, then from the Kelyey model of the recording process (cf. Section 7.1.5), the *effective* exposing intensity distribution is

$$\mathcal{I}_{\text{eff}}(x, y) = A^2 + |a|^2 + 2A|a|M(f_x, f_y - \alpha) \cos\{2\pi[f_x x + (f_y - \alpha)y] + \phi\} \quad (9-114)$$

where we assumed that the MTF of the recording medium is entirely real. We conclude that the particular Fourier component of the object with spatial frequency (f_x, f_y) is *attenuated* by the factor $M(f_x, f_y - \alpha)$, relative to the amplitude that would be obtained with a perfect recording medium.

Figure 9.45 illustrates this result pictorially. The spectrum of the object is assumed to be centered at the origin in the frequency domain. The MTF of the recording medium

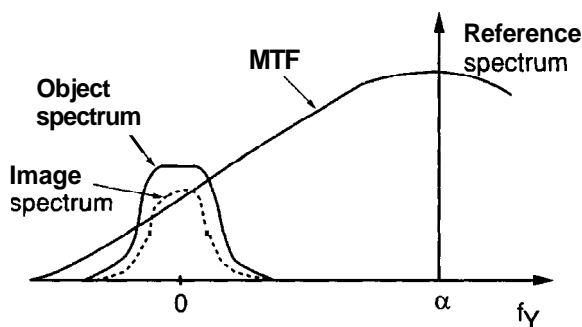


FIGURE 9.45
Effects of the MTF of the recording medium when the reference wave is collimated.

is erected about the frequency ($f_X = 0, f_Y = a$); the product of the object amplitude spectrum and the film MTF yields the image amplitude spectrum.

From the preceding analysis it is clear that, when a collimated reference wave is used, the effect of the recording medium frequency response may be represented by a simple multiplicative factor $M(f_X, f_Y - a)$ applied to the object spectrum. This result directly implies an amplitude transfer function

$$H(f_X, f_Y) = M(f_X, f_Y - \alpha)$$

might be associated with the imaging process. However, the spatial frequency response of the system is limited not only by the effects of the recording medium, but also by the effects of diffraction. Therefore, as mentioned previously the finite pupil function of the hologram aperture must be included. Let

$$P(x, y) = \begin{cases} 1 & (x, y) \text{ in the hologram aperture} \\ 0 & \text{otherwise.} \end{cases}$$

Then the complete amplitude transfer function of the imaging system may be written

$$H(f_X, f_Y) = P(\lambda_2 z_i f_X, \lambda_2 z_i f_Y) M(f_X, f_Y - \alpha), \quad (9-115)$$

where the wavelength used in reconstruction has been assumed to be λ_2 . Note that the effect of the recording medium MTF may be regarded as entirely equivalent to that of inserting an attenuating mask with amplitude transmittance

$$t_A(x, y) = M\left(\frac{x}{\lambda_2 z_i}, \frac{y}{\lambda_2 z_i} - \alpha\right) \quad (9-116)$$

across the pupil of the imaging system.

Fourier transform and lensless Fourier transform holograms

A second type of hologram we shall consider here is one in which each object point is encoded as a fringe pattern of a unique and constant spatial frequency. Such is the case for the Fourier transform hologram and the **lensless** Fourier transform hologram discussed earlier. For both of these types of holograms the reference wave originates from a point that is co-located with the object, and each object point is encoded in a fringe with a spatial frequency that is proportional to the distance of that point from the reference point.

For both types of holograms, the intensity distribution falling upon the recording medium when the object is a point source at coordinates (x_o, y_o) and the reference is at coordinates (x_r, y_r) is

$$\mathcal{I}(x, y) = |a|^2 + 2A|a| \cos \left[2\pi \frac{(x_o - x_r)x}{\lambda_1 z} + 2\pi \frac{(y_o - y_r)y}{\lambda_1 z} + \phi \right], \quad (9-117)$$

where z is the focal length of the lens in the case of the Fourier transform hologram, or the common perpendicular distance of the object and reference points from the recording plane in the case of the **lensless** Fourier transform hologram. Here ϕ is a phase angle that, for the **lensless** Fourier transform hologram, depends on the locations of the reference and the object points but not on the coordinates in the recording plane. For the true Fourier transform hologram, ϕ depends only on the relative phases of the

object and reference points. Thus the object point with coordinates (x_o, y_o) is encoded as a sinusoidal fringe of spatial frequency

$$f_x = \frac{x_o - x_r}{\lambda_1 z}$$

$$f_y = \frac{y_o - y_r}{\lambda_1 z}.$$
(9-118)

To find the effects of the MTF of the recording medium on the image obtained from such a hologram, we find the effective recording intensity by applying the MTF to the sinusoidal fringe in the interference pattern,

$$\mathcal{I}_{\text{eff}}(x, y) = A^2 + |a|^2 + 2A|a|M \left(\frac{x_o - x_r}{\lambda_1 z}, \frac{y_o - y_r}{\lambda_1 z} \right) \cos \left[2\pi \frac{(x_o - x_r)x}{\lambda_1 z} + 2\pi \frac{(y_o - y_r)y}{\lambda_1 z} + \phi \right].$$
(9-119)

If the factor M in this expression is less than unity, then the amplitude of the fringe generated by this object point will be reduced, the diffraction efficiency will be lower, and the light amplitude incident on the image of this particular point will have been reduced by the MTF of the recording medium. Since object points furthest from the reference point generate the highest spatial frequencies, these points will be attenuated the most.

While in the first case examined, i.e. a collimated reference wave, the effect of the MTF was found to be representable by a mask in the pupil plane of the imaging system, in the cases of the Fourier transform and **lensless** Fourier transform holograms, the effect of the MTF is seen to be representable by an attenuating mask placed over the **object** (or equivalently, over the image if magnification is taken into account). The amplitude transmittance of this mask is given by

$$t_A(x_o, y_o) = M \left(\frac{x_o - x_r}{\lambda_1 z}, \frac{y_o - y_r}{\lambda_1 z} \right)$$

and its intensity transmittance by the square of this quantity. Thus for these cases, the effect of the MTF of the recording medium is seen to restrict **the field of view** about the reference point, but not to affect the resolution attained within that field of view.

Since each object point receives a different amplitude weighting in these geometries, the imaging system is seen to be space-variant and there is no amplitude transfer function that can describe the effects of the MTF of the recording medium.

More general recording geometries

Van Ligten's analysis [288], [289] shows that the effects of the MTF of the recording medium are in all cases equivalent to those of an attenuating mask, again with amplitude transmittance proportional to a scaled version of the MTF, and placed at a certain position between the object and the recording plane. The particular position of this attenuating mask depends on the recording geometry used. The two types of geometries examined above represent limiting cases for which the location of the effective

mask is against the pupil of the system on the one hand, and against the object itself on the other hand.

Some further consideration of the effects of the MTF on holographic recordings is found in the problems.

9.10.2 Effects of Film Nonlinearities

Throughout our discussions of holography, we have repeatedly assumed that, at least for thin holograms, the recording medium is exposed in such a way as to assure operation within a linear region of the amplitude transmittance vs. exposure curve. However, real recording media are never perfectly linear in this respect, the deviation from linearity depending to a large degree on the magnitude of the variations of exposure to which the medium is subjected and the exposure bias point selected. In this section we present a brief discussion of the effects of recording medium nonlinearities on the reconstructed image. It should be emphasized that, when the average exposure produced by the object is comparable with that produced by the reference, nonlinear effects can present a serious limitation to image quality. This case may be contrasted with that of a very weak object, for which film-grain noise or other scattered light is generally the limiting factor.

In what follows we omit any detailed analysis of nonlinear effects, preferring to give the reader a set of references that will provide an excellent overview of previous work. Almost all previous work has been devoted to *thin* holograms.

In discussing the effects of nonlinearities on the reconstructed images, it is important to distinguish between two different classes of object. One class consists of objects that contain a collection of isolated point sources; another class consists of diffuse objects, such as a transparency illuminated through a diffuser or a three-dimensional object with an optically rough surface. For both types of objects, the analyses use a model of the nonlinearities first introduced by Kozma [176]. A simpler model was subsequently introduced by Bryngdahl and Lohmann [43].

For objects consisting of collections of point sources, the first analysis was that of Friesem and Zelenka [105], who demonstrated both analytically and experimentally several important effects. First, a phenomenon found with all types of objects, nonlinearities introduce higher-order images, i.e. images in the second, third, or higher diffraction orders. Such extraneous images are not of great concern since they generally do not overlap the first-order image. More important effects occur in the first-order image itself. If the object consists of two point sources, one of greater amplitude than the other, small-signal suppression effects are anticipated and observed. That is, the image of the weaker object point source is suppressed relative to that of the stronger point source. In addition, owing to intermodulation effects, false images may be generated within the first-order image by nonlinear interaction of the two point sources, yielding apparent images of point sources that are not actually present on the object itself.

The effects of film nonlinearities for diffuse objects have also been investigated [127], [178]. In this case the exposure is most properly treated as a random process, employing techniques that are somewhat more complex than required for points-source objects. For further details, the reader may wish to consult the previous references. In this case it is found that the effects of nonlinearities are primarily to

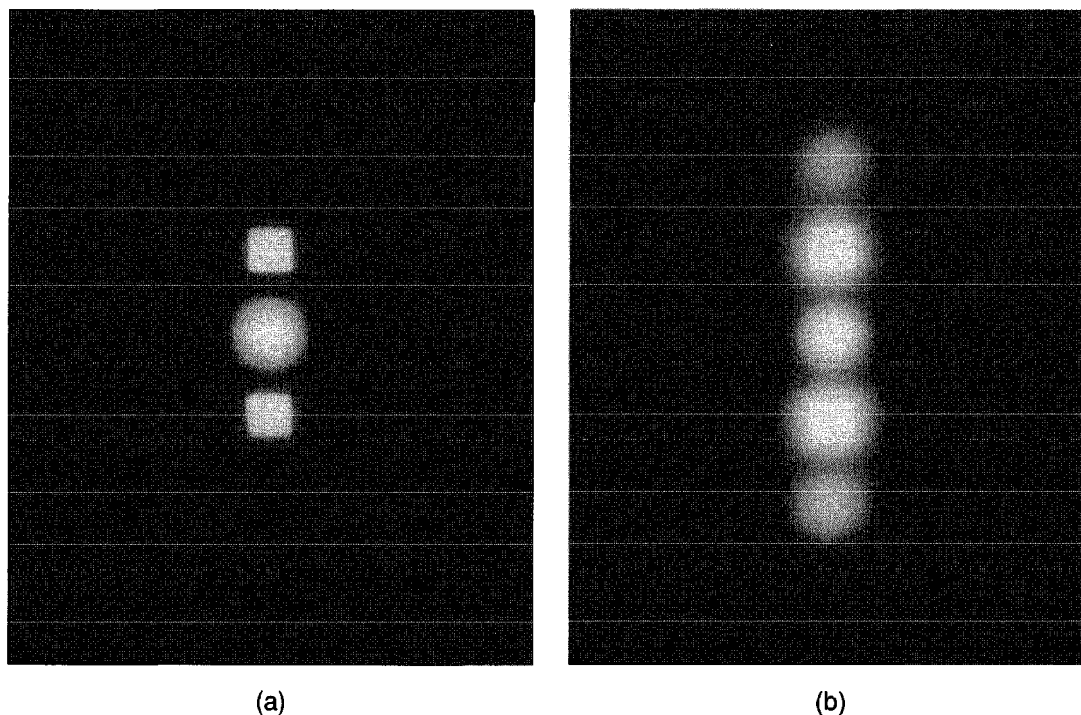


FIGURE 9.46

Effects of film nonlinearities on reconstructed images. The object is a diffuse uniform square patch and the hologram is recorded in the lenless Fourier transform geometry. (a) Twin images obtained when the recording is nearly linear; (b) images obtained under highly nonlinear recording conditions.

introduce a diffuse halo on and around the images of the object. If the diffuse object has fine structure, then the diffuse halo will have related structure. The effects can be quite severe, as illustrated in Fig. 9.46.

9.10.3 Effects of Film-Grain Noise

When the object wave is very weak compared with the reference wave, the primary source of image degradations is often grain noise arising from the film or plate on which the hologram is recorded. The effects of film-grain noise have been analyzed by Goodman [121] and by Kozma [177].

The effects of finite grain size in the photographic emulsions used for holography manifest themselves in a spatial power spectrum of scattered light which overlaps the locations of the desired images. This noise spectrum reduces the contrast of the images obtained, and because it is coherent with respect to the images, it interferes with them to cause unwanted fluctuations of image brightness. Such effects are most noticeable with low resolution films, for which the grain size is relatively large and the scattered light spectrum correspondingly strong. They are also most important when the light arriving from the object during recording is very weak, in which case the statistics of the detection process become quite important.

A particularly important fact, first recognized by **Gabor**, is that holographic recording can in some circumstances provide greater signal detectability than can

conventional photography of the same coherent object. This enhancement comes about from the interference of a strong reference wave with the weak object wave, and the resulting enhancement of the strength of the fringes, a phenomenon analogous to "heterodyne conversion gain" observed in heterodyne detection. Experimental work has shown that such advantages exist in practice [128].

9.10.4 Speckle Noise

For diffuse objects illuminated with coherent light, granularity arising from speckle can be expected in the resulting images, regardless of the imaging method. Since holography is nearly always a process that depends on coherent light, speckle is of special concern in holographic imaging.

When viewing a virtual image, the pupil of the eye is the limiting aperture, and a speckle pattern appears on the retina of the observer. When detecting a real image, using either film or an electronic detector, the aperture of the hologram is the limiting aperture, and together with the distance of the image from the hologram, defines the size of the speckles. If D is the size of the hologram, and z_i is the image distance, then the speckle size is of the order of the diffraction limit, $\lambda z_i/D$. Speckle has been found to reduce the detectability of image detail by a significant factor, particularly when the size of that detail is comparable with the speckle size. Its effects can be suppressed only by smoothing or averaging the intensity over several speckle sizes, for example with detector elements that are several times larger than the diffraction limit. However, if such smoothing is performed, the resolution in the image is reduced accordingly, so whether the reduction of detail comes from smoothing in the detection process or from the inability of the human visual system to extract details when their size is comparable with a speckle, the results are essentially the same.

For a more detailed discussion of the effects of speckle on the ability of the human observer to resolve image detail, see [10] and the references contained therein.

9.11

HOLOGRAPHY WITH SPATIALLY INCOHERENT LIGHT

While holography was originally conceived as a means for coherent image formation, certain techniques exist by means of which holograms of incoherently illuminated objects can be recorded. The extension of holographic techniques to the incoherent case was first suggested by Mertz and Young [210]. The theory and practice of incoherent holography were later extended by Lohmann [198], Stroke and Restrick [277], and Cochran [68]. For additional relevant information, see the book by Rogers [246].

The light from any one point on a spatially incoherent object will not interfere with the light from any other point. Nonetheless, if by means of some suitable optical trick the light from each object point is split into two parts, then it is possible for each pair of waves of common origin to interfere and form a fringe pattern. Thus each object point may be encoded in a suitable pattern of fringes, and if the encoding is a unique one, with no two object points generating identical fringe patterns, then in principle an image of the object can be obtained.

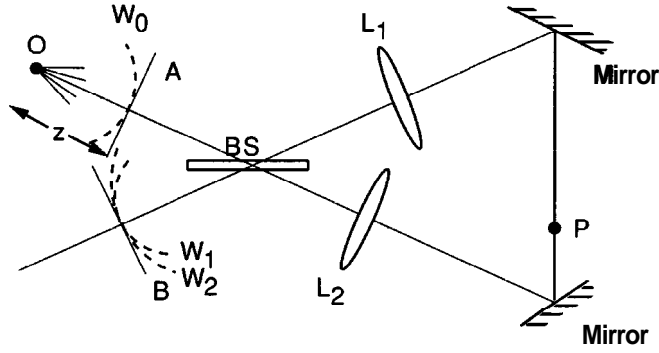


FIGURE 9.47
Triangular interferometer for
incoherent holography.

While many optical systems for achieving the required splitting of the object waves are known, we illustrate here with one particular system suggested by Cochran [68]. As shown in Fig. 9.47, the system consists of a triangular interferometer, in which are placed two lenses L_1 and L_2 with different focal lengths f_1 and f_2 . We assume that both lenses are positive, although a combination of one positive and one negative lens may also be used. The lenses are separated by a path length $f_1 + f_2$, their focal points coinciding at the point P in the figure. Plane A and plane B both lie at path length f_1 from lens L_1 and path length f_2 from L_2 .

Light may travel from plane A to plane B along either of two paths, one clockwise around the interferometer and the second counterclockwise. Considering first the clockwise path, light travels a distance f_1 from plane A to lens L_1 by means of a reflection at the beam splitter BS . From L_1 to L_2 the path length is $f_1 + f_2$, and from L_2 to plane B (again by means of reflection at BS) the path length is f_2 . Because of the particular choice of the path lengths in relation to the focal lengths f_1 and f_2 , plane A is *imaged* onto plane B ; due to the particular sequence in which L_1 and L_2 are encountered on this path, the imaging is performed with magnification $M_1 = -f_2/f_1$.

For the counterclockwise path, light is in each case transmitted (rather than reflected) by the beam splitter. Again plane A is imaged onto plane B , but for this path the lenses are encountered in opposite sequence, and the magnification is $M_2 = -f_1/f_2$.

Consider now the single point O (see Fig. 9.47) of an incoherent object located at distance z from plane A . Regarding the light from that one point as providing a phase reference, we may express the resulting spherical wave (wavefront W_0 in the figure) incident on plane A as the complex function

$$U_a(x, y) = U_o \exp \left[j \frac{\pi}{\lambda z} (x^2 + y^2) \right], \quad (9-120)$$

where a **paraxial** approximation has been used. At plane B we find two spherical waves (wavefronts W_1 and W_2 in the figure), one magnified by M_1 and the second by M_2 . Thus the total amplitude is

$$U_b(x, y) = U_1 \exp \left\{ j \frac{\pi}{\lambda z} \left[\left(\frac{x}{M_1} \right)^2 + \left(\frac{y}{M_1} \right)^2 \right] \right\} \\ + U_2 \exp \left\{ j \frac{\pi}{\lambda z} \left[\left(\frac{x}{M_2} \right)^2 + \left(\frac{y}{M_2} \right)^2 \right] \right\}. \quad (9-121)$$

The corresponding intensity distribution is

$$\mathcal{I}(x, y) = |U_1|^2 + |U_2|^2 + 2 \cos \left[\frac{\pi}{\lambda z} \frac{f_1^4 - f_2^4}{f_1^2 f_2^2} (x^2 + y^2) \right], \quad (9-122)$$

where we have used the relation

$$\frac{1}{M_1^2} - \frac{1}{M_2^2} = \frac{f_1^4 - f_2^4}{f_1^2 f_2^2}.$$

If a photographic plate is exposed by the intensity pattern of Eq. (9-122), and processed to produce a positive transparency with amplitude transmittance linearly proportional to exposure, the resulting transmittance may be written

$$t_A(x, y) = t_b + \beta' U_1 U_2^* \exp \left\{ j \left[\frac{\pi}{\lambda z} \left(\frac{f_1^4 - f_2^4}{f_1^2 f_2^2} \right) (x^2 + y^2) \right] \right\} \\ + \beta' U_1^* U_2 \exp \left\{ -j \left[\frac{\pi}{\lambda z} \left(\frac{f_1^4 - f_2^4}{f_1^2 f_2^2} \right) (x^2 + y^2) \right] \right\}. \quad (9-123)$$

We recognize the second and third terms as the transmittance functions of a negative and positive lens, respectively (cf. Eq. (5-10)), each of focal length

$$f = \frac{f_1^2 f_2^2}{f_1^4 - f_2^4} z. \quad (9-124)$$

Thus if the transparency is illuminated by a coherent source, both a virtual and a real image of the original object will be formed.

Generalizing now to an object consisting of a multitude of mutually incoherent point sources, each point source generates its own fringe pattern on the recording medium. Since the various sources are not coherent, the total intensity is found simply by adding the various intensity patterns so generated. The (x, y) coordinates of each point source determine the center of the corresponding pattern of fringes, and therefore fix the (x, y) coordinates of the real and virtual images. Similarly, the z coordinate of the point source influences the focal length of its contribution to the transmittance function, as seen in Eq. (9-124), and the image formed is thus a three-dimensional one.

Although the possibility of using incoherent illumination, rather than coherent illumination, is an attractive one in many applications, there exists one serious problem that limits the usefulness of incoherent holography. The problem arises because each elementary fringe pattern is formed by two extremely tiny portions of the light incident on the recording medium. Whereas for *coherent* holography light from each object point interferes with all the light contributed by the reference wave, for *incoherent* holography, the interfering waves represent only a minute fraction of the total light. The summation of many weak interference patterns, each with its own bias level of exposure, results in a very large bias level, in general much larger than for a hologram of a similar object formed with coherent light. As a consequence of this bias problem, incoherent holography has been successfully applied only to objects composed of small numbers of resolution elements. This limitation restricts its use significantly.

9.12 APPLICATIONS OF HOLOGRAPHY

Holography is a mature scientific field: most of the basic science has been done, and the techniques have undergone a great deal of refinement. During this process, a multitude of applications have been explored, some leading to highly successful businesses, others to important diagnostic tools that are widely used in some branches of both science and engineering. In this section we present a brief summary of the major applications to date.

9.12.1 Microscopy and High-Resolution Volume Imagery

From an historical perspective, microscopy has been the application of holography which has motivated much of the early work on wavefront reconstruction; it was certainly the chief motivating force behind the early works of **Gabor** [106], [107], [108] and **El-Sum** [93]. Interest in applications to electron microscopy has remained (cf. [285]), and interest in extending holographic microscopy to the X-ray region of the spectrum remains strong as well [207]. Interest in both electron and X-ray holography is motivated by the potential for achieving extremely high resolutions, comparable with the wavelength in each case.

In the visible region of the spectrum, holography is not a serious competitor with the conventional microscope in ordinary, run-of-the-mill microscopy. Nonetheless there does exist one area in which holography offers a unique potential to microscopy, namely in high-resolution volume imagery. In conventional microscopy, high lateral resolution is achieved only at the price of a limited depth of focus. As seen in Chapter 6, the best lateral resolution achievable by an imaging system is of the order of $\lambda/(NA)$ (cf. Eq. (6-46)), where (NA) is the numerical aperture. It can be shown that with this lateral resolution comes a depth of focus that is limited to an axial distance on the order of $\lambda/(NA)^2$. Note that for numerical apertures approaching unity, the depth of focus becomes as shallow as one wavelength! Thus there is a limited volume that can be brought into focus at one time.

It is, of course, possible to explore a large volume in sequence, by continuously refocusing to explore new regions of the object volume, but such an approach is often unsatisfactory if the object is a dynamic one, continuously in motion.

A solution to these problems can be obtained by recording a hologram of the object using a pulsed laser to obtain a very short exposure time. The dynamic object is then "frozen" in time, but the recording retains all the information necessary to explore the full object volume. If the hologram is illuminated, the real or virtual image can be explored in depth with an auxiliary optical system. Sequential observation of the image volume is now acceptable because the object (i.e. the holographic image) is no longer dynamic.

This approach was fruitfully applied by **C. Knox** in the microscopy of three-dimensional volumes of living biological specimens [167], and by **Thompson, Ward, and Zinky** in measurement of the particle-size distributions in aerosols [281]. The reader may consult these references for further details.

9.12.2 Interferometry

Some of the most important scientific applications of holography have proven to arise from the unique modalities of interferometry that it offers. Holographic interferometry can **take** many different forms, but all are dependent on the ability of a hologram to store two or more separate complex wave fields on the same recording medium, and the subsequent interference of those fields when they are reconstructed together. More detailed treatments of holographic interferometry can be found, for example, in the books by Vest [294] and Schumann [256].

Multiple-exposure holographic interferometry

The most powerful holographic interferometry techniques are based on a property, emphasized by **Gabor** et al. [110], that, by means of multiple exposures of holograms, coherent additions of complex wavefronts can be achieved. This property can easily be demonstrated as follows: let a holographic recording material be exposed sequentially by N different intensity distributions $\mathcal{I}_1, \mathcal{I}_2, \dots, \mathcal{I}_N$. The total exposure to which the medium has been subjected can be written

$$E = \sum_{k=1}^N T_k \mathcal{I}_k, \quad (9-125)$$

where T_1, T_2, \dots, T_N are the N individual exposure times. Now suppose that during each individual exposure interval the incident radiation is the sum of a reference wave $A(x, y)$ (the same for all exposures) and an object wave $a_k(x, y)$ which changes from exposure interval to exposure interval. The total exposure becomes

$$E = \sum_{k=1}^N T_k |A|^2 + \sum_{k=1}^N T_k |a_k|^2 + \sum_{k=1}^N T_k A^* a_k + \sum_{k=1}^N T_k A a_k^*. \quad (9-126)$$

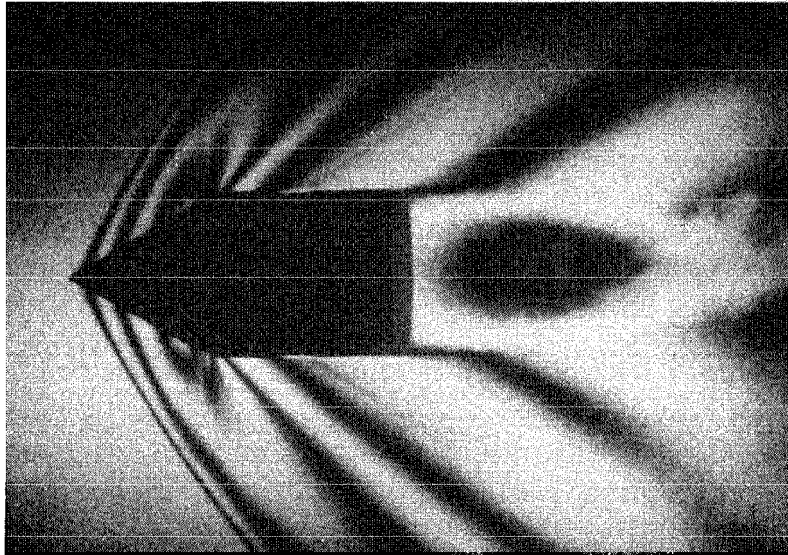
Assuming linear operation in the t_A vs. E characteristic of the recording medium, we find components of transmittance

$$t_\alpha = \beta \sum_{k=1}^N T_k A^* a_k \quad (9-127)$$

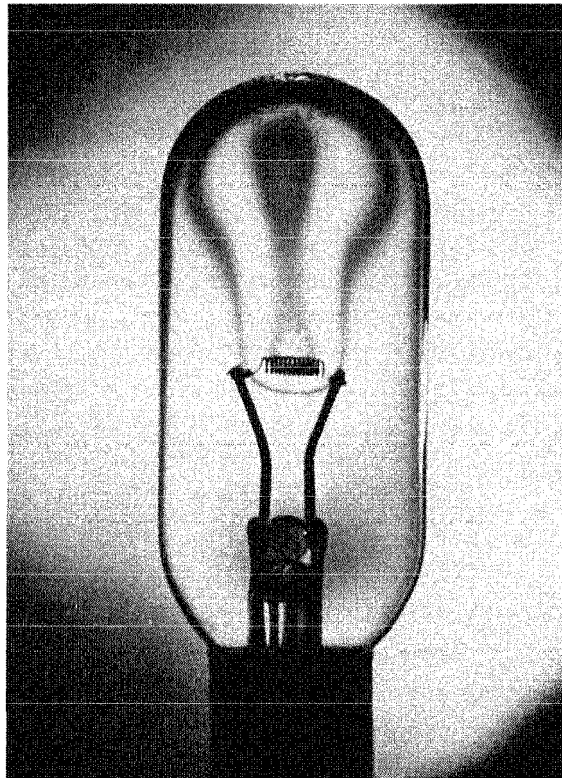
$$t_\beta = \beta \sum_{k=1}^N T_k A a_k^*.$$

From Eq. (9-127) it is clear that illumination of the processed hologram with a wavefront A will generate a transmitted field component proportional to the product of $|A|^2$ and the sum of the complex wavefronts a_1, a_2, \dots, a_N . As a consequence, N coherent virtual images of the objects that gave rise to the N wavefronts will be linearly superimposed and will mutually interfere. In a similar fashion, illumination of the transparency by a wavefront A^* will generate N coherent real images which likewise interfere.

The earliest dramatic demonstrations of the potential of this type of interferometry were performed by Brooks et al. [40] using a Q-switched ruby laser. Figure 9.48 shows two photographs obtained in each case by double exposure of a hologram with two laser pulses. In the case of part (a) of the figure, the first pulse records a hologram of only a



(a)



(b)

FIGURE 9.48
Double-exposure holographic interferometry with a Q-switched
ruby laser. [By permission of R.E. Brooks, L.O. Heflinger, and R.F.
Wuerker.]

diffuse background, while the second pulse records a hologram of a bullet in flight in front of the same diffuse background. The shock waves generated by the bullet produce changes in the local refractive index of the air. As a consequence, the two images of the diffuse background, one recorded in the absence of the bullet and the other recorded through the refractive-index perturbations of the air, will mutually interfere, producing

interference fringes that outline the shock waves generated by the bullet. These fringes have the appearance of being fixed in three-dimensional space around the bullet.

Part (b) of the same figure is a similarly obtained image of an incandescent bulb. During the first exposure the filament is off, and again a hologram of a diffuse background is recorded, this time through the imperfect glass envelope of the bulb. The filament is then turned on, and a second laser pulse exposes the hologram. The incoherent light generated by the lamp does not interfere with the laser light, so the filament does not appear lighted in the final image. However, the heating of the gases within the envelope has resulted in changes of the local index of refraction, which again generate fringes of interference in the final image, outlining the patterns of gas expansion. It should be emphasized that these interference fringes have been obtained in the presence of the optically imperfect glass envelope, a feat which would be impossible by other classical methods of interferometry.

Real-time holographic interferometry

Another important type of holographic interferometry depends on interference between a prerecorded, holographically produced wavefront and the coherent waves reflected from or transmitted by the same object in real time [39]. The holographically produced wavefront can be regarded as a reference, representing the reflected or transmitted light when the object is in a "relaxed" state. If the same object, located in the same position relative to the hologram that it occupied when the reference wavefront was recorded, is now perturbed, perhaps by placing it under stress with some form of loading, then the complex fields intercepted from the object change, and upon interference with the reference wavefront, produce fringes that are visible on the image of the object seen through the hologram. Two slightly different coherent images are being superimposed by this process, one the image of the object in its original state, and the second the image of the object in its stressed or modified state. The fringes observed can reveal quantitative information about the nature of the object deformations that have taken place.

Note that we are causing the coherent image of the object now to interfere with the coherent image of the object that existed sometime in *the past* (or perhaps, using a computer-generated hologram, with an object that never actually existed previously), a feat that would be impossible to accomplish with conventional interferometry.

Contour generation

The interference of multiple coherent images described previously has also led to the development of techniques for obtaining three-dimensional images with superimposed constant-range contours. These techniques are applicable to the problems of cross-section tracing and contour mapping. Two distinctly different techniques have been demonstrated by Hildebrand and Haines [144].

In the first of these techniques, the object is illuminated by two mutually coherent but spatially separated point sources. The two object illuminations may be applied simultaneously or the hologram may be double-exposed, with a different position of the object illumination source during each exposure. If the pattern of interference between the two object illumination sources is considered, it is found to consist of interference fringes that follow hyperbolas of constant path-length difference, as shown in Fig. 9.49. If the object is illuminated from the side and the hologram is recorded from above, then

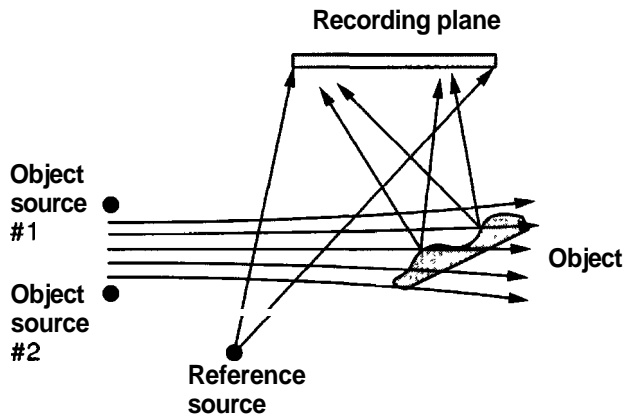


FIGURE 9.49
Contour generation by the two-source method.

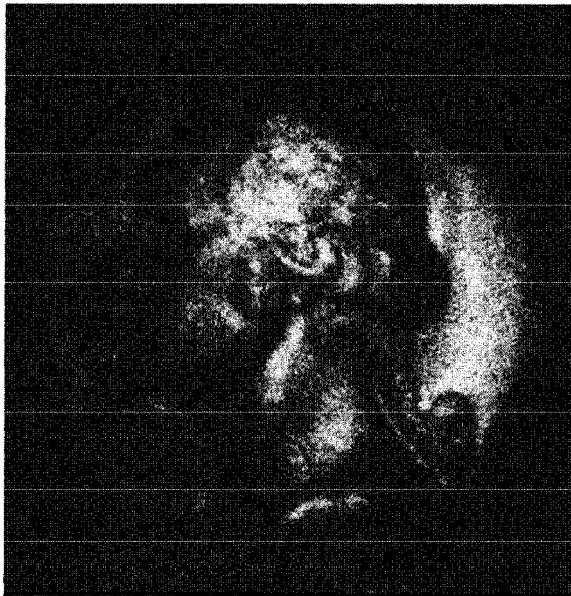
depth contours (i.e. the intersection of the object with the hyperbolic fringes) are readily seen on the reconstructed image. Identical results are obtained whether the two object illumination sources were used simultaneously in a single exposure or separately in individual exposures, for in either case the two illuminations add coherently.

The two-source method of contour generation suffers from the requirement that the directions of illumination and observation must differ by nearly 90° . Thus if the object has significant relieving, shadows will be cast and parts of the object will simply not be illuminated. This deficiency is overcome by the two-frequency or two-wavelength method of contour generation. In this case the object and reference illuminations both contain the same two distinct wavelengths, say λ_1 and λ_2 . In effect, each wavelength records a separate and independent hologram on the same recording medium. When the resulting hologram is illuminated by light of a single wavelength, two images with slightly different positions and magnifications are produced. These two images will interfere, and for certain geometries the resulting image contours will be accurate indications of depth. We do not dwell on a detailed analysis of this case; the interested reader may consult the original reference for further details [144]. Figure 9.50 shows the results of contour mapping by the two-wavelength method. In part (a) we see a holographic image of a coin, illuminated in the usual manner with single-wavelength light. When two-wavelength light is used to record the hologram, the image of part (b) is obtained. In this case the two wavelengths were obtained from two different lines of an argon laser. The two lines are separated by 650 nm, and the resulting contours on the image are spaced by $20 \mu\text{m}$.

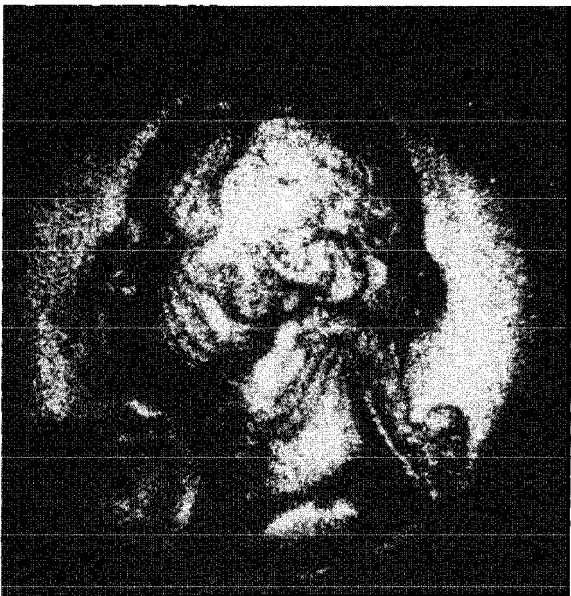
Vibration analysis

A holographic technique for vibration analysis, first proposed by Powell and Stetson [234], may be regarded as a generalization of multiple-exposure holographic interferometry to the case of a continuous time exposure of a vibrating object.

With reference to the geometry of Fig. 9.51, we consider a point at coordinates (x_o, y_o) on a planar object which is vibrating sinusoidally with angular frequency Ω . The peak amplitude of the vibration at that point is represented by $m(x_o, y_o)$, and the fixed phase of the vibration is $\mu(x_o, y_o)$. The light incident at the hologram recording plane coordinates (x, y) from that particular point may be regarded as having a time-varying phase modulation



(a)



(b)

FIGURE 9.50
Contour generation by the two-wavelength method. [By permission of B.P. Hildebrand and K.A. Haines.]

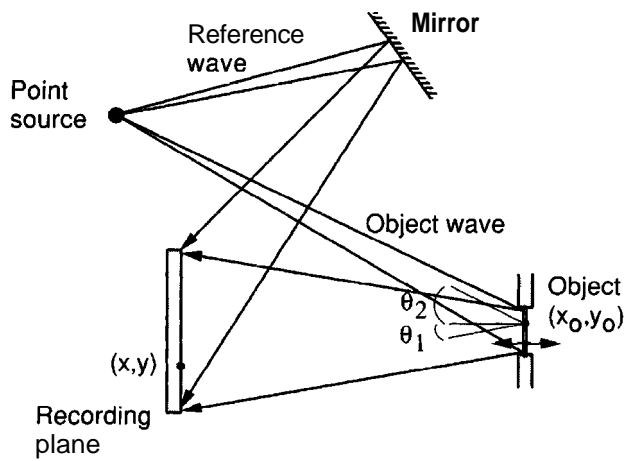


FIGURE 9.51
Recording a hologram of a vibrating object.

$$\phi(x, y; t) = \frac{2\pi}{\lambda} (\cos \theta_1 + \cos \theta_2) m(x_o, y_o) \cos[\Omega t + \mu(x_o, y_o)], \quad (9-128)$$

where λ is the optical wavelength of the illuminating source, θ_1 is the angle between the vector displacement of the object at (x_o, y_o) and the line joining that point to (x, y) , and θ_2 is the angle between the vector displacement and the direction of propagation of the incident light at (x_o, y_o) .

Using what by now should be a familiar expansion into Bessel functions, the temporal spectrum of the time-varying phaser representing the modulated light incident at (x, y) can be written

$$\begin{aligned} F(\nu) &= \mathcal{F}\{\exp[-j\phi(x, y; t)]\} \\ &= \sum_{k=-\infty}^{\infty} J_k \left[2\pi \frac{\cos \theta_1 + \cos \theta_2}{\lambda} m(x_o, y_o) \right] \delta\left(\nu - \frac{k\Omega}{2\pi}\right). \end{aligned} \quad (9-129)$$

When the exposure time is much longer than the vibration period (that is, when $T \gg 2\pi/\Omega$), only the $k = 0$ term, which is at the same optical frequency as the reference wave, will cause stable interference fringes to be formed. All other terms will fail to produce such fringes. If the variations of the modulation depth introduced by the term $\cos \theta_1$ are nearly independent of (x, y) (that is, if the angle subtended by the film at (x_o, y_o) is small), then the amplitude of the image at (x, y) will be suppressed by the factor

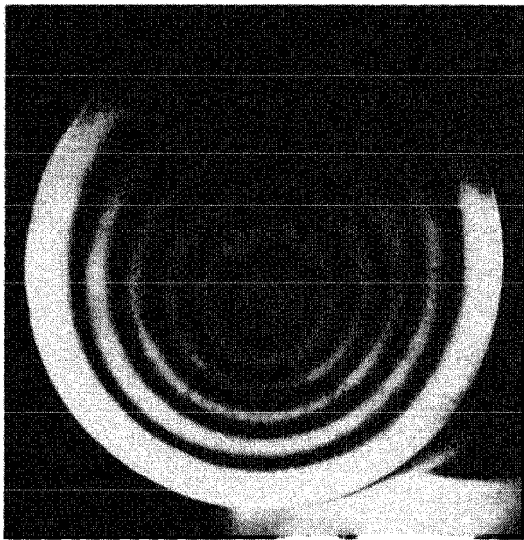
$$J_0 \left[\frac{2\pi}{\lambda} (\cos \theta_1 + \cos \theta_2) m(x_o, y_o) \right], \quad (9-130)$$

and the intensity will be suppressed by the square of this factor. Thus the intensity of the image depends at each point on the depth of vibration of the corresponding object point.

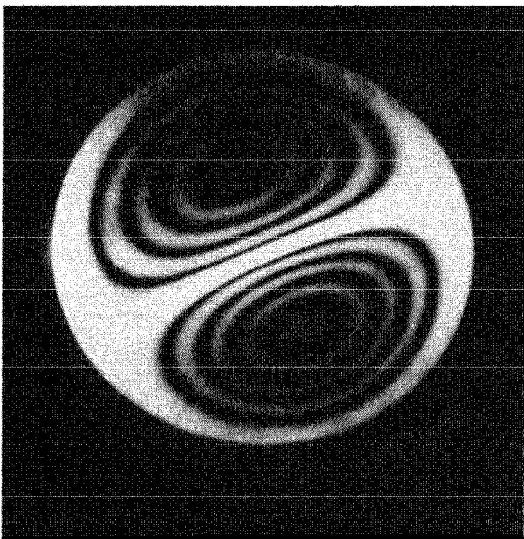
Figure 9.52 shows images of a vibrating diaphragm obtained experimentally by Powell and Stetson. In part (a) of the figure, the diaphragm is vibrating in its **lowest-order mode**, with a single vibration maximum at the center of the diaphragm. In part (b), the diaphragm is vibrating in a higher-order mode, with two vibration maxima. By counting the number of fringes from the center of the diaphragm to any point in question, it is possible, with the help of Eq. (9-130), to determine the vibration amplitude at that point.

9.12.3 Imaging Through Distorting Media

In many cases of practical interest, an optical system may be required to form images in the presence of uncontrollable aberrations. These aberrations may result from imperfections of the image-forming components themselves, or they may be introduced by an external medium, such as the Earth's atmosphere. The techniques of holography offer several unique advantages for problems of imaging in the presence of such aberrations. We discuss here three distinctly different holographic techniques for obtaining high resolution in the presence of severe aberrations.



(a)



(b)

FIGURE 9.52
Holographic images of a diaphragm vibrating in two different modes. [By permission of R.L. Powell and K.A. Stetson.]

The first technique ([191], [168]) of interest is applicable only when the distorting medium is constant in time. As illustrated in Fig. 9.53, a hologram of the distorted object waves is recorded with an undistorted reference wave. The processed hologram is then illuminated with an "anti-reference" wave, i.e. a reconstruction wave that duplicates the reference wave but propagates in the reverse direction. A real, conjugate image of the distorting medium will form precisely at the location of the medium itself, between the hologram and the image plane. If the object wave incident on the distorting medium during the recording process is represented by $U_o(\xi, \eta)$ and if the amplitude of transmittance of the distorting medium is $\exp[jW(\xi, \eta)]$, then the wave falling on the distorting medium during reconstruction is $U_o^*(\xi, \eta) \exp[-jW(\xi, \eta)]$. Note that when this conjugate wave passes back through the identically same distorting medium that was originally present, the aberrations entirely cancel,

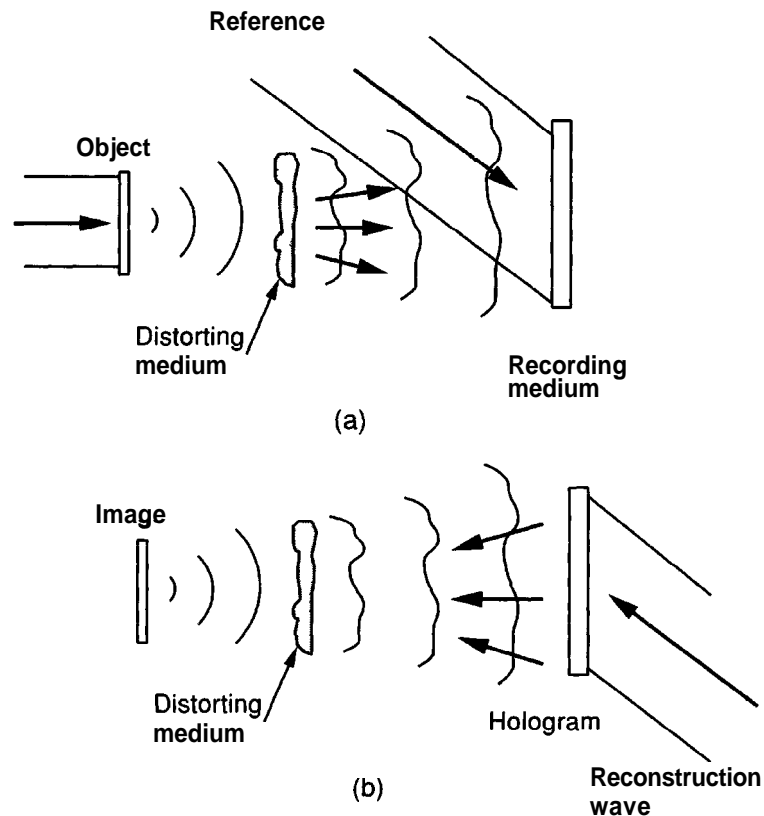


FIGURE 9.53

Use of the original distorting medium for compensating aberrations. (a) Recording the hologram; (b) reconstructing the image.

$$U_o^*(\xi, \eta) \exp[-jW(\xi, \eta)] \exp[jW(\xi, \eta)] = U_o^*(\xi, \eta),$$

leaving a wave $U_o^*(\xi, \eta)$ to propagate on to the image plane, where an aberration-free image appears.

A limitation of the technique in some applications is that the image must appear where the object originally was situated, whereas in practice it is often desired to obtain an image on the other side of the distorting medium (i.e. to the right of the distorting medium in Fig. 9.53). If the distorting medium is movable, then this difficulty can be overcome.

A second technique of interest is illustrated in Fig. 9.54. Again the distorting medium should be unchanging in time. In this case we record a hologram of the distorted waves transmitted by the medium when it is illuminated by a simple point source (i.e. a record of the point response of the medium). This hologram may now be used as a "compensating plate" to enable a more conventional optical system to form an aberration-free image. Let the waves incident on the recording medium due to the point source be represented by $\exp[jW(x, y)]$. We have assumed that the distorting medium is such that only phase distortions are introduced. The portion of the hologram amplitude transmittance that normally contributes the real image is proportional to $\exp[-jW(x, y)]$. Thus if we replace the point source by a more general object, and

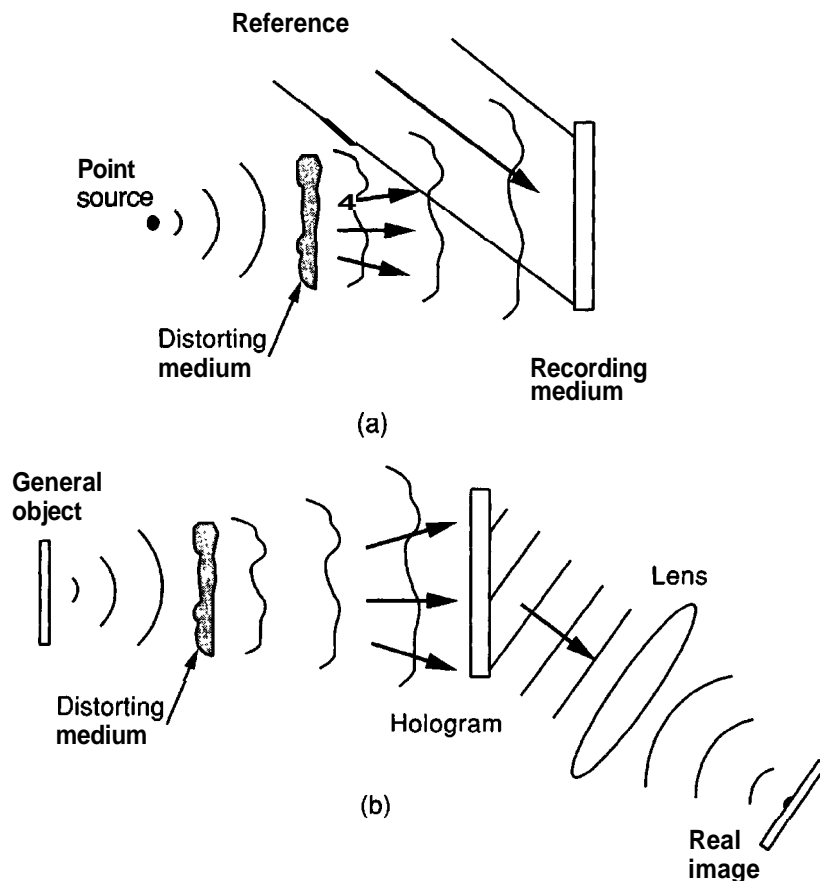


FIGURE 9.54 Use of a hologram compensating plate. (a) Recording the compensating plate; (b) cancellation of the aberrations.

reinsert the hologram in the same position it originally occupied, we find that the curvatures of the object waves reaching the hologram are canceled on passage through the hologram, with the waves from different object points producing plane waves traveling at different angles. The lens then forms a distortion-free image in its focal plane.

This technique will work well over only a restricted field of view, for if an object point is too far from the position of the original point source used in recording the hologram, the aberrations imparted to its wave may differ from those recorded on the hologram. This restriction is less severe if the hologram is recorded very close to the distorting medium. Upatnieks et al. [287] have successfully applied this technique to the compensation of lens aberrations, an application to which it is well suited.

A third technique, which may be applied to imaging through media that are **time-varying** or time-invariant, is accomplished by passing **both** the reference wave and the object wave through the same distorting medium [125]. As indicated in Fig. 9.55 the **lensless** Fourier transform recording geometry is most commonly used, with a reference point source existing in the same plane or nearly the same plane as the object of interest. For simplicity it is assumed that the distorting medium is located immediately in front of the recording plane, although this restriction can be relaxed with some loss of the field of view over which compensation is effective. The reference and

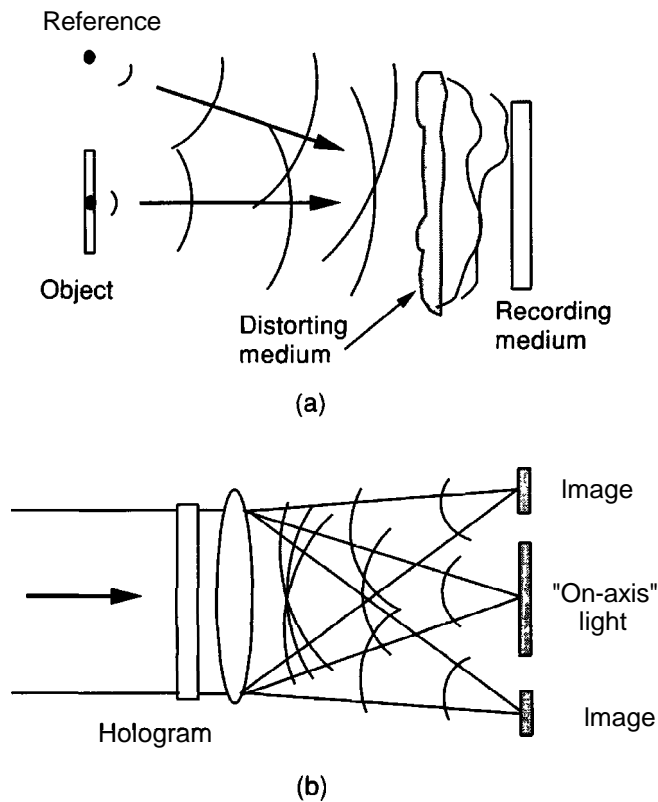


FIGURE 9.55
Aberration-free imaging when the object and reference waves are identically distorted. (a) Recording the hologram; (b) obtaining the image.

object waves reaching the recording medium can be written as $A(x, y) \exp[jW(x, y)]$ and $a(x, y) \exp[jW(x, y)]$, where A and a are the waves that would have been present in the absence of a distorting medium. Interference of the two distorted waves yields a pattern of intensity that is unaffected by the presence of the distorting medium,

$$\begin{aligned} \mathcal{I}(x, y) &= |A(x, y) \exp[jW(x, y)] + a(x, y) \exp[jW(x, y)]|^2 \\ &= |A|^2 + |a|^2 + A^*a + Aa^*, \end{aligned}$$

and distortion-free twin images can be obtained from the hologram.

Again the technique will work over only a limited object field, since points too far from the reference may produce waves with aberrations that differ significantly from those of the reference wave. The working field is largest when the aberrations are introduced close to the recording plane.

This method is an example of a more general set of techniques in optics known as "common path interferometry". It has been applied to the problem of obtaining high-resolution images of distant objects through the Earth's atmosphere [112], [113], [126].

9.12.4 Holographic Data Storage

There are many attractive properties of holography as a potential data storage technique, and as a consequence, much attention has been given over the years to this application. Most obvious, perhaps, is the highly diffused nature of holographic storage, in the sense that a single pixel of an analog image or a single bit in a binary data array is stored in

a distributed fashion over a considerable area of the hologram. Nonlocalization is most complete for a Fourier transform hologram, and least complete for an image hologram, with Fresnel holograms falling in between these two extremes. When there is a large amount of nonlocalization, a dust speck or a defect in the recording medium that obscures or destroys a very localized area on the hologram will not create a localized defect in the image, and therefore there will not be a localized loss of stored data.

A second advantage, associated particularly with the Fourier transform recording geometry, arises from the fact that a shift in the hologram domain results in only a linear phase tilt in the Fourier domain, and therefore has no effect on the location of the image intensity distribution. As a consequence, Fourier holograms are extremely tolerant to misalignment or registration errors. This property is extremely important for high-density memories, especially those that have high magnification in the sense that a small hologram produces a much larger image.

A third attraction of holography as a storage method comes from our ability to use the third dimension of a three-dimensional recording material, such as a thick recording film or a photorefractive crystal, for recording. Thus holography offers one method of three-dimensional optical storage, and by utilizing the third dimension the volume storage density that can be achieved is quite high.

Early work on holographic storage concentrated on thin holograms and storage of two-dimensional binary arrays [5]. Figure 9.56 shows a typical arrangement. Separate two-dimensional pages of binary data are stored in a two-dimensional array of holograms. The light from a CW laser is deflected by a pair of acousto-optic beam deflectors to a particular hologram in the array. The particular hologram selected generates an array of binary spots on a two-dimensional detector array. Thus to determine the state of one particular binary element in the memory, a combination of the right hologram and the right detector element must be interrogated.

More recent emphasis has been on three-dimensional storage media, such as photorefractive crystals (see, for example, [141]), which are capable of providing Bragg selectivity. Multiplexing of holograms within the crystal and selective recall of the data recorded in those holograms can be achieved by means of angle multiplexing, wavelength multiplexing, or multiplexing with phase-coded reference beams. A typical geometry for such a system (from [141]) is shown in Fig. 9.57. A spatial light modulator serves to generate an array of binary data, and the reference wave is introduced at a particular angle designated for that page of data. The reference beams are introduced from the side of the crystal, an orientation that maximizes angular selectivity. The hologram is recorded, and the data can be read out onto a CCD detector array by illumination of

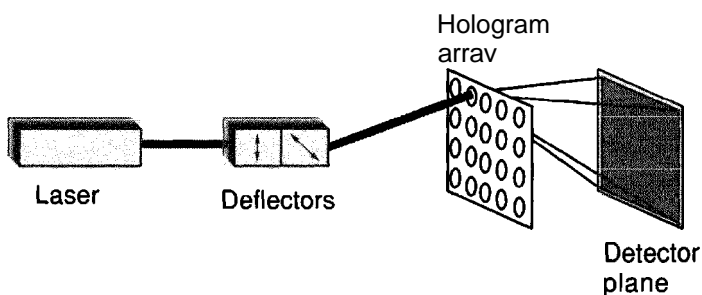


FIGURE 9.56
Page-oriented holographic storage.

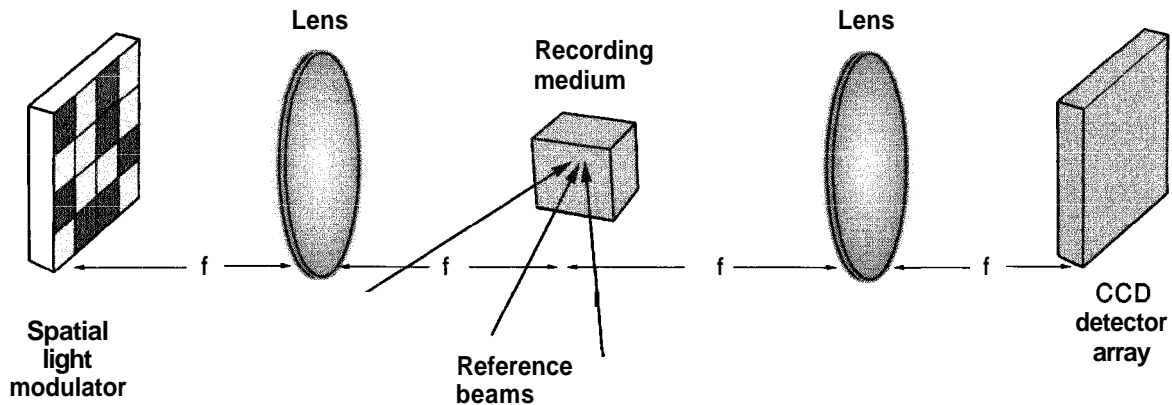


FIGURE 9.57

A volume holographic storage system. The case of angle multiplexing is illustrated.

the crystal with a duplicate of the reference beam. Other holograms are superimposed in the crystal by use of other reference angles, and must be read out with duplicates of those reference beams. The diffraction efficiency associated with a single bit falls as $1/N^2$ when N holograms are superimposed, due to the partial erasure of early holograms caused by the recording of later holograms. Superposition of several thousand holograms by angular multiplexing has been demonstrated experimentally [45].

Finally, mention should be made of the use of holography for associative memories, an idea first described by Gabor [109]. Discussion of related ideas can be found in [115], [164], [165], and [166].

9.12.5 Holographic Weights for Artificial Neural Networks

Neural network models provide an interesting and powerful approach to many pattern recognition and associative memory problems. One approach to constructing an artificial "neural" processor is through the use of volume holography. In this section we provide the briefest introduction to this subject, together with references that will allow the reader to pursue the ideas further. The terminology used to describe networks of this type is borrowed from the neurological sciences, but it is important to understand that the models used in artificial neural networks contain only the simplest extraction of the essence of the types of processing that are believed to take place in real biological neural systems. An introduction to the field of neural computing can be found in, for example, [142].

Model of a neuron

Neural networks consist of a multitude of nonlinear elements referred to as *neurons*, highly interconnected with each other. A simple model of a neuron is illustrated in Fig. 9.58(a). The summation of a multitude of different *weighted* binary inputs is applied to the input of a nonlinear element, usually taken to have a "sigmoid" nonlinear characteristic described by the input-output relation

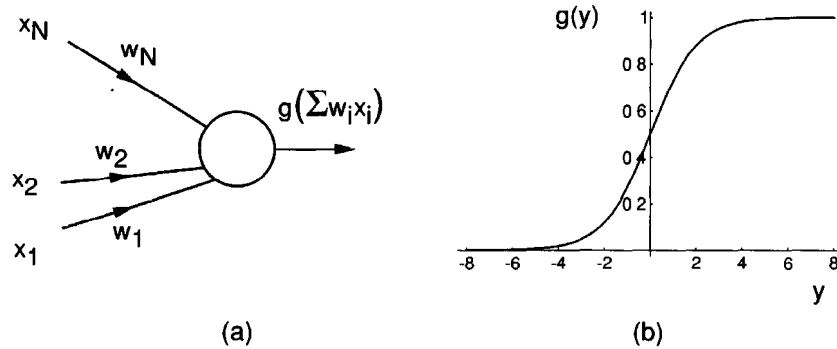


FIGURE 9.58
(a) Model of a single neuron; (b) sigmoidal nonlinearity.

$$z = g(y) = \frac{1}{1 + e^{-y}}, \quad (9-131)$$

which is illustrated in Fig. 9.58(b).

The input y to the nonlinearity is the sum of N weighted inputs x_i , as described by the relation

$$y = \sum_{i=1}^N w_i x_i = \vec{w} \cdot \vec{x}, \quad (9-132)$$

where the w_i are the weights applied to those inputs.

A single neuron can be trained to produce either a 1 or a 0 in response to a particular input vector \vec{x} by adjusting the weight vector so that it is either co-directional with or orthogonal to that input vector, respectively. In the former case a large positive input to the sigmoid nonlinearity drives the output to a result very close to unity, and in the latter case a large negative input drives the output result very close to zero. In this way, by adjusting the weights, it is possible to "train" the neuron to recognize a particular input vector. Extending this idea, one finds if the neuron is to be presented with an entire set of vectors, each of which is to be classified into one of two possible sets, by training the neuron with examples of the two classes, it can be taught to separate the input vectors by means of a simple hyperplane in the N -dimensional space of the vectors. Classes of vectors that are separable with a hyperplane will then be distinguished by the neuron, while those that are not separable with a hyperplane cannot be distinguished.

Networks of neurons

To obtain functionality that is more complex than that possible with a single neuron, collections of such elements are joined together to form a *neural network*. An example of such a network having four layers of interconnected neurons is shown in Fig. 9.59. The layer of neurons on the far left can be thought of as containing input neurons. In an image-recognition problem, for example, each such neuron might receive a single input representing the value of one pixel of an image that is to be classified by the network. The layer on the right can be thought of as containing output neurons. Each such neuron represents one of the possible classes into which the image is to be classified. Ideally, when a single image is presented at the input of the network, with appropriate training,

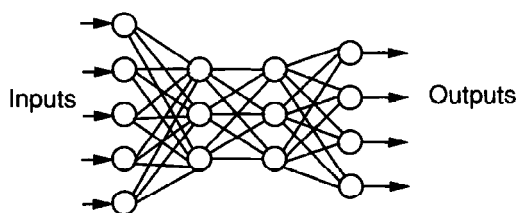


FIGURE 9.59
A four-layer neural network.

the network will cause a single output neuron to produce a value at or near unity and the rest to produce values at or near zero. The particular neuron that has unity output indicates the class of which the input image is a member. The middle layers of neurons are referred to as "hidden" layers. The number of hidden layers determines the complexity of the dividing surfaces in N -dimensional space that can separate input images into classes.

The neural network must be trained by presenting it with samples from the various classes of input images, and adjusting all of the weights according to some predetermined algorithm. A variety of training algorithms exist, all of which involve the minimization of an error metric. We mention in particular the LMS algorithm [299] for single-layer networks and the backpropagation algorithm [251] for multilayer networks, but must refer the reader to the references for details.

Optical neural networks based on volume holographic weights

One popular implementation of neural networks using optics is based upon storage of weights in an erasable, thick holographic medium. Photorefractive crystals are most commonly used. Figure 9.60 illustrates one manner in which a hologram can introduce a weighted interconnection. We assume that the input to the neural network consists of a spatial light modulator that generates a coherent amplitude distribution proportional to the input to be processed. The lens L_1 Fourier transforms the input, and thus each pixel in the spatial light modulator generates a plane wave with a unique k vector at the crystal. We assume that a collection of sinusoidal volume gratings has been written into the crystal; the exposure times or the strengths of the waves used in recording these

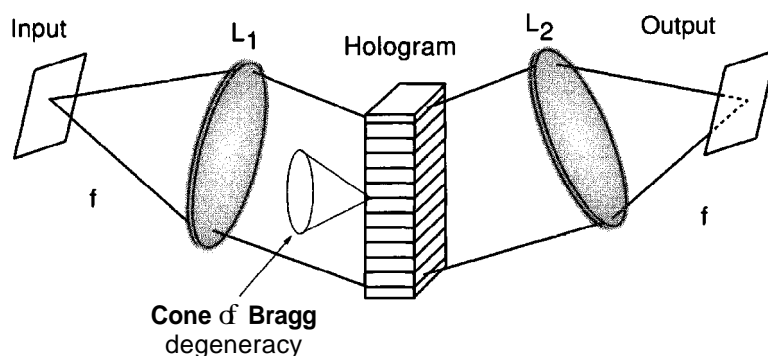


FIGURE 9.60
Illustration of a single weighted interconnection using a hologram. In practice, many such interconnections would be realized simultaneously.

gratings determine their diffraction efficiencies, and therefore control the weights that they will apply to incident Bragg-aligned plane waves. By means of a second Fourier transforming lens, all plane waves traveling in a common direction are summed onto a single output pixel. There are many output pixels, each corresponding to a different direction of light diffracted from the crystal.

Thus a multitude of volume gratings are written into the crystal, each grating representing one weight. A weighted sum of input pixel values then results at each output pixel. A training procedure can be implemented under computer control that changes the strengths of the volume gratings in accord with a chosen training algorithm.

The attraction of optics, and in particular volume holography, in this application comes from the very large number of gratings (or weights) that can be superimposed in a single crystal, coupled with the fact that large numbers of pixels (or neurons) can be realized with SLM technology. The thickness of the recording material is important if a multitude of different gratings are to be angularly multiplexed (using the Bragg effect) in the medium. The goal of achieving large numbers of weights is hindered by two phenomena. One, known as Bragg degeneracy, refers to the fact that the Bragg condition can be satisfied by an entire cone of angles, rather than just a single angle, and therefore there is the potential for significant crosstalk to exist between weighted interconnections. This problem can be combatted by utilizing only a properly chosen subset of the possible gratings, such that light can pass from one input pixel to one output pixel by means of one and only one volume grating [238]. A second solution is to break the Bragg degeneracy by forcing a single path from an input pixel to an output pixel to diffract from more than one volume grating [224].

A second limitation arises from the fact that, for photorefractive crystals subjected to a sequence of exposures, the later exposures partially erase the early exposures. This limits the total number of gratings that can be superimposed; however, experiments have demonstrated storage of several thousands of exposures [215]. Actually, the tendency of the photorefractive medium to "forget" early exposures can be used to advantage in some learning processes.

We have only touched the subject of optical neural networks in the above discussion. Other approaches to utilizing optics for neural-like computation exist. We mention in particular the realization of **Hopfield** neural networks using the matrix-vector architecture of Section 8.11.3 [237] and the use of competitive and cooperative phenomena in systems utilizing nonlinear optical elements [4]. For additional references, see the March 1993 issue of *Applied Optics*, which was devoted to the subject of optical neural networks.

9.12.6 Other Applications

Many other applications of holography exist, but space limitations prevent us from reviewing them all here. In this section we give brief mention of several areas that are particularly important, and present some references for further study.

Holographic optical elements

Holography has found considerable application to the construction of waveshaping optical elements, which are referred to as holographic optical elements (HOEs). Such

elements are one further example of diffractive optical elements. Again it is the light weight and compact volume associated with such elements that make them particularly attractive. Holographic optical elements have been used, for example, for optical scanning [18], [181], for heads-up displays in aircraft cockpits [66], and in many other applications.

The reader is referred to Section 7.3 for further discussion of diffractive optical elements. References [59], [60], [61], [62], and [63] all contain examples of applications.

Holographic display and holographic art

The striking character of three-dimensional holographic images has been the factor most responsible for interest in holography on the part of the nontechnical public. Holography has been applied to advertising, and a multitude of artists have taken up holography as a medium of choice. A Museum of Holography was established in New York City and recently moved to the Massachusetts Institute of Technology. Holographic jewelry can be found in many shops around the world.

Holograms for security applications

The application that brings holography into direct contact with the largest number of people is the use of holograms for prevention of counterfeiting and fraud. The ubiquitous embossed hologram on the credit card is the most common example in the United States, although in Europe its use has been extended even further. Holography is used in such applications to provide a deterrent to counterfeiting, since the presence of a hologram as an integral part of a credit card or a bank note makes the unauthorized duplication of that item considerably more difficult than would otherwise be the case.

To gain a better appreciation for the variety of applications of holography in the security field, the reader may wish to consult [96], which contains many papers on the subject.

PROBLEMS—CHAPTER 9

- 9-1. A hologram is recorded using a spherical reference wave that is diverging from the point (x_r, y_r, z_r) , and the images from that hologram are played back with a reconstruction beam that is diverging from the point (x, y, z_p) . The wavelength used for both recording and reconstruction is λ_1 . The hologram is taken to be circular, with diameter D . See Fig. P9.1(a) below. It is claimed that the image of an arbitrary three-dimensional object obtained by this method is entirely equivalent to that obtained by a lens of the same diameter and at the same distance from the object, as shown in part (b) of the figure, and a prism (for simplicity, not shown), where again the wavelength is λ_1 . What are the two possible focal lengths for the lens that will produce equivalence?
- 9-2. A hologram is recorded with light from an argon laser at 488 nm wavelength, and the images are reconstructed with light from a HeNe laser with wavelength 632.8 nm. There is no scaling of the hologram.
- (a) Assuming $z_p = \infty$, $z_r = \infty$, and $z_o = -10$ cm, what are the axial distances z_i of the twin images? What are the transverse and axial magnifications of the images?

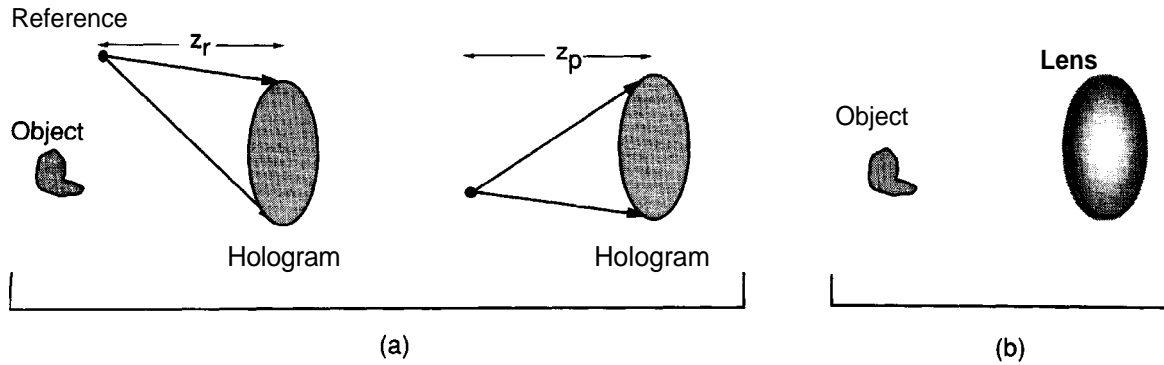


FIGURE P9.1

(b) Assuming $z_p = \infty$, $z_r = 2z_o$, and $z_o = -10$ cm, what are the axial distances and the transverse and axial magnifications of the twin images?

9-3. A hologram is recorded, and its images reconstructed with the same wavelength λ . Assuming $z_o < 0$, show that when $z_p = z_r$ there results a virtual image with unity transverse magnification, whereas with $z_p = -z_r$ there results a real image with unity transverse magnification. What is the transverse magnification of the twin image in each case?

9-4. The **lensless** Fourier transform geometry (see Fig. 9.14) is used to record a hologram of an object consisting of a square transparency of width L . The amplitude transmittance of the object is $t_A(x_o, y_o)$, and the distance of the object from the recording plane is $|z|$. The reconstruction wavelength is the same as the recording wavelength. The images are obtained by illuminating the hologram with a plane wave, followed by a positive lens of focal length f . For simplicity, both the object illumination and the reconstruction wave may be taken to have amplitude unity.

(a) What is the transverse magnification M_t of the first-order images?

(b) Show that the amplitude of the zero-order (i.e. on-axis) image term can be expressed as

$$U_f(u, v) = \frac{1}{\lambda f} \iint_{-\infty}^{\infty} U'_o(x_o, y_o) U_o^* \left(x_o + \frac{u}{M_t}, y_o + \frac{v}{M_t} \right) dx_o dy_o$$

(plus a central diffraction-limited spot), where

$$U'_o(x_o, y_o) = t_A(x_o, y_o) e^{j \frac{\pi}{\lambda |z|} (x_o^2 + y_o^2)}.$$

(c) How far from the center of the object transparency should the reference point source be placed in order to assure no overlap of the zero-order light with the first-order images?

9-5. We wish to make a holographic display that will project a real image of a planar transparency object. The recording and reconstruction geometries are shown in Fig. P9.5. The reference and object are constrained to lie to the left of the hologram during recording. The reconstruction source must lie to the left of the hologram, and the projected image must lie to the right of the hologram. The hologram is not turned around or changed in size between recording and reconstruction. The recording wavelength is 632.8 nm, the reconstruction

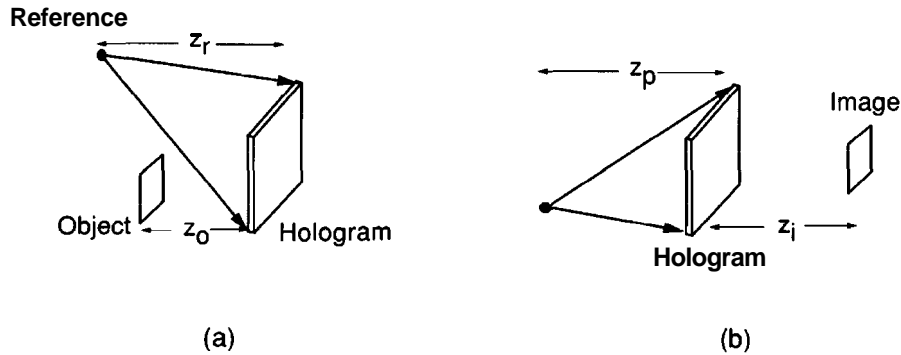


FIGURE P9.5

wavelength is 488 nm, the object transparency size is 2×2 cm, the desired image size is 4×4 cm, the axial distance from the hologram to the image must be 1 m, and the axial distance of the reconstruction source to the hologram is constrained to be 0.5 m.

- (a) Subject to the above constraints, specify all possible axial object and reference distances z_o and z_r that will together yield the desired image.
- (b) Repeat part (a), but with the hologram rotated 180° left to right (i.e. back and front interchanged) between the recording and reconstruction steps.
- 9-6.** It is proposed to record an X-ray hologram using coherent radiation of wavelength 0.1 nm and to reconstruct images optically using light of wavelength 600 nm. The object is a square transparency with a pattern of absorption at the X-ray wavelength. The lensless Fourier transform recording geometry is chosen. The width of the object is $100 \mu\text{m}$, and the minimum distance between the object and the reference is to be $200 \mu\text{m}$ to assure that the twin images will be separated from the "on-axis" interference. The X-ray film is placed 2 cm from the object.
- (a) What is the maximum spatial frequency (cycles/mm) in the interference pattern falling on the film?
- (b) Assume that the film has sufficient resolution to record all of the incident intensity variations. It is proposed to reconstruct the images in the usual manner, i.e. by looking in the rear focal plane of a Fourier transforming lens. Why will this experiment fail?
- 9-7.** A thick unslanted transmission phase grating is to be produced by bleaching the recording that results from interfering two plane waves in a photographic emulsion. The wavelength of the exposing radiation in air is 488 nm and the angle between the two interfering beams, also in air, is 60° . The thickness of the emulsion is $15 \mu\text{m}$. The average refractive index of the emulsion, both before exposure and after bleaching, is 1.52 (the index of gelatin). The same wavelength is used for reconstruction as for recording.
- (a) What are the wavelength and the angle between the two beams *inside* the emulsion during recording? How does the period of the grating predicted by the angle and wavelength outside the emulsion compare with the period predicted using the wavelength and angle inside the emulsion?

- (b) Assuming Bragg matched conditions, what peak refractive index modulation n_1 is required in order to reach the first 100% peak of the diffraction efficiency curve for a thick transmission phase grating?
- (c) Assuming operation at this same first maximum of diffraction efficiency, and assuming no error $\Delta\theta$ in the illumination angle, what wavelength error $\Delta\lambda$ (external to the emulsion) will result in the diffraction efficiency dropping to 50%?
- (d) Again assuming operation at the first maximum of the diffraction efficiency, and assuming no error in the reconstruction wavelength, what angular error $\Delta\theta$ (external to the emulsion) will result in a reduction of the diffraction efficiency to 50%?
- 9-8.** A holographic plate of thickness $15\ \mu\text{m}$ records a hologram by interference of two plane waves with equal but opposite angles to the emulsion normal. The wavelength for both recording and reconstruction is 633 nm, and the refractive index of the emulsion is 1.52 before and after development. For what angle (in air) between the two interfering beams will the thickness parameter Q of Eq. (9-44) have value 2π ?
- 9-9.** Consider a thick transmission, unslanted sinusoidal absorption grating. Let the absorption modulation have its maximum possible value. Assume that the interfering plane waves are separated in angle by 60° . Under Bragg matched conditions, what average density D of the transparency yields the maximum possible diffraction efficiency of 3.7%?
- 9-10.** Using Eq. (9-66), show that, in the absence of wavelength mismatch, the angular selectivity of a volume grating is maximized when the object and reference waves are separated by an angle of 90° . Hint: remember that K depends on θ .
- 9-11.** Show that for a Fourier transform computer-generated hologram, the number of samples required for the object amplitude equals the number of samples required for the hologram wave amplitude.
- 9-12.** Consider the problem of constructing a computer-generated hologram when the geometry is as illustrated in Fig. 9.39(a), but with the object plane moved by distance Δz to the left of the front focal plane. Determine the approximate bandwidth of the hologram field and the minimum allowed spacing of the samples in the hologram plane, as a function of Δz and any other needed parameters.
- 9-13.** For a certain binary detour-phase hologram, square cells (size $L \times L$) are allocated to each Fourier coefficient, and the amplitudes $|a_{pq}|$ of those coefficients are represented by opening a rectangular subcell within each cell. The width w_x for all transparent subcells is constrained to be 1/10th of the cell width to satisfy the approximations used. The width w_y can range from 0 to the full size of a cell, depending on the amplitude to be represented. The hologram is uniformly illuminated by a normally incident plane wave, and no light is lost by the Fourier transforming lens that follows the hologram. For the purposes of this problem, the object is taken to be a point source located at the center of the object space, yielding Fourier coefficients that are all the same constant, say of value a , which we shall take to be a number somewhere between 0 and 1. When the Fourier amplitude is to be a , the vertical height of all subcells is set to $w_y = aL$.
- (a) For a given value of a , find the coefficients of the two-dimensional Fourier series representation of the amplitude transmittance of the binary hologram.

- (b) Calculate the fraction of the total light intensity incident on the hologram that ends up in the zero-frequency spot on the axis in the image plane.
- (c) Calculate the fraction of total incident light that is blocked by the opaque portions of the hologram.
- (d) Find the diffraction efficiency for both of the two first-order images.

9-14. The following table lists approximate cutoff frequencies (in line-pairs/mm) for several different types of film:

Kodak Tri-X	25
Kodak High-Contrast Copy	30
Kodak SO-243	150
Agfa Agepan FF	300

Assume illumination at 632.8 nm and a lensless Fourier transform geometry with reference source and object both 10 cm from the film. For each film above, estimate the radius of the circle about the reference point outside of which object points will be at the respective cutoff frequencies.

9-15. A certain film has a nonlinear t_A vs. E curve which, over its region of operation, may be described by

$$t_A = t_b + \beta E_1^3,$$

where E_1 represents the variations of exposure about the reference exposure.

- (a) Assuming a reference wave $A \exp(-j2\pi\alpha x)$ and an object wave

$$a(x, y) \exp[-j\phi(x, y)]$$

at the film, find an expression for that portion of the transmitted field that generates the twin first-order images.

- (b) To what does this expression reduce if $A \gg |a|$?
- (c) How do the amplitude and phase modulations obtained in the previous parts of the problem compare with the ideal amplitude and phase modulations present when the film has a linear t_A vs. E curve?

APPENDIX A

Delta Functions and Fourier Transform Theorems

A.1 DELTA FUNCTIONS

The one-dimensional Dirac delta function, widely used in systems analysis, is in fact not a function at all, but rather is a more general entity, often called a "functional" or a "distribution". While a function is an entity that maps a number (the argument of the function), into a number (the value of the function), a functional maps a function into a number. A simple example of a functional is a definite integral, for example

$$\int_{-\infty}^{\infty} h(\xi) d\xi,$$

which maps any given function $h(\xi)$ into the value of its area.

In this spirit, the defining characteristic of the delta function¹ is its so-called "sifting" property under integration, namely

$$\int_{-\infty}^{\infty} \delta(\xi - b)h(\xi) d\xi = \begin{cases} h(b) & b \text{ a point of continuity of } h \\ \frac{1}{2}[h(b^+) + h(b^-)] & b \text{ a point of discontinuity of } h. \end{cases} \quad (\text{A-1})$$

In this equation, the symbols $h(b^+)$ and $h(b^-)$ represent the limiting values of h as its argument approaches the discontinuity from above and from below, respectively. The mapping of a function h into the values on the right of the above equation defines the functional we call the delta function. The integral is a convenient representation for this mapping, but should not be interpreted literally as an integral. It can be viewed as the limit of a set of integrals, i.e.

$$\int_{-\infty}^{\infty} \delta(\xi - b)h(\xi) d\xi \equiv \lim_{N \rightarrow \infty} \int_{-\infty}^{\infty} g_N(\xi - b)h(\xi) d\xi \quad (\text{A-2})$$

¹We continue to use the word *function* due to its common use, even though it is not strictly correct.

where g_N is a sequence of functions that in the limit $N \rightarrow \infty$ exhibit the required sifting property. Such functions must all have unit area, and must *in some sense* become narrower and narrower as N grows large.

It has become fairly common practice in the engineering literature to represent the delta function by the limit of the sequence of functions g_N in Eq. (A-2), i.e. to write

$$\delta(x) = \lim_{N \rightarrow \infty} g_N(x). \quad (\text{A-3})$$

Although this representation is not strictly correct, the limit of the sequence of integrals being the proper representation, nonetheless we use it here with the understanding that it really means what is expressed in Eq. (A-2). Thus we write, for example that

$$\delta(x) = \lim_{N \rightarrow \infty} N \exp(-N^2 \pi x^2)$$

$$\delta(x) = \lim_{N \rightarrow \infty} N \text{rect}(Nx)$$

$$\delta(x) = \lim_{N \rightarrow \infty} N \text{sinc}(Nx).$$

A plot of the last of the above functions shows that $N \text{sinc}(Nx)$ does not become a very narrow pulse as $N \rightarrow \infty$, but rather it retains a finite spread and develops ever more rapid oscillations everywhere except at the origin. Such oscillations in the limit assure that under an integral sign the value of h at the location of the center of the function sequence will be sifted out. Thus it is not necessary that the functions g_N vanish in the limit everywhere except the origin. A somewhat more bizarre example is the function sequence

$$g_N(x) = N e^{j\pi/4} \exp[j\pi(Nx)^2],$$

each member of which has unit area and magnitude N everywhere, but still exhibits a sifting property in the limit.

While the δ function is used in electrical systems analysis to represent a sharp, intense pulse of current or voltage, the analogous concept in optics is a point source of light, or a *spatial* pulse of unit volume. The definition of the δ function in two dimensions is a simple extension of the one-dimensional case, although there is even greater latitude in the possible functional forms of the pulse sequences used. Many possible definitions use separable pulse sequences, e.g.

$$\delta(x, y) = \lim_{N \rightarrow \infty} N^2 \exp[-N^2 \pi(x^2 + y^2)]$$

$$\delta(x, y) = \lim_{N \rightarrow \infty} N^2 \text{rect}(Nx) \text{rect}(Ny)$$

$$\delta(x, y) = \lim_{N \rightarrow \infty} N^2 \text{sinc}(Nx) \text{sinc}(Ny).$$

Other possible definitions use circularly symmetric functions, e.g.

$$\delta(x, y) = \lim_{N \rightarrow \infty} \frac{N^2}{\pi} \text{circ}(N \sqrt{x^2 + y^2}) \quad (\text{A-4})$$

$$\delta(x, y) = \lim_{N \rightarrow \infty} N \frac{J_1(2\pi N \sqrt{x^2 + y^2})}{\sqrt{x^2 + y^2}} \quad (\text{A-4 cont.})$$

In some applications one definition may be more convenient than others, and the definition best suited for the problem can be chosen.

A property of all two-dimensional delta functions that can be easily proved (see Prob. 2-1(a)) is

$$\delta(ax, by) = \frac{1}{|ab|} \delta(x, y), \quad (\text{A-5})$$

which describes how such entities behave under scaling of coordinates. Again this statement has meaning only under integral signs.

A.2

DERIVATION OF FOURIER TRANSFORM THEOREMS

In this section, brief proofs of basic Fourier transform theorems are presented. For more complete derivations, see [32], [226], and [131].

1. **Linearity theorem.** $\mathcal{F}\{\alpha g + \beta h\} = \alpha \mathcal{F}\{g\} + \beta \mathcal{F}\{h\}$

Proof: This theorem follows directly from the linearity of the integrals that define the Fourier transform.

2. **Similarity theorem.** If $\mathcal{F}\{g(x, y)\} = G(f_x, f_y)$, then

$$\mathcal{F}\{g(ax, by)\} = \frac{1}{|ab|} G\left(\frac{f_x}{a}, \frac{f_y}{b}\right).$$

Proof:

$$\begin{aligned} \mathcal{F}\{g(ax, by)\} &= \iint_{-\infty}^{\infty} g(ax, by) \exp[-j2\pi(f_x x + f_y y)] dx dy \\ &= \iint_{-\infty}^{\infty} g(ax, by) \exp\left[-j2\pi\left(\frac{f_x}{a} ax + \frac{f_y}{b} by\right)\right] \frac{d(ax)}{|a|} \frac{d(by)}{|b|} \\ &= \frac{1}{|ab|} G\left(\frac{f_x}{a}, \frac{f_y}{b}\right). \end{aligned}$$

3. **Shift theorem.** If $\mathcal{F}\{g(x, y)\} = G(f_x, f_y)$, then

$$\mathcal{F}\{g(x - a, y - b)\} = G(f_x, f_y) \exp[-j2\pi(f_x a + f_y b)].$$

Proof:

$$\begin{aligned}
\mathcal{F}\{g(x-a, y-b)\} &= \iint_{-\infty}^{\infty} g(x-a, y-b) \exp[-j2\pi(f_X x + f_Y y)] dx dy \\
&= \iint_{-\infty}^{\infty} g(x', y') \exp\{-j2\pi[f_X(x'+a) + f_Y(y'+b)]\} dx' dy' \\
&= G(f_X, f_Y) \exp[-j2\pi(f_X a + f_Y b)].
\end{aligned}$$

4. Rayleigh's (Parseval's) theorem. If $\mathcal{F}\{g(x, y)\} = G(f_X, f_Y)$, then

$$\iint_{-\infty}^{\infty} |g(x, y)|^2 dx dy = \iint_{-\infty}^{\infty} |G(f_X, f_Y)|^2 df_X df_Y.$$

Proof:

$$\begin{aligned}
\iint_{-\infty}^{\infty} |g(x, y)|^2 dx dy &= \iint_{-\infty}^{\infty} g(x, y) g^*(x, y) dx dy \\
&= \iint_{-\infty}^{\infty} dx dy \left[\iint_{-\infty}^{\infty} d\xi d\eta G(\xi, \eta) \exp[j2\pi(x\xi + y\eta)] \right] \\
&\quad \times \left[\iint_{-\infty}^{\infty} d\alpha d\beta G^*(\alpha, \beta) \exp[-j2\pi(x\alpha + y\beta)] \right] \\
&= \iint_{-\infty}^{\infty} d\xi d\eta G(\xi, \eta) \iint_{-\infty}^{\infty} d\alpha d\beta G^*(\alpha, \beta) \\
&\quad \times \left[\iint_{-\infty}^{\infty} \exp\{j2\pi[x(\xi - \alpha) + y(\eta - \beta)]\} dx dy \right] \\
&= \iint_{-\infty}^{\infty} d\xi d\eta G(\xi, \eta) \iint_{-\infty}^{\infty} d\alpha d\beta G^*(\alpha, \beta) \delta(\xi - \alpha, \eta - \beta) \\
&= \iint_{-\infty}^{\infty} |G(\xi, \eta)|^2 d\xi d\eta.
\end{aligned}$$

5. Convolution theorem. If $\mathcal{F}\{G(x, y)\} = G(f_X, f_Y)$ and $\mathcal{F}\{h(x, y)\} = H(f_X, f_Y)$, then

$$\mathcal{F} \left\{ \iint_{-\infty}^{\infty} g(\xi, \eta) h(x - \xi, y - \eta) d\xi d\eta \right\} = G(f_X, f_Y) H(f_X, f_Y).$$

Proof:

$$\begin{aligned} \mathcal{F} \left\{ \iint_{-\infty}^{\infty} g(\xi, \eta) h(x - \xi, y - \eta) d\xi d\eta \right\} &= \iint_{-\infty}^{\infty} g(\xi, \eta) \mathcal{F} \{ h(x - \xi, y - \eta) \} d\xi d\eta \\ &= \iint_{-\infty}^{\infty} g(\xi, \eta) \exp[-j2\pi(f_X \xi + f_Y \eta)] d\xi d\eta H(f_X, f_Y) \\ &= G(f_X, f_Y) H(f_X, f_Y). \end{aligned}$$

6. Autocorrelation theorem. If $\mathcal{F}\{g(x, y)\} = G(f_X, f_Y)$, then

$$\mathcal{F} \left\{ \iint_{-\infty}^{\infty} g(\xi, \eta) g^*(\xi - x, \eta - y) d\xi d\eta \right\} = |G(f_X, f_Y)|^2.$$

Proof:

$$\begin{aligned} \mathcal{F} \left\{ \iint_{-\infty}^{\infty} g(\xi, \eta) g^*(\xi - x, \eta - y) d\xi d\eta \right\} &= \mathcal{F} \left\{ \iint_{-\infty}^{\infty} g(\xi' + x, \eta' + y) g^*(\xi', \eta') d\xi' d\eta' \right\} \\ &= \iint_{-\infty}^{\infty} d\xi' d\eta' g^*(\xi', \eta') \mathcal{F} \{ g(\xi' + x, \eta' + y) \} \\ &= \iint_{-\infty}^{\infty} d\xi' d\eta' g^*(\xi', \eta') \exp[j2\pi(f_X \xi' + f_Y \eta')] G(f_X, f_Y) \\ &= G^*(f_X, f_Y) G(f_X, f_Y) = |G(f_X, f_Y)|^2. \end{aligned}$$

7. Fourier integral theorem. At each point of continuity of g ,

$$\mathcal{F}\mathcal{F}^{-1}\{g(x, y)\} = \mathcal{F}^{-1}\mathcal{F}\{g(x, y)\} = g(x, y).$$

At each point of discontinuity of g , the two successive transformations yield the angular average of the value of g in a small neighborhood of that point.

Proof: Let the function $g_R(x, y)$ be defined by

$$g_R(x, y) = \iint G(f_X, f_Y) \exp[j2\pi(f_X x + f_Y y)] df_X df_Y$$

where A_R is a circle of radius R , centered at the origin of the (f_x, f_y) plane. To prove the theorem, it suffices to show that, at each point of continuity of g ,

$$\lim_{R \rightarrow \infty} g_R(x, y) = g(x, y),$$

and that, at each point of discontinuity of g ,

$$\lim_{R \rightarrow \infty} g_R(x, y) = \frac{1}{2\pi} \int_0^{2\pi} g_o(\theta) d\theta,$$

where $g_o(\theta)$ is the angular dependence of g in a small neighborhood about the point in question. Some initial straightforward manipulation can be performed as follows:

$$\begin{aligned} g_R(x, y) &= \iint_{A_R} \left\{ \iint_{-\infty}^{\infty} d\xi d\eta g(\xi, \eta) e^{-j2\pi(f_x \xi + f_y \eta)} \right\} e^{j2\pi(f_x x + f_y y)} df_x df_y \\ &= \iint_{-\infty}^{\infty} d\xi d\eta g(\xi, \eta) \iint_{A_R} df_x df_y \exp\{j2\pi[f_x(x - \xi) + f_y(y - \eta)]\}. \end{aligned}$$

Noting that

$$\iint_{A_R} df_x df_y \exp\{j2\pi[f_x(x - \xi) + f_y(y - \eta)]\} = R \left[\frac{J_1(2\pi Rr)}{r} \right],$$

where $r = \sqrt{(x - \xi)^2 + (y - \eta)^2}$, we have

$$g_R(x, y) = \iint_{-\infty}^{\infty} d\xi d\eta g(\xi, \eta) R \left[\frac{J_1(2\pi Rr)}{r} \right].$$

Suppose initially that (x, y) is a point of continuity of g . Then

$$\begin{aligned} \lim_{R \rightarrow \infty} g_R(x, y) &= \iint_{-\infty}^{\infty} d\xi d\eta g(\xi, \eta) \lim_{R \rightarrow \infty} R \left[\frac{J_1(2\pi Rr)}{r} \right] \\ &= \iint_{-\infty}^{\infty} d\xi d\eta g(\xi, \eta) \delta(x - \xi, y - \eta) = g(x, y), \end{aligned}$$

where Eq. (A-4) has been used in the second step. Thus the first part of the theorem has been proved.

Consider next a point of discontinuity of g . Without loss of generality that point can be taken to be the origin. Thus we write

$$g_R(0, 0) = \iint_{-\infty}^{\infty} d\xi d\eta g(\xi, \eta) R \left[\frac{J_1(2\pi Rr)}{r} \right],$$

where $r = \sqrt{\xi^2 + \eta^2}$. But for sufficiently large R , the quantity in brackets has significant value only in a small neighborhood of the origin. In addition, in this small neighborhood the function g depends (approximately) only on the angle θ about that point, and therefore

$$g_R(0, 0) \approx \int_0^{2\pi} g_o(\theta) d\theta \int_0^\infty rR \left[\frac{J_1(2\pi Rr)}{r} \right] dr$$

where $g_o(\theta)$ represents the θ dependence of g about the origin. Finally, noting that

$$\int_0^\infty rR \left[\frac{J_1(2\pi Rr)}{r} \right] dr = \frac{1}{2\pi}$$

we conclude that

$$\lim_{R \rightarrow \infty} g_R(0, 0) = \frac{1}{2\pi} \int_0^{2\pi} g_o(\theta) d\theta,$$

and the proof is thus complete.

APPENDIX B

Introduction to Paraxial Geometrical Optics

B.1 THE DOMAIN OF GEOMETRICAL OPTICS

If the wavelength of light is imagined to become vanishingly small, we enter a domain in which the concepts of geometrical optics suffice to analyze optical systems. While the actual wavelength of light is always finite, nonetheless provided all variations or changes of the amplitude and phase of a wavefield take place on spatial scales that are very large compared with a wavelength, the predictions of geometrical optics will be accurate. Examples of situations for which geometrical optics does not yield accurate predictions occur when we insert a sharp edge or a sharply defined aperture in a beam of light, or when we change the phase of a wave by a significant fraction of 2π radians over spatial scales that are comparable with a wavelength.

Thus if we imagine a periodic phase grating for which a "smooth" change of phase by 2π radians takes place only over a distance of many wavelengths, the predictions of geometrical optics for the amplitude distribution behind the grating will be reasonably accurate. On the other hand, if the changes of 2π radians take place in only a few wavelengths, or take place very abruptly, then diffraction effects can not be ignored, and a full wave-optics (or "physical-optics") treatment of the problem is needed.

This appendix is not a complete introduction to the subject of geometrical optics. Rather, we have selected several topics that will help the reader better understand the relationship between geometrical optics and physical optics. In addition, several geometrical concepts that are needed in formulating the physical-optics description of imaging and spatial filtering systems are introduced.

The Concept of a Ray

Consider a monochromatic disturbance traveling in a medium with refractive index that varies slowly on the scale of an optical wavelength. Such a disturbance can be described

by an amplitude and phase distribution

$$U(\vec{r}) = A(\vec{r}) \exp[-jk_o S(\vec{r})], \quad (\text{B-1})$$

where $A(\vec{r})$ is the amplitude and $k_o S(\vec{r})$ is the phase of the wave. Here k_o is the free-space wavenumber $2\pi/\lambda_o$; the refractive index n of the medium is contained in the definition of S . $S(\vec{r})$ is called the Eikonal function. We follow the argument presented in [253] (p. 52) in finding the equation that must be satisfied by the Eikonal function.

Surfaces defined by

$$S(\vec{r}) = \text{constant}$$

are called wavefronts of the disturbance. The direction of power flow and the direction of the wave vector \vec{k} are both normal to the wavefronts at each point \vec{r} in an isotropic medium. A ray is defined as a trajectory or a path through space that starts at any particular point on a wavefront and moves through space with the wave, always remaining perpendicular to the wavefront at every point on the trajectory. Thus a ray traces out the path of power flow in an isotropic medium. Substitution of (B-1) in the Helmholtz equation of Eq. (3-12) yields the following equation that must be satisfied by both $A(\vec{r})$ and $S(\vec{r})$:

$$k_o^2 \left[n^2 - |\nabla S|^2 \right] A + \nabla^2 A - jk_o \left[2\nabla S \cdot \nabla A + A \nabla^2 S \right] = 0.$$

The real and imaginary parts of this equation must vanish independently. For the real part to vanish, we require

$$|\nabla S|^2 = n^2 + \left(\frac{\lambda_o}{2\pi} \right)^2 \frac{\nabla^2 A}{A}. \quad (\text{B-2})$$

Using the artifice of allowing the wavelength to approach zero to recover the geometrical-optics limit of this equation, the last term is seen to vanish, leaving the so-called Eikonal equation, which is perhaps the most fundamental description of the behavior of light under the approximations of geometrical optics,

$$|\nabla S(\vec{r})|^2 = n^2(\vec{r}). \quad (\text{B-3})$$

This equation serves to define the wavefront S . Once the wavefronts are known, the trajectories defining rays can be determined.

Rays and Local Spatial Frequency

Consider a monochromatic wave propagating in three dimensional space defined by an (x, y, z) coordinate system, with propagation being in the positive z direction. At each point on a plane of constant z , there is a well defined direction of the ray through that point, a direction that coincides with the direction of the wave vector \vec{k} at that point.

We have seen previously that an arbitrary distribution of complex field across a plane can be decomposed by means of a Fourier transform into a collection of plane-wave components traveling in different directions. Each such plane wave component

has a unique wave vector with direction cosines (α, β, γ) defined by Eq. (3-58), and can be regarded as one spatial frequency associated with the wave.

The spatial frequencies defined through the Fourier decomposition exist everywhere in space and cannot be regarded as being localized. However, for complex functions with a phase that does not vary too rapidly, the concept of a local spatial frequency can be introduced, as was done in Section 2.2. The definitions of the local spatial frequencies (f_{lX}, f_{lY}) given there can also be viewed as defining the local direction cosines $(\alpha_l, \beta_l, \gamma_l)$ of the wavefront through the relations

$$\alpha_l = \lambda f_{lX} \quad \beta_l = \lambda f_{lY} \quad \gamma_l = \sqrt{1 - \alpha_l^2 - \beta_l^2}. \quad (\text{B-4})$$

These local direction cosines are in fact the direction cosines of the ray through the (x, y) plane at each point. This leads us to the following important observation:

The description of the local spatial frequencies of a wavefront is identical with the description of that wavefront in terms of the rays of geometrical optics. Ray direction cosines are found from local spatial frequencies simply by multiplication by the wavelength.

B.2 REFRACTION, SNELL'S LAW, AND THE PARAXIAL APPROXIMATION

Rays traveling in a medium with constant index of refraction always travel in straight lines, as can be derived from the Eikonal equation. However, when the wave travels through a medium having an index of refraction that changes in space (i.e. an **inhomogeneous** medium), the ray directions will undergo changes that depend on the changes of refractive index. When the changes of refractive index are gradual, the ray trajectories will be smoothly changing curves in space. Such bending of the rays is called refraction.

However, when a wave encounters an abrupt boundary between two media having different refractive indices, the ray directions are changed suddenly as they pass through the interface. The angles of incidence θ_1 and refraction θ_2 , as shown in Fig. 3.1, are related by Snell's law,

$$n_1 \sin \theta_1 = n_2 \sin \theta_2, \quad (\text{B-5})$$

where n_1 and n_2 are the refractive indices of the first and second media, respectively. In the problems of interest here, the changes of refractive index, as encountered, for example, on passage through a lens, will always be abrupt, so Snell's law will form the basis for our analyses.

A further simplifying approximation can be made if we restrict attention to rays that are traveling close to the optical axis and at small angles to that axis, the geometrical optics version of the *paraxial* approximation. In such a case, Snell's law reduces to a simple linear relationship between the angle of incidence and the angle of refraction,

$$n_1 \theta_1 = n_2 \theta_2, \quad (\text{B-6})$$

and in addition the cosines of these angles can be replaced by unity.

The product $\mathcal{S} = n\theta$ of the refractive index n and an angle θ within that medium is called a *reduced angle*. Thus the **paraxial** version of **Snell's** law states that the reduced angle remains constant as light passes through a sharp interface between media of different refractive indices,

$$\hat{\theta}_1 = \hat{\theta}_2. \quad (\text{B-7})$$

B.3 THE RAY-TRANSFER MATRIX

Under **paraxial** conditions, the properties of rays in optical systems can be treated with an elegant matrix formalism, which in many respects is the geometrical-optics equivalent of the operator methods of wave optics introduced in Section 5.4. Additional references for this material are [163], [253], and [261]. To apply this methodology, it is necessary to consider only *meridional* rays, which are rays traveling in paths that are completely contained in a single plane containing the z axis. We call the transverse axis in this plane the y axis, and therefore the plane of interest is the (y, z) plane.

Figure B.1 shows the typical kind of ray propagation problem that must be solved in order to understand the effects of an optical system. On the left, at axial coordinate z_1 , is an input plane of an optical system, and on the right, at axial coordinate z_2 , is an output plane. A ray with transverse coordinate y_1 enters the system at angle θ_1 , and the same ray, now with transverse coordinate y_2 , leaves the system with angle θ_2 . The goal is to determine the position y_2 and angle θ_2 of the output ray for every possible y_1 and θ_1 associated with an input ray.

Under the **paraxial** condition, the relationships between (y_2, θ_2) and (y_1, θ_1) are linear and can be written explicitly as

$$\begin{aligned} y_2 &= Ay_1 + B\hat{\theta}_1 \\ \hat{\theta}_2 &= Cy_1 + D\hat{\theta}_1, \end{aligned} \quad (\text{B-8})$$

where for reasons that will become evident, we use reduced angles rather than just angles. The above equation can be expressed more compactly in matrix notation,

$$\begin{bmatrix} y_2 \\ \hat{\theta}_2 \end{bmatrix} = \begin{bmatrix} A & B \\ C & D \end{bmatrix} \begin{bmatrix} y_1 \\ \hat{\theta}_1 \end{bmatrix}. \quad (\text{B-9})$$

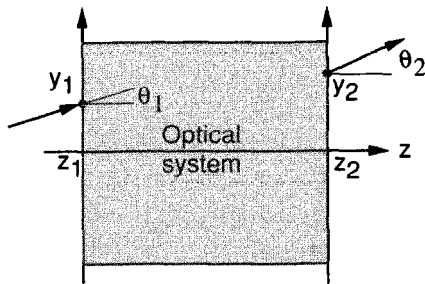


FIGURE B.1
Input and output of an optical system.

The matrix

$$\mathbf{M} = \begin{bmatrix} A & B \\ C & D \end{bmatrix}$$

is called the ray-transfer matrix or the ABCD matrix.

The ray-transfer matrix has an interesting interpretation in terms of local spatial frequencies. In the (y, z) plane under paraxial conditions, the reduced ray angle $\hat{\theta}$ with respect to the z axis is related to local spatial frequency f_l through

$$f_l = \frac{\theta}{\lambda} = \frac{\hat{\theta}}{\lambda_0}.$$

Therefore the ray-transfer matrix can be regarded as specifying a transformation between the spatial distribution of local spatial frequency at the input and the corresponding distribution at the output.

Elementary Ray-Transfer Matrices

Certain simple structures are commonly encountered in ray tracing problems. Here we specify the ray-transfer matrices for the most important of these structures. They are all illustrated in Fig. B.2.

- 1. Propagation through free space of index n .** Geometrical rays travel in straight lines in a medium with constant refractive index. Therefore the effect of propagation through free space is to translate the location of the ray in proportion to the angle at which it travels and to leave the angle of the ray unchanged. The ray-transfer matrix describing propagation over distance d is therefore

$$\mathbf{M} = \begin{bmatrix} 1 & d/n \\ 0 & 1 \end{bmatrix} \tag{B-10}$$

- 2. Refraction at a planar interface.** At a planar interface the position of the ray is unchanged but the angle of the ray is transformed according to Snell's law; the reduced

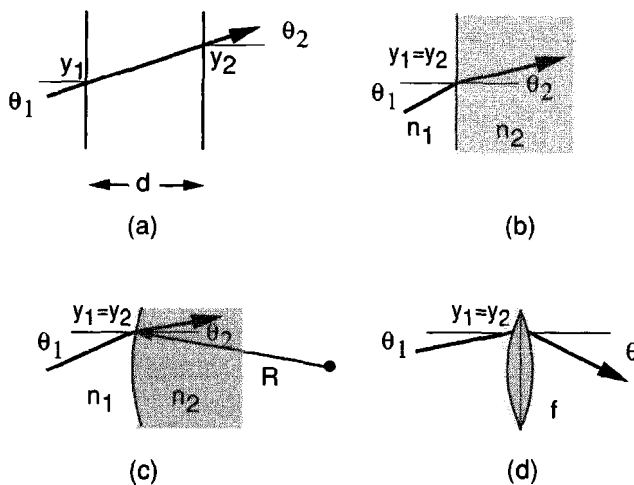


FIGURE B.2
Elementary structures for ray-transfer matrix calculations. (a) Free space, (b) a planar interface, (c) a spherical interface, and (d) a thin lens.

angle remains unchanged. Therefore the ray-transfer matrix for a planar interface between a medium of refractive index n_1 and a medium of refractive index n_2 is

$$\mathbf{M} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}. \quad (\text{B-11})$$

- 3. Refraction at a spherical interface.** At a spherical interface between an initial medium with refractive index n_1 and a final medium with refractive index n_2 , the position of a ray is again not changed, but the angle is changed. However at a point on the interface at distance y from the optical axis, the normal to the interface is not parallel to the optical axis, but rather is inclined with respect to the optical axis by angle

$$\psi = \arctan \frac{Y}{R} \approx \frac{Y}{R}$$

where R is the radius of the spherical surface. Therefore if θ_1 and θ_2 are measured with respect to the optical axis, Snell's law at transverse coordinate y becomes

$$n_1 \theta_1 + n_1 \frac{y}{R} = n_2 \theta_2 + n_2 \frac{y}{R},$$

or, using reduced angles,

$$\hat{\theta}_1 + n_1 \frac{y}{R} = \hat{\theta}_2 + n_2 \frac{y}{R}.$$

Solving for $\hat{\theta}_2$ yields

$$\hat{\theta}_2 = \hat{\theta}_1 + \frac{n_1 - n_2}{R} y.$$

The ray-transfer matrix for a spherical interface can now be written as

$$\mathbf{M} = \begin{bmatrix} 1 & 0 \\ \frac{n_1 - n_2}{R} & 1 \end{bmatrix}. \quad (\text{B-12})$$

Note that a positive value for R signifies a convex surface encountered from left to right, while a negative value for R signifies a concave surface.

- 4. Passage through a thin lens.** A thin lens (index n_2 embedded in a medium of index n_1) can be treated by cascading two spherical interfaces. The roles of n_1 and n_2 are interchanged for the two surfaces. Representing the ray-transfer matrices of the surfaces on the left and the right by \mathbf{M}_1 and \mathbf{M}_2 , respectively, the ray-transfer matrix for the sequence of two surfaces is

$$\begin{aligned} \mathbf{M} &= \mathbf{M}_2 \mathbf{M}_1 \\ &= \begin{bmatrix} 1 & 0 \\ \frac{n_2 - n_1}{R_2} & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ \frac{n_1 - n_2}{R_1} & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ -(n_2 - n_1) \left(\frac{1}{R_2} - \frac{1}{R_1} \right) & 1 \end{bmatrix}. \end{aligned}$$

We define the focal length of the lens by

$$\frac{1}{f} = \frac{n_2 - n_1}{n_1} \left(\frac{1}{R_1} - \frac{1}{R_2} \right). \quad (\text{B-13})$$

in which case the ray-transfer matrix for a thin lens becomes

$$\mathbf{M} = \begin{bmatrix} 1 & 0 \\ -\frac{n_1}{f} & 1 \end{bmatrix}. \quad (\text{B-14})$$

The most useful elementary ray-transfer matrices have now been presented. Propagation through a system consisting of regions of free space separated by thin lenses can be treated with these matrices. Note that, just as with the wave-optics operators presented in Chapter 5, the ray-transfer matrices should be applied in the sequence in which the structures are encountered. If light propagates first through a structure with ray-transfer matrix \mathbf{M}_1 , then through a structure with ray-transfer matrix \mathbf{M}_2 , etc., with a final structure having ray-transfer matrix \mathbf{M}_n , then the overall ray-transfer matrix for the entire system is

$$\mathbf{M} = \mathbf{M}_n \cdots \mathbf{M}_2 \mathbf{M}_1. \quad (\text{B-15})$$

We note also that, because we have chosen to use *reduced* angles, rather than the angles themselves in the definition of the ray-transfer matrix, all of the elementary matrices presented have a determinant that is unity.

B.4

CONJUGATE PLANES, FOCAL PLANES, AND PRINCIPAL PLANES

There exist certain planes within an optical system that play important conceptual and practical roles. In this section we explain the three most important of these types of planes.

Conjugate Planes

Two planes within an optical system are said to be *conjugate planes* if the intensity distribution across one plane is an image (generally magnified or demagnified) of the intensity distribution across the other plane. Likewise, two points are said to be conjugate points if one is the image of the other.

The properties that must be satisfied by the ray-transfer matrix between two conjugate planes can be deduced by considering the relation between two conjugate points y_1 and y_2 , as implied by Eq. (B-8). The position of the point y_2 that is conjugate to y_1 should be independent of the reduced angle $\hat{\theta}_1$ of a ray through y_1 , implying that the matrix element B should be zero. The position y_2 should be related to the position y_1 only through the transverse magnification m_t , which is the scale factor between coordinates in the two planes. We conclude that the matrix element A must equal m_t . In addition, the angles of the rays passing through y_2 will generally be magnified or demagnified with respect to the angles of the same rays passing through y_1 . The magnification for reduced angles is represented by m_α , and we conclude that the matrix element D must satisfy $D = m_\alpha$. There is no general restriction on the matrix element C, so the ray-transfer matrix between conjugate planes takes the general form

$$\mathbf{M} = \begin{bmatrix} m_t & 0 \\ C & m_\alpha \end{bmatrix}.$$

Recalling that angles and positions are conjugate Fourier variables, the scaling theorem of Fourier analysis implies that the transverse magnification and the angular magnification must be related in a reciprocal fashion. The magnifications m_t and m_α are in fact related by

$$m_t m_\alpha = 1. \quad (\text{B-16})$$

Thus the form of the ray-transfer matrix for conjugate planes is

$$\mathbf{M} = \begin{bmatrix} m_t & 0 \\ C & m_t^{-1} \end{bmatrix}.$$

Note that both m_t and m_α can be positive or negative (signifying image inversion), but they must be of the same sign.

The **paraxial** relation (B-16) has a more general nonparaxial form, known as the sine condition, which states that for conjugate points y_1 and y_2 the following equation must be satisfied:

$$n_1 y_1 \sin \theta_1 = n_2 y_2 \sin \theta_2. \quad (\text{B-17})$$

Focal Planes

Consider a parallel bundle of rays traveling parallel to the optical axis and entering a lens. Whether that lens is thick or thin, for **paraxial** rays there will exist a point on the optical axis toward which that ray bundle will converge (positive lens) or from which it will appear to diverge (negative lens). See Fig. B.3 for an illustration. Considering a positive lens for the moment, the point behind the lens at which this originally parallel ray bundle crosses in a focused point is called the rearfocal point or the secondfocal point of the lens. A plane constructed through that point perpendicular to the optical axis is called the rearfocal plane or the secondfocal plane. It has the property that a **paraxial** parallel bundle of rays traveling into the lens at any angle with respect to the optical axis will be brought to a focus at a point in the focal plane that depends on the initial angle of the bundle.

In a similar fashion, consider a point source on the optical axis in front of a positive lens, thick or thin. The particular point in front of the lens for which the diverging bundle of rays is made to emerge as a parallel bundle traveling parallel to the optical axis behind the lens is called the front focal point (or **the firstfocal point**) of the lens. A plane erected through the front focal point normal to the optical axis is called the front focal plane (or **the firstfocal plane**) of the lens.

For a negative lens, the roles of the front and rear focal points and planes are reversed. The front focal point is now the point from which a bundle of rays, originally parallel to the optical axis, appears to be diverging when viewed from the exit side of the lens. The rear focal point is defined by the point of convergence of an incident bundle of rays that emerges parallel or collimated after passage through the lens.

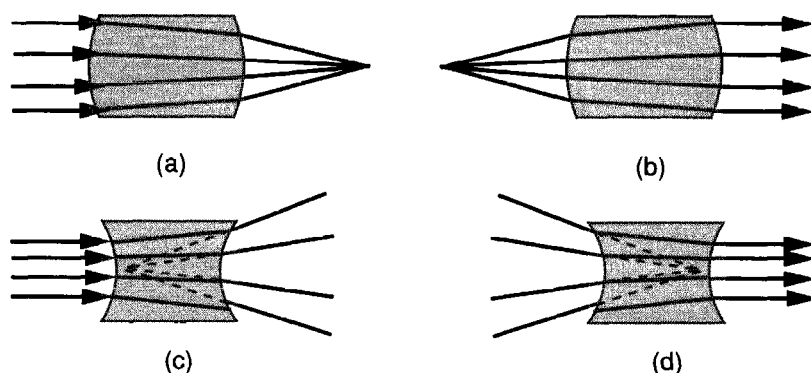


FIGURE B.3
 Definition of focal points. (a) Rear focal point of a positive lens, (b) front focal point of a positive lens, (c) front focal point of a negative lens, and (d) rear focal point of a negative lens.

The mapping from the front focal plane to the rear focal plane is one that maps angles into positions, and positions into angles. If f is the focal length of the lens, then the ray-transfer matrix between focal planes takes the form

$$\mathbf{M} = \begin{bmatrix} 0 & f/n_1 \\ -n_1/f & 0 \end{bmatrix},$$

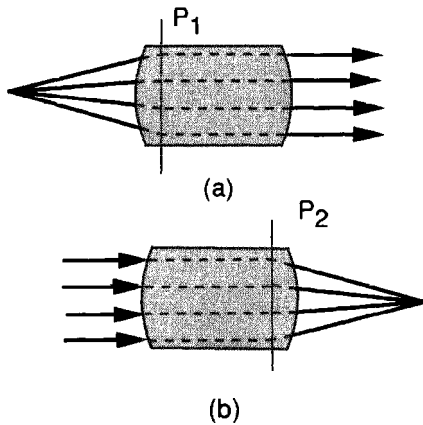
as can be readily verified by multiplying together three matrices representing propagation over distance f , passage through a thin lens with focal length f , and propagation over a second distance f .

Principal Planes

By the definition of a thin lens, a ray incident at input coordinate y_1 exits that lens at the same coordinate $y_2 = y_1$. For a thick lens this simple idealization is no longer valid. A ray entering the first spherical surface at coordinate y_1 will in general leave the second spherical surface at a different coordinate $y_2 \neq y_1$, as can be seen in Fig. B.3.

Much of the simplicity of a thin lens can be retained for a thick lens by introducing the concept of principal planes. Principal planes are planes where the focusing power of the lens can be imagined to be concentrated.

To find the first principal plane of a lens, trace a ray from the front focal point to the first lens surface, as shown in Fig. B.4. By definition of the focal point, that ray will exit the second surface of the lens parallel to the optical axis, i.e. in a collimated beam. If we project the incident ray forwards and the exiting ray backwards into the lens, retaining their original angles, they will intersect at a point. A plane through this point normal to the optical axis defines the *first* principal plane. For this geometry it is possible to imagine that all the refraction associated with the lens takes place in this principal plane.

**FIGURE B.4**

Definitions of principal planes. (a) First principal plane P_1 , (b) second principal plane P_2 .

In the most general case, different rays diverging from the front focal point might define different planes, which would be an indication that the principal plane is not a plane at all, but rather is a curved surface. Such can be the case for lenses with very large aperture or for special lenses such as wide-angle lenses, but for the lenses of interest to us in this book the principal planes are indeed flat to an excellent approximation.

The second principal plane is found by starting with a ray that is parallel to the optical axis, and tracing it through the rear focal point of the lens. The extension of the incident ray and the exiting ray intersect in a point, which in turn defines the *second principal plane* of the lens, again normal to the optical axis. For this geometry it is possible to imagine that all of the power of the lens is concentrated in the second principal plane.

For more general geometries, ray bending can be imagined to take place in both of the principal planes. As will be seen shortly, the two planes are in fact conjugate to one another with unit magnification. A ray incident at particular transverse coordinates on the first principal plane will exit from the second principal plane at those same coordinates, but in general with a change of angle.

In general, the first and second principal planes are separate planes. However, the definition of a thin lens implies that for such a lens the distinguishing characteristic is that the first and second principal planes coincide, and all the focusing power can be imagined to be concentrated in a single plane.

The relationship between the principal planes can be more fully understood if we derive the ray-transfer matrix that holds for propagation between the two principal planes. The derivation is based on the two geometries already introduced, namely that of a point source at the front focal point that yields a collimated ray bundle leaving the second principal plane, and that of a collimated bundle incident on the first principal plane that yields a ray bundle converging from the second principal plane toward a focus at the rear focal point. Considering the case of collimated input light passing through the rear focal point, we find that the matrix element A must be unity, and the matrix element C must be $-n_1/f$. Consideration of the case of input rays diverging from the front focal point shows that $B = 0$ and $D = 1$. Thus the ray-transfer matrix for the passage between principal planes is

$$\mathbf{M} = \begin{bmatrix} 1 & 0 \\ -\frac{n_1}{f} & 1 \end{bmatrix}.$$

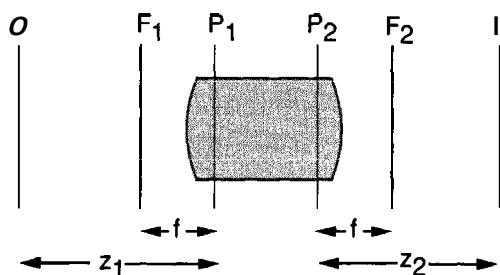


FIGURE B.5
Relations between principal planes, focal lengths, and object/image distances.

This matrix is identical with the ray-transfer matrix describing passage through a thin lens. Thus by constructing the principal planes, and by tracing **rays only** up to the first principal plane and away from the second principal plane, we are able to treat a complex lens system as if it were a simpler thin lens. Note that the ray-transfer matrix above implies that the two principal planes are conjugate to one another, and the magnification between them is unity.

The **focal length** of a lens is by definition the distance of a principal plane from the corresponding focal point that was used in its definition. Assuming that the refractive indices of the media in front of and behind the lens are the same, the distance of the front focal plane from the first principal plane is identical with the distance of the rear focal point from the second principal plane. That is, the two focal lengths of the lens are the same. Note that for some lenses the second principal plane may lie to the left of the first principal plane. Such an occurrence does not change the definition of the focal length. It can also be shown that the distances z_1 and z_2 in the lens law

$$\frac{1}{z_1} + \frac{1}{z_2} = \frac{1}{f}$$

are measured from the first and second principal planes. These various relations are illustrated in Fig. B.5.

B.5 ENTRANCE AND EXIT PUPILS

Until now, we have not considered the effects of pupils (*i.e.* finite apertures) in optical systems. Apertures, of course, give rise to diffraction effects. The concepts of entrance and exit apertures are of great importance in calculations of the effects of diffraction on optical systems.

A system of lenses may contain several or many different apertures, but one such aperture always provides the severest limitation to the extent of the optical wavefront captured at the input of the system, and to the extent of the optical wavefront leaving the system. That aperture may lie deep within the system of lenses, but the single aperture that most severely restricts the bundle of rays passing through the system is in effect the aperture that limits the extent of the wavefront at both the input and at the output.

The **entrance pupil** of the optical system is defined as the **image of the most severely limiting aperture**, when viewed from the object space, looking through any optical elements that may precede the physical aperture. The **exit pupil** of the system is also

defined as the image of the physical aperture, but this time looking from the image space through any optical elements that may lie between that aperture and the image plane.

Figure B.6 illustrates the entrance and exit pupils for a very simple system consisting of a single lens, for three cases: a limiting pupil (1) in the plane of the lens, (2) following the lens, and (3) preceding the lens. In the first case, the entrance and exit apertures coincide with the real physical aperture in the plane of the lens. In the second case, the exit pupil coincides with the physical pupil (which is assumed to limit the angle of the bundle of rays more severely than does the lens aperture), and the entrance pupil is a virtual image of the physical aperture, lying to the right of the lens. In the third case, the entrance pupil is the real physical aperture lying to the left of the lens. In this case, the exit pupil is a virtual image of the physical aperture, lying in a plane to the left of the lens.

In a more complex optical system, containing many lenses and many apertures, it is in general necessary to trace rays through the entire system in order to determine which aperture constitutes the most severe restriction on the ray bundles and therefore which aperture must be imaged to find the entrance and exit pupils.

Once the location and extent of the exit pupil are known, the effects of diffraction on the image of a point-source object can be calculated. For an object point source,

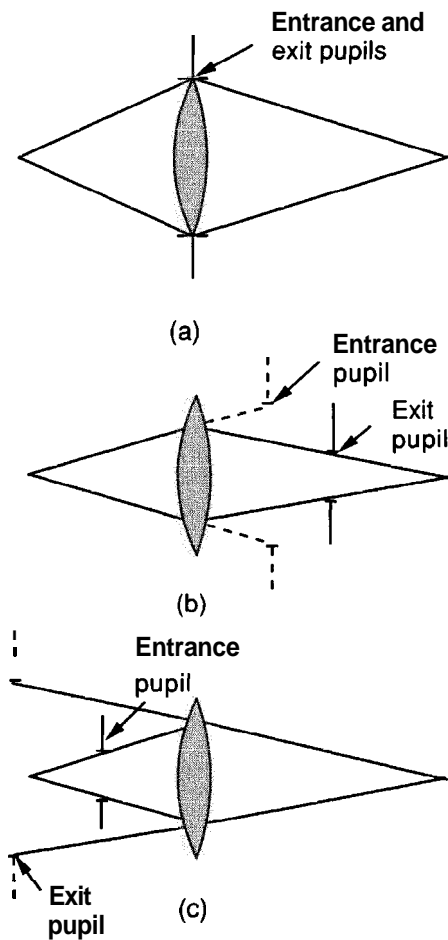


FIGURE B.6

Entrance and exit pupils. (a) Entrance and exit pupils coincide with the physical pupil, (b) the exit pupil coincides with the physical pupil, and (c) the entrance pupil coincides with the physical pupil.

a converging bundle of rays fills the exit pupil on its way to a geometrical image. If the optical system has no aberrations, the geometrical image is an ideal point and the converging bundle defines a perfect spherical wave. The exit pupil limits the angular extent of the converging bundle. The Fraunhofer diffraction formula can now be applied at the exit pupil, using the distance from that pupil to the image as the distance appearing in the formula.

APPENDIX C

Polarization and Jones Matrices

Birefringent media play an important role in the analysis of spatial light modulators of various kinds, as described in Chapter 7. In this appendix we introduce a tool for analyzing polarization-based devices, the so-called Jones calculus, first introduced by R.C. Jones. For an alternative discussion, together with references, see Ref. [123], Section 4.3.

For simplicity, we restrict attention here to monochromatic light, since the problems of interest here arise primarily in coherent optical systems. However, the theory is more general, and can be extended to both narrowband and broadband optical signals with appropriate modifications.

C.1 DEFINITION OF THE JONES MATRIX

Consider a monochromatic light wave, polarized in the (x, y) plane, but with an arbitrary state of polarization in that plane. Let the polarization state be defined by a vector \vec{U} formed from the complex amplitudes (phasor amplitudes) of the x and y components of polarization, U_X and U_Y as follows:

$$\vec{U} = \begin{bmatrix} U_X \\ U_Y \end{bmatrix}. \quad (\text{C-1})$$

We will refer to \vec{U} as the *polarization vector* of the light. Some examples of unit-length polarization vectors describing light with different states of polarization are as follows:

$$\begin{aligned}
&\text{Linearly polarized in the } x \text{ direction: } \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \\
&\text{Linearly polarized in the } y \text{ direction: } \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \\
&\text{Linearly polarized at } +45 \text{ degrees: } \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \\
&\text{Right-hand circularly polarized: } \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ -j \end{bmatrix}, \\
&\text{Left-hand circularly polarized: } \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ j \end{bmatrix}.
\end{aligned} \tag{C-2}$$

As an aside, the convention adopted in optics is to define left-hand and right-hand circular polarization as follows. The observer always looks "head-on" into the wave as it approaches, i.e. towards the source of the light. If from such a perspective the polarization vector is rotating (with a period equal to the optical period and without change of length) in the *clockwise* sense, then the wave is said to be *right-hand circularly polarized*. This is because if you point the thumb of your right hand towards the source, the direction your fingers curl is clockwise, which in this case is the direction of rotation of the polarization vector. If, on the other hand, the direction of rotation is *counter-clockwise*, then for reasons that are probably now obvious we call this wave *left-hand circularly polarized*.

Left-hand and right-hand elliptical polarizations are similar to circular polarizations except that the length of the polarization vector changes periodically as the vector rotates.

When light passes through a polarization-sensitive device, the state of polarization of the wave will in general change, and it is of interest to find a simple representation of the new state of polarization, described by the vector \vec{U}' , in terms of the initial state of polarization described by the vector \vec{U} . All of the polarization devices of interest here are *linear*, and for such devices the initial and final polarization vectors can be related through a 2×2 matrix \mathbf{L} , known as the *Jones matrix*,

$$\vec{U}' = \mathbf{L}\vec{U} = \begin{bmatrix} l_{11} & l_{12} \\ l_{21} & l_{22} \end{bmatrix} \vec{U}. \tag{C-3}$$

The four elements of the Jones matrix fully describe the effects of a linear device on the state of polarization of the wave.

When light passes through a sequence of linear polarization devices, the Jones matrices of the various transformations can be chained together, defining a single new Jones matrix for the sequence of devices through the relation

$$\mathbf{L} = \mathbf{L}_N \cdots \mathbf{L}_2 \mathbf{L}_1, \tag{C-4}$$

where \mathbf{L}_1 is the Jones matrix of the first device encountered, \mathbf{L}_2 that of the second device, etc.

C.2 EXAMPLES OF SIMPLE POLARIZATION TRANSFORMATIONS

Perhaps the simplest transformation of the state of polarization of a wave is that defined by a rotation of the coordinate system within which the wave is described (the wave itself does not change under such a rotation, only our mathematical description of it). If the (x, y) coordinate system is rotated by angle θ in the counterclockwise direction (as illustrated in Fig. C.1), simple geometry shows that

$$\begin{aligned} U'_X &= \cos \theta U_X + \sin \theta U_Y \\ U'_Y &= -\sin \theta U_X + \cos \theta U_Y, \end{aligned} \quad (\text{C-5})$$

and therefore that the Jones matrix for a coordinate rotation is given by

$$\mathbf{L}_{\text{rotate}}(\theta) = \begin{bmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{bmatrix}. \quad (\text{C-6})$$

Closely related to the Jones matrix of a coordinate rotation is the Jones matrix of a polarization device that transforms the polarization of a linearly polarized wave, initially polarized in direction θ_1 with respect to the x axis, into a linearly polarized wave with new polarization direction $\theta_2 = \theta_1 + \theta$. Such a device is called a *polarization rotator*. Since the polarization vectors before and after rotation are given, respectively, by $\begin{bmatrix} \cos \theta_1 \\ \sin \theta_1 \end{bmatrix}$ and $\begin{bmatrix} \cos \theta_2 \\ \sin \theta_2 \end{bmatrix}$, the Jones matrix of a device that rotates the polarization counterclockwise by angle θ must be given by

$$\mathbf{L}_R(\theta) = \mathbf{L}_{\text{rotate}}(-\theta) = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix} \quad (\text{C-7})$$

A second simple case is one in which the X and Y components of the wave undergo different phase delays. A device introducing such a polarization transformation is called a *wave retarder*. For example, a transparent birefringent plate of thickness d having refractive index n_X for the polarization component in the x direction and refractive index n_Y for the polarization component in the y direction, will introduce phase delays $\phi_X = 2\pi n_X d / \lambda_0$ and $\phi_Y = 2\pi n_Y d / \lambda_0$, respectively, in those two components. The Jones matrix for such a transformation can be written

$$\mathbf{L}_{\text{retard}}(\Delta) = \begin{bmatrix} 1 & 0 \\ 0 & e^{-j\Delta} \end{bmatrix}, \quad (\text{C-8})$$

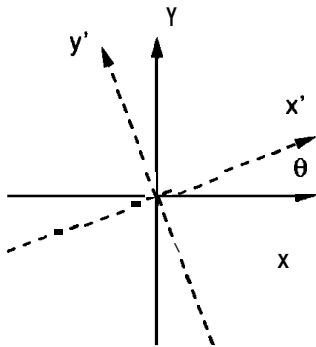


FIGURE C.1

Coordinate rotation. The direction of wave propagation is out of the page.

where λ_o is the vacuum wavelength of light, a common phase delay suffered by both components has been dropped, and the *relative* phase shift A is given by

$$\Delta = \frac{2\pi(n_X - n_Y)d}{\lambda_o} \quad (\text{C-9})$$

A wave retarder of special interest is a *quarter wave plate*, for which $A = d/2$. The Jones matrix for such a device is

$$\mathbf{L}_{\text{retard}}(\pi/2) = \begin{bmatrix} 1 & 0 \\ 0 & -j \end{bmatrix} \quad (\text{C-10})$$

It is easily seen to convert linearly polarized light with polarization direction at 45° to the x axis, described by the polarization vector $\frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 1 \end{bmatrix}$, into right-hand circularly polarized light described by polarization vector $\begin{bmatrix} 1 \\ -j \end{bmatrix}$. Equivalently, this device converts left-hand circularly polarized light $\frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ j \end{bmatrix}$ into linearly polarized light $\frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 1 \end{bmatrix}$.

Another wave retarder of special interest is a *half-wave plate*, for which $A = \pi$ and

$$\mathbf{L}_{\text{retard}}(\pi) = \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix}. \quad (\text{C-11})$$

Comparison of the Jones matrix for such a device with Eq. (C-7) shows that a half-wave plate is a device that rotates the polarization of a wave, initially linearly polarized at 45° to the x axis, by 90° .

As a final example of a polarization device we consider a *polarizer* (or equivalently a *polarization analyzer*) which passes only the wave component that is linearly polarized at angle a to the x axis. With a small amount of work it can be shown that the Jones matrix for such a device is given by

$$\mathbf{L}(a) = \begin{bmatrix} \cos^2 a & \sin a \cos a \\ \sin a \cos a & \sin^2 a \end{bmatrix}. \quad (\text{C-12})$$

C.3 REFLECTIVE POLARIZATION DEVICES

Until this point we have considered only polarization devices used in transmission. Since many spatial light modulators operate in a reflective mode, we turn attention to such a geometry.

Consider a reflective polarization device as illustrated in Fig. C.2. Light enters the device from the left, with normal incidence assumed. It passes through a polarization element having Jones matrix L , is normally incident on a *lossless* mirror, reflects from the mirror, and passes a second time through the same polarization element. We wish to specify the Jones matrix for an equivalent *transmissive* device that will function in the same way as this reflective device.

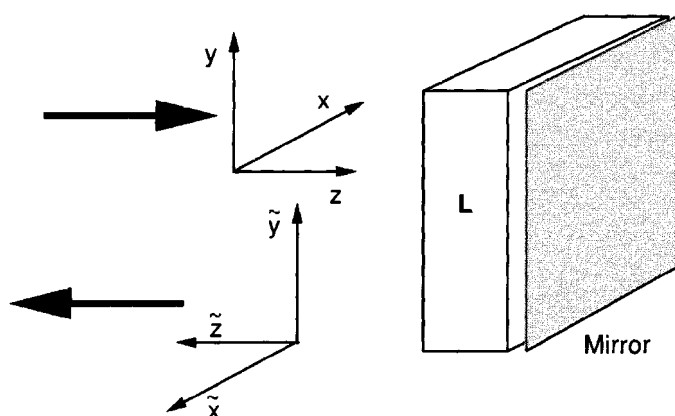


FIGURE C.2
Reflective polarization device.

An important point to consider at the start is that we will consider only *reciprocal* polarization elements before the mirror. For a reciprocal element, the coupling from, say, the x component of polarization to the y component of polarization on the forward pass through the device must equal the coupling from the y component back to the x component on the reverse pass. In addition the forward coupling from the y component to the x component must be the same as the backward coupling from x to y . For a reciprocal element, the Jones matrix for backward passage of light is exactly equal to the *transpose* of the Jones matrix for forward passage of light. Most polarization elements are reciprocal, the most important exceptions being devices that operate by the Faraday effect in the presence of a magnetic field. For such devices the dependence on the direction of the magnetic field destroys reciprocity, and the Jones matrix for reverse propagation is identical with the Jones matrix for forward propagation.

It is also important to note several geometrical factors at the start. First, we specify the polarization vectors of waves by examining the polarization state from a “heads-on” geometry, looking towards the source and using x and y axes that form a right-hand coordinate system, with the z axis pointing in the direction of propagation. This is a convention that we must consistently apply. Note that for a transmissive device, the coordinate system both before and after passage through the device is right-handed. We attempt to retain this convention even with the reflective device.

As shown in Fig. C.2, the z axis is taken to reverse direction after reflection to become \tilde{z} . We have also shown the x axis reversed to obtain a right-hand system, with x being changed to \tilde{x} . However, for the time being, we allow the coordinate system to be left-handed, converting to a right-handed system shown only at the very end.

Consider now the progress of a wave as it travels through the reflective device. It begins with a polarization state described by a vector \vec{U} . This polarization state is modified by passage through the polarization element, yielding a polarization state

$$\vec{U}' = \mathbf{L}\vec{U}.$$

Next the light reflects from the mirror. Since the tangential components of the electric field must be zero at a perfectly conducting boundary, the electric field components after reflection are the negative of their values before reflection. However, we regularly drop constant phase factors, and a negation of the two components of the electric field is just a common phase factor of 180° that we drop. So with this understanding, after

reflection, the field components U_X and U_Y are the same as they were before reflection, when measured in the original (x, y) coordinate system.

The wave now proceeds back through the polarization element. As argued above, for a reciprocal device, the Jones matrix under reverse propagation is \mathbf{L}^t , where the superscript t indicates the ordinary matrix transpose operation.

Finally, if we wish to specify the polarization vector leaving the element in a *right-hand* coordinate system, rather than a left-hand system, we must reverse either the direction of the x axis or the direction of the y axis. We choose to reverse the direction of the x axis. Such a reversal is accounted for by a Jones matrix of the form

$$\mathbf{R} = \begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix}.$$

Thus the transmission equivalent of the reflective device has a Jones matrix of the form

$$\mathbf{L}_{\text{reflect}} = \mathbf{R}\mathbf{L}^t\mathbf{L}. \quad (\text{C-13})$$

As an example, consider a polarization device that consists of a simple coordinate rotation by angle $+\theta$, followed by reflection from a mirror. On passage through the coordinate rotation the second time, in the backwards direction, the coordinate system is once again rotated, but this time back to its original orientation. Utilizing Eq. (C-13), the Jones matrix for the entire device, expressed in a right-hand coordinate system at the output, is

$$\mathbf{L} = \begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix} \begin{bmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{bmatrix} = \begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix}.$$

Thus the only effect of passage through this simple device is a reversal of the direction of the x axis, a reversal we intentionally introduced to assure a right-handed system. Note that in this case the transpose operation was critical to obtaining the right result.

Bibliography

- [1] E. Abbe. **Beitrage zur Theorie des Mikroskops und der Mikroskopischen wahrnehmung.** *Archiv. Mikroskopische Anat.*, 9:413–468, 1873.
- [2] J.J. Amodei. Analysis of transport processes during hologram recording in crystals. *RCA Rev.*, 32:185–198, 1971.
- [3] J.J. Amodei. Electron diffusion effects during hologram recording in crystals. *Appl. Phys. Lett.*, 18:22–24, 1971.
- [4] D.Z. Anderson. Competitive and cooperative dynamics in nonlinear optical circuits. In S.F. Zornetzer, J.L. Davis, and C. Lau, editors, *An Introduction to Neural and Electronic Networks*. Academic Press, San Diego, CA, 1990.
- [5] L.K. Anderson. Holographic optical memory for bulk data storage. *Bell Lab. Record*, 46:318–325, 1968.
- [6] M. Arm, L. Lambert, and I. Weissman. Optical correlation techniques for radar pulse compression. *Proc. I.E.E.E.*, 52:842, 1964.
- [7] J.D. Armitage and A.W. Lohmann. Character recognition by incoherent spatial filtering. *Appl. Opt.*, 4:461, 1965.
- [8] E.H. Armstrong. A method for reducing disturbances in radio signalling by a system of frequency modulation. *Proc. IRE*, 24:689, 1936.
- [9] H.H. Arsenault. Distortion-invariant pattern recognition using circular harmonic matched filters. In H.H. Arsenault, T. Szoplik, and B. Macukow, editors, *Optical Processing and Computing*. Academic Press, San Diego, CA, 1989.
- [10] J.M. Artigas, M.J. Buades, and A. Filipe. Contrast sensitivity of the visual system in speckle imaging. *J. Opt. Soc. Am. A*, 11:2345–2349, 1994.
- [11] R.A. Athale and W.C. Collins. Optical matrix-matrix multiplier based on outer product decomposition. *Appl. Opt.*, 21:2089–2090, 1982.
- [12] B.F. Aull, B.E. Burke, K.B. Nichols, and W.B. Goodhue. Multiple-quantum-well CCD spatial light modulators. *Proc. S.P.I.E.*, 825:2–7, 1987.
- [13] B.B. Baker and E.T. Copson. *The Mathematical Theory of Huygen's Principle*. Clarendon Press, Oxford, second edition, 1949.
- [14] P.R. Barbier, L. Wang, and G. Moddel. Thin-film photosensor design for liquid crystal spatial light modulators. *Opt. Engin.*, 33:1322–1329, 1994.

- [15] C.W. Barnes. Object restoration in a diffraction-limited imaging system. *J. Opt. Soc. Am.*, 56:575, 1966.
- [16] G. Barton. *Elements of Green's Functions and Propagation*. Oxford University Press, New York, NY, 1989.
- [17] M. Bass, editor. *Handbook of Optics*. McGraw-Hill, Inc., New York, NY, 1995.
- [18] L. Beiser. *Holographic Scanning*. John Wiley & Sons, New York, NY, 1988.
- [19] S.A. Benton. On a method for reducing the information content of holograms. *J. Opt. Soc. Am.*, 59:1545, 1969.
- [20] M.J. Beran and G.B. Parrent, Jr. *Theory of Partial Coherence*. Prentice-Hall, Inc., Englewood Cliffs, NJ, 1964.
- [21] N.J. Berg and J.N. Lee. *Acousto-Optic Signal Processing: Theory & Applications*. Marcel Dekker, New York, NY 1983.
- [22] H.I. Bjelkhagen. *Silver-Halide Recording Materials*. Springer-Verlag, Berlin, 1993.
- [23] R.P. Bocker. Matrix multiplication using incoherent optical techniques. *Appl. Opt.*, 13:1670–1676, 1974.
- [24] G. Bonnet. Introduction to metaxial optics. I. *Ann. des Télécom.*, 33:143–165, 1978.
- [25] G. Bonnet. Introduction to metaxial optics. II. *Ann. des Télécom.*, 33:225–243, 1978.
- [26] B.L. Booth. Photopolymer material for holography. *Appl. Opt.*, 11:2994–2995, 1972.
- [27] B.L. Booth. Photopolymer material for holography. *Appl. Opt.*, 14:593–601, 1975.
- [28] M. Born and E. Wolf. *Principles of Optics*. Pergamon Press, New York, second revised edition, 1964.
- [29] C.J. Bouwkamp. Diffraction theory. In A.C. Strickland, editor, *Reports on Progress in Physics*, volume XVII. The Physical Society, London, 1954.
- [30] S.R. Bowman, W.S. Rabinovich, C.S. Kyono, D.S. Katzer, and K. Ikossi-Anastasiou. High-resolution spatial light modulators using GaAs/AlGaAs multiple quantum wells. *Appl. Phys. Lett.*, 65:956–958, 1994.
- [31] R.N. Bracewell. Two-dimensional aerial smoothing in radio astronomy. *Australia J. Phys.*, 9:297, 1956.
- [32] R.N. Bracewell. *The Fourier Transform and Its Applications*. McGraw-Hill Book Company, Inc., New York, second revised edition, 1965.
- [33] R.N. Bracewell. *Two Dimensional Imaging*. Prentice-Hall, Inc., Englewood Cliffs, NJ, 1994.
- [34] W.L. Bragg. The X-ray microscope. *Nature*, 149:470, 1942.
- [35] B. Braunecker, R. Hauck, and A.W. Lohmann. Optical character recognition based on nonredundant correlator measurements. *Appl. Opt.*, 18:2746–2753, 1979.
- [36] B. Braunecker and A.W. Lohmann. Character recognition by digital holography. *Optics Commun.*, 11:141, 1974.
- [37] E.O. Brigham. *The Fast Fourier Transform*. Prentice-Hall, Englewood Cliffs, NJ, 1974.
- [38] K. Bromley. An incoherent optical correlator. *Optica Acta*, 21:35–41, 1974.
- [39] R.E. Brooks, L.O. Heflinger, and R.F. Wuerker. Interferometry with a holographically reconstructed reference beam. *Appl. Phys. Lett.*, 7:248, 1965.
- [40] R.E. Brooks, L.O. Heflinger, and R.F. Wuerker. Pulsed laser holograms. *I.E.E.E. J. Quant. Electr.*, QE-2:275, 1966.
- [41] B.R. Brown and A.W. Lohmann. Complex spatial filter. *Appl. Opt.*, 5:967, 1966.
- [42] O. Bryngdahl. Geometrical transformations in optics. *J. Opt. Soc. Am.*, 64:1092–1099, 1974.
- [43] O. Bryngdahl and A. Lohmann. Nonlinear effects in holography. *J. Opt. Soc. Am.*, 58:1325–1334, 1968.
- [44] C.B. Burckhardt. A simplification of Lee's method of generating holograms by computer. *Appl. Opt.*, 9:1949, 1970.

- [45] G.W. Bun; F. Mok, and D. Psaltis. Large-scale holographic memory: experimental results. *Proc. S.P.I.E.*, 2026:630–641, 1993.
- [46] H.J. Butterweck. General theory of linear, coherent optical data processing systems. *J. Opt. Soc. Am.*, 67:60–70, 1977.
- [47] D. Casasent, editor. *Optical Data Processing—Applications*. Springer-Verlag, Berlin, 1978.
- [48] D. Casasent. Unified synthetic discriminant function computational formalism. *Appl. Opt.*, 23:1620–1627, 1984.
- [49] D. Casasent and W-T. Chang. Correlation synthetic discriminant functions. *Appl. Opt.*, 25:2343–2350, 1986.
- [50] D. Casasent and D. Psaltis. New optical transforms for pattern recognition. *Proc. I.E.E.E.*, 65:77–84, 1977.
- [51] S.K. Case and R. Alferness. Index modulation and spatial harmonic generation in dichromated gelatin films. *Appl. Phys.*, 10:41–51, 1976.
- [52] W.T. Cathey, B.R. Frieden, W.T. Rhodes, and C.K. Rushforth. Image gathering and processing for enhanced resolution. *J. Opt. Soc. Am. A*, 1:241–250, 1984.
- [53] H.J. Caulfield, editor. *Handbook of Holography*. Academic Press, New York, NY, 1979.
- [54] H.J. Caulfield and R. Haimes. Generalized matched filtering. *Appl. Opt.*, 19:181–183, 1980.
- [55] H.J. Caulfield, W.T. Rhodes, M.J. Foster, and S. Horvitz. Optical implementation of systolic array processing. *Optics Commun.*, 40:86–90, 1982.
- [56] M. Chang. Dichromated gelatin of improved optical quality. *Appl. Opt.*, 10:2550–2551, 1971.
- [57] F.S. Chen, J.T. LaMacchia, and D.B. Fraser. Holographic storage in lithium niobate. *Appl. Phys. Lett.*, 13:223–225, 1968.
- [58] D.C. Chu, J.R. Fienup, and J.W. Goodman. Multi-emulsion, on-axis, computer generated hologram. *Appl. Opt.*, 12:1386–1388, 1973.
- [59] I. Cindrich, editor. *Holographic Optics: Design and Application*, volume 883 of *Proceedings of the S.P.I.E.*, 1988.
- [60] I. Cindrich and S.H. Lee, editors. *Holographic Optics: Optically and Computer Generated*, volume 1052 of *Proceedings of the S.P.I.E.*, 1989.
- [61] I. Cindrich and S.H. Lee, editors. *Computer and Optically Formed Holographic Optics*, volume 1211 of *Proceedings of the S.P.I.E.*, 1990.
- [62] I. Cindrich and S.H. Lee, editors. *Computer and Optically Generated Holographic Optics*, volume 1555 of *Proceedings of the S.P.I.E.*, 1990.
- [63] I. Cindrich and S.H. Lee, editors. *Diffraction and Holographic Optics Technology*, volume 2152 of *Proceedings of the S.P.I.E.*, 1994.
- [64] N.A. Clark, M.A. Handschy, and S.T. Lagerwall. Ferroelectric liquid crystal electro-optics using the surface-stabilized structure. *Mol. Cryst. Liq. Cryst.*, 94:213–234, 1983.
- [65] N.A. Clark and S.T. Lagerwall. Submicrosecond bistable electrooptic switching in liquid crystals. *Appl. Phys. Lett.*, 36:899–901, 1980.
- [66] D.H. Close. Holographic optical elements. *Opt. Engin.*, 14:408–419, 1975.
- [67] D.H. Close, A.D. Jacobson, J.D. Margerum, R.G. Brault, and F.J. McClung. Hologram recording on photopolymer materials. *Appl. Phys. Lett.*, 14:159–160, 1969.
- [68] G. Cochran. New method of making Fresnel transforms with incoherent light. *J. Opt. Soc. Am.*, 56:1513, 1966.
- [69] W.S. Colburn and K.A. Haines. Volume hologram formation in photopolymer materials. *Appl. Opt.*, 10:1636–1641, 1971.
- [70] R.J. Collier, C.B. Burckhardt, and L.H. Lin. *Optical Holography*. Academic Press, New York, NY, 1971.

- [71] P.S. Considine. Effects of coherence on imaging systems. *J. Opt. Soc. Am.*, **56**:1001, 1966.
- [72] L. Cross. Multiplex holography. Presented at the S.P.I.E. Seminar on Three Dimensional Imaging, but unpublished, August 1977.
- [73] L.J. Cutrona et al. Optical data processing and filtering systems. *IRE Trans. Inform. Theory*, **IT-6**:386, 1960.
- [74] L.J. Cutrona et al. On the application of coherent optical processing techniques to synthetic-aperture radar. *Proc. I.E.E.E.*, **54**:1026–1032, 1966.
- [75] J.C. Dainty, editor. *Laser Speckle and Related Phenomena*. Springer-Verlag, New York, NY, second edition, 1984.
- [76] J.C. Dainty and R. Shaw. *Image Science*. Academic Press, London, 1974.
- [77] W.J. Dallas. Phase quantization—a compact derivation. *Appl. Opt.*, **10**:673–674, 1971.
- [78] W.J. Dallas. Phase quantization in holograms—a few illustrations. *Appl. Opt.*, **10**:674–676, 1971.
- [79] H. Dammann. Spectral characteristics of stepped-phase gratings. *Optik*, **53**:409–417, 1979.
- [80] J.A. Davis and J.M. Waas. Current status of the magneto-optic spatial light modulator. *Proc. S.P.I.E.*, **1150**:27–43, 1990.
- [81] D.J. De Bitteto. Holographic panoramic stereograms synthesized from white light recordings. *Appl. Opt.*, **8**:1740–1741, 1970.
- [82] Y.N. Denisyuk. Photographic reconstruction of the optical properties of an object in its own scattered radiation field. *Sov. Phys.—Dokl.*, **7**:543, 1962.
- [83] J.B. Develis and G.O. Reynolds. *Theory and Applications of Holography*. Addison-Wesley Publishing Company, Reading, MA, 1967.
- [84] A.R. Dias, R.F. Kalman, J.W. Goodman, and A.A. Sawchuk. Fiber-optic crossbar switch with broadcast capability. *Opt. Engin.*, **27**:955–960, 1988.
- [85] D.E. Dudgeon and R.M. Mersereau. *Multidimensional Digital Signal Processing*. Prentice-Hall, Inc., Englewood Cliffs, NJ, 1984.
- [86] P.M. Duffieux. *L'Intégrale de Fourier et ses Applications a l'Optique*. Faculté des Sciences, Besançon, 1946.
- [87] P.M. Duffieux. *The Fourier Transform and Its Applications to Optics*. John Wiley & Sons, New York, NY, second edition, 1983.
- [88] U. Efron, editor. *Spatial Light Modulators and Applications I*, volume **465** of *Proceedings of the S.P.I.E.*, 1984.
- [89] U. Efron, editor. *Spatial Light Modulators and Applications II*, volume **825** of *Proceedings of the S.P.I.E.*, 1988.
- [90] U. Efron, editor. *Spatial Light Modulators and Applications III*, volume **1150** of *Proceedings of the S.P.I.E.*, 1990.
- [91] U. Efron, editor. *Spatial Light Modulator Technology*. Marcel Dekker, New York, NY, 1994.
- [92] U. Efron, T.Y. Hsu, J.N. Schulman, W.Y. Wu, I. Rouse, and I.J. D'Haenena. Multiple quantum well-based spatial light modulators. *Proc. S.P.I.E.*, **825**:8–18, 1988.
- [93] H.M.A. El-Sum. *Reconstructed Wavefront Microscopy*. PhD thesis, Stanford University, Dept. of Physics, 1952.
- [94] P. Elias. Optics and communication theory. *J. Opt. Soc. Am.*, **43**:229, 1953.
- [95] P. Elias, D.S. Gray, and D.Z. Robinson. Fourier treatment of optical processes. *J. Opt. Soc. Am.*, **42**:127, 1952.
- [96] W.F. Fagan, editor. *Holographic Optical Security Systems*, volume **1509** of *Proceedings of the S.P.I.E.*, 1991.
- [97] N.H. Farhat, D. Psaltis, A. Prata, and E. Paek. Optical implementation of the Hopfield model. *Appl. Opt.*, **24**:1469–1475, 1985.

- [98] M.W. Farn. Binary optics. In R. Stern, editor, *Handbook of Photonics*. CRC Press, Boca Raton, FL, 1995.
- [99] M.W. Farn and J.W. Goodman. Diffractive doublets corrected at two wavelengths. *J. Opt. Soc. Am. A*, 8:860–867, 1991.
- [100] D. Feitelson. *Optical Computing*. MIT Press, Cambridge, MA, 1988.
- [101] H.A. Ferwerda. Frits Zernike — life and achievements. *Opt. Engin.*, 32:3176–3181, 1993.
- [102] J.P. Fitch. *Synthetic Aperture Radar*. Springer-Verlag, Berlin, 1987.
- [103] M. Françon. *Modern Applications of Physical Optics*. John Wiley & Sons, New York, NY, 1963.
- [104] G. Toraldo di Francia. Degrees of freedom of an image. *J. Opt. Soc. Am.*, 59:799–804, 1969.
- [105] A.A. Friesem and J.S. Zelenka. Effects of film nonlinearities in holography. *Appl. Opt.*, 6:1755–1759, 1967.
- [106] D. Gabor. A new microscope principle. *Nature*, 161:777, 1948.
- [107] D. Gabor. Microscopy by reconstructed wavefronts. *Proc. Roy. Soc.*, A197:454, 1949.
- [108] D. Gabor. Microscopy by reconstructed wavefronts II. *Proc. Phys. Soc.*, B64:449, 1951.
- [109] D. Gabor. Associative holographic memories. *IBM J. Res. Dev.*, 13:156–159, 1969.
- [110] D. Gabor et al. Optical image synthesis (complex addition and subtraction) by holographic Fourier transformation. *Phys. Lett.*, 18:116, 1965.
- [111] D. Gabor and W.P. Goss. Interference microscope with total wavefront reconstruction. *J. Opt. Soc. Am.*, 56:849, 1966.
- [112] J.D. Gaskill. Imaging through a randomly inhomogeneous medium by wavefront reconstruction. *J. Opt. Soc. Am.*, 58:600–608, 1968.
- [113] J.D. Gaskill. Atmospheric degradation of holographic images. *J. Opt. Soc. Am.*, 59:308–318, 1969.
- [114] T.K. Gaylord and M.G. Moharam. Analysis and applications of optical diffraction by gratings. *Proc. I.E.E.E.*, 73:894–937, 1985.
- [115] S.A. Gerasimova and V.M. Zakharchenko. Holographic processor for associative information retrieval. *Soviet J. Opt. Techn.*, 48:404–406, 1981.
- [116] R.W. Gerchberg. Super-resolution through error energy reduction. *Optica Acta*, 21:709–720, 1974.
- [117] W.W. Goj. *Synthetic-Aperture Radar & Electronic Warfare*. Artec House, New York, NY 1992.
- [118] E. Goldberg. Statistical Machine: U.S. Patent No. 1,838,389, Dec. 29, 1931.
- [119] R.C. Gonzalez and R.E. Woods. *Digital Image Processing*. Addison-Wesley Publishing Company, Reading, MA, 1992.
- [120] J.W. Goodman. Some effects of target-induced scintillation on optical radar performance. *Proc. I.E.E.E.*, 53:1688, 1965.
- [121] J.W. Goodman. Film grain noise in wavefront reconstruction imaging. *J. Opt. Soc. Am.*, 57:493–502, 1967.
- [122] J.W. Goodman. Noise in coherent optical processing. In Y.E. Nesterikhin, G.W. Stroke, and W.E. Kock, editors, *Optical Information Processing*. Plenum Press, New York, NY, 1976.
- [123] J.W. Goodman. *Statistical Optics*. John Wiley & Sons, New York, NY, 1985.
- [124] J.W. Goodman, A.R. Dias, and L.M. Woody. Fully parallel, high-speed incoherent optical method for performing discrete Fourier transforms. *Optics Letters*, 2:1–3, 1978.
- [125] J.W. Goodman, W.H. Huntley, Jr., D.W. Jackson, and M. Lehmann. Wavefront-reconstruction imaging through random media. *Appl. Phys. Lett.*, 8:311, 1966.
- [126] J.W. Goodman, D.W. Jackson, M. Lehmann, and J. Knotts. Experiments in long-distance holographic imagery. *Appl. Opt.*, 8:1581, 1969.

- [127] J.W. Goodman and G.R. Knight. Effects of film nonlinearities on wavefront-reconstruction images of diffuse objects. *J. Opt. Soc. Am.*, **58**:1276–1283, 1967.
- [128] J.W. Goodman, R.B. Miles, and R.B. Kimball. Comparative noise performance of photographic emulsions in conventional and holographic imagery. *J. Opt. Soc. Am.*, **58**:609–614, 1968.
- [129] J.W. Goodman and A.M. Silvestri. Some effects of Fourier domain phase quantization. *IBM J. Res. and Dev.*, **14**:478–484, 1970.
- [130] J.W. Goodman and L.M. Woody. Method for performing complex-valued linear operations on complex-valued data using incoherent light. *Appl. Opt.*, **16**:2611–2612, 1977.
- [131] R.M. Gray and J.W. Goodman. *Fourier Transforms: An Introduction for Engineers*. Kluwer Academic Publishers, Norwell, MA, 1995.
- [132] D.A. Gregory. Real-time pattern recognition using a modified liquid crystal television in a coherent optical correlator. *Appl. Opt.*, **25**:467–469, 1986.
- [133] J. Grinberg, A. Jacobson, W. Bleha, L. Lewis, L. Fraas, D. Boswell, and G. Myer. A new real-time non-coherent to coherent light image converter: the hybrid field effect liquid crystal light valve. *Opt. Engin.*, **14**:217–225, 1975.
- [134] R.D. Guenther. *Modern Optics*. John Wiley & Sons, New York, NY, 1990.
- [135] P.S. Guilfoyle. Systolic acousto-optic binary convolver. *Opt. Engin.*, **23**:20–25, 1984.
- [136] E.A. Guillemin. *The Mathematics of Circuit Analysis*. John Wiley & Sons, New York, NY, 1949.
- [137] P. Gunter and J.-P. Huignard, editors. *Photorefractive Materials and Their Applications I*. Springer-Verlag, Berlin, 1988.
- [138] P. Gunter and J.-P. Huignard, editors. *Photorefractive Materials and Their Applications II*. Springer-Verlag, Berlin, 1989.
- [139] P. Hariharan. *Optical Holography: Principles, Techniques & Applications*. Cambridge University Press, Cambridge, U.K., 1984.
- [140] J.L. Harris. Diffraction and resolving power. *J. Opt. Soc. Am.*, **54**:931, 1964.
- [141] J.F. Heanue, M.C. Bashaw, and L. Hesselink. Volume holographic storage and retrieval of digital data. *Science*, **265**:749–752, 1994.
- [142] R. Hecht-Nelson. *Neurocomputing*. Addison-Wesley Publishing Company, Reading, MA, 1980.
- [143] J.C. Heurtley. Scalar Rayleigh-Sommerfeld and Kirchhoff diffraction integrals: a comparison of exact evaluations for axial points. *J. Opt. Soc. Am.*, **63**:1003, 1973.
- [144] B.P. Hildebrand and K.A. Haines. Multiple-wavelength and multiple-source holography applied to contour generation. *J. Opt. Soc. Am.*, **57**:155, 1967.
- [145] H.A. Hoenl, A.W. Maue, and K. Westpfahl. Theorie der Beugung. In S. Fluegge, editor, *Handbuch der Physik*, volume 25. Springer-Verlag, Berlin, 1961.
- [146] H.H. Hopkins. *Wave Theory of Aberrations*. Oxford University Press, Oxford, 1950.
- [147] L.J. Hornbeck. Deformable-mirror spatial light modulators. *Proc. S.P.I.E.*, **1150**:86–102, 1990.
- [148] J. Horner, editor. *Optical Signal Processing*. Academic Press, Inc., Orlando, FL, 1988.
- [149] J.R. Homer and J.R. Leger. Pattern recognition with binary phase-only filters. *Appl. Opt.*, **24**:609–611, 1985.
- [150] Y-N. Hsu and H.H. Arsenault. Optical pattern recognition using circular harmonic expansion. *Appl. Opt.*, **21**:4016–4019, 1982.
- [151] Y-N. Hsu and H.H. Arsenault. Rotation invariant digital pattern recognition using circular harmonic expansion. *Appl. Opt.*, **21**:4012–4015, 1982.
- [152] A.L. Ingalls. The effects of film thickness variations on coherent light. *Phot. Sci. Eng.*, **4**:135, 1960.
- [153] P.L. Jackson. Diffractive processing of geophysical data. *Appl. Opt.*, **4**:419, 1965.

- [154] B. Javidi, J. Ruiz, and C. Ruiz. Image enhancement by nonlinear signal processing. *Appl. Opt.*, 29:4812–4818, 1990.
- [155] B.M. Javidi. Nonlinear joint transform correlators. In B. Javidi and J.L. Horner, editors, *Real-Time Optical Information Processing*. Academic Press, San Diego, CA, 1994.
- [156] J. Jahns and S.H. Lee. *Optical Computing Hardware*. Academic Press, San Diego, CA, 1993.
- [157] K.M. Johnson, D.J. McKnight, and I. Underwood. Smart spatial light modulators using liquid crystals on silicon. *I.E.E.E.J. Quant. Electr.*, 29:699–714, 1993.
- [158] T. Kailath. Channel characterization: Time varying dispersive channels. In E.J. Baghdady, editor, *Lectures on Communications System Theory*. McGraw-Hill Book Company, New York, NY, 1960.
- [159] E. Kaneko. *Liquid Crystal TV Displays: Principles and Applications of Liquid Crystal Displays*. KTK Scientific Publishers, Tokyo, Japan, 1987.
- [160] M.A. Karim and A.A.S. Awwal. *Optical Computing: An Introduction*. John Wiley & Sons, New York, NY, 1992.
- [161] J.B. Keller. Geometrical theory of diffraction. *J. Opt. Soc. Am.*, 52:116, 1962.
- [162] G. Kirchhoff. Zur Theorie der Lichtstrahlen. *Weidemann Ann. (2)*, 18:663, 1883.
- [163] M.V. Klein and T.E. Furtak. *Optics*. John Wiley & Sons, New York, NY, second edition, 1986.
- [164] G. Knight. Page-oriented associative holographic memory. *Appl. Opt.*, 13:904–912, 1974.
- [165] G. Knight. Holographic associative memory and processor. *Appl. Opt.*, 14:1088–1092, 1975.
- [166] G.R. Knight. Holographic memories. *Opt. Engin.*, 14:453–459, 1975.
- [167] C. Knox. Holographic microscopy as a technique for recording dynamic microscopic subjects. *Science*, 153:989, 1966.
- [168] H. Kogelnik. Holographic image projection through inhomogeneous media. *Bell Syst. Tech. J.*, 44:2451, 1965.
- [169] H. Kogelnik. Reconstructing response and efficiency of hologram gratings. In J. Fox, editor, *Proc. Symp. Modern Opt.*, pages 605–617. Polytechnic Press, Brooklyn, NY, 1967.
- [170] H. Kogelnik. Coupled wave theory for thick hologram gratings. *Bell Syst. Tech. J.*, 48:2909–2947, 1969.
- [171] A. Korpel. *Acousto-Optics*. Marcel Dekker, New York, NY, 1988.
- [172] F. Kottler. Electromagnetische Theorie der Beugung an schwarzen Schirmen. *Ann. Physik*, (4) 71:457, 1923.
- [173] F. Kottler. Zur Theorie der Beugung an schwarzen Schirmen. *Ann. Physik*, (4) 70:405, 1923.
- [174] F. Kottler. Diffraction at a black screen. In E. Wolf, editor, *Progress in Optics*, volume IV. North Holland Publishing Company, Amsterdam, 1965.
- [175] L.S.G. Kovanay and A. Arman. Optical autocorrelation measurement of two-dimensional random patterns. *Rev. Sci. Instr.*, 28:793, 1957.
- [176] A. Kozma. Photographic recording of spatially modulated coherent light. *J. Opt. Soc. Am.*, 56:428, 1966.
- [177] A. Kozma. Effects of film grain noise in holography. *J. Opt. Soc. Am.*, 58:436–438, 1968.
- [178] A. Kozma, G.W. Jull, and K.O. Hill. An analytical and experimental study of nonlinearities in hologram recording. *Appl. Opt.*, 9:721–731, 1970.
- [179] A. Kozma and D.L. Kelly. Spatial filtering for detection of signals submerged in noise. *Appl. Opt.*, 4:387, 1965.
- [180] A. Kozma, E.N. Leith, and N.G. Massey. Tilted-plane optical processor. *Appl. Opt.*, 11:1766–1777, 1972.

- [181] C.J. Kramer. Holographic laser scanners for nonimpact printing. *Laser Focus*, 17:70–82, 1981.
- [182] S.H. Lee, editor. *Optical Information Processing—Fundamentals*. Springer-Verlag, Berlin, 1981.
- [183] W.H. Lee. Sampled Fourier transform hologram generated by computer. *Appl. Opt.*, 9:639–643, 1970.
- [184] W.H. Lee. Binary synthetic holograms. *Appl. Opt.*, 13:1677–1682, 1974.
- [185] W.H. Lee. Computer-generated holograms: Techniques and applications. In E. Wolf, editor, *Progress in Optics*, volume 16, pages 121–232. North-Holland Publishing Company, Amsterdam, 1978.
- [186] W.H. Lee. Binary computer-generated holograms. *Appl. Opt.*, 18:3661–3669, 1979.
- [187] E.N. Leith. Photographic film as an element of a coherent optical system. *Phot. Sci. Eng.*, 6:75, 1962.
- [188] E.N. Leith and J. Upatnieks. Wavefront reconstruction and communication theory. *J. Opt. Soc. Am.*, 52:1123, 1962.
- [189] E.N. Leith and J. Upatnieks. Wavefront reconstruction with continuous-tone objects. *J. Opt. Soc. Am.*, 53:1377, 1963.
- [190] E.N. Leith and J. Upatnieks. Wavefront reconstruction with diffused illumination and three-dimensional objects. *J. Opt. Soc. Am.*, 54:1295, 1964.
- [191] E.N. Leith and J. Upatnieks. Holograms: their properties and uses. *S.P.I.E. J.*, 4:3–6, 1965.
- [192] A.L. Lentine, L.M.F. Chirovsky, L.A. D’Asaro, E.J. Laskowski, S-S. Pei, M.W. Focht, J.M. Freund, G.D. Guth, R.E. Leibenguth, R.E. Smith, and T.K. Woodward. Field-effect-transistor self-electro-optic-effect-device (FET-SEED) electrically addressed differential modulator array. *Appl. Opt.*, 33:2849–2855, 1994.
- [193] L.B. Lesem, P.M. Hirsch, and J.A. Jordan, Jr. The kinoform: a new wavefront reconstruction device. *IBM J. Res. and Dev.*, 13:150–155, 1969.
- [194] M.J. Lighthill. *Introduction to Fourier Analysis and Generalized Functions*. Cambridge University Press, New York, NY, 1960.
- [195] L.H. Lin. Hologram formation in hardened dichromated gelatin films. *Appl. Opt.*, 8:963–966, 1969.
- [196] E.H. Linfoot. *Fourier Methods in Optical Image Evaluation*. Focal Press, Ltd., London, 1964.
- [197] A.W. Lohmann. Optical single-sideband transmission applied to the Gabor microscope. *Optica Acta*, 3:97, 1956.
- [198] A.W. Lohmann. Wavefront reconstruction for incoherent objects. *J. Opt. Soc. Am.*, 55:1555, 1965.
- [199] A.W. Lohmann and D.P. Paris. Space-variant image formation. *J. Opt. Soc. Am.*, 55:1007, 1965.
- [200] A.W. Lohmann and D.P. Paris. Binary Fraunhofer holograms generated by computer. *Appl. Opt.*, 6:1739–1748, 1967.
- [201] X.J. Lu, F.T.S. Yu, and D.A. Gregory. Comparison of Vander Lugt and joint transform correlators. *Appl. Phys. B*, 51:153–164, 1990.
- [202] G.A. Maggi. Sulla propagazione libera e perturbata della onde lurninose in un mezzo isotropo. *Ann. Mathematica*, 16:21, 1888.
- [203] A.S. Marathay. *Elements of Optical Coherence Theory*. John Wiley & Sons, New York, NY, 1982.
- [204] A.D. McAulay. *Optical Computer Architectures*. John Wiley & Sons, New York, NY, 1991.
- [205] J.T. McCrickerd and N. George. Holographic stereogram from sequential component photographs. *Appl. Phys. Lett.*, 12:10–12, 1968.

- [206] D.J. McKnight, K.M. Johnson, and R.A. Serati. Electrically addressed 256 by 256 liquid-crystal-on-silicon spatial light modulator. *Optics Letters*, 18:2159–2161, 1993.
- [207] I. McNulty et al. Experimental demonstration of high-resolution three-dimensional X-ray holography. *Proc. S.P.I.E.*, 1741:78–84, 1993.
- [208] C.E.K. Mees and T.H. James, editors. *The Theory of the Photographic Process*. The Macmillan Company, New York, NY, third edition, 1966.
- [209] R.W. Meier. Magnification and third-order aberrations in holography. *J. Opt. Soc. Am.*, 55:987–992, 1965.
- [210] L. Mertz and N.O. Young. Fresnel transformations of images. In K.J. Habell, editor, *Proc. Conf. Optical Instruments and Techniques*, page 305. John Wiley and Sons, New York, NY, 1963.
- [211] D. Meyerhofer. Phase holograms in dichromated gelatin. *RCA Rev.*, 33:110–130, 1972.
- [212] D.A.B. Miller. Quantum-well self-electro-optic effect devices. *Opt. and Quant. Electr.*, 22:S61–S98, 1990.
- [213] D.A.B. Miller. Novel analog self-electrooptic effect devices. *I.E.E.E. J. Quant. Electr.*, 29:678–698, 1993.
- [214] D.A.B. Miller, D.S. Chemla, T.C. Damen, A.C. Gossard, W. Wiegmann, T.H. Wood, and C.A. Burrus. Bandedge electroabsorption in quantum well structures: the quantum confined Stark effect. *Phys. Rev. Lett.*, 53:2173–2177, 1984.
- [215] F.H. Mok and H.M. Stoll. Holographic inner-product processor for pattern recognition. *Proc. S.P.I.E.*, 1701:312–322, 1992.
- [216] R.M. Montgomery. Acousto-optical signal processing system, U.S. Patent 3,634,749, 1972.
- [217] G.M. Morris and D.L. Zweig. White light Fourier transformations. In J.L. Horner, editor, *Optical Signal Processing*. Academic Press, Orlando, FL, 1987.
- [218] M. Murdocca. *A Digital Design Methodology for Optical Computing*. MIT Press, Cambridge, MA, 1990.
- [219] M. Nazarathy and J. Shamir. Fourier optics described by operator algebra. *J. Opt. Soc. Am.*, 70:150–159, 1980.
- [220] J.A. Neff, R.A. Athale, and S.H. Lee. Two-dimensional spatial light modulators: a tutorial. *Proc. I.E.E.E.*, 78:826–855, 1990.
- [221] B.M. Oliver. Sparkling spots and random diffraction. *Proc. I.E.E.E.*, 51:220, 1963.
- [222] E.L. O’Neill. Spatial filtering in optics. *IRE Trans. Inform. Theory*, IT-256–65, 1956.
- [223] E.L. O’Neill. *Introduction to Statistical Optics*. Addison-Wesley Publishing Company, Reading, MA, 1963.
- [224] Y. Owechko and B.H. Soffer. Holographic neural networks based on multi-grating processes. In B. Javidi and J.L. Horner, editors, *Real-Time Optical Information Processing*. Academic Press, San Diego, CA, 1994.
- [225] H.M. Ozaktas and D. Mendlovic. Fractional Fourier optics. *J. Opt. Soc. Am. A*, 12:743–751, 1995.
- [226] A. Papoulis. *The Fourier Integral and Its Applications*. McGraw-Hill Book Company, New York, NY, 1962.
- [227] A. Papoulis. *Systems and Transforms with Applications to Optics*. McGraw-Hill Book Company, New York, NY, 1968.
- [228] A. Papoulis. A new algorithm in spectral analysis and band-limited extrapolation. *I.E.E.E. Trans. on Circuits and Systems*, CAS-22:735–742, 1975.
- [229] A. Papoulis. Pulse compression, fiber communications, and diffraction: a unified approach. *J. Opt. Soc. Am. A*, 11:3–13, 1994.
- [230] D.P. Peterson and D. Middleton. Sampling and reconstruction of wave-number-limited functions in n-dimensional space. *Information and Control*, 5:279, 1962.

- [231] D.K. Pollack, C.J. Koester, and J.T. Tippett, editors. *Optical Processing of Information*. Spartan Books, Baltimore, MD, 1963.
- [232] D.A. Pomet, M.G. Moharam, and E.B. Grann. Limits of scalar diffraction theory for diffractive phase elements. *J. Opt. Soc. Am. A*, 11:1827–1834, 1994.
- [233] A.B. Porter. On the diffraction theory of microscope vision. *Phil. Mag. (6)*, 11:154, 1906.
- [234] R.L. Powell and K.A. Stetson. Interferometric vibration analysis by wavefront reconstruction. *J. Opt. Soc. Am.*, 55:1593, 1965.
- [235] K. Preston, Jr. Use of the Fourier transformable properties of lenses for signal spectrum analysis. In J.T. Tippett et al., editors, *Optical and Electro-Optical Information Processing*. M.I.T. Press, Cambridge, MA, 1965.
- [236] D. Psaltis, D. Casasent, and M. Carlotto. Iterative color-multiplexed electro-optical processor. *Optics Letters*, 4:348–350, 1979.
- [237] D. Psaltis and N. Farhat. Optical information processing based on an associative-memory model of neural nets with thresholding and feedback. *Optics Letters*, 10:98–100, 1985.
- [238] D. Psaltis, X. Gu, and D. Brady. Holographic implementations of neural networks. In S.F. Zornetzer, J.L. Davis, and C. Lau, editors, *An Introduction to Neural and Electronic Networks*. Academic Press, San Diego, CA 1990.
- [239] S.I. Ragnarsson. A new holographic method of generating a high efficiency, extended range spatial filter with application to restoration of defocussed images. *Physica Scripta*, 2:145–153, 1970.
- [240] J.A. Ratcliffe. Aspects of diffraction theory and their application to the ionosphere. In A.C. Strickland, editor, *Reports on Progress in Physics*, volume XIX. The Physical Society, London, 1956.
- [241] Lord Rayleigh. On the theory of optical images, with special references to the microscope. *Phil. Mag. (5)*, 42:167, 1896.
- [242] J.D. Redman, W.P. Wolton, and E. Shuttleworth. Use of holography to make truly three-dimensional X-ray images. *Nature*, 220:58–60, 1968.
- [243] J. Rhodes. Analysis and synthesis of optical images. *Am. J. Phys.*, 21:337, 1953.
- [244] W.T. Rhodes. Two pupil synthesis of optical transfer functions. *Appl. Opt.*, 17:1141–1151, 1978.
- [245] G.L. Rogers. Gabor diffraction microscope: the hologram as a generalized zone plate. *Nature*, 166:237, 1950.
- [246] G.L. Rogers. *Noncoherent Optical Processing*. John Wiley & Sons, New York, NY, 1977.
- [247] W.E. Ross. Advanced magneto-optic spatial light modulator device development update. *Proc. S.P.I.E.*, 1704:222–229, 1992.
- [248] W.E. Ross, D. Psaltis, and R.H. Anderson. Two-dimensional magneto-optic spatial light modulator. *Opt. Engin.*, 22:485–489, 1983.
- [249] A. Rubinowicz. Die Beugungswelle in der Kirchhoffschen Theorie der Beugungserscheinungen. *Ann. Physik*, 53:257, 1917.
- [250] A. Rubinowicz. The Miyamoto-Wolf diffraction wave. In E. Wolf, editor, *Progress in Optics*, volume IV. North Holland Publishing Company, Amsterdam, 1965.
- [251] D.E. Rumelhart, G.E. Hinton, and R.J. Williams. Learning internal representations by error propagation. In D.E. Rumelhart and J.L. McClelland, editors, *Parallel Distributed Processing*. MIT Press, Cambridge, MA, 1986.
- [252] C.K. Rushforth and R.W. Harris. Restoration, resolution and noise. *J. Opt. Soc. Am.*, 58:539–545, 1968.
- [253] B.E.A. Saleh and M.C. Teich. *Fundamentals of Photonics*. John Wiley & Sons, New York, NY, 1991.
- [254] G. Saxby. *Practical Holography*. Prentice-Hall, Englewood Cliffs, NJ, 1988.

- [255] O.H. Schade. Electro-optical characteristics of television systems. *RCA Review*, IX:5 (Part I), 245 (Part 11), 490 (Part 111), 653 (Part IV), 1948.
- [256] W. Schumann. *Holography and Deformation Analysis*. Springer-Verlag, Berlin, 1985.
- [257] P.J. Sementilli, B.R. Hunt, and M.S. Nadar. Analysis of the limit to superresolution in incoherent imaging. *J. Opt. Soc. Am. A*, 10:2265–2276, 1993.
- [258] T.A. Shankoff. Phase holograms in dichromated gelatin. *Appl. Opt.*, 7:2101–2105, 1968.
- [259] C.E. Shannon. Communication in the presence of noise. *Proc. IRE*, 37:10, 1949.
- [260] G.C. Sherman. Application of the convolution theorem to Rayleigh's integral formulas. *J. Opt. Soc. Am.*, 57:546, 1967.
- [261] A.E. Siegman. *Lasers*. University Science Books, Mill Valley, CA, 1986.
- [262] S. Silver. Microwave aperture antennas and diffraction theory. *J. Opt. Soc. Am.*, 52:131, 1962.
- [263] T.J. Skinner. Surface texture effects in coherent imaging. *J. Opt. Soc. Am.*, 53:1350A, 1963.
- [264] L. Slobodin. Optical correlation techniques. *Proc. I.E.E.E.*, 51:1782, 1963.
- [265] H.M. Smith. *Principles of Holography*. John Wiley & Sons, New York, NY, second edition, 1975.
- [266] H.M. Smith, editor. *Holographic Recording Materials*. Springer-Verlag, Berlin, 1977.
- [267] L. Solymar and D.J. Cook. *Volume Holography and Volume Gratings*. Academic Press, New York, NY, 1981.
- [268] A. Sommerfeld. Mathematische Theorie der Diffraction. *Math. Ann.*, 47:317, 1896.
- [269] A. Sommerfeld. Die Greensche Funktion der Schwingungsgleichung. *Jahresber. Deut. Math. Vex*, 21:309, 1912.
- [270] A. Sommerfeld. *Optics*, volume IV of *Lectures on Theoretical Physics*. Academic Press, New York, NY, 1954.
- [271] W.H. Southwell. Validity of the Fresnel approximation in the near field. *J. Opt. Soc. Am.*, 71:7, 1981.
- [272] R.A. Sprague and C.L. Koliopoulos. Time integrating acousto-optic correlator. *Appl. Opt.*, 15:89–92, 1975.
- [273] H. Stark. *Image Recovery: Theory and Application*. Academic Press, Orlando, FL, 1987.
- [274] T. Stone and N. George. Hybrid diffractive-refractive lenses and achromats. *Appl. Opt.*, 27:2960–2971, 1988.
- [275] G.W. Stroke. Image **deblurring** and aperture synthesis using a posteriori processing by Fourier-transform holography. *Optica Acta*, 16:401–422, 1969.
- [276] G.W. Stroke and A.T. Funkhouser. Fourier-transform spectroscopy using holographic imaging without computing and with stationary interferometers. *Phys. Lett.*, 16:272, 1965.
- [277] G.W. Stroke and R.C. Restruck III. Holography with spatially incoherent light. *Appl. Phys. Lett.*, 7:229, 1965.
- [278] K.J. Strozewski, C-Y. Wang, Jr., G.C. Wetsel, R.M. Boysel, and J.M. Florence. Characterization of a micromechanical spatial light modulator. *J. Appl. Phys.*, 73:7125–7128, 1993.
- [279] M.R. Taghizadeh and J. Turunen. Synthetic diffractive elements for optical interconnection. *Opt. Comp. and Proc.*, 2:221–242, 1992.
- [280] H.F. Talbot. *Philos. Mag.*, 9:401, 1836.
- [281] B.J. Thompson, J.H. Ward, and W.R. Zinky. Application of hologram techniques for particle size analysis. *Appl. Opt.*, 6:519, 1967.
- [282] D.A. Tichenor and J.W. Goodman. Coherent transfer function. *J. Opt. Soc. Am.*, 62:293–295, 1972.

- [283] D.A. Tichenor and J.W. Goodman. Restored impulse response of finite-range image deblurring filter. *Appl. Optics*, 14:1059–1060, 1975.
- [284] J.T. Tippett et al., editors. *Optical and Electro-Optical Information Processing*. MIT Press, Cambridge, MA, 1965.
- [285] A. Tonomura. *Electron Holography*. Springer-Verlag, Berlin, 1994.
- [286] G.L. Turin. An introduction to matched filters. *IRE Trans. Info. Theory*, IT-6:311, 1960.
- [287] J. Upatnieks, A. Vander Lugt, and E. Leith. Correction of lens aberrations by means of holograms. *Appl. Opt.*, 5:589, 1966.
- [288] R.F. van Ligten. Influence of photographic film on wavefront reconstruction. I: Plane wavefronts. *J. Opt. Soc. Am.*, 56:1, 1966.
- [289] R.F. van Ligten. Influence of photographic film on wavefront reconstruction. II: 'Cylindrical' wavefronts. *J. Opt. Soc. Am.*, 56:1009–1014, 1966.
- [290] A.B. Vander Lugt. Signal detection by complex spatial filtering. Technical report, Institute of Science and Technology, University of Michigan, Ann Arbor, MI, 1963.
- [291] A.B. Vander Lugt. Signal detection by complex spatial filtering. *I.E.E.E. Trans. Info. Theory*, IT-10:139–145, 1964.
- [292] A.B. Vander Lugt. Operational notation for analysis and synthesis of optical data processing systems. *Proc. I.E.E.E.*, 54:1055, 1966.
- [293] A.B. VanderLugt. *Optical Signal Processing*. John Wiley & Sons, New York, NY, 1992.
- [294] C.M. Vest. *Holographic Interferometry*. John Wiley & Sons, New York, NY, 1979.
- [295] C.J. Weaver and J.W. Goodman. A technique for optically convolving two functions. *Appl. Opt.*, 5:1248–1249, 1966.
- [296] W.T. Welford. *Aberrations of the Symmetrical Optical System*. Academic Press, New York, NY, 1974.
- [297] B.S. Wherrett and F.A.P. Tooley, editors. *Optical Computing*. Edinburgh University Press, Edinburgh, U.K., 1989.
- [298] E.T. Whittaker. On the functions which are represented by the expansions of the interpolation theory. *Proc. Roy. Soc. Edinburgh, Sect. A*, 35:181, 1915.
- [299] B. Widrow and S.D. Stearns. *Adaptive Signal Processing*. Prentice-Hall, Englewood Cliffs, NJ, 1985.
- [300] C.S. Williams and O.A. Becklund. *Introduction to the Optical Transfer Function*. John Wiley & Sons, New York, NY, 1989.
- [301] E. Wolf and E.W. Marchand. Comparison of the Kirchhoff and Rayleigh-Sommerfeld theories of diffraction at an aperture. *J. Opt. Soc. Am.*, 54:587, 1964.
- [302] S. Wolfram. *Mathematica*. Addison-Wesley Publishing Company, Reading, MA, second edition, 1991.
- [303] T.H. Wood. Multiple quantum well (MQW) waveguide modulators. *J. Lightwave Techn.*, 6:743, 1988.
- [304] L.P. Yaroslavskii and N.S. Merzlyakov. *Methods of Digital Holography*. Consultants Bureau, Plenum Publishing Company, New York, NY, 1980.
- [305] F. Zernike. Das Phasenkontrastverfahren bei der Mikroskopischen beobachtung. *Z. Tech. Phys.*, 16:454, 1935.

INDEX

- Abbe, E., 126, 129, 218
Abbe-Porter experiments, 218–220
Abbe theory of image formation, 129, 218
ABCD matrix (see ray-transfer matrix)
Aberrations, 100, 128, 145–151
 and amplitude transfer function,
 effects on, 146
 in holography, 363
 and optical transfer function, effects on,
 146–151
Acoustic holography, 319
Acousto-optic effect, 206
Acousto-optic signal processing, 276–281
Acousto-optic spatial light modulators, 205–209
Adjacency effect, 183
Agfa-Gevaert Scientia series, 346
Airy, G.B., 77
Airy pattern, 77
Alignment direction, for liquid crystal
 molecule, 186
Alignment layer, for liquid crystal cell, 188
Amplitude impulse response, 133
Amplitude mapping, by a photographic
 emulsion, 178–180
Amplitude transfer function, 135–137
 for hologram, 363
Amplitude transmittance:
 definition, 59
 vs. exposure curve, 179–180
 of photographic emulsion, 178–180
Analytic continuation, 161
Anamorphic processor, 234
Angle mismatch, for thick hologram, 339
Angle multiplexing, in holography, 383
Angular spectrum
 of plane waves, 55–61
 propagation of, 57–58
Anisotropy, of liquid crystals, 190
Antenna designer's formula, 74
Anti-reference wave, 310, 324, 325, 379
Apodization, 151–154
 Gaussian, 152
 inverse, 154
Arago, F., 35
ASA speed, of a photographic emulsion, 313n
Aspheric wavefront, 211
Autocorrelation theorem, 9
 proof, 397
Axial magnification, in holography, 317

Axicon (see conical lens)
Azimuthal image, in synthetic aperture
 radar, 267

Backpropagation algorithm, 386
Bandlimited function, 23
Bandwidth extrapolation, 160–165
 intuitive explanation, 161–162
 iterative method, 164–165
 practical limitations, 165
 and sampling theorem, 162–164
Barium titanate, 348
Benton, S., 324
Besinc function, 16
Bessel function, 11, 12, 297
Binary intensity modulator, 193
Binary optics, 210–214
 fabrication, 213–214
Binomial expansion, 66
Bipolar data, processing with incoherent systems,
 287–289
Birefringence, of liquid crystals, 190
Bismuth germanium oxide, 348
Bismuth silicon oxide, 348
Bistability, in ferroelectric liquid crystal, 189
Bleaching, of photographic emulsion,
 183–184
Blur point-spread function, 264
Bocker, R.P., 283
Booth, B.L., 348
Boundary conditions, 37
Boundary diffraction, 54–55
Bracewell, **R.N.**, 5, 26
Bragg, W.L., 295
Bragg angle, 207
Bragg cell spectrum analyzer, 276–278
Bragg condition, 332
Bragg degeneracy, 332, 387
Bragg diffraction, 330
Bragg effect, 207
Bragg regime, of acoustic diffraction,
 206, 208, 209
Braunecker, B., 256
Brooks, R.E., 373
Brown, B.R., 355
Bryngdahl, O., 367
Burckhardt, C.B., 295, 360
Butterweck, H.J., 114

434 Introduction to Fourier Optics

- Cadmium sulfide, 194
- Cadmium telluride, 195
- Carrier frequency, 241
- Casasent, D., 253, 256
- Caulfield, H.J., 256, 287, 295
- Central dark ground imaging, 220, 290
- Characteristic impedance, 64
- Character recognition, with coherent optics, 246–257
- Charge-coupled detector array, 281
- Charge transport, in photorefractive materials, 349, 350
- Chemical bleaching, 183–184
- Chirp function, 17*n*, 30, 267
- Cholesteric liquid crystal, 186
- Chromatic blur, in holography, 321
- Circle function:
 - definition, 14, 15
 - Fourier transform of, 15, 16
- Circuit theory, approximations of, 38
- Circular harmonic correlation, 254–256
- Circular harmonic decomposition, 254
- Circularly polarized light, 416
- Circular symmetry, 11
- CMOS (see complementary metal oxide silicon)
- Cochran, G., 369
- Coherence, 4
- Coherence length, 313
- Coherent imaging, vs. incoherent imaging, 154–160
- Coherent optical processing architectures, 232–236
- Colburn, W.S., 348
- Collier, R.J., 295
- Comb function, definition, 13
- Common path interferometry, 382
- Compensating filter, 222
- Compensating plate, 380
- Complementary metal oxide silicon (CMOS), 202
- Complex data, processing with incoherent systems, 289
- Complex exponential functions, 7, 22
- Computer-generated hologram, 351–363
 - computational problem of, 354–355
 - devices for plotting, 351
 - representational problem of, 355–363
 - sampling problem of, 352–354
- Confocal spherical surfaces, diffraction
 - between, 72
- Conical lens, 121, 271
- Conjugate planes, 407–408
- Contour generation, with holography, 375–376
- Contrast ratio, of liquid crystal cell, 196
- Contrast reversal, 89, 147
- Converging lens (see positive lens)
- Conversion gain, in holography, 369
- Convolution:
 - with incoherent processing, 226–229
 - with joint transform correlator, 245
 - with VanderLugt filter, 241
- Convolution integral, 22
- Convolution theorem, 9
 - proof, 396
- Coordinate rotation, Jones matrix
 - representation, 191
- Cornu's spiral, 85*n*
- Corpuscular theory, 33
- Coupled mode theory, 336–346
- Coupling constant, for a thick hologram, 338
- Cross, L., 326
- Crossbar switch, based on matrix-vector multiplier, 286
- Crosscorrelation:
 - with joint transform correlator, 245
 - with VanderLugt filter, 241
- Cross-linking:
 - of gelatin molecules, 184
 - of photopolymer molecules, 347
- Cutrona, L.J., 224
- Cylindrical lens, 120
- Data storage, holographic, 382–384
- Deformable mirror, spatial light modulators, 200–202
 - and cantilever beam, 201
 - and membrane, 200
 - and torsion beam, 201
- Degradation of holographic images, 363–369
- Delta functions, 7, 393–395
- Denisyuk, Y.N., 296, 321
- Density, photographic (see photographic density, definition)
- Depth distortion, in holography, 318
- Detour phase, 355–360
- Detuning parameter, for thick hologram, 338
- Develis, J.B., 295
- Developer, 174
- Development speck, 174
- Diamond turning, 214
- Dichromate, 348
- Dichromated gelatin, 348
- Diffraction rays, 55
- Diffraction, 32
 - by a conducting half-plane, 55
- Diffraction efficiency, 330
 - of binary optic element, 211, 214

- Diffraction efficiency (continued)
 - definition, 81
 - of general periodic grating, 92
 - of sawtooth grating, 211
 - of sinusoidal amplitude grating, 81
 - of sinusoidal phase grating, 82–83
 - of thick holograms, 336–346
 - and volume gratings, 346
- Diffraction-limited system, definition, 128
- Diffraction optical elements, 209–214
- Diffuse object, 367
- Diffuser, 160
- Diffusion, in a photographic emulsion, 181
- Direction cosine, 56
- Discrete analog optical processors, 282–289
- Discrete Fourier transform, 355
- Distorting media, imaging through, 378–382
- Diverging lens (see negative lens)
- Doppler history, 267
- Doppler shifts:
 - from acoustic grating, 207
 - in synthetic aperture radar, 267
- Driffield, V.C., 175
- Dudgeon, D.E., 26
- Duffieux, P.M., 126
- Dynamic range:
 - of inverse filter, 258, 261
 - of Wiener filter, 263
- Edge images, in coherent and incoherent light, 158–159
- Eigenfunction:
 - Bessel function as, 29
 - complex exponentials as, 22
 - definition, 22
- Eigenvalue, 22
- Eikonal equation, 402
- Eikonal function, 402
- Electric dipole, in a liquid crystal molecule, 188
- Electroforming process, 329
- Electron holography, 295
- Electronic charge, 64
- Electro-optic effect, in photorefractive materials, 349
- Elementary function, 20, 22
- Elias, P., 223
- El-Sum, H.M.A., 296, 372
- Embossed hologram, 328–329
- Embossing, 210
- Emulsion, photographic, 173
- Energy spectrum, 103
- Entrance pupil, 127, 128, 411–413
- Equivalent area, 28
- Equivalent bandwidth, 28
- Erasure, in photorefractive materials, 350, 387
- Evanescent waves, 58
- Exciton, 202, 203
- Exit pupil, 127, 128, 411–413
- Exposure, definition, 174
- Extraordinary refractive index, 190
- Fabry-Perot *étalon*, 204
- False images, in holography, 367
- Faraday effect, 198
 - rotation angle, 200
- Faraday rotation, 215, 216
- Far field, 74
- Fast Fourier transform, 355
- Ferroelectric liquid crystal, 186–187, 189, 190–194
 - and spatial light modulator, 197–198
- FET-SEED, 205
- Film grain noise, effects in holography, 368–369
- Film MTF, effects in holography, 363–367
- Film **nonlinearities**, effects in holography, 367–368
- Filter realization, constraints of, 236–237
- Finite thickness, effects with volume grating, 333–336
- Fixing:
 - of photographic image, 174
 - for photorefractive materials, 350
- F-number, of a lens, 170
- Focal length, 99
- Focal plane, 101, 104, 408–409
- Focal properties, of synthetic aperture radar data, 268–271
- Focusing error, 148–151
- Fourier-Bessel transform, 12, 28
- Fourier coefficient, 356
- Fourier hologram, computer-generated, 352–353
- Fourier integral, 5
 - three-dimensional, 335
- Fourier integral theorem, 9
 - proof, 397–399
- Fourier transform hologram, 320, 365
- Fourier transform:
 - as decomposition, 7
 - definition, 5
 - existence conditions of, 5
 - generalized, 6
 - inverse, 5
 - with lens, 104–106
 - optical, 101–107
 - optical, of array of one-dimensional functions, 122
 - optical, exact, 104
 - optical, example of, 107

- Fourier transform (continued)
 optical, geometries, 101–107
 optical, input against lens, 102–104
 optical, input behind lens, 106–107
 optical, input in front of lens,
 optical, location of transform plane, 119
 optical, vignetting effect, 105–106
 pairs, table of, 14
 theorems, 8, 395–399
 two-dimensional, 7
- Fraunhofer approximation, 63, 73–75
- Fraunhofer diffraction, 74
 by circular aperture, 77–78
 distance required for, 74
 by rectangular aperture, 75–77
 by sinusoidal amplitude grating, 78–81
 by sinusoidal phase grating, 81–83
- Fraunhofer diffraction pattern, obtained with
 lens, 103
- Fraunhofer hologram, 319, 321
- Frequency spectrum:
 of coherent image intensity, 155
 of incoherent image intensity, 155
- Fresnel, **A.J.**, 34
- Fresnel approximation, 63, 66–67
 accuracy of, 69–71
- Fresnel diffraction:
 between confocal spherical surfaces, 72
 by square aperture, 84–87
- Fresnel diffraction integral, 67
 as convolution, 67
 as Fourier transform, 67
- Fresnel hologram, 319, 321
 computer-generated, 353–354
- Fresnel integral, 18, 69, 84
- Fresnel-Kirchhoff diffraction formula, 45–46
- Fresnel number, 85
- Fresnel zone plate, 124
- Friesem, **A.A.**, 367
- Fringe orientation, in thick holograms, 332–334
- Fringe period, in holography, 321
- Gabor**, D., 295, 368, 372, 373, 384
- Gabor** hologram, 302–304
- Gamma, photographic, 176
- Gaussian reference sphere, 145–146
- Geometrical optics, 19, 32
- Geometrical theory of diffraction, 55
- Gerchberg, R.W., 164
- Goldberg, E., 225
- Goodman, J.W., 5, 111, 243, 263, 368
- Grain size, in photographic emulsions, 368
- Grating vector, 331
- Grating vector cloud, 335
- Grating vector spectrum, 335
- Gray, D.S., 223
- Gray, R.M., 5
- Green's function, 35, 40–41, 47–49
- Green's theorem, 39
- Gross fog, 176
- Guilfoyle, P., 287
- Haines, R., 256
- Haines, K.A., 348, 375
- Halftone process, 223
- Half-wave plate, 418
- Hankel** transform, 10–12
- Hard-clipped filter, 237n
- Hariharan, P., 295
- H&D** curve (see Hurter-Driffield curve)
- Helmholtz equation, 38–39, 41, 57
 paraxial, 62
- Heurtley, **J.C.**, 51
- Hidden layer, in a neural network, 386
- High-contrast film, 176
- High-definition television, 202
- High-gamma film (see high-contrast film)
- High-order images, in holography, 367
- Hildebrand, B.P., 375
- HOE (see holographic optical element)
- Holographic art, 388
- Holographic data storage, 382–384
- Holographic display, 388
- Holographic interferometry, 373–378
 contour generation, 375–376
 multiple exposure, 373–375
 real-time, 375
- Holographic memory (see holographic data
 storage)
- Holographic optical element (HOE), 387–388
- Holographic stereograms, 322–324
- Holographic weights, for neural networks,
 386–387
- Holography:
 applications of, 372–388
 history of, 295–296
 with incoherent light, 369–371
 linearity of imaging process, 299–300
 practical problems in, 313–314
- Homogeneous medium, 36
- Hopfield** neural network, 286, 387
- Hopkins, H.H., 126, 130n
- Hughes liquid crystal light valve, 194–196
- Hurter, F., 175
- Hurter-Driffield curve, 175
- Huygens, C., 33
- Huygens envelope construction, 34
- Huygens-Fresnel principle, 35, 52–53, 65

- Huygens-Fresnel principle (*continued*)
 as convolution integral, 53
- Hybrid-field-effect mode, of liquid crystal cell, **195**
- Hydrogenated amorphous silicon, 197
- Hyperplane decision surface, 385
- Ideal image, definition, 130
- Ilford, Ltd., holographic emulsions from, 347
- Image amplification, with SLM, 185
- Image casting, 225
- Image deblurring, 222–223
- Image formation:
 by holography, 301–302
 in monochromatic light, 108–114
 with polychromatic illumination, **130–134**
- Image hologram, 319
- Image location, in holography, 314–317
- Image magnification, in holography, 317–319
- Image restoration, fundamentals, 257–260
- Imaging system, generalized model, 127–128
- Impedance, of liquid crystal cell, 195
- Impulse response, 20
 of imaging system, 112
 of positive lens, 108–112
- Incoherent imaging:
 vs. coherent imaging, 154–160
 conditions for, 134
- Incoherent optical processing, 224–231
 and geometrical optics, 225–229
 limitations of, 229
- Incoherent-to-coherent converter, with SLM, **185**
- Integral theorem, of Helmholtz and Kirchhoff, **40–42**
- Intensity:
 definition, 63–65
 instantaneous, 65
- Intensity impulse response, 134
- Intensity mapping, by photographic emulsion, 177–178
- Intensity modulator, and liquid crystals, 192–193
- Intensity transmittance, definition, 175
- Interference gain, 160
- Interferometry:
 with holography, 373–378
 holography as, 297
- Intermodulation effects, in holography, 367
- Invariant pattern recognition, optical approaches, 252–257
- Invariant system, 21
- Inverse filter, 258–259
 realization of, 261
- Inverse Fourier transform, 5, 12
- Inverse **Hankel** transform, 12
- Isoplanatic patch, 21
- Isoplanatic system, 21
- Isotropic medium, 36
- Jinc function, 16
- Joint transform correlator, 243–246
- Jones calculus, 190
- Jones matrix, definition, 190, **415–416**
- Keller, J.**, 55
- Kelley, D.H., 181
- Kelley model, 181
- Kinoform, 360–361
- Kirchhoff, G., 35
- Kirchhoff boundary conditions, **44–45, 49**
- Kirchhoff diffraction theory, 42–46
- Knox, C., 372
- Kodak 649F spectroscopic plate, 183, 313, 346
- Kogelnik, H., 336
- Koliopoulos, C.L., 279
- Kottler, F., 35
- Kozma, A., 179, 367, 368
- k vector diagram (*see* wave vector diagram)
- Laser ablation, 214
- Latent image, 174
- Least-mean-square-error filter (*see* Wiener filter)
- Lee, W.H., 360, 361
- Leith, E.N., **237n**, 296, 313
- Leith-Upatnieks hologram, **304–314**
 with diffused illumination, 313
 minimum reference angle, 308–309
 obtaining images from, **306–307**
 recording, 305–306
 for three-dimensional scenes, 309–312
- Lens, thin, 96–101
- Lens law, **110**
- Lensless** Fourier transform hologram, 320, 365
- Light emitting diode, 283
- Lighthill, M.J., 6
- Light-mod (*see* magneto-optic spatial light modulator)
- Lin, L.H., 295
- Linear system, 4
 definition, 7
- Linearity, 4
- Linearity theorem, 9
 proof, 395
- Line-spread function, 166
- Lippmann, G., 296
- Liquid crystals, 185–198
 electrical properties, 188–190

- Liquid crystals (continued)
 mechanical properties, 186–188
 optical properties, 190
 Liquid gate, 178–179
 Lithium niobate, 348
 LMS algorithm, 386
 Local spatial frequency, 17
 relation to ray optics, 402–403
 Lohmann, **A.W.**, 256, 296, 355, 367, 369
 Low contrast film, 176
- Mach-Zehnder interferometer, 239, 240
 Maggi, **G.A.**, 55
 Magneto-optic spatial light modulator
 (MOSLM), 198–200
 Magnification, definition, 112
 Marchand, E.W., 51
Maréchal, A., 179, 222, 236
 Matched filters, 246–251
 bank of, 250
 Matrix-matrix multiplier (see outer product
 processor)
 Matrix-vector multiplier:
 parallel, 284–286
 serial, 283–284
 Matrix-vector product, 282–283
 Maxwell, J.C., 35
 Maxwell's equations, 36
 Meier, **R.W.**, 363
Mellin transform, 252
 insensitivity to scale size, 253
 Meridional ray, 404
 Mersereau, R.M., 26
 Mertz, L., 369
 Metal master hologram, 329
 Micromachining, in fabrication of binary
 optics, 210
 Microscopy, with holography, 372
 Microwave holography, 315, 317–319
 Middleton, D., 26
 Minimum reference angle, for Leith-Upatnieks
 hologram, 308–309
 Modulation transfer function (MTF), 139
 definition, 182
 measurement, 183
 of photographic emulsion, 180–183
 Molecular beam epitaxy, 202, 204
 Montgomery, R.M., 279
 MOSLM (see magneto-optic spatial light
 modulator)
 MTF (see modulation transfer function)
 Multilayer neural network, 385–386
 Multiple exposure holographic interferometry,
 373–375
- Multiple quantum well structure, 203
 Multiplex hologram, 326–328
 Multiplexing, in holography, 383
 Museum of Holography, 388
 Mutual intensity, 133
 Mylar-base film, 173
- Narrowband light, 131
 Nazarathy, M., 114
 Near field, 67
 Negative lens, 99–101
 Nematic liquid crystal, 186–187, 190–194
 Networks of neurons, 385–386
 Neural network, 384
 Neuron, 384–385
 Newton, I., 34
 Nonlocalization of data, in Fourier
 hologram, 383
 Nonmagnetic medium, 36
 Nonmonochromatic wave, 53
 Nontanning bleach, 183–184, 262
 Numerical aperture, 157
 Nyquist sampling, 281
- Obliquity factor, 51
 Offset reference hologram (see Leith-Upatnieks
 hologram)
 O'Neill, E.L., 223
 Operator notation, 114–120
 applications of, 116–120
 operator relations, table of, 117
 Optical transfer function, 138–145
 of aberration-free system, 140
 and apodization, effects of, 151–154
 as autocorrelation function, 139
 examples of, 142–144, 148–150
 and geometrical calculation, 141
 Ordinary refractive index, 190
 Orthoscopic image, 311
 Outer product processor, 286–287
- Palermo, C., 237n
Papoulis, A., 5, 164
Paraxial approximation, 403
Paraxial diffraction, 72, 73
Paraxial geometrical optics, 401–413
 Paris, D.P., 355
 Parseval's theorem, 9
 proof, 396
 Partial coherence, 131
 Penumbra effect, 33
 Permeability, 36

- Permittivity, 36, 37
 Peterson, D.P., 26
 Phase-coded reference beams, 383
 Phase contour interferogram, 361–363
 Phase-contrast microscope, 220
 Phase matching, 360
 Phase modulation, by photographic emulsion, 183–184
 Phase-only filter, 237n
 Phase shift, by photographic emulsion, 178
 Phase-shifting plate, 222
 Phasor, 39
 time-varying, 131
 Photo-elastic effect, 206
 Photographic density, definition, 175
 Photographic film, 173–184
 Photolithography, in fabrication of binary optics, 210
 Photopolymer films, 347–349
 Photorefractive crystal, 386
 Photorefractive effect, 349
 Photorefractive materials, 348–351
 Photoresist, for recording hologram, 329
 Pinhole camera, 168
Planck's constant, 64
 Plane wave, 56
 Plus-X film, 183
 Point-spread function, 20
 Poisson, S., 34
 Poisson's spot, 35
 Polarization analyzer, 418
 Polarization rotator, 191–193, 199, 417
 Polarization transformations, 417–418
 Polarization vector, 190, 415
 Polymerization, 347
 Porter, A.B., 218
 Positive lens, 99–101
 Potassium tantalum niobate, 348
 Potential function, 46
 Powell, R.L., 376
 Power spectral densities, 259
 Power spectrum, 103
 Principal plane, 409–411
 Prism, 120
Psaltis, D., 253
 Pseudoscopic image, 311
 Pupil function, 102
 generalized, 145
- Q factor:
 for acousto-optic diffraction, 208
 for holograms, 329–330
 Q parameter (see Q factor)
 Quadratic phase dispersion, 72
 Quadratic phase factors, in imaging equation, 109–112
 Quality criterion, 154
 Quantum-confined Stark effect, 203
 Quantum confinement, 203
 Quantum efficiency, of detector, 64
 Quantum well, 202
 Quarter-wave plate, 418
 Quasi-monochromatic light, 127
- Ragnarsson, S.I.**, 262, 263
 Rainbow hologram, 324–326
Raman-Nath regime, 206, 208, 209, 276
 Ratcliffe, J.A., 55
Ray, 401–402
 Ray directions, 19
 Rayleigh, Lord, 126, 129
 Rayleigh interferometer, 239, 240
 Rayleigh resolution criterion, 157
 Rayleigh resolution distance, 157
 Rayleigh-Sommerfeld diffraction formula, 49–50
 Rayleigh-Sommerfeld diffraction theory, 46–50
 Rayleigh's theorem, 9, 396
 Ray-transfer matrix, 404–407
 Real image, 108
 in holography, 301
 Real-time holographic interferometry, 375
 Reciprocal polarization element, 419
 Reciprocity theorem, 46
 Reconstruction wave, 297
 Recording materials, for holography, 346–351
 Rectangle function, 13
 Rectangular sampling lattice, 23, 26
 Reduced coordinates, 130
 Referenceless on-axis complex hologram (ROACH), 355, 360–361
 Reference wave, 297
 in **VanderLugt** filter, 239
 Reflection, 33
 Reflection hologram, 296, 321–322, 333
 Reflection intensity modulator, 193
 Reflective polarization devices, 418–420
 Refraction, 32
 Refractive index, 32, 36
 Registration errors, in holographic memories, 383
 Relief image, photographic, 184
 Replication of holograms, 328–329
 Resistor-biased SEED (R-SEED), 204
 Resolution, beyond classical limit, 160–165
 Responsivity, of detector, 64
 Restrick, R.C., 369
 Retardation, Jones matrix representation of, 191
 Retarded time, 54

- Reynolds, G.O., 295
 Rhodes, W.T., 230
 ROACH (see referenceless on-axis complex hologram)
 Robinson, D.Z., 223
 Rogers, G.L., 296, 369
 Rotation sensitivity, of matched filter, 251
 R-SEED (see resistor-biased SEED)
 Rubinowicz, A., 55
- Saltus** problem, 35
 Sampled function, 25
 Sampling theory, 22–27
 SAR (see synthetic aperture radar)
 Sawtooth grating, 211
Saxby, G., 295
 Scalar diffraction theory, 35, 36–38
 limitations of, 214
 Scale size sensitivity, of matched filter, 251
 Scaling of hologram, 317, 319
 Scattering, in photographic emulsion, 181
 Schade, O., 126
 Schlieren method, in microscopy, 220
 Schumann, W., 373
Schwarz's inequality, 140, 147
 Secondary source, 46
 Secondary wavelets, 34
 Security applications, of holography, 388
 SEED (see self-electro-optic effect device)
 Self-electro-optic effect, 204–205
 Self-electro-optic effect device (SEED), 204–205
 Self images, 87–90
 Separable functions, 10
 Shamir, J., 114
 Shannon, C., 23
 Sherman, G., 61
 Shift theorem, 8
 proof, 395–396
 Shoulder, of **H&D** curve, 176
 Side lobes, 152
 Side-looking radar, 159
 Sifting property, of delta function, 20, 393
 Sight-mod (see magneto-optic spatial light modulator)
 Sigmoid nonlinearity, 384
 Sign convention, for radii of curvature, 97
 Signum function, 13
 Silver halide, 173
 Similarity theorem, 8
 proof, 395
Sinc function, 13
 SLM (see spatial light modulator)
 Small angle diffraction, 72
 Small signal suppression, in holography, 367
Smectic-A phase, of liquid crystal, 188*n*
 Smectic-C* phase, of liquid crystal, 186,
 187, 188*n*
 Smectic liquid crystal, 186
 Smith, H.M., 295
Snell's law, 32, 403
 Sommerfeld, A., 33, 35
 Sommerfeld radiation condition, 44
 Space-bandwidth product, 27, 229
 Space-integrating correlator, 278–279
 Space invariance, 21
 Space-variant impulse response, 113
 Sparrow resolution criterion, 170
 Spatial coherence, 133
 Spatial frequency, 5
 local, 17
 Spatial light modulator (SLM), 101, 184–209
 acousto-optic, 205–209
 and liquid crystals, 194–198
 and multiple quantum wells, 202–205
 Spatially coherent illumination, 131, 133
 Spatially incoherent illumination, 131, 134
 Speckle, effects in holography, 369
 Speckle effect, 159–160, 224
 Speckle size, 159
 Sprague, R.A., 279
 S-SEED (see symmetric SEED)
 Stanford matrix-vector multiplier, 284–286
 Stationary phase, principle of, 71
 Stereo effect, 323, 328
Stetson, K.A., 376
 Strehl definition, 168
 Stroke, G.W., 369
 Strontium barium nitrate, 348
 Superposition, 20
 Superposition integral, 20, 52
 for imaging system, 129
 Superposition property, 20
 Super-resolution, 160–165
 Symmetric SEED (S-SEED), 204
 Synthetic aperture radar (SAR), 264–275
 data processing for, 268–275
 data recording format, 265–268
 image from, 275
 Synthetic discriminant functions, 256–257
 System, 19
 Systolic processor, 287
- Talbot image, 87–90
 Talbot subimage, 89
 Tanning bleach, 183–184
 Telescopic system, imaging properties of, 273

- Television display, and liquid crystals, 197
 Thick amplitude reflection grating, 344
 Thick amplitude transmission grating, 341
 Thick hologram, 329–346
 Thick phase reflection grating, 343
 Thick phase transmission grating, 339
 Thin lens, 96–101
 definition, 96
 as phase transformation, 99–101
 Thompson, B.J., 372
 Three-dimensional imaging, with holography, 309–312
 Three-dimensional interference pattern, 330
 Three-dimensional optical storage, 383
 Tichenor, D.A., 111, 263, 264
 Tilt angle, of ferroelectric liquid crystal, 192
 Tilted plane processor, 272–275
 Time constant, for liquid crystals, 190
 Time-bandwidth product, 281
 Time-integrating correlator, 279–281
 Time invariance, 21
 Toe, of H&D curve, 176
 Training of neural network, 385
 Transfer function, 22
 of wave propagation, 59–61
 Transmission hologram, 321–322
 Transverse magnification, 273
 in holography, 317
 Triangle function, 13
 Triangular interferometer, 370
 Twin image problem, 303–304
 Twisted nematic liquid crystal, 186, 191
 Two-point resolution, 157
 Two-pupil OTF synthesis, 230–231
- Ultrasound imaging, 159
 Ultraviolet holography, 319
 Upatnieks, J., 296, 313, 381
- VanderLugt, A.B., 114, 237
 VanderLugt filter, 237–243
 advantages of, 242–243
 impulse response of, 250
 and synthesis of frequency-plane mask, 238–240
- Van Ligten, R.F., 366
 Vectorial diffraction theory, 35, 36–38
 Velocity of propagation, 37
 Vertical parallax, 328
 Vest, C.M., 373
 Vibration analysis, with holography, 376–378
 Virtual image, 108
 in holography, 301
 Volume grating (**see** thick hologram)
 Volume imagery, with holography, 372
- Ward, J.H., 372
 Wave equation, scalar, 37, 39
 Wave equation, vector, 36
 Wavefront, 402
 Wavefront reconstruction, 295, **296–302**
 Wavelength conversion, with SLM, 185
 Wavelength mismatch, for thick hologram, 339
 Wavelength multiplexing, in holography, 383
 Wave number, 39
 Wave retarder, 417–418
 Wave vector diagram, 208
 Weaver, C.S., 243
 White light illumination, in holography, 32
 Whittaker, E.T., 23
 Whittaker-Shannon sampling theorem, 23–26
 Wiener filter, 259–260
 optical realization, 261–264
 Wigner distribution function, 30
 Windowing, 151
 Wolf, E., 51
- X-ray crystallography, 295
 X-ray holography, 319, 372
- Young, N.O., 369
 Young, T., 34, 54
- Zelenka, J.S., 367
 Zernike, F., 2, 220
 Zernike phase contrast microscope, 220
 Zero-spread nonlinearity, 181
 Zinky, W.R., 372