# CONTENTS

---

`T2`  Track Two; see page xv

## 10   Nonlinear Optics   513

# Geometric Optics

Solar rays parallel to $OB$ and passing through this solid are refracted at the hyperbolic surface,
and the refracted rays converge at $A$.

IBN SAHL (984)

## 7.1 Overview

Geometric optics, the study of "rays," is the oldest approach to optics. It is an accurate
description of wave propagation when the wavelengths and periods of the waves are
far smaller than the lengthscales and timescales on which the wave amplitude and the
medium supporting the waves vary.

After reviewing wave propagation in a homogeneous medium (Sec. 7.2), we begin
our study of geometric optics in Sec. 7.3. There we derive the geometric-optics prop-
agation equations with the aid of the eikonal approximation, and we elucidate the
connection to Hamilton-Jacobi theory (which we assume the reader has already en-
countered). This connection is made more explicit by demonstrating that a classical,
geometric-optics wave can be interpreted as a flux of quanta. In Sec. 7.4, we special-
ize the geometric-optics formalism to any situation where a bundle of nearly parallel
rays is being guided and manipulated by some sort of apparatus. This is the parax-
ial approximation, and we illustrate it with a magnetically focused beam of charged
particles and show how matrix methods can be used to describe the particle (i.e., ray)
trajectories. In Sec. 7.5, we explore how imperfect optics can produce multiple images
of a distant source, and that as one moves from one location to another, the images
appear and disappear in pairs. Locations where this happens are called "caustics" and
are governed by catastrophe theory, a topic we explore briefly. In Sec. 7.6, we describe
gravitational lenses, remarkable astronomical phenomena that illustrate the forma-
tion of multiple images and caustics. Finally, in Sec. 7.7, we turn from scalar waves to
the vector waves of electromagnetic radiation. We deduce the geometric-optics prop-
agation law for the waves' polarization vector and explore the classical version of a
phenomenon called geometric phase.

> ## BOX 7.1.   READERS' GUIDE
>
> - This chapter does not depend substantially on any previous chapter, but it does assume familiarity with classical mechanics, quantum mechanics, and classical electromagnetism.
> - Secs. 7.1–7.4 are foundations for the remaining optics chapters, 8, 9, and 10.
> - The discussion of caustics in Sec. 7.5 is a foundation for Sec. 8.6 on diffraction at a caustic.
> - Secs. 7.2 and 7.3 (monochromatic plane waves and wave packets in a homogeneous, time-independent medium; the dispersion relation; and the geometric-optics equations) are used extensively in subsequent parts of this book, including
>     - Chap. 12 for elastodynamic waves,
>     - Chap. 16 for waves in fluids,
>     - Sec. 19.7 and Chaps. 21–23 for waves in plasmas, and
>     - Chap. 27 for gravitational waves.
>     - Sec. 28.6.2 for weak gravitational lensing.

**7.2**

## 7.2  Waves in a Homogeneous Medium

**7.2.1**

### 7.2.1  Monochromatic Plane Waves; Dispersion Relation

Consider a monochromatic plane wave propagating through a homogeneous medium. Independently of the physical nature of the wave, it can be described mathematically by

$$\psi = Ae^{i(\mathbf{k}\cdot\mathbf{x}-\omega t)} \equiv Ae^{i\varphi},\tag{7.1}$$

**plane wave: complex amplitude, phase, angular frequency, wave vector, wavelength, and propagation direction**

where $\psi$ is any oscillatory physical quantity associated with the wave, for example, the $y$ component of the magnetic field associated with an electromagnetic wave. If, as is usually the case, the physical quantity is real (not complex), then we must take the real part of Eq. (7.1). In Eq. (7.1), $A$ is the wave's *complex amplitude*; $\varphi = \mathbf{k}\cdot\mathbf{x} - \omega t$ is the wave's *phase*; $t$ and $\mathbf{x}$ are time and location in space; $\omega = 2\pi f$ is the wave's *angular frequency*; and $\mathbf{k}$ is its *wave vector* (with $k \equiv |\mathbf{k}|$ its *wave number*, $\lambda = 2\pi/k$ its *wavelength*, $\lambdabar = \lambda/(2\pi)$ its *reduced wavelength*, and $\hat{\mathbf{k}} \equiv \mathbf{k}/k$ its *propagation direction*). Surfaces of constant phase $\varphi$ are orthogonal to the propagation direction $\hat{\mathbf{k}}$ and move in the $\hat{\mathbf{k}}$ direction with the *phase velocity*

**phase velocity**

$$\mathbf{V}_{\mathrm{ph}} \equiv \left(\frac{\partial \mathbf{x}}{\partial t}\right)_{\varphi} = \frac{\omega}{k}\hat{\mathbf{k}}\tag{7.2}$$

(cf. Fig. 7.1). The frequency $\omega$ is determined by the wave vector $\mathbf{k}$ in a manner that depends on the wave's physical nature; the functional relationship

**FIGURE 7.1** A monochromatic plane wave in a homogeneous medium.

$$\omega = \Omega(\mathbf{k})$$

(7.3)

**dispersion relation**

is called the wave's *dispersion relation,* because (as we shall see in Ex. 7.2) it governs the dispersion (spreading) of a wave packet that is constructed by superposing plane waves.

Some examples of plane waves that we study in this book are:

**examples:**

1. Electromagnetic waves propagating through an isotropic dielectric medium with index of refraction $\mathfrak{n}$ [Eq. 10.20)], for which $\psi$ could be any Cartesian component of the electric or magnetic field or vector potential and the dispersion relation is

**electromagnetic waves**

$$\omega = \Omega(\mathbf{k}) = Ck \equiv C|\mathbf{k}|,$$

(7.4)

with $C = c/\mathfrak{n}$ the phase speed and $c$ the speed of light in vacuum.

2. Sound waves propagating through a solid (Sec. 12.2.3) or fluid (liquid or vapor; Secs. 7.3.1 and 16.5), for which $\psi$ could be the pressure or density perturbation produced by the sound wave (or it could be a potential whose gradient is the velocity perturbation), and the dispersion relation is the same as for electromagnetic waves, Eq. (7.4), but with $C$ now the sound speed.

**sound waves**

3. Waves on the surface of a deep body of water (depth $\gg \lambda$; Sec. 16.2.1), for which $\psi$ could be the height of the water above equilibrium, and the dispersion relation is [Eq. (16.9)]:

**water waves**

$$\omega = \Omega(\mathbf{k}) = \sqrt{gk} = \sqrt{g|\mathbf{k}|},$$

(7.5)

with $g$ the acceleration of gravity.

4. Flexural waves on a stiff beam or rod (Sec. 12.3.4), for which $\psi$ could be the transverse displacement of the beam from equilibrium, and the dispersion relation is

**flexural waves**

$$\omega = \Omega(\mathbf{k}) = \sqrt{\frac{D}{\Lambda}k^2} = \sqrt{\frac{D}{\Lambda}\mathbf{k}\cdot\mathbf{k}},$$

(7.6)

with $\Lambda$ the rod's mass per unit length and $D$ its "flexural rigidity" [Eq. (12.33)].

**Alfvén waves**

5. Alfvén waves in a magnetized, nonrelativistic plasma (bending waves of the plasma-laden magnetic field lines; Sec. 19.7.2), for which $\psi$ could be the transverse displacement of the field and plasma, and the dispersion relation is [Eq. (19.75)]

$$\omega = \Omega(\mathbf{k}) = \mathbf{a} \cdot \mathbf{k}, \tag{7.7}$$

with $\mathbf{a} = \mathbf{B}/\sqrt{\mu_o \rho}$, $[= \mathbf{B}/\sqrt{4\pi\rho}]^1$ the Alfvén speed, $\mathbf{B}$ the (homogeneous) magnetic field, $\mu_o$ the magnetic permittivity of the vacuum, and $\rho$ the plasma mass density.

**gravitational waves**

6. Gravitational waves propagating across the universe, for which $\psi$ can be a component of the waves' metric perturbation which describes the waves' stretching and squeezing of space; these waves propagate nondispersively at the speed of light, so their dispersion relation is Eq. (7.4) with $C$ replaced by the vacuum light speed $c$.

In general, one can derive the dispersion relation $\omega = \Omega(\mathbf{k})$ by inserting the plane-wave ansatz (7.1) into the dynamical equations that govern one's physical system [e.g., Maxwell's equations, the equations of elastodynamics (Chap. 12), or the equations for a magnetized plasma (Part VI)]. We shall do so time and again in this book.

7.2.2

### 7.2.2 Wave Packets

Waves in the real world are not precisely monochromatic and planar. Instead, they occupy wave packets that are somewhat localized in space and time. Such wave packets can be constructed as superpositions of plane waves:

**wave packet**

$$\psi(\mathbf{x}, t) = \int A(\mathbf{k}) e^{i\alpha(\mathbf{k})} e^{i(\mathbf{k}\cdot\mathbf{x} - \omega t)} \frac{d^3 k}{(2\pi)^3}, \tag{7.8a}$$
where $A(\mathbf{k})$ is concentrated around some $\mathbf{k} = \mathbf{k}_o$.

Here $A$ and $\alpha$ (both real) are the modulus and phase of the complex amplitude $A e^{i\alpha}$, and the integration element is $d^3 k \equiv dV_k \equiv dk_x dk_y dk_z$ in terms of components of $\mathbf{k}$ along Cartesian axes $x$, $y$, and $z$. In the integral (7.8a), the contributions from adjacent $\mathbf{k}$s will tend to cancel each other except in that region of space and time where the oscillatory phase factor changes little with changing $\mathbf{k}$ (when $\mathbf{k}$ is near $\mathbf{k}_o$). This is the spacetime region in which the wave packet is concentrated, and its center is where $\mathbf{V_k}$(phase factor) $= 0$:

$$\left( \frac{\partial \alpha}{\partial k_j} + \frac{\partial}{\partial k_j}(\mathbf{k} \cdot \mathbf{x} - \omega t) \right)_{\mathbf{k}=\mathbf{k}_o} = 0. \tag{7.8b}$$

---

1. Gaussian unit equivalents will be given with square brackets.

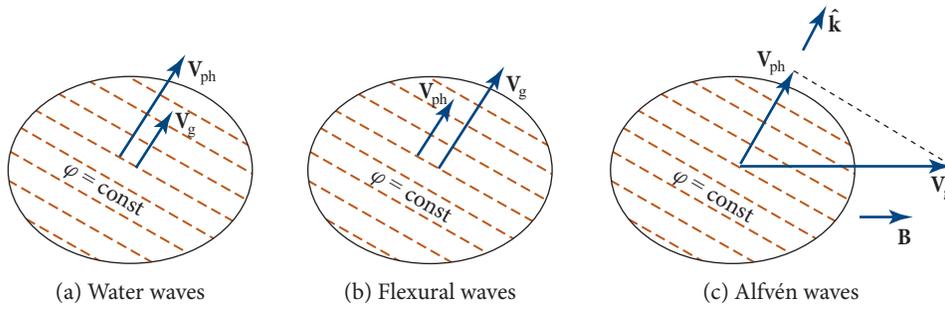(a) Water waves        (b) Flexural waves        (c) Alfvén waves

**FIGURE 7.2** (a) A wave packet of waves on a deep body of water. The packet is localized in the spatial region bounded by the ellipse. The packet's (ellipse's) center moves with the group velocity $\mathbf{V}_g$. The ellipse expands slowly due to wave-packet dispersion (spreading; Ex. 7.2). The surfaces of constant phase (the wave's oscillations) move twice as fast as the ellipse and in the same direction, $\mathbf{V}_{ph} = 2\mathbf{V}_g$ [Eq. (7.11)]. This means that the wave's oscillations arise at the back of the packet and move forward through the packet, disappearing at the front. The wavelength of these oscillations is $\lambda = 2\pi/k_o$, where $k_o = |\mathbf{k}_o|$ is the wave number about which the wave packet is concentrated [Eq. (7.8a) and associated discussion]. (b) A flexural wave packet on a beam, for which $\mathbf{V}_{ph} = \frac{1}{2}\mathbf{V}_g$ [Eq. (7.12)], so the wave's oscillations arise at the packet's front and, traveling more slowly than the packet, disappear at its back. (c) An Alfvén wave packet. Its center moves with a group velocity $\mathbf{V}_g$ that points along the direction of the background magnetic field [Eq. (7.13)], and its surfaces of constant phase (the wave's oscillations) move with a phase velocity $\mathbf{V}_{ph}$ that can be in any direction $\hat{\mathbf{k}}$. The phase speed is the projection of the group velocity onto the phase propagation direction, $|\mathbf{V}_{ph}| = \mathbf{V}_g \cdot \hat{\mathbf{k}}$ [Eq. (7.13)], which implies that the wave's oscillations remain fixed inside the packet as the packet moves; their pattern inside the ellipse does not change. (An even more striking example is provided by the Rossby wave, discussed in Sec. 16.4, in which the group velocity is equal and oppositely directed to the phase velocity.)

Evaluating the derivative with the aid of the wave's dispersion relation $\omega = \Omega(\mathbf{k})$, we obtain for the location of the wave packet's center

$$x_j - \left(\frac{\partial\Omega}{\partial k_j}\right)_{\mathbf{k}=\mathbf{k}_o} t = -\left(\frac{\partial\alpha}{\partial k_j}\right)_{\mathbf{k}=\mathbf{k}_o} = \text{const.} \tag{7.8c}$$

This tells us that the *wave packet* moves with the *group velocity*

$$\mathbf{V}_g = \nabla_{\mathbf{k}}\Omega, \quad \text{i.e.,} \quad V_{g\,j} = \left(\frac{\partial\Omega}{\partial k_j}\right)_{\mathbf{k}=\mathbf{k}_o}. \tag{7.9}$$

group velocity

When, as for electromagnetic waves in a dielectric medium or sound waves in a solid or fluid, the dispersion relation has the simple form of Eq. (7.4), $\omega = \Omega(\mathbf{k}) = Ck$ with $k \equiv |\mathbf{k}|$, then the group and phase velocities are the same,

$$\mathbf{V}_g = \mathbf{V}_{ph} = C\hat{\mathbf{k}}, \tag{7.10}$$

and the waves are said to be *dispersionless*. If the dispersion relation has any other form, then the group and phase velocities are different, and the wave is said to exhibit

*dispersion;* cf. Ex. 7.2. Examples are (see Fig. 7.2 and the list in Sec. 7.2.1, from which our numbering is taken):

3. Waves on a deep body of water [dispersion relation (7.5); Fig. 7.2a], for which

$$\mathbf{V}_g = \frac{1}{2}\mathbf{V}_{ph} = \frac{1}{2}\sqrt{\frac{g}{k}}\,\hat{\mathbf{k}};\tag{7.11}$$

4. Flexural waves on a stiff beam or rod [dispersion relation (7.6); Fig. 7.2b], for which

$$\mathbf{V}_g = 2\mathbf{V}_{ph} = 2\sqrt{\frac{D}{\Lambda}}\,k\hat{\mathbf{k}};\tag{7.12}$$

5. Alfvén waves in a magnetized and nonrelativistic plasma [dispersion relation (7.7); Fig. 7.2c], for which

$$\mathbf{V}_g = \mathbf{a}, \qquad \mathbf{V}_{ph} = (\mathbf{a}\cdot\hat{\mathbf{k}})\hat{\mathbf{k}}.\tag{7.13}$$

Notice that, depending on the dispersion relation, the group speed $|\mathbf{V}_g|$ can be less than or greater than the phase speed, and if the homogeneous medium is anisotropic (e.g., for a magnetized plasma), the group velocity can point in a different direction than the phase velocity.

Physically, it should be obvious that the energy contained in a wave packet must remain always with the packet and cannot move into the region outside the packet where the wave amplitude vanishes. Correspondingly, the wave packet's energy must propagate with the group velocity $\mathbf{V}_g$ and not with the phase velocity $\mathbf{V}_{ph}$. When one examines the wave packet from a quantum mechanical viewpoint, its quanta must move with the group velocity $\mathbf{V}_g$. Since we have required that the wave packet have its wave vectors concentrated around $\mathbf{k}_o$, the energy and momentum of each of the packet's quanta are given by the standard quantum mechanical relations:

$$\boxed{\mathcal{E} = \hbar\Omega(\mathbf{k}_o), \quad \text{and} \quad \mathbf{p} = \hbar\mathbf{k}_o.}\tag{7.14}$$

## EXERCISES

**Exercise 7.1** *Practice: Group and Phase Velocities*
Derive the group and phase velocities (7.10)–(7.13) from the dispersion relations (7.4)–(7.7).

**Exercise 7.2** *\*\*Example: Gaussian Wave Packet and Its Dispersion*
Consider a 1-dimensional wave packet, $\psi(x,t) = \int A(k)e^{i\alpha(k)}e^{i(kx-\omega t)}dk/(2\pi)$, with dispersion relation $\omega = \Omega(k)$. For concreteness, let $A(k)$ be a narrow Gaussian peaked around $k_o$: $A \propto \exp[-\kappa^2/(2(\Delta k)^2)]$, where $\kappa = k - k_o$.

(a) Expand $\alpha$ as $\alpha(k) = \alpha_o - x_o\kappa$ with $x_o$ a constant, and assume for simplicity that higher order terms are negligible. Similarly, expand $\omega \equiv \Omega(k)$ to quadratic order,

and explain why the coefficients are related to the group velocity $V_g$ at $k = k_o$ by $\Omega = \omega_o + V_g\kappa + (dV_g/dk)\kappa^2/2$.

(b) Show that the wave packet is given by

$$\psi \propto \exp[i(\alpha_o + k_o x - \omega_o t)] \int_{-\infty}^{+\infty} \exp[i\kappa(x - x_o - V_g t)] \tag{7.15a}$$

$$\times \exp\left[-\frac{\kappa^2}{2}\left(\frac{1}{(\Delta k)^2} + i\frac{dV_g}{dk}t\right)\right] d\kappa.$$

The term in front of the integral describes the phase evolution of the waves inside the packet; cf. Fig. 7.2.

(c) Evaluate the integral analytically (with the help of a computer, if you wish). From your answer, show that the modulus of $\psi$ satisfies

$$|\psi| \propto \frac{1}{L^{1/2}} \exp\left[-\frac{(x - x_o - V_g t)^2}{2L^2}\right], \quad \text{where } L = \frac{1}{\Delta k}\sqrt{1 + \left(\frac{dV_g}{dk}(\Delta k)^2 \, t\right)^2} \tag{7.15b}$$

is the packet's half-width.

(d) Discuss the relationship of this result at time $t = 0$ to the uncertainty principle for the localization of the packet's quanta.

(e) Equation (7.15b) shows that the wave packet spreads (i.e., disperses) due to its containing a range of group velocities [Eq. (7.11)]. How long does it take for the packet to enlarge by a factor 2? For what range of initial half-widths can a water wave on the ocean spread by less than a factor 2 while traveling from Hawaii to California?

## 7.3 Waves in an Inhomogeneous, Time-Varying Medium: The Eikonal Approximation and Geometric Optics

**7.3**

Suppose that the medium in which the waves propagate is spatially inhomogeneous and varies with time. If the lengthscale $\mathcal{L}$ and timescale $\mathcal{T}$ for substantial variations are long compared to the waves' reduced wavelength and period,

$$\mathcal{L} \gg \lambda = 1/k, \qquad \mathcal{T} \gg 1/\omega, \tag{7.16}$$

then the waves can be regarded locally as planar and monochromatic. The medium's inhomogeneities and time variations may produce variations in the wave vector $\mathbf{k}$ and frequency $\omega$, but those variations should be substantial only on scales $\gtrsim \mathcal{L} \gg 1/k$ and $\gtrsim \mathcal{T} \gg 1/\omega$. This intuitively obvious fact can be proved rigorously using a two-lengthscale expansion (i.e., an expansion of the wave equation in powers of $\lambda/\mathcal{L} = 1/k\mathcal{L}$ and $1/\omega\mathcal{T}$). Such an expansion, in this context of wave propagation, is called

**two-lengthscale expansion**

**eikonal approximation**

**WKB approximation**

the *geometric-optics approximation* or the *eikonal approximation* (after the Greek word $\epsilon\iota\kappa\omega\nu$, meaning image). When the waves are those of elementary quantum mechanics, it is called the *WKB approximation.*[2] The eikonal approximation converts the laws of wave propagation into a remarkably simple form, in which the waves' amplitude is transported along trajectories in spacetime called *rays.* In the language of quantum mechanics, these rays are the world lines of the wave's quanta (photons for light, phonons for sound, plasmons for Alfvén waves, and gravitons for gravitational waves), and the law by which the wave amplitude is transported along the rays is one that conserves quanta. These ray-based propagation laws are called the laws of *geometric optics.*

In this section we develop and study the eikonal approximation and its resulting laws of geometric optics. We begin in Sec. 7.3.1 with a full development of the eikonal approximation and its geometric-optics consequences for a prototypical dispersion-free wave equation that represents, for example, sound waves in a weakly inhomogeneous fluid. In Sec. 7.3.3, we extend our analysis to cover all other types of waves. In Sec. 7.3.4 and a number of exercises we explore examples of geometric-optics waves, and in Sec. 7.3.5 we discuss conditions under which the eikonal approximation breaks down and some non-geometric-optics phenomena that result from the breakdown. Finally, in Sec. 7.3.6 we return to nondispersive light and sound waves, deduce Fermat's principle, and explore some of its consequences.

**7.3.1**

### 7.3.1 Geometric Optics for a Prototypical Wave Equation

Our prototypical wave equation is

**prototypical wave equation in a slowly varying medium**

$$\frac{\partial}{\partial t}\left(W\frac{\partial\psi}{\partial t}\right) - \nabla \cdot (WC^2\nabla\psi) = 0. \tag{7.17}$$

Here $\psi(\mathbf{x}, t)$ is the quantity that oscillates (the *wave field*), $C(\mathbf{x}, t)$ will turn out to be the wave's slowly varying *propagation speed,* and $W(\mathbf{x}, t)$ is a slowly varying *weighting function* that depends on the properties of the medium through which the wave propagates. As we shall see, $W$ has no influence on the wave's dispersion relation or on its geometric-optics rays, but it does influence the law of transport for the waves' amplitude.

The wave equation (7.17) describes sound waves propagating through a static, isentropic, inhomogeneous fluid (Ex. 16.13), in which case $\psi$ is the wave's pressure perturbation $\delta P$, $C(\mathbf{x}) = \sqrt{(\partial P/\partial\rho)_s}$ is the adiabatic sound speed, and the weighting function is $W(\mathbf{x}) = 1/(\rho C^2)$, with $\rho$ the fluid's unperturbed density. This wave equation also describes waves on the surface of a lake or pond or the ocean, in the limit that the slowly varying depth of the undisturbed water $h_o(\mathbf{x})$ is small compared

---

2. Sometimes called "JWKB," adding Jeffreys to the attribution, though Carlini, Liouville, and Green used it a century earlier.

to the wavelength (shallow-water waves; e.g., tsunamis); see Ex. 16.3. In this case $\psi$ is the perturbation of the water's depth, $W = 1$, and $C = \sqrt{gh_o}$ with $g$ the acceleration of gravity. In both cases—sound waves in a fluid and shallow-water waves—if we turn on a slow time dependence in the unperturbed fluid, then additional terms enter the wave equation (7.17). For pedagogical simplicity we leave those terms out, but in the analysis below we do allow $W$ and $C$ to be slowly varying in time, as well as in space: $W = W(\mathbf{x}, t)$ and $C = C(\mathbf{x}, t)$.

Associated with the wave equation (7.17) are an energy density $U(\mathbf{x}, t)$ and energy flux $\mathbf{F}(\mathbf{x}, t)$ given by

$$U = W \left[ \frac{1}{2} \left( \frac{\partial \psi}{\partial t} \right)^2 + \frac{1}{2} C^2 (\nabla \psi)^2 \right], \qquad \mathbf{F} = -W C^2 \frac{\partial \psi}{\partial t} \nabla \psi; \qquad (7.18)$$

**energy density and flux**

see Ex. 7.4. It is straightforward to verify that, if $C$ and $W$ are independent of time $t$, then the scalar wave equation (7.17) guarantees that the $U$ and $\mathbf{F}$ of Eq. (7.18) satisfy the law of energy conservation:

$$\frac{\partial U}{\partial t} + \nabla \cdot \mathbf{F} = 0; \qquad (7.19)$$

cf. Ex. 7.4.[3]

We now specialize to a weakly inhomogeneous and slowly time-varying fluid and to nearly plane waves, and we seek a solution of the wave equation (7.17) that locally has approximately the plane-wave form $\psi \simeq A e^{i\mathbf{k}\cdot\mathbf{x}-\omega t}$. Motivated by this plane-wave form, (i) we express the waves in the eikonal approximation as the product of a real amplitude $A(\mathbf{x}, t)$ that varies slowly on the length- and timescales $\mathcal{L}$ and $\mathcal{T}$, and the exponential of a complex phase $\varphi(\mathbf{x}, t)$ that varies rapidly on the timescale $1/\omega$ and lengthscale $\lambdabar$:

**eikonal approximated wave: amplitude, phase, wave vector, and angular frequency**

$$\psi(\mathbf{x}, t) = A(\mathbf{x}, t) e^{i\varphi(\mathbf{x},t)}; \qquad (7.20)$$

and (ii) we define the wave vector (field) and angular frequency (field) by

$$\mathbf{k}(\mathbf{x}, t) \equiv \nabla \varphi, \qquad \omega(\mathbf{x}, t) \equiv -\partial\varphi/\partial t. \qquad (7.21)$$

In addition to our two-lengthscale requirement, $\mathcal{L} \gg 1/k$ and $\mathcal{T} \gg 1/\omega$, we also require that $A$, $\mathbf{k}$, and $\omega$ vary slowly (i.e., vary on lengthscales $\mathcal{R}$ and timescales $\mathcal{T}'$ long compared to $\lambdabar = 1/k$ and $1/\omega$).[4] This requirement guarantees that the waves are locally planar, $\varphi \simeq \mathbf{k} \cdot x - \omega t + \text{constant}$.

---

3. Alternatively, one can observe that a stationary medium will not perform work.
4. Note that these variations can arise both (i) from the influence of the medium's inhomogeneity (which puts limits $\mathcal{R} \lesssim \mathcal{L}$ and $\mathcal{T}' \lesssim \mathcal{T}$ on the wave's variations) and (ii) from the chosen form of the wave. For example, the wave might be traveling outward from a source and so have nearly spherical phase fronts with radii of curvature $r \simeq$ (distance from source); then $\mathcal{R} = \min(r, \mathcal{L})$.

## BOX 7.2. BOOKKEEPING PARAMETER IN TWO-LENGTHSCALE EXPANSIONS

When developing a two-lengthscale expansion, it is sometimes helpful to introduce a bookkeeping parameter $\sigma$ and rewrite the ansatz (7.20) in a fleshed-out form:

$$\psi = (A + \sigma B + \ldots)e^{i\varphi/\sigma}. \tag{1}$$

The numerical value of $\sigma$ is unity, so it can be dropped when the analysis is finished. We use $\sigma$ to tell us how various terms scale when $\lambdabar$ is reduced at fixed $\mathcal{L}$ and $\mathcal{R}$. The amplitude $A$ has no attached $\sigma$ and so scales as $\lambdabar^0$, $B$ is multiplied by $\sigma$ and so scales proportional to $\lambdabar$, and $\varphi$ is multiplied by $\sigma^{-1}$ and so scales as $\lambdabar^{-1}$. When one uses these factors of $\sigma$ in the evaluation of the wave equation, the first term on the right-hand side of Eq. (7.22) gets multiplied by $\sigma^{-2}$, the second term by $\sigma^{-1}$, and the omitted terms by $\sigma^0$. These factors of $\sigma$ help us to quickly group together all terms that scale in a similar manner and to identify which of the groupings is leading order, and which subleading, in the two-lengthscale expansion. In Eq. (7.22) the omitted $\sigma^0$ terms are the first ones in which $B$ appears; they produce a propagation law for $B$, which can be regarded as a post-geometric-optics correction.

Occasionally the wave equation itself will contain terms that scale with $\lambdabar$ differently from one another (e.g., Ex. 7.9). One should always look out for this possibility.

We now insert the eikonal-approximated wave field (7.20) into the wave equation (7.17), perform the differentiations with the aid of Eqs. (7.21), and collect terms in a manner dictated by a two-lengthscale expansion (see Box 7.2):

$$0 = \frac{\partial}{\partial t}\left(W\frac{\partial\psi}{\partial t}\right) - \boldsymbol{\nabla}\cdot(WC^2\boldsymbol{\nabla}\psi) \tag{7.22}$$

$$= \left(-\omega^2 + C^2 k^2\right)W\psi$$

$$+ i\left[-2\left(\omega\frac{\partial A}{\partial t} + C^2 k_j A_{,j}\right)W - \frac{\partial(W\omega)}{\partial t}A - (WC^2 k_j)_{,j}A\right]e^{i\varphi} + \cdots.$$

The first term on the right-hand side, $(-\omega^2 + C^2 k^2)W\psi$, scales as $\lambdabar^{-2}$ when we make the reduced wavelength $\lambdabar$ shorter and shorter while holding the macroscopic lengthscales $\mathcal{L}$ and $\mathcal{R}$ fixed; the second term (in square brackets) scales as $\lambdabar^{-1}$; and the omitted terms scale as $\lambdabar^0$. This is what we mean by "collecting terms in a manner dictated by a two-lengthscale expansion." Because of their different scaling, the first,

second, and omitted terms must vanish separately; they cannot possibly cancel one another.

The vanishing of the first term in the eikonal-approximated wave equation (7.22) implies that the waves' frequency field $\omega(\mathbf{x}, t) \equiv -\partial\varphi/\partial t$ and wave-vector field $\mathbf{k} \equiv \nabla\varphi$ satisfy the dispersionless dispersion relation,

$$\omega = \Omega(\mathbf{k}, \mathbf{x}, t) \equiv C(\mathbf{x}, t)k, \tag{7.23}$$

**dispersion relation**

where (as throughout this chapter) $k \equiv |\mathbf{k}|$. Notice that, as promised, this dispersion relation is independent of the weighting function $W$ in the wave equation. Notice further that this dispersion relation is identical to that for a precisely plane wave in a homogeneous medium, Eq. (7.4), except that the propagation speed $C$ is now a slowly varying function of space and time. This will always be so.

*One can always deduce the geometric-optics dispersion relation by (i) considering a precisely plane, monochromatic wave in a precisely homogeneous, time-independent medium and deducing $\omega = \Omega(\mathbf{k})$ in a functional form that involves the medium's properties (e.g., density) and then (ii) allowing the properties to be slowly varying functions of $\mathbf{x}$ and $t$. The resulting dispersion relation [e.g., Eq. (7.23)] then acquires its $\mathbf{x}$ and $t$ dependence from the properties of the medium.*

The vanishing of the second term in the eikonal-approximated wave equation (7.22) dictates that the wave's real amplitude $A$ is transported with the group velocity $\mathbf{V}_\mathrm{g} = C\hat{\mathbf{k}}$ in the following manner:

$$\frac{dA}{dt} \equiv \left(\frac{\partial}{\partial t} + \mathbf{V}_\mathrm{g} \cdot \nabla\right) A = -\frac{1}{2W\omega}\left[\frac{\partial(W\omega)}{\partial t} + \nabla \cdot (WC^2\mathbf{k})\right] A. \tag{7.24}$$

**propagation law for amplitude**

This propagation law, by contrast with the dispersion relation, does depend on the weighting function $W$. We return to this propagation law shortly and shall understand more deeply its dependence on $W$, but first we must investigate in detail the directions along which $A$ is transported.

The time derivative $d/dt = \partial/\partial t + \mathbf{V}_\mathrm{g} \cdot \nabla$ appearing in the propagation law (7.24) is similar to the derivative with respect to proper time along a world line in special relativity, $d/d\tau = u^0\partial/\partial t + \mathbf{u} \cdot \nabla$ (with $u^\alpha$ the world line's 4-velocity). This analogy tells us that the waves' amplitude $A$ is being propagated along some sort of world lines (trajectories). Those world lines (the waves' rays), in fact, are governed by Hamilton's equations of particle mechanics with the dispersion relation $\Omega(\mathbf{x}, t, \mathbf{k})$ playing the role of the hamiltonian and $\mathbf{k}$ playing the role of momentum:

**rays**

$$\boxed{\frac{dx_j}{dt} = \left(\frac{\partial\Omega}{\partial k_j}\right)_{\mathbf{x},t} \equiv V_{\mathrm{g}\,j},} \quad \boxed{\frac{dk_j}{dt} = -\left(\frac{\partial\Omega}{\partial x_j}\right)_{\mathbf{k},t},} \quad \boxed{\frac{d\omega}{dt} = \left(\frac{\partial\Omega}{\partial t}\right)_{\mathbf{x},\mathbf{k}}.} \tag{7.25}$$

**Hamilton's equations for rays**

The first of these Hamilton equations is just our definition of the group velocity, with which [according to Eq. (7.24)] the amplitude is transported. The second tells us how

the wave vector $\mathbf{k}$ changes along a ray, and together with our knowledge of $C(\mathbf{x}, t)$, it tells us how the group velocity $\mathbf{V}_g = C\hat{\mathbf{k}}$ for our dispersionless waves changes along a ray, and thence defines the ray itself. The third tells us how the waves' frequency changes along a ray.

To deduce the second and third of these Hamilton equations, we begin by inserting the definitions $\omega = -\partial\varphi/\partial t$ and $\mathbf{k} = \nabla\varphi$ [Eqs. (7.21)] into the dispersion relation $\omega = \Omega(\mathbf{x}, t; \mathbf{k})$ for an arbitrary wave, thereby obtaining

$$\boxed{\frac{\partial\varphi}{\partial t} + \Omega(\mathbf{x}, t; \nabla\varphi) = 0.} \tag{7.26a}$$

**eikonal equation and Hamilton-Jacobi equation**

This equation is known in optics as the *eikonal equation*. It is formally the same as the Hamilton-Jacobi equation of classical mechanics (see, e.g., Goldstein, Poole, and Safko, 2002), if we identify $\Omega$ with the hamiltonian and $\varphi$ with Hamilton's principal function (cf. Ex. 7.9). This suggests that, to derive the second and third of Eqs. (7.25), we can follow the same procedure as is used to derive Hamilton's equations of motion. We take the gradient of Eq. (7.26a) to obtain

$$\frac{\partial^2\varphi}{\partial t \partial x_j} + \frac{\partial\Omega}{\partial k_l}\frac{\partial^2\varphi}{\partial x_l \partial x_j} + \frac{\partial\Omega}{\partial x_j} = 0, \tag{7.26b}$$

where the partial derivatives of $\Omega$ are with respect to its arguments $(\mathbf{x}, t; \mathbf{k})$; we then use $\partial\varphi/\partial x_j = k_j$ and $\partial\Omega/\partial k_l = V_{g\,l}$ to write Eq. (7.26b) as $dk_j/dt = -\partial\Omega/\partial x_j$. This is the second of Hamilton's equations (7.25), and it tells us how the wave vector changes along a ray. The third Hamilton equation, $d\omega/dt = \partial\Omega/\partial t$ [Eq. (7.25)], is obtained by taking the time derivative of the eikonal equation (7.26a).

Not only is the waves' amplitude $A$ propagated along the rays, so also is their phase:

**propagation equation for phase**

$$\frac{d\varphi}{dt} = \frac{\partial\varphi}{\partial t} + \mathbf{V}_g \cdot \nabla\varphi = -\omega + \mathbf{V}_g \cdot \mathbf{k}. \tag{7.27}$$

Since our dispersionless waves have $\omega = Ck$ and $\mathbf{V}_g = C\hat{\mathbf{k}}$, this vanishes. Therefore, for the special case of dispersionless waves (e.g., sound waves in a fluid and electromagnetic waves in an isotropic dielectric medium), the phase is constant along each ray:

$$\boxed{d\varphi/dt = 0.} \tag{7.28}$$

### 7.3.2  7.3.2 Connection of Geometric Optics to Quantum Theory

Although the waves $\psi = Ae^{i\varphi}$ are classical and our analysis is classical, their propagation laws in the eikonal approximation can be described most nicely in quantum mechanical language.[5] Quantum mechanics insists that, associated with any wave in

---

5. This is intimately related to the fact that quantum mechanics underlies classical mechanics; the classical world is an approximation to the quantum world, often a very good approximation.

the geometric-optics regime, there are real quanta: the wave's quantum mechanical particles. If the wave is electromagnetic, the quanta are photons; if it is gravitational, they are gravitons; if it is sound, they are phonons; if it is a plasma wave (e.g., Alfvén), they are plasmons. When we multiply the wave's $\mathbf{k}$ and $\omega$ by $\hbar$, we obtain the particles' momentum and energy:

$$\mathbf{p} = \hbar\mathbf{k}, \qquad \mathcal{E} = \hbar\omega. \tag{7.29}$$

**quanta**

**momentum and energy of quanta**

Although the originators of the nineteenth-century theory of classical waves were unaware of these quanta, once quantum mechanics had been formulated, the quanta became a powerful conceptual tool for thinking about classical waves.

In particular, we can regard the rays as the world lines of the quanta, and by multiplying the dispersion relation by $\hbar$, we can obtain the hamiltonian for the quanta's world lines:

$$H(\mathbf{x}, t; \mathbf{p}) = \hbar\Omega(\mathbf{x}, t; \mathbf{k} = \mathbf{p}/\hbar). \tag{7.30}$$

**hamiltonian for quanta**

Hamilton's equations (7.25) for the rays then immediately become Hamilton's equations for the quanta: $dx_j/dt = \partial H/\partial p_j$, $dp_j/dt = -\partial H/\partial x_j$, and $d\mathcal{E}/dt = \partial H/\partial t$.

Return now to the propagation law (7.24) for the waves' amplitude, and examine its consequences for the waves' energy. By inserting the ansatz $\psi = \Re(Ae^{i\varphi}) = A\cos(\varphi)$ into Eqs. (7.18) for the energy density $U$ and energy flux $\mathbf{F}$ and averaging over a wavelength and wave period (so $\overline{\cos^2\varphi} = \overline{\sin^2\varphi} = 1/2$), we find that

$$U = \frac{1}{2}WC^2k^2A^2 = \frac{1}{2}W\omega^2A^2, \qquad \mathbf{F} = U(C\hat{\mathbf{k}}) = U\mathbf{V}_{\mathrm{g}}. \tag{7.31}$$

Inserting these into the expression $\partial U/\partial t + \mathbf{\nabla} \cdot \mathbf{F}$ for the rate at which energy (per unit volume) fails to be conserved and using the propagation law (7.24) for $A$, we obtain

$$\frac{\partial U}{\partial t} + \mathbf{\nabla} \cdot \mathbf{F} = U\frac{\partial \ln C}{\partial t}. \tag{7.32}$$

Thus, as the propagation speed $C$ slowly changes at a fixed location in space due to a slow change in the medium's properties, the medium slowly pumps energy into the waves or removes it from them at a rate per unit volume of $U\partial \ln C/\partial t$.

This slow energy change can be understood more deeply using quantum concepts. The number density and number flux of quanta are

$$n = \frac{U}{\hbar\omega}, \qquad \mathbf{S} = \frac{\mathbf{F}}{\hbar\omega} = n\mathbf{V}_{\mathrm{g}}. \tag{7.33}$$

**number density and flux for quanta**

By combining these equations with the energy (non)conservation equation (7.32), we obtain

$$\frac{\partial n}{\partial t} + \mathbf{\nabla} \cdot \mathbf{S} = n\left[\frac{\partial \ln C}{\partial t} - \frac{d \ln \omega}{dt}\right]. \tag{7.34}$$

The third Hamilton equation (7.25) tells us that

$$d\omega/dt = (\partial\Omega/\partial t)_{x,k} = [\partial(Ck)/\partial t]_{x,k} = k\partial C/\partial t,$$

whence $d\ln\omega/dt = \partial\ln C/\partial t$, which, when inserted into Eq. (7.34), implies that the quanta are conserved:

**conservation of quanta**

$$\boxed{\frac{\partial n}{\partial t} + \nabla\cdot\mathbf{S} = 0.}$$

(7.35a)

Since $\mathbf{S} = n\mathbf{V}_g$ and $d/dt = \partial/\partial t + \mathbf{V}_g\cdot\nabla$, we can rewrite this conservation law as a propagation law for the number density of quanta:

$$\boxed{\frac{dn}{dt} + n\nabla\cdot\mathbf{V}_g = 0.}$$

(7.35b)

The propagation law for the waves' amplitude, Eq. (7.24), can now be understood much more deeply: *The amplitude propagation law is nothing but the law of conservation of quanta in a slowly varying medium, rewritten in terms of the amplitude. This is true quite generally, for any kind of wave (Sec. 7.3.3); and the quickest route to the amplitude propagation law is often to express the wave's energy density U in terms of the amplitude and then invoke conservation of quanta, Eq. (7.35b).*

In Ex. 7.3 we show that the conservation law (7.35b) is equivalent to

$$\boxed{\frac{d(nC\mathcal{A})}{dt} = 0, \quad \text{i.e., } nC\mathcal{A} \text{ is a constant along each ray.}}$$

(7.35c)

Here $\mathcal{A}$ is the cross sectional area of a bundle of rays surrounding the ray along which the wave is propagating. Equivalently, by virtue of Eqs. (7.33) and (7.31) for the number density of quanta in terms of the wave amplitude $A$, we have

$$\frac{d}{dt}A\sqrt{CW\omega\mathcal{A}} = 0, \quad \text{i.e., } A\sqrt{CW\omega\mathcal{A}} \text{ is a constant along each ray.} \quad (7.35d)$$

In Eqs. (7.33) and (7.35), we have boxed those equations that are completely general (because they embody conservation of quanta) and have not boxed those that are specialized to our prototypical wave equation.

**EXERCISES**

**Exercise 7.3** ** *Derivation and Example: Amplitude Propagation for Dispersionless Waves Expressed as Constancy of Something along a Ray*

(a) In connection with Eq. (7.35b), explain why $\nabla\cdot\mathbf{V}_g = d\ln\mathcal{V}/dt$, where $\mathcal{V}$ is the tiny volume occupied by a collection of the wave's quanta.

(b) Choose for the collection of quanta those that occupy a cross sectional area $\mathcal{A}$ orthogonal to a chosen ray, and a longitudinal length $\Delta s$ along the ray, so $\mathcal{V} = \mathcal{A}\Delta s$. Show that $d\ln\Delta s/dt = d\ln C/dt$ and correspondingly, $d\ln\mathcal{V}/dt = d\ln(C\mathcal{A})/dt$.

(c) Given part (b), show that the conservation law (7.35b) is equivalent to the constancy of $nC\mathcal{A}$ along a ray, Eq. (7.35c).

(d) From the results of part (c), derive the constancy of $A\sqrt{CW\omega\mathcal{A}}$ along a ray (where $A$ is the wave's amplitude), Eq. (7.35d).

**Exercise 7.4** **\*\*Example: Energy Density and Flux, and Adiabatic Invariant,** *for a Dispersionless Wave*

(a) Show that the prototypical scalar wave equation (7.17) follows from the variational principle

$$\delta \int \mathcal{L}dt d^3x = 0, \tag{7.36a}$$

where $\mathcal{L}$ is the lagrangian density

$$\mathcal{L} = W\left[\frac{1}{2}\left(\frac{\partial\psi}{\partial t}\right)^2 - \frac{1}{2}C^2(\nabla\psi)^2\right] \tag{7.36b}$$

(not to be confused with the lengthscale $\mathcal{L}$ of inhomogeneities in the medium).

(b) For any scalar-field lagrangian density $\mathcal{L}(\psi, \partial\psi/\partial t, \nabla\psi, \mathbf{x}, t)$, the energy density and energy flux can be expressed in terms of the lagrangian, in Cartesian coordinates, as

$$U(\mathbf{x}, t) = \frac{\partial\psi}{\partial t}\frac{\partial\mathcal{L}}{\partial\psi/\partial t} - \mathcal{L}, \qquad F_j = \frac{\partial\psi}{\partial t}\frac{\partial\mathcal{L}}{\partial\psi/\partial x_j} \tag{7.36c}$$

(Goldstein, Poole, and Safko, 2002, Sec. 13.3). Show, from the Euler-Lagrange equations for $\mathcal{L}$, that these expressions satisfy energy conservation, $\partial U/\partial t + \nabla \cdot \mathbf{F} = 0$, *if* $\mathcal{L}$ has no explicit time dependence [e.g., for the lagrangian (7.36b) if $C = C(\mathbf{x})$ and $W = W(\mathbf{x})$ do not depend on time $t$].

(c) Show that expression (7.36c) for the field's energy density $U$ and its energy flux $F_j$ agree with Eqs. (7.18).

(d) Now, regard the wave amplitude $\psi$ as a generalized (field) coordinate. Use the lagrangian $L = \int \mathcal{L}d^3x$ to define a field momentum $\Pi$ conjugate to this $\psi$, and then compute a *wave action,*

$$J \equiv \int_0^{2\pi/\omega}\int \Pi(\partial\psi/\partial t)d^3x\, dt, \tag{7.36d}$$

which is the continuum analog of Eq. (7.43) in Sec. 7.3.6. The temporal integral is over one wave period. Show that this $J$ is proportional to the wave energy divided by the frequency and thence to the number of quanta in the wave.

It is shown in standard texts on classical mechanics that, for approximately periodic oscillations, the particle action (7.43), with the integral limited to one period of oscillation of $q$, is an *adiabatic invariant.* By the extension of that proof to continuum physics, the wave action (7.36d) is also an adiabatic invariant. This

means that the wave action and hence the number of quanta in the waves are conserved when the medium [in our case the index of refraction $\mathfrak{n}(\mathbf{x})$] changes very slowly in time—a result asserted in the text, and one that also follows from quantum mechanics. We study the particle version (7.43) of this adiabatic invariant in detail when we analyze charged-particle motion in a slowly varying magnetic field in Sec. 20.7.4.

**Exercise 7.5** *Problem: Propagation of Sound Waves in a Wind*
Consider sound waves propagating in an atmosphere with a horizontal wind. Assume that the sound speed $C$, as measured in the air's local rest frame, is constant. Let the wind velocity $\mathbf{u} = u_x \mathbf{e}_x$ increase linearly with height $z$ above the ground: $u_x = Sz$, where $S$ is the constant shearing rate. Consider only rays in the $x$-$z$ plane.

(a) Give an expression for the dispersion relation $\omega = \Omega(\mathbf{x}, t; \mathbf{k})$. [Hint: In the local rest frame of the air, $\Omega$ should have its standard sound-wave form.]

(b) Show that $k_x$ is constant along a ray path, and then demonstrate that sound waves will not propagate when

$$\left| \frac{\omega}{k_x} - u_x(z) \right| < C. \tag{7.37}$$

(c) Consider sound rays generated on the ground that make an angle $\theta$ to the horizontal initially. Derive the equations describing the rays, and use them to sketch the rays, distinguishing values of $\theta$ both less than and greater than $\pi/2$. (You might like to perform this exercise numerically.)

### 7.3.3 Geometric Optics for a General Wave

With the simple case of nondispersive sound waves (Secs. 7.3.1 and 7.3.2) as our model, we now study an arbitrary kind of wave in a weakly inhomogeneous and slowly time varying medium (e.g., any of the examples in Sec. 7.2.1: light waves in a dielectric medium, deep water waves, flexural waves on a stiff beam, or Alfvén waves). Whatever the wave may be, we seek a solution to its wave equation using the eikonal approximation $\psi = Ae^{i\varphi}$ with slowly varying amplitude $A$ and rapidly varying phase $\varphi$. Depending on the nature of the wave, $\psi$ and $A$ might be a scalar (e.g., sound waves), a vector (e.g., light waves), or a tensor (e.g., gravitational waves).

When we insert the ansatz $\psi = Ae^{i\varphi}$ into the wave equation and collect terms in the manner dictated by our two-lengthscale expansion [as in Eq. (7.22) and Box 7.2], the leading-order term will arise from letting every temporal or spatial derivative act on the $e^{i\varphi}$. This is precisely where the derivatives would operate in the case of a plane wave in a homogeneous medium, and here, as there, the result of each differentiation is $\partial e^{i\varphi}/\partial t = -i\omega e^{i\varphi}$ or $\partial e^{i\varphi}/\partial x_j = ik_j e^{i\varphi}$. Correspondingly, the leading-order terms in the wave equation here will be identical to those in the homogeneous plane wave

case: they will be the dispersion relation multiplied by something times the wave,

$$[-\omega^2 + \Omega^2(\mathbf{x}, t; \mathbf{k})] \times (\text{something}) A e^{i\varphi} = 0, \qquad (7.38a)$$

with the spatial and temporal dependence of $\Omega^2$ entering through the medium's properties. This guarantees that (as we claimed in Sec. 7.3.1) the dispersion relation can be obtained by analyzing a plane, monochromatic wave in a homogeneous, time-independent medium and then letting the medium's properties, in the dispersion relation, vary slowly with $\mathbf{x}$ and $t$.

**dispersion relation for general wave**

Each next-order ("subleading") term in the wave equation will entail just one of the wave operator's derivatives acting on a slowly varying quantity ($A$, a medium property, $\omega$, or $\mathbf{k}$) and all the other derivatives acting on $e^{i\varphi}$. The subleading terms that interest us, for the moment, are those in which the one derivative acts on $A$, thereby propagating it. Therefore, the subleading terms can be deduced from the leading-order terms (7.38a) by replacing just one $i\omega A e^{i\varphi} = -A(e^{i\varphi})_{,t}$ by $-A_{,t} e^{i\varphi}$, and replacing just one $ik_j A e^{i\varphi} = A(e^{i\varphi})_{,j}$ by $A_{,j} e^{i\varphi}$ (where the subscript commas denote partial derivatives in Cartesian coordinates). A little thought then reveals that the equation for the vanishing of the subleading terms must take the form [deducible from the leading terms (7.38a)]:

$$-2i\omega\frac{\partial A}{\partial t} - 2i\Omega(\mathbf{k}, \mathbf{x}, t)\frac{\partial\Omega(\mathbf{k}, \mathbf{x}, t)}{\partial k_j}\frac{\partial A}{\partial x_j} = \text{terms proportional to } A. \quad (7.38b)$$

Using the dispersion relation $\omega = \Omega(\mathbf{x}, t; \mathbf{k})$ and the group velocity (first Hamilton equation) $\partial\Omega/\partial k_j = V_{g\,j}$, we bring this into the "propagate $A$ along a ray" form:

$$\frac{dA}{dt} \equiv \frac{\partial A}{\partial t} + \mathbf{V}_g \cdot \nabla A = \text{terms proportional to } A. \qquad (7.38c)$$

Let us return to the leading-order terms (7.38a) in the wave equation [i.e., to the dispersion relation $\omega = \Omega(\mathbf{x}, t; k)$]. For our general wave, as for the prototypical dispersionless wave of the previous two sections, the argument embodied in Eqs. (7.26) shows that the rays are determined by Hamilton's equations (7.25),

$$\boxed{\frac{dx_j}{dt} = \left(\frac{\partial\Omega}{\partial k_j}\right)_{\mathbf{x},t} \equiv V_{g\,j}, \quad \frac{dk_j}{dt} = -\left(\frac{\partial\Omega}{\partial x_j}\right)_{\mathbf{k},t}, \quad \frac{d\omega}{dt} = \left(\frac{\partial\Omega}{\partial t}\right)_{\mathbf{x},\mathbf{k}},} \quad (7.39)$$

**Hamilton's equations for general wave**

but using the general wave's dispersion relation $\Omega(\mathbf{k}, \mathbf{x}, t)$ rather than $\Omega = C(\mathbf{x}, t)k$. These Hamilton equations include propagation laws for $\omega = -\partial\varphi/\partial t$ and $k_j = \partial\varphi/\partial x_j$, from which we can deduce the propagation law (7.27) for $\varphi$ along the rays:

$$\boxed{\frac{d\varphi}{dt} = -\omega + \mathbf{V}_g \cdot \mathbf{k}.} \qquad (7.40)$$

**propagation law for phase of general wave**

For waves with dispersion, by contrast with sound in a fluid and other waves that have $\Omega = Ck$, $\varphi$ will not be constant along a ray.

For our general wave, as for dispersionless waves, the Hamilton equations for the rays can be reinterpreted as Hamilton's equations for the world lines of the waves' quanta [Eq. (7.30) and associated discussion]. And for our general wave, as for dispersionless waves, the medium's slow variations are incapable of creating or destroying wave quanta.[6] Correspondingly, if one knows the relationship between the waves' energy density $U$ and their amplitude $A$, and thence the relationship between the waves' quantum number density $n = U/\hbar\omega$ and $A$, then *from the quantum conservation law* [boxed Eqs. (7.35)]

conservation of quanta and propagation of amplitude for general wave

$$\frac{\partial n}{\partial t} + \boldsymbol{\nabla} \cdot (n\mathbf{V}_{\mathrm{g}}) = 0, \quad \frac{dn}{dt} + n\boldsymbol{\nabla} \cdot \mathbf{V}_{\mathrm{g}} = 0, \quad \text{or} \quad \frac{d(nC\mathcal{A})}{dt} = 0, \tag{7.41}$$

*one can deduce the propagation law for $A$—and the result must be the same propagation law as one obtains from the subleading terms in the eikonal approximation.*

### 7.3.4 Examples of Geometric-Optics Wave Propagation

#### SPHERICAL SOUND WAVES

As a simple example of these geometric-optics propagation laws, consider a sound wave propagating radially outward through a homogeneous fluid from a spherical source (e.g., a radially oscillating ball; cf. Sec. 16.5.3). The dispersion relation is Eq. (7.4): $\Omega = Ck$. It is straightforward (Ex. 7.6) to integrate Hamilton's equations and learn that the rays have the simple form $\{r = Ct + \text{constant}, \theta = \text{constant}, \phi = \text{constant}, \mathbf{k} = (\omega/C)\mathbf{e}_r\}$ in spherical polar coordinates, with $\mathbf{e}_r$ the unit radial vector. Because the wave is dispersionless, its phase $\varphi$ must be conserved along a ray [Eq. (7.28)], so $\varphi$ must be a function of $Ct - r, \theta$, and $\phi$. For the waves to propagate radially, it is essential that $\mathbf{k} = \boldsymbol{\nabla}\varphi$ point very nearly radially, which implies that $\varphi$ must be a rapidly varying function of $Ct - r$ and a slowly varying one of $\theta$ and $\phi$. The law of conservation of quanta in this case reduces to the propagation law $d(rA)/dt = 0$ (Ex. 7.6), so $rA$ is also a constant along the ray; we call it $\mathcal{B}$. Putting this all together, we conclude that the sound waves' pressure perturbation $\psi = \delta P$ has the form

$$\psi = \frac{\mathcal{B}(Ct - r, \theta, \phi)}{r} e^{i\varphi(Ct-r,\theta,\phi)}, \tag{7.42}$$

where the phase $\varphi$ is rapidly varying in $Ct - r$ and slowly varying in the angles, and the amplitude $\mathcal{B}$ is slowly varying in $Ct - r$ and the angles.

#### FLEXURAL WAVES

As another example of the geometric-optics propagation laws, consider flexural waves on a spacecraft's tapering antenna. The dispersion relation is $\Omega = k^2\sqrt{D/\Lambda}$ [Eq. (7.6)] with $D/\Lambda \propto h^2$, where $h$ is the antenna's thickness in its direction of bend (or the

---

6. This is a general feature of quantum theory; creation and destruction of quanta require imposed oscillations at the high frequency and short wavelength of the waves themselves, or at some submultiple of them (in the case of nonlinear creation and annihilation processes; Chap. 10).

antenna's diameter, if it has a circular cross section); cf. Eq. (12.33). Since $\Omega$ is independent of $t$, as the waves propagate from the spacecraft to the antenna's tip, their frequency $\omega$ is conserved [third of Eqs. (7.39)], which implies by the dispersion relation that $k = (D/\Lambda)^{-1/4}\omega^{1/2} \propto h^{-1/2}$; hence the wavelength decreases as $h^{1/2}$. The group velocity is $V_g = 2(D/\Lambda)^{1/4}\omega^{1/2} \propto h^{1/2}$. Since the energy per quantum $\hbar\omega$ is constant, particle conservation implies that the waves' energy must be conserved, which in this 1-dimensional problem means that the energy flowing through a segment of the antenna per unit time must be constant along the antenna. On physical grounds this constant energy flow rate must be proportional to $A^2 V_g h^2$, which means that the amplitude $A$ must increase $\propto h^{-5/4}$ as the flexural waves approach the antenna's end. A qualitatively similar phenomenon is seen in the cracking of a bullwhip (where the speed of the end can become supersonic).

## LIGHT THROUGH A LENS AND ALFVÉN WAVES

Figure 7.3 sketches two other examples: light propagating through a lens and Alfvén waves propagating in the magnetosphere of a planet. In Sec. 7.3.6 and the exercises we explore a variety of other applications, but first we describe how the geometric-optics propagation laws can fail (Sec. 7.3.5).

---

**EXERCISES**

**Exercise 7.6** *Derivation and Practice: Quasi-Spherical Solution to Vacuum Scalar Wave Equation*
Derive the quasi-spherical solution (7.42) of the vacuum scalar wave equation $-\partial^2\psi/\partial t^2 + \nabla^2\psi = 0$ from the geometric-optics laws by the procedure sketched in the text.

---

### 7.3.5 Relation to Wave Packets; Limitations of the Eikonal Approximation and Geometric Optics

7.3.5

The form $\psi = Ae^{i\varphi}$ of the waves in the eikonal approximation is remarkably general. At some initial moment of time, $A$ and $\varphi$ can have any form whatsoever, so long as the two-lengthscale constraints are satisfied [$A$, $\omega \equiv -\partial\varphi/\partial t$, $\mathbf{k} \equiv \nabla\varphi$, and dispersion relation $\Omega(\mathbf{k}; \mathbf{x}, t)$ all vary on lengthscales long compared to $\lambda = 1/k$ and on timescales long compared to $1/\omega$]. For example, $\psi$ could be as nearly planar as is allowed by the inhomogeneities of the dispersion relation. At the other extreme, $\psi$ could be a moderately narrow wave packet, confined initially to a small region of space (though not too small; its size must be large compared to its mean reduced wavelength). In either case, the evolution will be governed by the above propagation laws.

Of course, the eikonal approximation is an approximation. Its propagation laws make errors, though when the two-lengthscale constraints are well satisfied, the errors will be small for sufficiently short propagation times. Wave packets provide an important example. Dispersion (different group velocities for different wave vectors) causes wave packets to spread (disperse) as they propagate; see Ex. 7.2. This spreading

**phenomena missed by geometric optics**

(a)



(b)

**FIGURE 7.3** (a) The rays and the surfaces of constant phase $\varphi$ at a fixed time for light passing through a converging lens [dispersion relation $\Omega = ck/\mathfrak{n}(\mathbf{x})$, where $\mathfrak{n}$ is the index of refraction]. In this case the rays (which always point along $\mathbf{V}_g$) are parallel to the wave vector $\mathbf{k} = \nabla\varphi$ and thus are also parallel to the phase velocity $\mathbf{V}_{ph}$, and the waves propagate along the rays with a speed $V_g = V_{ph} = c/\mathfrak{n}$ that is independent of wavelength. The strange self-intersecting shape of the last phase front is due to caustics; see Sec. 7.5. (b) The rays and surfaces of constant phase for Alfvén waves in the magnetosphere of a planet [dispersion relation $\Omega = \mathbf{a}(\mathbf{x}) \cdot \mathbf{k}$]. In this case, because $\mathbf{V}_g = \mathbf{a} \equiv \mathbf{B}/\sqrt{\mu_0\rho}$, the rays are parallel to the magnetic field lines and are not parallel to the wave vector, and the waves propagate along the field lines with speeds $V_g$ that are independent of wavelength; cf. Fig. 7.2c. As a consequence, if some electric discharge excites Alfvén waves on the planetary surface, then they will be observable by a spacecraft when it passes magnetic field lines on which the discharge occurred. As the waves propagate, because $\mathbf{B}$ and $\rho$ are time independent and hence $\partial\Omega/\partial t = 0$, the frequency $\omega$ and energy $\hbar\omega$ of each quantum is conserved, and conservation of quanta implies conservation of wave energy. Because the Alfvén speed generally diminishes with increasing distance from the planet, conservation of wave energy typically requires the waves' energy density and amplitude to increase as they climb upward.

is not included in the geometric-optics propagation laws; it is a fundamentally wave-based phenomenon and is lost when one goes to the particle-motion regime. In the limit that the wave packet becomes very large compared to its wavelength or that the packet propagates for only a short time, the spreading is small (Ex. 7.2). This is the geometric-optics regime, and geometric optics ignores the spreading.

Many other wave phenomena are missed by geometric optics. Examples are diffraction (e.g., at a geometric-optics caustic; Secs. 7.5 and 8.6), nonlinear wave-wave coupling (Chaps. 10 and 23, and Sec. 16.3), and parametric amplification of waves by

rapid time variations of the medium (Sec. 10.7.3)—which shows up in quantum mechanics as particle production (i.e., a breakdown of the law of conservation of quanta). In Sec. 28.7.1 , we will encounter such particle production in inflationary models of the early universe.

### 7.3.6 Fermat's Principle

Hamilton's equations of optics allow us to solve for the paths of rays in media that vary both spatially and temporally. When the medium is time independent, the rays $\mathbf{x}(t)$ can be computed from a variational principle due to Fermat. This is the optical analog of the classical dynamics principle of least action,[7] which states that, when a particle moves from one point to another through a time-independent potential (so its energy, the hamiltonian, is conserved), then the path $\mathbf{q}(t)$ that it follows is one that extremizes the action

**principle of least action**

$$J = \int \mathbf{p} \cdot d\mathbf{q} \tag{7.43}$$

(where $\mathbf{q}$ and $\mathbf{p}$ are the particle's generalized coordinates and momentum), subject to the constraint that the paths have a fixed starting point, a fixed endpoint, and constant energy. The proof (e.g., Goldstein, Poole, and Safko, 2002, Sec. 8.6) carries over directly to optics when we replace the hamiltonian by $\Omega$, $\mathbf{q}$ by $\mathbf{x}$, and $\mathbf{p}$ by $\mathbf{k}$. The resulting Fermat principle, stated with some care, has the following form.

Consider waves whose hamiltonian $\Omega(\mathbf{k}, \mathbf{x})$ is independent of time. Choose an initial location $\mathbf{x}_{\text{initial}}$ and a final location $\mathbf{x}_{\text{final}}$ in space, and consider the rays $\mathbf{x}(t)$ that connect these two points. The rays (usually only one) are those paths that satisfy the variational principle

**Fermat's principle**

$$\boxed{\delta \int \mathbf{k} \cdot d\mathbf{x} = 0.} \tag{7.44}$$

In this variational principle, $\mathbf{k}$ must be expressed in terms of the trial path $\mathbf{x}(t)$ using Hamilton's equation $dx^j/dt = -\partial\Omega/\partial k_j$; the rate that the trial path is traversed (i.e., the magnitude of the group velocity) must be adjusted to keep $\Omega$ constant along the trial path (which means that the total time taken to go from $\mathbf{x}_{\text{initial}}$ to $\mathbf{x}_{\text{final}}$ can differ from one trial path to another). And of course, the trial paths must all begin at $\mathbf{x}_{\text{initial}}$ and end at $\mathbf{x}_{\text{final}}$.

#### PATH INTEGRALS

Notice that, once a ray has been identified by this action principle, it has $\mathbf{k} = \nabla\varphi$, and therefore the extremal value of the action $\int \mathbf{k} \cdot d\mathbf{x}$ along the ray is equal to the waves'

---

7. This is commonly attributed to Maupertuis, though others, including Leibniz and Euler, understood it earlier or better. This "action" and the rules for its variation are different from those in play in Hamilton's principle.

phase difference $\Delta\varphi$ between $\mathbf{x}_{\text{initial}}$ and $\mathbf{x}_{\text{final}}$. Correspondingly, for any trial path, we can think of the action as a phase difference along that path,

$$\Delta\varphi = \int \mathbf{k} \cdot d\mathbf{x}, \tag{7.45a}$$

and we can think of Fermat's principle as saying that the particle travels along a path of extremal phase difference $\Delta\varphi$. This can be reexpressed in a form closely related to *Feynman's path-integral formulation of quantum mechanics* (Feynman, 1966). We can regard all the trial paths as being followed with equal probability. For each path, we are to construct a probability amplitude $e^{i\Delta\varphi}$, and we must then add together these amplitudes,

$$\sum_{\text{all paths}} e^{i\Delta\varphi}, \tag{7.45b}$$

to get the net complex amplitude for quanta associated with the waves to travel from $\mathbf{x}_{\text{initial}}$ to $\mathbf{x}_{\text{final}}$. The contributions from almost all neighboring paths will interfere destructively. The only exceptions are those paths whose neighbors have the same values of $\Delta\varphi$, to first order in the path difference. These are the paths that extremize the action (7.44): they are the wave's rays, the actual paths of the quanta.

SPECIALIZATION TO $\Omega = C(\mathbf{x})k$

Fermat's principle takes on an especially simple form when not only is the hamiltonian $\Omega(\mathbf{k}, \mathbf{x})$ time independent, but it also has the simple dispersion-free form $\Omega = C(\mathbf{x})k$—a form valid for the propagation of light through a time-independent dielectric, and sound waves through a time-independent, inhomogeneous fluid, and electromagnetic or gravitational waves through a time-independent, Newtonian gravitational field (Sec. 7.6). In this $\Omega = C(\mathbf{x})k$ case, the hamiltonian dictates that for each trial path, $\mathbf{k}$ is parallel to $d\mathbf{x}$, and therefore $\mathbf{k} \cdot d\mathbf{x} = kds$, where $s$ is distance along the path. Using the dispersion relation $k = \Omega/C$ and noting that Hamilton's equation $dx^j/dt = \partial\Omega/\partial k_j$ implies $ds/dt = C$ for the rate of traversal of the trial path, we see that $\mathbf{k} \cdot d\mathbf{x} = kds = \Omega dt$. Since the trial paths are constrained to have $\Omega$ constant, Fermat's principle (7.44) becomes a principle of extremal time: The rays between $\mathbf{x}_{\text{initial}}$ and $\mathbf{x}_{\text{final}}$ are those paths along which

**principle of extreme time for dispersionless wave**

$$\boxed{\int dt = \int \frac{ds}{C(\mathbf{x})} = \int \frac{n(\mathbf{x})}{c}ds} \tag{7.46}$$

is extremal. In the last expression we have adopted the convention used for light in a dielectric medium, that $C(\mathbf{x}) = c/n(\mathbf{x})$, where $c$ is the speed of light in vacuum, and $n$

**index of refraction**

is the medium's index of refraction. Since $c$ is constant, the rays are paths of extremal optical path length $\int n(\mathbf{x})ds$.

We can use Fermat's principle to demonstrate that, if the medium contains no opaque objects, then there will always be at least one ray connecting any two

points. This is because there is a lower bound on the optical path between any two points, given by $n_{min} L$, where $n_{min}$ is the lowest value of the refractive index anywhere in the medium, and $L$ is the distance between the two points. This means that for some path the optical path length must be a minimum, and that path is then a ray connecting the two points.

From the principle of extremal time, we can derive the Euler-Lagrange differential equation for the ray. For ease of derivation, we write the action principle in the form

$$\delta \int n(\mathbf{x}) \sqrt{\frac{d\mathbf{x}}{d\mathbf{s}} \cdot \frac{d\mathbf{x}}{d\mathbf{s}}} \, ds, \tag{7.47}$$

where the quantity in the square root is identically one. Performing a variation in the usual manner then gives

$$\boxed{\frac{d}{ds}\left(n\frac{d\mathbf{x}}{ds}\right) = \nabla n, \quad \text{i.e.,} \quad \frac{d}{ds}\left(\frac{1}{C}\frac{d\mathbf{x}}{ds}\right) = \nabla\left(\frac{1}{C}\right).} \tag{7.48}$$

**ray equation for dispersionless wave**

This is equivalent to Hamilton's equations for the ray, as one can readily verify using the hamiltonian $\Omega = kc/n$ (Ex. 7.7).

Equation (7.48) is a second-order differential equation requiring two boundary conditions to define a solution. We can either choose these to be the location of the start of the ray and its starting direction, or the start and end of the ray. A simple case arises when the medium is stratified [i.e., when $n = n(z)$, where $(x, y, z)$ are Cartesian coordinates]. Projecting Eq. (7.48) perpendicular to $\mathbf{e}_z$, we discover that $ndy/ds$ and $ndx/ds$ are constant, which implies

$$\boxed{n \sin\theta = \text{constant},} \tag{7.49}$$

**Snell's law**

where $\theta$ is the angle between the ray and $\mathbf{e}_z$. This is a variant of Snell's law of refraction. Snell's law is just a mathematical statement that the rays are normal to surfaces (wavefronts) on which the eikonal (phase) $\varphi$ is constant (cf. Fig. 7.4).[8] Snell's law is valid not only when $n(\mathbf{x})$ varies slowly but also when it jumps discontinuously, despite the assumptions underlying geometric optics failing at a discontinuity.

---

**EXERCISES**

**Exercise 7.7** *Derivation: Hamilton's Equations for Dispersionless Waves; Fermat's Principle*
Show that Hamilton's equations for the standard dispersionless dispersion relation (7.4) imply the same ray equation (7.48) as we derived using Fermat's principle.

---

8. Another important application of this general principle is to the design of optical instruments, where it is known as the *Abbé condition*. See, e.g., Born and Wolf (1999).

**FIGURE 7.4** Illustration of Snell's law of refraction at the interface between two media, for which the refractive indices are $n_1$ and $n_2$ (assumed less than $n_1$). As the wavefronts must be continuous across the interface, simple geometry tells us that $\lambda_1/\sin\theta_1 = \lambda_2/\sin\theta_2$. This and the fact that the wavelengths are inversely proportional to the refractive index, $\lambda_j \propto 1/n_j$, imply that $n_1\sin\theta_1 = n_2\sin\theta_2$, in agreement with Eq. (7.49).

**Exercise 7.8** *Example: Self-Focusing Optical Fibers*

Optical fibers in which the refractive index varies with radius are commonly used to transport optical signals. When the diameter of the fiber is many wavelengths, we can use geometric optics. Let the refractive index be

$$n = n_0(1 - \alpha^2 r^2)^{1/2}, \tag{7.50a}$$

where $n_0$ and $\alpha$ are constants, and $r$ is radial distance from the fiber's axis.

(a) Consider a ray that leaves the axis of the fiber along a direction that makes an angle $\beta$ to the axis. Solve the ray-transport equation (7.48) to show that the radius of the ray is given by

$$r = \frac{\sin\beta}{\alpha}\left|\sin\left(\frac{\alpha z}{\cos\beta}\right)\right|, \tag{7.50b}$$

where $z$ measures distance along the fiber.

(b) Next consider the propagation time $T$ for a light pulse propagating along the ray with $\beta \ll 1$, down a long length $L$ of fiber. Show that

$$T = \frac{n_0 L}{C}[1 + O(\beta^4)], \tag{7.50c}$$

and comment on the implications of this result for the use of fiber optics for communication.

**Exercise 7.9** **Example: Geometric Optics for the Schrödinger Equation*

Consider the nonrelativistic Schrödinger equation for a particle moving in a time-dependent, 3-dimensional potential well:

$$-\frac{\hbar}{i}\frac{\partial \psi}{\partial t} = \left[ \frac{1}{2m}\left( \frac{\hbar}{i}\nabla \right)^2 + V(\mathbf{x}, t) \right]\psi. \tag{7.51}$$

(a) Seek a geometric-optics solution to this equation with the form $\psi = A e^{iS/\hbar}$, where $A$ and $V$ are assumed to vary on a lengthscale $\mathcal{L}$ and timescale $\mathcal{T}$ long compared to those, $1/k$ and $1/\omega$, on which $S$ varies. Show that the leading-order terms in the two-lengthscale expansion of the Schrödinger equation give the Hamilton-Jacobi equation

$$\frac{\partial S}{\partial t} + \frac{1}{2m}(\nabla S)^2 + V = 0. \tag{7.52a}$$

Our notation $\varphi \equiv S/\hbar$ for the phase $\varphi$ of the wave function $\psi$ is motivated by the fact that the geometric-optics limit of quantum mechanics is classical mechanics, and the function $S = \hbar\varphi$ becomes, in that limit, "Hamilton's principal function," which obeys the Hamilton-Jacobi equation (see, e.g., Goldstein, Poole, and Safko, 2002, Chap. 10). [Hint: Use a formal parameter $\sigma$ to keep track of orders (Box 7.2), and argue that terms proportional to $\hbar^n$ are of order $\sigma^n$. This means there must be factors of $\sigma$ in the Schrödinger equation (7.51) itself.]

(b) From Eq. (7.52a) derive the equation of motion for the rays (which of course is identical to the equation of motion for a wave packet and therefore is also the equation of motion for a classical particle):

$$\frac{d\mathbf{x}}{dt} = \frac{\mathbf{p}}{m}, \qquad \frac{d\mathbf{p}}{dt} = -\nabla V, \tag{7.52b}$$

where $\mathbf{p} = \nabla S$.

(c) Derive the propagation equation for the wave amplitude $A$ and show that it implies

$$\frac{d|A|^2}{dt} + |A|^2 \frac{\nabla \cdot \mathbf{p}}{m} = 0. \tag{7.52c}$$

Interpret this equation quantum mechanically.

## 7.4 Paraxial Optics

It is quite common in optics to be concerned with a bundle of rays that are almost parallel (i.e., for which the angle the rays make with some reference ray can be treated as small). This approximation is called *paraxial optics,* and it permits one to linearize the geometric-optics equations and use matrix methods to trace their

**FIGURE 7.5** A reference ray (the $z$-axis) and an adjacent ray identified by its transverse distances $x(z)$ and $y(z)$, from the reference ray.

rays. The resulting matrix formalism underlies the first-order theory of simple optical instruments (e.g., the telescope and the microscope).

We develop the paraxial optics formalism for waves whose dispersion relation has the simple, time-independent, nondispersive form $\Omega = kc/\mathfrak{n}(\mathbf{x})$. This applies to light in a dielectric medium—the usual application. As we shall see, it also applies to charged particles in a storage ring or electron microscope (Sec. 7.4.2) and to light being lensed by a weak gravitational field (Sec. 7.6).

We restrict ourselves to a situation where there exists a ray that is a straight line, except when it reflects off a mirror or other surface. We choose this as a *reference ray* (also called the *optic axis*) for our formalism, and we orient the $z$-axis of a Cartesian coordinate system along it (Fig. 7.5). Let the 2-dimensional vector $\mathbf{x}(z)$ be the transverse displacement of some other ray from this reference ray, and denote by $(x, y) = (x_1, x_2)$ the Cartesian components of $\mathbf{x}$.

Under paraxial conditions, $|\mathbf{x}|$ is small compared to the $z$ lengthscales of the propagation, so we can Taylor expand the refractive index $\mathfrak{n}(\mathbf{x}, z)$ in $(x_1, x_2)$:

$$\mathfrak{n}(\mathbf{x}, z) = \mathfrak{n}(0, z) + x_i \mathfrak{n}_{,i}(0, z) + \frac{1}{2} x_i x_j \mathfrak{n}_{,ij}(0, z) + \ldots. \tag{7.53a}$$

Here the subscript commas denote partial derivatives with respect to the transverse coordinates, $\mathfrak{n}_{,i} \equiv \partial \mathfrak{n}/\partial x_i$. The linearized form of the ray-propagation equation (7.48) is then given by

$$\frac{d}{dz}\left(\mathfrak{n}(0, z)\frac{dx_i}{dz}\right) = \mathfrak{n}_{,i}(0, z) + x_j \mathfrak{n}_{,ij}(0, z). \tag{7.53b}$$

In order for the reference ray $x_i = 0$ to satisfy this equation, $\mathfrak{n}_{,i}(0, z)$ must vanish, so Eq. (7.53b) becomes a linear, homogeneous, second-order equation for the path of a nearby ray, $\mathbf{x}(z)$:

**paraxial ray equation**

$$\boxed{\left(\frac{d}{dz}\right)\left(\frac{\mathfrak{n}\,dx_i}{dz}\right) = x_j \mathfrak{n}_{,ij}.} \tag{7.54}$$

Here $\mathfrak{n}$ and $\mathfrak{n}_{,ij}$ are evaluated on the reference ray. It is helpful to regard $z$ as "time" and think of Eq. (7.54) as an equation for the 2-dimensional motion of a particle (the

ray) in a quadratic potential well. We can solve Eq. (7.54) given starting values $\mathbf{x}(z')$ and $\dot{\mathbf{x}}(z')$, where the dot denotes differentiation with respect to $z$, and $z'$ is the starting location. The solution at some later point $z$ is linearly related to the starting values. We can capitalize on this linearity by treating $\{\mathbf{x}(z), \dot{\mathbf{x}}(z)\}$ as a 4-dimensional vector $V_i(z)$, with

$$V_1 = x, \quad V_2 = \dot{x}, \quad V_3 = y, \quad \text{and} \quad V_4 = \dot{y}, \tag{7.55a}$$

and embodying the linear transformation [linear solution of Eq. (7.54)] from location $z'$ to location $z$ in a *transfer matrix* $J_{ab}(z, z')$:

$$V_a(z) = J_{ab}(z, z') \cdot V_b(z'), \tag{7.55b}$$

**paraxial transfer matrix**

where there is an implied sum over the repeated index $b$. The transfer matrix contains full information about the change of position and direction of all rays that propagate from $z'$ to $z$. As is always the case for linear systems, the transfer matrix for propagation over a large interval, from $z'$ to $z$, can be written as the product of the matrices for two subintervals, from $z'$ to $z''$ and from $z''$ to $z$:

$$J_{ac}(z, z') = J_{ab}(z, z'') J_{bc}(z'', z'). \tag{7.55c}$$

### 7.4.1  Axisymmetric, Paraxial Systems: Lenses, Mirrors, Telescopes, Microscopes, and Optical Cavities

7.4.1

If the index of refraction is everywhere axisymmetric, so $\mathfrak{n} = \mathfrak{n}(\sqrt{x^2 + y^2}, z)$, then there is no coupling between the motions of rays along the $x$ and $y$ directions, and the equations of motion along $x$ are identical to those along $y$. In other words, $J_{11} = J_{33}$, $J_{12} = J_{34}$, $J_{21} = J_{43}$, and $J_{22} = J_{44}$ are the only nonzero components of the transfer matrix. This reduces the dimensionality of the propagation problem from 4 dimensions to 2: $V_a$ can be regarded as either $\{x(z), \dot{x}(z)\}$ or $\{y(z), \dot{y}(z)\}$, and in both cases the $2 \times 2$ transfer matrix $J_{ab}$ is the same.

Let us illustrate the paraxial formalism by deriving the transfer matrices of a few simple, axisymmetric optical elements. In our derivations it is helpful conceptually to focus on rays that move in the $x$-$z$ plane (i.e., that have $y = \dot{y} = 0$). We write the 2-dimensional $V_i$ as a column vector:

**axisymmetric transfer matrices**

$$V_a = \begin{pmatrix} x \\ \dot{x} \end{pmatrix}. \tag{7.56a}$$

The simplest case is a straight section of length $d$ extending from $z'$ to $z = z' + d$. The components of $V$ will change according to

$$x = x' + \dot{x}'d,$$
$$\dot{x} = \dot{x}',$$

so

**for straight section**

$$J_{ab} = \begin{pmatrix} 1 & d \\ 0 & 1 \end{pmatrix} \text{ for a straight section of length } d, \qquad (7.56b)$$

where $x' = x(z')$, and so forth. Next, consider a thin lens with focal length $f$. The usual convention in optics is to give $f$ a positive sign when the lens is converging and a negative sign when diverging. A thin lens gives a deflection to the ray that is linearly proportional to its displacement from the optic axis, but does not change its transverse location. Correspondingly, the transfer matrix in crossing the lens (ignoring its thickness) is

**for thin lens**

$$J_{ab} = \begin{pmatrix} 1 & 0 \\ -f^{-1} & 1 \end{pmatrix} \text{ for a thin lens with focal length } f. \qquad (7.56c)$$

Similarly, a spherical mirror with radius of curvature $R$ (again adopting a positive sign for a converging mirror and a negative sign for a diverging mirror) has a transfer matrix

**for spherical mirror**

$$J_{ab} = \begin{pmatrix} 1 & 0 \\ -2R^{-1} & 1 \end{pmatrix} \text{ for a spherical mirror with radius of curvature } R.$$

$$(7.56d)$$

(Recall our convention that $z$ always increases along a ray, even when the ray reflects off a mirror.)

As a simple illustration, we consider rays that leave a point source located a distance $u$ in front of a converging lens of focal length $f$, and we solve for the ray positions a distance $v$ behind the lens (Fig. 7.6). The total transfer matrix is the transfer matrix (7.56b) for a straight section, multiplied by the product of the lens transfer matrix (7.56c) and a second straight-section transfer matrix:

$$J_{ab} = \begin{pmatrix} 1 & v \\ 0 & 1 \end{pmatrix} \begin{pmatrix} 1 & 0 \\ -f^{-1} & 1 \end{pmatrix} \begin{pmatrix} 1 & u \\ 0 & 1 \end{pmatrix} = \begin{pmatrix} 1 - vf^{-1} & u + v - uvf^{-1} \\ -f^{-1} & 1 - uf^{-1} \end{pmatrix}. \quad (7.57)$$

When the 1-2 element (upper right entry) of this transfer matrix vanishes, the position of the ray after traversing the optical system is independent of the starting direction. In other words, rays from the point source form a point image. When this happens, the planes containing the source and the image are said to be conjugate. The condition for this to occur is

**conjugate planes**

**thin-lens equations**

$$\frac{1}{u} + \frac{1}{v} = \frac{1}{f}. \qquad (7.58)$$

**FIGURE 7.6** Simple converging lens used to illustrate the use of transfer matrices. The total transfer matrix is formed by taking the product of the straight-section transfer matrix with the lens matrix and another straight-section matrix.



**FIGURE 7.7** Simple refracting telescope. By convention $\theta > 0$ and $\mathcal{M}\theta < 0$, so the image is inverted.

This is the standard thin-lens equation. The linear magnification of the image is given by $\mathcal{M} = J_{11} = 1 - v/f$, that is,

$$\mathcal{M} = -\frac{v}{u}, \tag{7.59}$$

where the negative sign means that the image is inverted. Note that, if a ray is reversed in direction, it remains a ray, but with the source and image planes interchanged; $u$ and $v$ are exchanged, Eq. (7.58) is unaffected, and the magnification (7.59) is inverted: $\mathcal{M} \to 1/\mathcal{M}$.

**EXERCISES**

**Exercise 7.10** *Problem: Matrix Optics for a Simple Refracting Telescope*
Consider a simple refracting telescope (Fig. 7.7) that comprises two converging lenses, the *objective* and the *eyepiece*. This telescope takes parallel rays of light from distant stars, which make an angle $\theta \ll 1$ with the optic axis, and converts them into parallel rays making a much larger angle $\mathcal{M}\theta$. Here $\mathcal{M}$ is the magnification with $\mathcal{M}$ negative,

**FIGURE 7.8** Simple microscope.

$|\mathcal{M}| \gg 1$, and $|\mathcal{M}\theta| \ll 1$. (The parallel output rays are then focused by the lens of a human's eye, to a point on the eye's retina.)

(a) Use matrix methods to investigate how the output rays depend on the separation of the two lenses, and hence find the condition that the output rays are parallel when the input rays are parallel.

(b) How does the magnification $\mathcal{M}$ depend on the ratio of the focal lengths of the two lenses?

(c) If, instead of looking through the telescope with one's eye, one wants to record the stars' image on a photographic plate or CCD, how should the optics be changed?

**Exercise 7.11** *Problem: Matrix Optics for a Simple Microscope*
A microscope takes light rays from a point on a microscopic object, very near the optic axis, and transforms them into parallel light rays that will be focused by a human eye's lens onto the eye's retina (Fig. 7.8). Use matrix methods to explore the operation of such a microscope. A single lens (magnifying glass) could do the same job (rays from a point converted to parallel rays). Why does a microscope need two lenses? What focal lengths and lens separations are appropriate for the eye to resolve a bacterium 100 $\mu$m in size?

**Exercise 7.12** *Example: Optical Cavity—Rays Bouncing between Two Mirrors*
Consider two spherical mirrors, each with radius of curvature $R$, separated by distance $d$ so as to form an optical cavity (Fig. 7.9). A laser beam bounces back and forth between the two mirrors. The center of the beam travels along a geometric-optics ray. (We study such beams, including their diffractive behavior, in Sec. 8.5.5.)

(a) Show, using matrix methods, that the central ray hits one of the mirrors (either one) at successive locations $\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \ldots$ (where $\mathbf{x} \equiv (x, y)$ is a 2-dimensional vector in the plane perpendicular to the optic axis), which satisfy the difference equation

$$\mathbf{x}_{k+2} - 2b\mathbf{x}_{k+1} + \mathbf{x}_k = 0, \qquad\qquad (7.60a)$$

**FIGURE 7.9** An optical cavity formed by two mirrors, and a light beam bouncing back and forth inside it.

where

$$b = 1 - \frac{4d}{R} + \frac{2d^2}{R^2}. \tag{7.60b}$$

Explain why this is a difference-equation analog of the simple-harmonic-oscillator equation.

(b) Show that this difference equation has the general solution

$$\mathbf{x}_k = \mathbf{A} \cos(k \cos^{-1} b) + \mathbf{B} \sin(k \cos^{-1} b). \tag{7.60c}$$

Obviously, $\mathbf{A}$ is the transverse position $\mathbf{x}_0$ of the ray at its 0th bounce. The ray's 0th position $\mathbf{x}_0$ and its 0th direction of motion $\dot{\mathbf{x}}_0$ together determine $\mathbf{B}$.

(c) Show that if $0 \leq d \leq 2R$, the mirror system is stable. In other words, all rays oscillate about the optic axis. Similarly, show that if $d > 2R$, the mirror system is unstable and the rays diverge from the optic axis.

(d) For an appropriate choice of initial conditions $\mathbf{x}_0$ and $\dot{\mathbf{x}}_0$, the laser beam's successive spots on the mirror lie on a circle centered on the optic axis. When operated in this manner, the cavity is called a *Harriet delay line.* How must $d/R$ be chosen so that the spots have an angular step size $\theta$? (There are two possible choices.)

## 7.4.2 Converging Magnetic Lens for Charged Particle Beam

Since geometric optics is the same as particle dynamics, matrix equations can be used to describe paraxial motions of electrons or ions in a storage ring. (Note, however, that the hamiltonian for such particles is dispersive, since it does not depend linearly on the particle momentum, and so for our simple matrix formalism to be valid, we must confine attention to a monoenergetic beam of particles.)

The simplest practical lens for charged particles is a quadrupolar magnet. Quadrupolar magnetic fields are used to guide particles around storage rings. If we orient our axes appropriately, the magnet's magnetic field can be expressed in the form

$$\mathbf{B} = \frac{B_0}{r_0}(y\mathbf{e}_x + x\mathbf{e}_y) \quad \text{independent of } z \text{ within the lens} \tag{7.61}$$

**quadrupolar magnetic field**

**FIGURE 7.10** Quadrupolar magnetic lens. The magnetic field lines lie in a plane perpendicular to the optic axis. Positively charged particles moving along $\mathbf{e}_z$ converge when $y = 0$ and diverge when $x = 0$.

(Fig. 7.10). Particles traversing this magnetic field will be subjected to a Lorentz force that curves their trajectories. In the paraxial approximation, a particle's coordinates satisfy the two differential equations

$$\ddot{x} = -\frac{x}{\lambda^2}, \qquad \ddot{y} = \frac{y}{\lambda^2}, \tag{7.62a}$$

where the dots (as above) mean $d/dz = v^{-1}d/dt$, and

$$\lambda = \left(\frac{pr_0}{qB_0}\right)^{1/2} \tag{7.62b}$$

[cf. Eq. (7.61)], with $q$ the particle's charge (assumed positive) and $p$ its momentum. The motions in the $x$ and $y$ directions are decoupled. It is convenient in this case to work with two 2-dimensional vectors, $\{V_{x1}, V_{x2}\} \equiv \{x, \dot{x}\}$ and $\{V_{y1}, V_{y2}\} = \{y, \dot{y}\}$. From the elementary solutions to the equations of motion (7.62a), we infer that the transfer matrices from the magnet's entrance to its exit are $J_{x\,ab}$, $J_{y\,ab}$, where

**transfer matrices for quadrupolar magnetic lens**

$$J_{x\,ab} = \begin{pmatrix} \cos\phi & \lambda\,\sin\phi \\ -\lambda^{-1}\,\sin\phi & \cos\phi \end{pmatrix}, \tag{7.63a}$$

$$J_{y\,ab} = \begin{pmatrix} \cosh\phi & \lambda\,\sinh\phi \\ \lambda^{-1}\,\sinh\phi & \cosh\phi \end{pmatrix}, \tag{7.63b}$$

and

$$\phi = L/\lambda, \tag{7.63c}$$

with $L$ the distance from entrance to exit (i.e., the lens thickness).

The matrices $J_{x\,ab}$ and $J_{y\,ab}$ can be decomposed as follows:

$$J_{x\,ab} = \begin{pmatrix} 1 & \lambda\tan\phi/2 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} 1 & 0 \\ -\lambda^{-1}\sin\phi & 1 \end{pmatrix} \begin{pmatrix} 1 & \lambda\tan\phi/2 \\ 0 & 1 \end{pmatrix} \qquad (7.63d)$$

$$J_{y\,ab} = \begin{pmatrix} 1 & \lambda\tanh\phi/2 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} 1 & 0 \\ \lambda^{-1}\sinh\phi & 1 \end{pmatrix} \begin{pmatrix} 1 & \lambda\tanh\phi/2 \\ 0 & 1 \end{pmatrix} \qquad (7.63e)$$

Comparing with Eqs. (7.56b) and (7.56c), we see that the action of a single magnet is equivalent to the action of a straight section, followed by a thin lens, followed by another straight section. Unfortunately, if the lens is focusing in the $x$ direction, it must be defocusing in the $y$ direction and vice versa. However, we can construct a lens that is focusing along both directions by combining two magnets that have opposite polarity but the same focusing strength $\phi = L/\lambda$.

**combining two magnets to make a converging lens**

Consider first the particles' motion in the $x$ direction. Let

$$f_+ = \lambda/\sin\phi \quad \text{and} \quad f_- = -\lambda/\sinh\phi \qquad (7.64)$$

be the equivalent focal lengths of the first converging lens and the second diverging lens. If we separate the magnets by a distance $s$, this must be added to the two effective lengths of the two magnets to give an equivalent separation of $d = \lambda\tan(\phi/2) + s + \lambda\tanh(\phi/2)$ for the two equivalent thin lenses. The combined transfer matrix for the two thin lenses separated by this distance $d$ is then

$$\begin{pmatrix} 1 & 0 \\ -f_-^{-1} & 1 \end{pmatrix} \begin{pmatrix} 1 & d \\ 0 & 1 \end{pmatrix} \begin{pmatrix} 1 & 0 \\ -f_+^{-1} & 1 \end{pmatrix} = \begin{pmatrix} 1 - df_+^{-1} & d \\ -f_*^{-1} & 1 - df_-^{-1} \end{pmatrix}, \qquad (7.65a)$$

**transfer matrix for converging magnetic lens**

where

$$\frac{1}{f_*} = \frac{1}{f_-} + \frac{1}{f_+} - \frac{d}{f_-f_+} = \frac{\sin\phi}{\lambda} - \frac{\sinh\phi}{\lambda} + \frac{d\sin\phi\sinh\phi}{\lambda^2}. \qquad (7.65b)$$

If we assume that $\phi \ll 1$ and $s \ll L$, then we can expand as a Taylor series in $\phi$ to obtain

$$f_* \simeq \frac{3\lambda}{2\phi^3} = \frac{3\lambda^4}{2L^3}. \qquad (7.66)$$

The effective focal length $f_*$ of the combined magnets is positive, and so the lens has a net focusing effect. From the symmetry of Eq. (7.65b) under interchange of $f_+$ and $f_-$, it should be clear that $f_*$ is independent of the order in which the magnets are encountered. Therefore, if we were to repeat the calculation for the motion in the $y$ direction, we would get the same focusing effect. (The diagonal elements of the transfer matrix are interchanged, but as they are both close to unity, this difference is rather small.)

The combination of two quadrupole lenses of opposite polarity can therefore imitate the action of a converging lens. Combinations of magnets like this are used to collimate particle beams in storage rings, particle accelerators, and electron microscopes.

## 7.5  Catastrophe Optics

### 7.5.1  Image Formation

CAUSTICS

Many simple optical instruments are carefully made to form point images from point sources. However, naturally occurring optical systems, and indeed precision optical instruments when examined in fine detail, bring light to a focus not at a point, but instead on a 2-dimensional surface—an envelope formed by the rays—called a *caustic*. Caustics are often seen in everyday life. For example, when bright sunlight is reflected by the inside of an empty coffee mug some of the rays are reflected *specularly* (angle of incidence equals angle of reflection) and some of the rays are reflected *diffusely* (in all directions due to surface irregularity and multiple reflections and refractions beneath the surface). The specular reflection by the walls—a cylindrical mirror—forms a caustic surface. The intersection of this surface with the bottom forms caustic lines that can be seen in diffuse reflection.[9] These caustic lines are observed to meet in a point. When the optical surfaces are quite irregular (e.g., the water surface in a swimming pool[10] or the type of glass used in bathrooms), then a *caustic network* forms. Caustic lines and points are seen, just the same as with the mug (Fig. 7.11).

What may be surprising is that caustics like these, formed under quite general conditions, can be classified into a rather small number of types, called *catastrophes,* possessing generic properties and scaling laws (Thom, 1994). The scaling laws are reminiscent of the renormalization group discussed in Sec. 5.8.3. Although we focus on catastrophes in the context of optics (e.g., Berry and Upstill, 1980), where they are caustics, the phenomenon is quite general and crops up in other subfields of physics, especially dynamics (e.g., Arnol'd, 1992) and thermodynamics (e.g., Ex. 7.16). It has also been invoked, often quite controversially, in fields outside physics (e.g., Poston and Stewart, 2012). Catastrophes can be found whenever we have a physical system whose states are determined by extremizing a function, such as energy. Our treatment will be quite heuristic, but the subject does have a formal mathematical foundation that connects it to bifurcations and *Morse theory* (see Sec. 11.6; see also, e.g., Petters et al., 2001).

STATE VARIABLES AND CONTROL PARAMETERS

Let us start with a specific, simple example. Suppose that there is a distant *source S* and a *detector D* separated by free space. If we consider all the paths from *S* to

---

9.  The curve that is formed is called a "nephroid."
10. The optics is quite complicated. Some rays from the Sun are reflected specularly by the surface of the water, creating multiple images of the Sun. As the Sun is half a degree in diameter, these produce thickened caustic lines. Some rays are refracted by the water, forming caustic surfaces that are intersected by the bottom of the pool to form a caustic pattern. Some light from this pattern is reflected diffusely before being refracted a second time on the surface of the water and ultimately detected by a retina or a CCD. Other rays are reflected multiple times.

(a)

(b)

**FIGURE 7.11** Photographs of caustics. (a) Simple caustic pattern formed by a coffee mug. (b) Caustic network formed in a swimming pool. The generic structure of these patterns comprises *fold* lines meeting at *cusp* points.



(a)

(b)

**FIGURE 7.12** (a) Alternative paths make up a sequence of straight segments from a source $S$ to a detector $D$. The path $S$-$\mathcal{P}_1$-$\mathcal{P}_2$-$\mathcal{P}_3$-$D$ can be simplified to a shorter path $S$-$\mathcal{P}_1$-$\mathcal{P}_3$-$D$, and this process can be continued until we have the minimum number of segments needed to exhibit the catastrophe. The (true) ray, with the smallest phase difference, is the axis $S$-$D$. (b) A single path from a distant source intersecting a screen at $\{a, b\}$ and ending at detector $D$ with coordinates $\{x, y, z\}$.

$D$ there is a single extremum—a minimum—in the phase difference, $\Delta \varphi = \omega t = \int \mathbf{k} \cdot d\mathbf{x}$. By Fermat's principle, this is the (true) ray—the *axis*—connecting the two points. There are an infinite number of alternative paths that could be defined by an infinite set of parameters, but wherever else the rays go, the phase difference is larger. Take one of these alternative paths connecting $S$ and $D$ and break it down into a sequence of connected segments (Fig. 7.12a). We can imagine replacing two successive segments with a single segment connecting their endpoints. This will reduce the phase difference. The operation can be repeated until we are left with the minimum number of segments, specified by the minimum number of necessary variables that we need to exhibit the catastrophe. The order in which we do this does not matter, and the final variables characterizing the path can be chosen for convenience. These variables are known as *state variables*.

state variables

Next, introduce a screen perpendicular to the $S$–$D$ axis and close to $D$ (Fig. 7.12b). Consider a path from $S$, nearly parallel to the axis and intersecting the screen at a point with Cartesian coordinates $\{a, b\}$ measured from the axis. There let it be deflected toward $D$. In this section and the next, introduce the *delay* $t \equiv \Delta\varphi/\omega$, subtracting off the constant travel time along the axis in the absence of the screen, to measure the phase. The additional geometric delay associated with this ray is given approximately by

$$t_{\text{geo}} = \frac{a^2 + b^2}{2zc}, \tag{7.67}$$

**control parameters**

where $z \gg \{a, b\}$ measures the distance from the screen to $D$, parallel to the axis, and $c$ is the speed of light. The coordinates $\{a, b\}$ act as state variables, and the true ray is determined by differentiating with respect to them. Next, move $D$ off the axis and give it Cartesian coordinates $\{x, y, z\}$ with the $x$-$y$ plane parallel to the $a$-$b$ plane, and the transverse coordinates measured from the original axis. As these coordinates specify one of the endpoints, they do not enter into the variation that determines the true ray, but they do change $t_{\text{geo}}$ to $[(a - x)^2 + (b - y)^2]/(2zc)$. These $\{x, y, z\}$ parameters are examples of *control parameters*. In general, the number of control parameters that we use is also the minimum needed to exhibit the catastrophe, and the choice is usually determined by algebraic convenience.

### FOLD CATASTROPHE

Now replace the screen with a thin lens of refractive index $\mathfrak{n}$ and thickness $w(a, b)$. This introduces an additional contribution to the delay, $t_{\text{lens}} = (\mathfrak{n} - 1)w/c$. The true ray will be fixed by the variation of the sum of the geometric and lens delays with respect to $a$ and $b$ plus any additional state variables that are needed. Suppose that the lens is cylindrical so rays are bent only in the $x$ direction and one state variable, $a$, suffices. Let us use an analytically tractable example, $t_{\text{lens}} = s^2(1 - 2a^2/s^2 + a^4/s^4)/(4fc)$, for $|a| < s$, where $f \gg s$ is the focal length (Fig. 7.13). Place the detector $D$ on the axis with $z < f$. The delay along a path is:

$$t \equiv t_{\text{geo}} + t_{\text{lens}} = \frac{a^2}{2fzc}\left(f - z + \frac{a^2 z}{2s^2}\right), \tag{7.68}$$

dropping a constant. This leaves the single minimum and the true ray at $a = 0$.

Now displace $D$ perpendicular to the axis a distance $x$ with $z$ fixed. $t_{\text{geo}}$ becomes $(a - x)^2/(2zc)$, and the true ray will be parameterized by the single real value of $a$ that minimizes $t$ and therefore solves

$$x = \frac{a}{f}\left(f - z + z\frac{a^2}{s^2}\right). \tag{7.69}$$

The first two terms, $f - z$, represent a perfect thin lens [cf. $J_{11}$ in Eq. (7.57)]; the third represents an imperfection.

A human eye at $D$ [coordinates $(x, y, z)$] focuses the ray through $D$ and adjacent rays onto its retina, producing there a point image of the point source at $S$. The

**FIGURE 7.13** Light from a distant source is normally incident on a thin, phase-changing lens. The phase change depends solely on the distance $a$ from the axis. The delay $t$ along paths encountering a detector $D$ located at $(x, z)$ can be calculated using an equation such as Eq. (7.68). The true rays are located where $t$ is extremized. The envelope created by these rays comprises two fold caustic curves (red) that meet at a cusp point. When $D$ lies outside the caustic, for example at $A$, there is only one true ray (cyan) where $t$ has its single minimum. When $D$ lies inside the caustic, for example at point $B$, there are three true rays associated with two minima (orange, cyan) and one maximum (purple) of $t$. When $D$ lies on the caustic, for example at $C$, there is one minimum (cyan) and a point of inflection (green). The magnification is formally infinite on the caustic. The cusp at the end of the caustic is the focus, a distance $f$ from the screen.

power $P = d\mathcal{E}/dt$ in that image is the energy flux at $D$ times the area inside the eye's iris, so the image power is proportional to the energy flux. Since the energy flux is proportional to the square of the field amplitude $A^2$ and $A \propto 1/\sqrt{\mathcal{A}}$, where $\mathcal{A}$ is the area of a bundle of rays [Eq. (7.35d)], the image power is $P \propto 1/\mathcal{A}$.

Consider a rectangular ray bundle that passes through $\{a, b\}$ on the screen with edges $da$ and $db$, and area $\mathcal{A}_{screen} = dadb$. The bundle's area when it arrives at $D$ is $\mathcal{A}_D = dxdy$. Because the lens is cylindrical, $dy = db$, so the ratio of power seen by an eye at $D$ to that seen by an eye at the screen (i.e., the lens's *magnification*) is $\mathcal{M} = dP_D/dP_{screen} = \mathcal{A}_{screen}/\mathcal{A}_D = da/dx$. Using Eq. (7.69), we find

$$\mathcal{H} \equiv \mathcal{M}^{-1} = \frac{dx}{da} = zc \left(\frac{d^2t}{da^2}\right) = \left(\frac{f - z}{f} + 3\frac{a^2z}{s^2f}\right). \qquad (7.70)$$

If the curvature $\mathcal{H}$ of the delay in the vicinity of a true ray is decreased, the magnification is increased. When $z < f$, the curvature is positive. When $z > f$ the curvature for $a = x = 0$ is negative, and the magnification is $\mathcal{M} = -f/(z - f)$, corresponding to an inverted image. However, there are now two additional images with $a = \pm z^{-1/2}(z - f)^{1/2}s$ at minima with associated magnifications $\mathcal{M} = \frac{1}{2}f(z - f)^{-1}$.

Next move the detector $D$ farther away from the axis. A maximum and a minimum in $t$ will become a point of inflection and, equivalently, two of the images will merge when $a = a_f = \pm 3^{-1/2}z^{-1/2}(z - f)^{1/2}s$ or

$$x = x_f \equiv \mp 2f^{-1}z^{-1/2}(z - f)^{3/2}s. \qquad (7.71)$$

**fold catastrophe**

This is the location of the caustic; in 3-dimensional space it is the surface shown in the first panel of Fig. 7.15 below.

The magnification of the two images will diverge at the caustic and, expanding $\mathcal{H}$ to linear order about zero, we find that $\mathcal{M} = \pm 2^{-1} 3^{-1/4} f^{1/2} s^{1/2} z^{-1/4} (z - f)^{-1/4} |x - x_f|^{-1/2}$ for each image. However, when $|x| > x_f$, the two images vanish. This abrupt change in the optics—two point images becoming infinitely magnified and then vanishing as the detector is moved—is an example of a *fold catastrophe* occurring at a caustic. (Note that the algebraic sum of the two magnifications is zero in the limit.)

It should be pointed out that the divergence of the magnification does not happen in practice for two reasons. The first is that a point source is only an idealization, and if we allow the source to have finite size, different parts will produce caustics at slightly different locations. The second is that geometric optics, on which our analysis is based, pretends that the wavelength of light is vanishingly small. In actuality, the wavelength is always nonzero, and near a caustic its finiteness leads to diffraction effects, which also limit the magnification to a finite value (Sec. 8.6).

Although we have examined one specific and stylized example, the algebraic details can be worked out for any configuration governed by geometric optics. However, they are less important than the scaling laws—for example, $\mathcal{M} \propto |x - x_f|^{-1/2}$—that become increasingly accurate as the catastrophe (caustic) is approached. For this reason, catastrophes are commonly given a *standard form* chosen to exhibit these features and only valid very close to the catastrophe. We discuss this here just in the context of geometrical optics, but the basic scalings are useful in other physics applications (see, e.g., Ex. 7.16).

**standard form for catastrophes**

First we measure the state variables $a, b, \ldots$ in units of some appropriate scale. Next we do likewise for the control parameters, $x, y, \ldots$. We call the new state variables and control parameters $\tilde{a}, \tilde{b}, \ldots$ and $\tilde{x}, \tilde{y} \ldots$, respectively. We then Taylor expand the delay about the catastrophe. In the case of the fold, we want to be able to find up to two extrema. This requires a cubic equation in $\tilde{a}$. The constant is clearly irrelevant, and an overall multiplying factor will not change the scalings. We are also free to change the origin of $\tilde{a}$, allowing us to drop either the linear or the quadratic term (we choose the latter, so that the coefficient is linearly related to $x$). If we adjust the scaled delay in the vicinity of a fold catastrophe, Eq. (7.68) can be written in the standard form:

**for fold catastrophe**

$$\tilde{t}_{\text{fold}} = \frac{1}{3}\tilde{a}^3 - \tilde{x}\tilde{a}, \tag{7.72}$$

where the coefficients are chosen for algebraic convenience. The maximum number of rays involved in the catastrophe is two, and the number of control parameters required is one, which we can think of as being used to adjust the difference in $\tilde{t}$ between two stationary points. The scaled magnifications are now given by $\widetilde{\mathcal{M}} \equiv (d\tilde{x}/d\tilde{a})^{-1} = \pm\frac{1}{2}\tilde{x}^{-1/2}$, and the combined, scaled magnification is $\widetilde{\mathcal{M}} = \tilde{x}^{-1/2}$ for $\tilde{x} > 0$.

## CUSP CATASTROPHE

So far, we have only allowed $D$ to move perpendicular to the axis along $x$. Now move it along the axis toward the screen. We find that $x_f$ decreases with decreasing $z$ until it vanishes at $z = f$. At this point, the central maximum in $t$ merges simultaneously with both minima, leaving a single image. This is an example of a *cusp catastrophe*. Working in 1-dimensional state-variable space with two control parameters and applying the same arguments as we just used with the fold, the standard form for the cusp can be written as

$$\tilde{t}_{\text{cusp}} = \frac{1}{4}\tilde{a}^4 - \frac{1}{2}\tilde{z}\tilde{a}^2 - \tilde{x}\tilde{a}. \qquad (7.73)$$

for cusp catastrophe

The parameter $\tilde{x}$ is still associated with a transverse displacement of $D$, and we can quickly persuade ourselves that $\tilde{z} \propto z - f$ by inspecting the quadratic term in Eq. (7.68).

The cusp then describes a transition between one and three images, one of which must be inverted with respect to the other two. The location of the image for a given $\tilde{a}$ and $\tilde{z}$ is

$$\tilde{x} = \tilde{a}^3 - \tilde{z}\tilde{a}. \qquad (7.74)$$

Conversely, for a given $\tilde{x}$ and $\tilde{z}$, there are one or three real solutions for $\tilde{a}(\tilde{x}, \tilde{z})$ and one or three images. The equation satisfied by the fold lines where the transition occurs is

$$\tilde{x} = \pm\frac{2}{3^{3/2}}\tilde{z}^{3/2}. \qquad (7.75)$$

These are the two branches of a semi-cubical parabola (the caustic surface in 3 dimensions depicted in the second panel of Fig. 7.15 below), and they meet at the cusp catastrophe where $\tilde{x} = \tilde{z} = 0$.

The scaled magnification at the cusp is

$$\widetilde{\mathcal{M}}(\tilde{x}, \tilde{z}) = \left(\frac{\partial\tilde{x}}{\partial\tilde{a}}\right)_{\tilde{z}}^{-1} = [3\tilde{a}(\tilde{x}, \tilde{z})^2 - z]^{-1}. \qquad (7.76)$$

## SWALLOWTAIL CATASTROPHE

Now let the rays propagate in 3 dimensions, so that there are three control variables, $x$, $y$, and $z$, where the $y$-axis is perpendicular to the $x$- and $z$-axes. The fold, which was a point in 1 dimension and a line in 2, becomes a surface in 3 dimensions, and the point cusp in 2 dimensions becomes a line in 3 (see Fig. 7.15 below). If there is still only one state variable, then $\tilde{t}$ should be a quintic with up to four extrema. (In general, a catastrophe involving as many as $N$ images requires $N - 1$ control parameters for its full description. These parameters can be thought of as independently changing the relative values of $\tilde{t}$ at the extrema.) The resulting, four-image catastrophe is called a *swallowtail*. (In practice, this catastrophe only arises when there are two state variables, and additional images are always present. However, these are not involved in the

catastrophe.) Again following our procedure, we can write the standard form of the swallowtail catastrophe as

**for swallowtail catastrophe**

$$\tilde{t}_{\text{swallowtail}} = \frac{1}{5}\tilde{a}^5 - \frac{1}{3}\tilde{z}\tilde{a}^3 - \frac{1}{2}\tilde{y}\tilde{a}^2 - \tilde{x}\tilde{a}. \tag{7.77}$$

There are two cusp lines in the half-space $\tilde{z} > 0$, and these meet at the catastrophe where $\tilde{x} = \tilde{y} = \tilde{z} = 0$ (see Fig. 7.15 below). The relationship between $\tilde{x}$, $\tilde{y}$, and $\tilde{z}$ and $x$, $y$, and $z$ in this or any other example is not simple, and so the variation of the magnification in the vicinity of the swallowtail catastrophe depends on the details.

### HYPERBOLIC UMBILIC CATASTROPHE

Next increase the number of essential state variables to two. We can choose these to be $\tilde{a}$ and $\tilde{b}$. To see what is possible, sketch contours of $\Delta t$ in the $\tilde{a}$-$\tilde{b}$ plane for fixed values of the control variables. The true rays will be associated with maxima, minima, or saddle points, and each distinct catastrophe corresponds to a different way to nest the contours (Fig. 7.14). The properties of the fold, cusp, and swallowtail are essentially unchanged by the extra dimension. We say that they are *structurally stable*. However, a little geometric experimentation uncovers a genuinely 2-dimensional nesting. The

*hyperbolic umbilic* catastrophe has two saddles, one maximum and one minimum. Further algebraic experiment produces a standard form:

**for hyperbolic umbilic catastrophe**

$$\tilde{t} = \frac{1}{3}(\tilde{a}^3 + \tilde{b}^3) - \tilde{z}\tilde{a}\tilde{b} - \tilde{x}\tilde{a} - \tilde{y}\tilde{b}. \tag{7.78}$$

This catastrophe can be exhibited by a simple generalization of our example. We replace the cylindrical lens described by Eq. (7.68) with a nearly circular lens where the focal length $f_a$ for rays in the $a$-$z$ plane differs from the focal length $f_b$ for rays in the $b$-$z$ plane.

$$t = \frac{a^2}{2f_a zc}\left(f_a - z + \frac{a^2 z}{2s_a^2}\right) + \frac{b^2}{2f_b zc}\left(f_b - z + \frac{b^2 z}{2s_b^2}\right). \tag{7.79}$$

**astigmatism**

This is an example of *astigmatism*. A pair of fold surfaces is associated with each of these foci. These surfaces can cross, and when this happens a cusp line associated with one fold surface can transfer onto the other fold surface. The point where this happens is the hyperbolic umbilic catastrophe.

This example also allows us to illustrate a simple feature of magnification. When the source and the detector are both treated as 2-dimensional, then we generalize the curvature to the *Hessian* matrix

**magnification matrix**

$$\widetilde{\mathcal{H}} = \widetilde{\mathcal{M}}^{-1} = \begin{pmatrix} \frac{\partial \tilde{x}}{\partial \tilde{a}} & \frac{\partial \tilde{y}}{\partial \tilde{a}} \\ \frac{\partial \tilde{x}}{\partial \tilde{b}} & \frac{\partial \tilde{y}}{\partial \tilde{b}} \end{pmatrix}. \tag{7.80}$$

The magnification matrix $\widetilde{\mathcal{M}}$, which describes the mapping from the source plane to the image plane, is simply the inverse of $\widetilde{\mathcal{H}}$. The four matrix elements also describe the deformation of the image. As we describe in more detail when discussing elastostatics (Sec. 11.2.2), the antisymmetric part of $\widetilde{\mathcal{M}}$ describes the rotation of the image, and

**FIGURE 7.14**  Distinct nestings of contours of $\tilde{t}$ in state-variable space. (a) A 1-dimensional arrangement of two saddle points and a maximum in the vicinity of a cusp. The locations of the extrema and their curvatures change as the control parameters change. (b) A cusp formed by two minima and a saddle. Although the nestings look different in 2 dimensions, this is essentially the same catastrophe when considered in 1 dimension, which is all that is necessary to determine its salient properties. These are the only contour nestings possible with two state variables and three extrema (or rays). (c) When we increase the number of extrema to four, two more nestings are possible. The swallowtail catastrophe is essentially a cusp with an additional extremum added to the end, requiring three control parameters to express. It, too, is essentially 1-dimensional. (d) The hyperbolic umbilic catastrophe is essentially 2-dimensional and is associated with a maximum, a minimum, and two saddles. A distinct nesting of contours with three saddle points and one extremum occurs in the elliptic umbilic catastrophe (Ex. 7.13).

7.5  Catastrophe Optics          **391**

the symmetric part its magnification and stretching or *shear*. Both eigenvalues of $\widetilde{\mathcal{M}}$
are positive at a minimum, and the image is a distorted version of the source. At a
saddle, one eigenvalue is positive, the other negative, and the image is inverted; at a
maximum, they are both negative, and the image is doubly inverted so that it appears
to have been rotated through a half-turn.

### ELLIPTIC UMBILIC CATASTROPHE

There is a second standard form that can describe the nesting of contours just
discussed—a distinct catastrophe called the *elliptic umbilic catastrophe* (Ex. 7.13b):

**standard form for elliptic umbilic catastrophe**

$$\tilde{t} = \frac{1}{3}\tilde{a}^3 - \tilde{a}\tilde{b}^2 - \tilde{z}(\tilde{a}^2 + \tilde{b}^2) - \tilde{x}\tilde{a} - \tilde{y}\tilde{b}. \tag{7.81}$$

The caustic surfaces in three dimensions $(\tilde{x}, \tilde{y}, \tilde{z})$ for the five elementary catastro-
phes discussed here are shown in Fig. 7.15. Additional types of catastrophe are found
with more control parameters, for example, time (e.g., Poston and Stewart, 2012). This
is relevant, for example, to the twinkling of starlight in the geometric-optics limit.

### EXERCISES

**Exercise 7.13** *Derivation and Problem: Cusps and Elliptic Umbilics* **T2**

(a) Work through the derivation of Eq. (7.73) for the scaled time delay in the vicinity
of the cusp caustic for our simple example [Eq. (7.68)], with the aid of a suitable
change of variables (Goodman, Romani, Blandford, and Narayan, 1987, Appen-
dix B).

(b) Sketch the nesting of the contours for the elliptic umbilic catastrophe as shown
for the other four catastrophes in Fig. 7.14. Verify that Eq. (7.81) describes this
catastrophe.

**Exercise 7.14** *Problem: Cusp Scaling Relations* **T2**

Consider a cusp catastrophe created by a screen as in the example and described by a
standard cusp potential, Eq. (7.73). Suppose that a detector lies between the folds, so
that there are three images of a single point source with state variables $\tilde{a}_i$.

(a) Explain how, in principle, it is possible to determine $\tilde{a}$ for a single image by
measurements at $D$.

(b) Make a 3-dimensional plot of the location of the image(s) in $\tilde{a}$-$\tilde{x}$-$\tilde{y}$ space and
explain why the names "fold" and "cusp" were chosen.

(c) Prove as many as you can of the following scaling relations, valid in the limit as
the cusp catastrophe is approached:

$$\sum_{i=1}^{3} \tilde{a}_i = 0, \quad \sum_{i=1}^{3} \frac{1}{\tilde{a}_i} = -\frac{\tilde{z}}{\tilde{x}}, \quad \sum_{i=1}^{3} \widetilde{\mathcal{M}}_i = 0, \quad \sum_{i=1}^{3} \tilde{a}_i\widetilde{\mathcal{M}}_i = 0,$$

$$\sum_{i=1}^{3} \tilde{a}_i^2\widetilde{\mathcal{M}}_i = 1, \quad \sum_{i=1}^{3} \tilde{a}_i^3\widetilde{\mathcal{M}}_i = 0 \quad \text{and} \quad \sum_{i=1}^{3} \tilde{a}_i^4\widetilde{\mathcal{M}} = \tilde{z}. \tag{7.82}$$

$$\tilde{t} = \tfrac{1}{3}\tilde{a}^3 - \tilde{x}\,\tilde{a}$$

$$\tilde{t} = \tfrac{1}{4}\tilde{a}^4 - \tfrac{1}{2}\tilde{z}\,\tilde{a}^2 - \tilde{x}\,\tilde{a}$$

fold

cusp

$$\tilde{t} = \tfrac{1}{5}\tilde{a}^5 - \tfrac{1}{3}\tilde{z}\,\tilde{a}^3 - \tfrac{1}{2}\tilde{y}\,\tilde{a}^2 - \tilde{x}\,\tilde{a}$$

$$\tilde{t} = \tfrac{1}{3}\tilde{a}^3 - \tilde{a}\,\tilde{b}^2 - \tilde{z}\,(\tilde{a}^2 + \tilde{b}^2) - \tilde{x}\,\tilde{a} - \tilde{y}\,\tilde{b}$$

swallowtail

elliptic umbilic

$$\tilde{t} = \tfrac{1}{3}(\tilde{a}^3 + \tilde{b}^3) - \tilde{z}\,\tilde{a}\,\tilde{b} - \tilde{x}\,\tilde{a} - \tilde{y}\,\tilde{b}$$

hyperbolic umbilic

**FIGURE 7.15** The five elementary catastrophes (caustic structures) that are possible for a set of light rays specified by one or two state varables $\{\tilde{a}, \tilde{b}\}$ in 3-dimensional space with coordinates (control parameters) $\{\tilde{x}, \tilde{y}, \tilde{z}\}$. The surfaces represent the loci of points of infinite magnification assuming a point source and geometric optics. The actual caustic surfaces will be deformed versions of these basic shapes. The hyperbolic umbilic surfaces are shown from two different viewpoints.

[Hint: You must retain the sign of the magnification.] Of course, not all of these are useful. However, relations like these exist for all catastrophes and are increasingly accurate as the separation of the images becomes much smaller than the scale of variation of $\tilde{t}$.

**Exercise 7.15** *Problem: Wavefronts* T2

As we have emphasized, representing light using wavefronts is complementary to treating it in terms of rays. Sketch the evolution of the wavefronts after they propagate through a phase-changing screen and eventually form caustics. Do this for a 2-dimensional cusp, and then consider the formation of a hyperbolic umbilic catastrophe by an astigmatic lens.

**Exercise 7.16** **Example: Van der Waals Catastrophe* T2

The van der Waals equation of state $(P + a/v^2)(v - b) = k_B T$ for $H_2O$ relates the pressure $P$ and specific volume (volume per molecule) $v$ to the temperature $T$; see Sec. 5.7. Figure 5.8 makes it clear that, at some temperatures $T$ and pressures $P$, there are three allowed volumes $v(T, P)$, one describing liquid water, one water vapor, and the third an unstable phase that cannot exist in Nature. At other values of $T$ and $P$, there is only one allowed $v$. The transition between three allowed $v$ values and one occurs along some curve in the $T$-$P$ plane—a catastrophe curve.

(a) This curve must correspond to one of the elementary catastrophes explored in the previous exercise. Based on the number of solutions for $v(T, P)$, which catastrophe must it be?

(b) Change variables in the van der Waals equation of state to $p = P/P_c - 1$, $\tau = T/T_c - 1$, and $\rho = v_c/v - 1$, where $T_c = 8a/(27bk_B)$, $P_c = a/(27b^2)$, and $v_c = 3b$ are the temperature, pressure, and specific volume at the critical point $C$ of Fig. 5.8. Show that this change of variables brings the van der Waals equation of state into the form

$$\rho^3 - z\rho - x = 0, \tag{7.83}$$

where $z = -(p/3 + 8\tau/3)$ and $x = 2p/3 - 8\tau/3$.

(c) This equation $\rho^3 - z\rho - x$ is the equilibrium surface associated with the catastrophe-theory potential $t(\rho; x, z) = \frac{1}{4}\rho^4 - \frac{1}{2}z\rho^2 - x\rho$ [Eq. (7.73)]. Correspondingly, the catastrophe [the boundary between three solutions $v(T, P)$ and one] has the universal cusp form $x = \pm 2(z/3)^{2/3}$ [Eq. (7.75)]. Plot this curve in the temperature-pressure plane.

Note that we were guaranteed by catastrophe theory that the catastrophe curve would have this form near its cusp point. However, it is a surprise and quite unusual that, for

the van der Waals case, the cusp shape $x = \pm 2(z/3)^{2/3}$ is not confined to the vicinity of the cusp point but remains accurate far from that point.

---

### 7.5.2 Aberrations of Optical Instruments T2

Much computational effort is expended in the design of expensive optical instruments prior to prototyping and fabrication. This is conventionally discussed in terms of *aberrations,* which provide a perturbative description of rays that complements the singularity-based approach of catastrophe theory. While it is possible to design instruments that take all the rays from a point source $S$ and focus them geometrically onto a point detector $D$,[11] this is not what is demanded of them in practice. Typically, they have to map an extended image onto an extended surface, for example, a CCD detector. Sometimes the source is large, and the instrument must achieve a large *field of view;* sometimes it is small, and image fidelity close to the axis matters. Sometimes light levels are low, and transmission losses must be minimized. Sometimes the bandwidth of the light is large, and the variation of the imaging with frequency must be minimized. Sometimes diffractive effects are important. The residual imperfections of an instrument are known as *aberrations.*

**aberrations**

As we have shown, any (geometric-optics) instrument will map, one to many, source points onto detector points. This mapping is usually expanded in terms of a set of basis functions, and several choices are in use, for example, those due to Seidel and Zernike (e.g., Born and Wolf, 1999, Secs. 5.3, 9.2). If we set aside effects caused by the variation of the refractive index with wavelength, known as *chromatic aberration,* there are five common types of geometrical aberration. *Spherical aberration* is the failure to bring a point on the optic axis to a single focus. Instead, an axisymmetric cusp/fold caustic is created. We have already exhibited *astigmatism* in our discussion of the hyperbolic umbilic catastrophe with a non-axisymmetric lens and an axial source (Sec. 7.5.1). It is not hard to make axisymmetric lenses and mirrors, so this does not happen much in practice. However, as soon as we consider off-axis surfaces, we break the symmetry, and astigmatism is unavoidable. *Curvature* arises when the surface on which the rays from point sources are best brought to a focus lies on a curved surface, not on a plane. It is sometimes advantageous to accept this aberration and to curve the detector surface.[12] To understand *coma,* consider a small pencil of rays from an off-axis source that passes through the center of an instrument and is brought to a focus. Now consider a cone of rays about this pencil that passes through the periphery of the lens. When there is coma, these rays will on average be displaced

**chromatic aberration**

**spherical aberration**

**curvature**

**coma**

---

11. A simple example is to make the interior of a prolate ellipsoidal detector perfectly reflecting and to place $S$ and $D$ at the two foci, as crudely implemented in whispering galleries.

12. For example, in a traditional Schmidt telescope.

**distortion**

radially. Coma can be ameliorated by reducing the aperture. Finally, there is *distortion,* in which the sides of a square in the source plane are pushed in (*pin cushion*) or out (*barrel*) in the image plane.

**7.6**

## 7.6 Gravitational Lenses  T2

**7.6.1**

### 7.6.1 Gravitational Deflection of Light  T2

Albert Einstein's general relativity theory predicts that light rays should be deflected by the gravitational field of the Sun (Ex. 27.3; Sec. 27.2.3). Newton's law of gravity combined with his corpuscular theory of light also predicts this deflection, but through an angle half as great as relativity predicts. A famous measurement, during a 1919 solar eclipse, confirmed the relativistic prediction, thereby making Einstein world famous.

The deflection of light by gravitational fields allows a cosmologically distant galaxy to behave like a crude lens and, in particular, to produce multiple images of a more distant quasar. Many examples of this phenomenon have been observed. The optics of these gravitational lenses provides an excellent illustration of the use of Fermat's principle (e.g., Blandford and Narayan, 1992; Schneider, Ehlers, and Falco, 1992). We explore these issues in this section.

The action of a gravitational lens can only be understood properly using general relativity. However, when the gravitational field is weak, there exists an equivalent Newtonian model, due to Eddington (1919), that is adequate for our purposes. In this model, curved spacetime behaves as if it were spatially flat and endowed with a refractive index given by

**refractive index model for gravitational lensing**

$$\mathfrak{n} = 1 - \frac{2\Phi}{c^2},$$  (7.84)

where $\Phi$ is the Newtonian gravitational potential, normalized to vanish far from the source of the gravitational field and chosen to have a negative sign (so, e.g., the field at a distance $r$ from a point mass $M$ is $\Phi = -GM/r$). Time is treated in the Newtonian manner in this model. In Sec. 27.2.3, we use a general relativistic version of Fermat's principle to show that for static gravitational fields this index-of-refraction model gives the same predictions as general relativity, up to fractional corrections of order $|\Phi|/c^2$, which are $\lesssim 10^{-5}$ for the lensing examples in this chapter.

**second refractive index model**

A second Newtonian model gives the same predictions as this index-of-refraction model to within its errors, $\sim |\Phi|/c^2$. We deduce it by rewriting the ray equation (7.48) in terms of Newtonian time $t$ using $ds/dt = C = c/\mathfrak{n}$. The resulting equation is $(\mathfrak{n}^3/c^2)d^2\mathbf{x}/dt^2 = \nabla\mathfrak{n} - 2(\mathfrak{n}/c)^2(d\mathfrak{n}/dt)d\mathbf{x}/dt$. The second term changes the length of the velocity vector $d\mathbf{x}/dt$ by a fractional amount of order $|\Phi|/c^2 \lesssim 10^{-5}$ (so as to keep the length of $d\mathbf{x}/ds$ unity). This is of no significance for our Newtonian model, so we drop this term. The factor $\mathfrak{n}^3 \simeq 1 - 6\Phi/c^2$ produces a fractional correction to $d^2\mathbf{x}/dt^2$ that is of the same magnitude as the fractional errors in our index of refraction

model, so we replace this factor by one. The resulting equation of motion for the ray is

$$\frac{d^2\mathbf{x}}{dt^2} = c^2 \nabla \mathfrak{n} = -2\nabla\Phi. \tag{7.85}$$

Equation (7.85) says that the photons that travel along rays feel a Newtonian gravitational potential that is twice as large as the potential felt by low-speed particles; the photons, moving at speed $c$ (aside from fractional changes of order $|\Phi|/c^2$), respond to that doubled Newtonian field in the same way as any Newtonian particle would. The extra deflection is attributable to the geometry of the spatial part of the metric being non-Euclidean (Sec. 27.2.3).

### 7.6.2 Optical Configuration  T2

To understand how gravitational lenses work, we adopt some key features from our discussion of optical catastrophes formed by an intervening screen. However, there are some essential differences.

- The source is not assumed to be distant from the screen (which we now call a lens, $L$).

- Instead of tracing rays emanating from a point source $S$, we consider a *congruence* of rays emanating from the observer $O$ (i.e., us) and propagating backward in time past the lens to the sources. This is because there are many stars and galaxies whose images will be distorted by the lens. The caustics envelop the sources.

- The universe is expanding, which makes the optics formally time-dependent. However, as we discuss in Sec. 28.6.2, we can work in comoving coordinates and still use Fermat's principle. For the moment, we introduce three distances: $d_{OL}$ for distance from the observer to the lens, $d_{OS}$ for the distance from the observer to the source, and $d_{LS}$ for the distance from the lens to the source. We evaluate these quantities cosmologically in Sec. 28.6.2.

- Instead of treating $a$ and $b$ as the state variables that describe rays (Sec. 7.5.1), we use a 2-dimensional (small) angular vector $\boldsymbol{\theta}$ measuring the image position on the sky. We also replace the control parameters $x$ and $y$ with the 2-dimensional angle $\boldsymbol{\beta}$, which measures the location that the image of the source would have in the absence of the lens. We can also treat the distance $d_{OS}$ as a third control parameter replacing $z$.

The Hessian matrix, replacing Eq. (7.80), is now the Jacobian of the vectorial angles that a small, finite source would subtend in the absence and in the presence of the lens:

$$\mathcal{H} = \mathcal{M}^{-1} = \frac{\partial\boldsymbol{\beta}}{\partial\boldsymbol{\theta}}. \tag{7.86}$$

As the specific intensity $I_\nu = dE/dAdtd\nu d\Omega$ is conserved along a ray (see Sec. 3.6), the determinant of $\mathcal{H}$ is just the ratio of the flux of energy per unit frequency without the lens to the flux with the lens, and correspondingly the determinant of $\mathcal{M}$ (the scalar magnification) is the ratio of flux with the lens to that without the lens.

7.6.3

### 7.6.3 Microlensing

Our first example of a gravitational lens is a point mass—specifically, a star. This phenomenon is known as *microlensing,* because the angles of deflection are typically microarcseconds.[13] The source is also usually another star, which we also treat as a point.

We first compute the deflection of a Newtonian particle with speed $v$ passing by a mass $M$ with impact parameter $b$. By computing the perpendicular impulse, it is straightforward to show that the deflection angle is $2GM/v^2$. Replacing $v$ by $c$ and doubling the answer gives the small deflection angle for light:

**microlensing deflection angle**

$$\boldsymbol{\alpha} = \frac{4GM}{bc^2} = 1.75 \left( \frac{M}{M_\odot} \right) \left( \frac{b}{R_\odot} \right)^{-1} \hat{\mathbf{b}} \text{ arcsec,} \tag{7.87}$$

where $\hat{\mathbf{b}}$ is a unit vector along the impact parameter, which allows us to treat the deflection as a 2-dimensional vector like $\boldsymbol{\theta}$ and $\boldsymbol{\beta}$. $M_\odot$ and $R_\odot$ are the solar mass and radius, respectively.

**microlensing lens equation**

The imaging geometry can be expressed as a simple vector equation called the *lens equation* (Fig. 7.16):

$$\boldsymbol{\theta} = \boldsymbol{\beta} + \frac{d_{LS}}{d_{OS}} \boldsymbol{\alpha}. \tag{7.88}$$

A point mass exhibits circular symmetry, so we can treat this equation as a scalar equation and rewrite it in the form

$$\theta = \beta + \frac{\theta_E^2}{\theta}, \tag{7.89}$$

where

$$\theta_E = \left( \frac{4GM}{d_{\text{eff}} c^2} \right)^{1/2} = 903 \left( \frac{M}{M_\odot} \right)^{1/2} \left( \frac{d_{\text{eff}}}{10 \text{ kpc}} \right)^{-1/2} \mu \text{ arcsec} \tag{7.90}$$

is the *Einstein radius,* and

$$d_{\text{eff}} = \frac{d_{OL} d_{OS}}{d_{LS}} \tag{7.91}$$

is the *effective distance.* (Here 10 kpc means 10 kiloparsecs, about 30,000 light years.)

---

13. Interestingly, Newton speculated that light rays could be deflected by gravity, and the underlying theory of microlensing was worked out correctly by Einstein in 1912, before he realized that the deflection was twice the Newtonian value.

**FIGURE 7.16** Geometry for microlensing of a stellar source $S$ by a stellar lens $L$ observed at $O$.

The solutions to this quadratic equation are

$$\theta_\pm = \frac{\beta}{2} \pm \sqrt{\theta_E{}^2 + \left(\frac{\beta}{2}\right)^2}. \tag{7.92}$$

**image locations**

The magnification of the two images can be computed directly by evaluating the reciprocal of the determinant of $\mathcal{H}$ from Eq. (7.86). However, it is quicker to exploit the circular symmetry and note that the element of source solid angle in polar coordinates is $\beta d\beta d\phi$, while the element of image solid angle is $\theta d\theta d\phi$, so that

$$\mathcal{M} = \frac{\theta \, d\theta}{\beta \, d\beta} = \frac{1}{1 - (\theta_E/\theta)^4}. \tag{7.93}$$

**image magnifications**

The eigenvalues of the magnification matrix are $[1 + (\theta_E/\theta)^2]^{-1}$ and $[1 - (\theta_E/\theta)^2]^{-1}$. The former describes the radial magnification, the latter the tangential magnification. As $\beta \to 0, \theta \to \theta_E$ for both images. If we consider a source of finite angular size, when $\beta$ approaches this size then the tangential stretching is pronounced and two nearly circular arcs are formed on opposite sides of the lens, one with $\theta > \theta_E$; the other, inverted, with $\theta < \theta_E$. When $\beta$ is reduced even more, the arcs join up to form an *Einstein ring* (Fig. 7.17).

**Einstein ring**

Astronomers routinely observe microlensing events when stellar sources pass behind stellar lenses. They are unable to distinguish the two images and so measure a combined magnification $\mathcal{M} = |\mathcal{M}_+| + |\mathcal{M}_-|$. If we substitute $\beta$ for $\theta$, then we have

$$\mathcal{M} = \frac{(\theta_E^2 + \frac{1}{2}\beta^2)}{(\theta_E^2 + \frac{1}{4}\beta^2)^{1/2}\beta}. \tag{7.94}$$

Note that in the limit of large magnifications, $\mathcal{M} \sim \theta_E/\beta$, and the probability that the magnification exceeds $\mathcal{M}$ is proportional to the cross sectional area $\pi\beta^2 \propto \mathcal{M}^{-2}$ (cf. Fig. 7.16).

If the speed of the source relative to the continuation of the $O$-$L$ line is $v$, and the closest approach to this line is $h$, which happens at time $t = 0$, then $\beta = (h^2 + v^2t^2)^{1/2}/d_{OS}$, and then there is a one-parameter family of magnification curves (shown in Fig. 7.18). The characteristic variation of the magnification can be used to distinguish this phenomenon from intrinsic stellar variation. This behavior is very

**FIGURE 7.17** The source LRG 3-757, as imaged by the Hubble Space Telescope. The blue Einstein ring is the image of two background galaxies formed by the gravitational field associated with the intervening central (yellow) lens galaxy. The accurate alignment of the lens and source galaxies is quite unusual. (ESA/Hubble and NASA.)



**FIGURE 7.18** Variation of the combined magnifications of a stellar source as it passes behind a lens star. The time $t$ is measured in units of $d_{OS}\theta_E/v$, and the parameter that labels the curves is $h/(d_{OS}\theta_E)$.

different from that at a generic caustic (Sec. 7.5.1) and is not structurally stable: if the axisymmetry is broken, then the behavior will change significantly. Nevertheless, for finite-sized sources of light and stars or nearly circular galaxies, stable Einstein rings are commonly formed.

**Exercise 7.17**  *Example: Microlensing Time Delay*  T2

An alternative derivation of the lens equation for a point-mass lens, Eq. (7.88), evaluates the time delay along a path from the source to the observer and finds the true ray by extremizing it with respect to variations of $\theta$ [cf. Eq. (7.68)].

(a) Show that the geometric time delay is given by

$$t_{\text{geo}} = \frac{1}{2c} d_{\text{eff}} (\boldsymbol{\theta} - \boldsymbol{\beta})^2. \tag{7.95}$$

(b) Next show that the lens time delay can be expressed as

$$t_{\text{lens}} = -(4GM/c^3) \ln b + \text{const},$$

where $b$ is the impact parameter. (It will be helpful to evaluate the difference in delays between two rays with differing impact parameters.) This is known as the *Shapiro delay* and is discussed further in Sec. 27.2.4.

(c) Show that the lens delay can also be written as

$$t_{\text{lens}} = -\frac{2}{c^3} \int dz \Phi = -\frac{2}{c^3} \Phi_2, \tag{7.96}$$

where $\Phi_2$ is the surface gravitational potential obtained by integrating the 3-dimensional potential $\Phi$ along the path. The surface potential is only determined up to an unimportant, divergent constant, which is acceptable because we are only interested in $dt_{\text{lens}}/db$ which is finite.

(d) By minimizing $t_{\text{geo}} + t_{\text{lens}}$, derive the lens equation (7.88).

**Exercise 7.18**  *Derivation: Microlensing Variation*  T2

Derive Eq. (7.94).

**Exercise 7.19**  *Problem: Magnification by a Massive Black Hole*  T2

Suppose that a large black hole forms two images of a background source separated by an angle $\theta$. Let the fluxes of the two images be $F_+$ and $F_- < F_+$. Show that the flux from the source would be $F_+ - F_-$ if there were no lens and that the black hole should be located an angular distance $[1 + (F_-/F_+)^{-1/2}]^{-1}\theta$ along the line from the brighter image to the fainter one. (Only consider small angle deflections.)

### 7.6.4  Lensing by Galaxies  T2

7.6.4

Most observed gravitational lenses are galaxies. Observing these systems brings out new features of the optics and proves useful for learning about galaxies and the universe. Galaxies comprise dark matter and stars, and the dispersion $\langle v_{\parallel}^2 \rangle$ in the

stars' velocities along the line of sight can be measured using spectroscopy. The virial
theorem (Goldstein, Poole, and Safko, 2002) tells us that the kinetic energy of the
matter in the galaxy is half the magnitude of its gravitational potential energy $\Phi$. We
can therefore make an order of magnitude estimate of the ray deflection angle caused
by a galaxy by using Eq. (7.87):

$$\alpha \sim \frac{4GM}{bc^2} \sim \frac{4|\Phi|}{c^2} \sim 4 \times 2 \times \frac{3}{2} \times \frac{\langle v_{||}^2 \rangle}{c^2} \sim \frac{12\langle v_{||}^2 \rangle}{c^2}. \tag{7.97}$$

This evaluates to $\alpha \sim 2$ arcsec for a typical galaxy velocity dispersion of $\sim 300$ km
$s^{-1}$. The images can typically be resolved using radio and optical telescopes, but their
separations are much less than the full angular sizes of distant galaxies and so the
imaging is sensitive to the lens galaxy's structure. As the lens is typically far more
complicated than a point mass, it is now convenient to measure the angle $\boldsymbol{\theta}$ relative to
the reference ray that connects us to the source.

   We can describe the optics of a galaxy gravitational lens by adapting the formal-
ism that we developed for a point-mass lens in Ex. 7.17. We assume that there is a
point source at $\boldsymbol{\beta}$ and consider paths designated by $\boldsymbol{\theta}$. The geometrical time delay is
unchanged. In Ex. 7.17, we showed that the lens time delay for a point mass was pro-
portional to the surface gravitational potential $\Phi_2$. A distributed lens is handled by
adding the potentials associated with all the point masses out of which it can be con-
sidered as being composed. In other words, we simply use the surface potential for
the distributed mass in the galaxy,

$$t = t_{\text{geo}} + t_{\text{lens}} = \frac{d_{\text{eff}}}{2c}(\boldsymbol{\theta} - \boldsymbol{\beta})^2 - \frac{2}{c^3}\Phi_2 = \frac{d_{\text{eff}}\tilde{t}}{c}, \tag{7.98}$$

where the *scaled time delay t* is defined by $\tilde{t}(\boldsymbol{\theta}; \boldsymbol{\beta}) = \frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\beta})^2 - \Psi(\boldsymbol{\theta})$, and $\Psi = 2\Phi_2/(c^2 d_{\text{eff}})$. The quantity $\Psi$ satisfies the 2-dimensional Poisson equation:

$$\nabla_{2,\theta}^2 \Psi = \frac{8\pi G \Sigma}{d_{\text{eff}} c^2}, \tag{7.99}$$

where $\Sigma$ is the density of matter per unit solid angle, and the 2-dimensional laplacian
describes differentiation with respect to the components of $\boldsymbol{\theta}$.[14]

   As written, Eq. (7.98) describes all paths for a given source position, only a small
number of which correspond to true rays. However, if, instead, we set $\boldsymbol{\beta} = 0$ and
choose any convenient origin for $\boldsymbol{\theta}$ so that

---

14. The minimum surface density (expressed as mass per area) of a cosmologically distant lens needed to
produce multiple images of a background source turns out to be $\sim 1$ g cm$^{-2}$. It is remarkable that such a
seemingly small surface density operating on these scales can make such a large difference to our view
of the universe. It is tempting to call this a "rule of thumb," because it is roughly the column density
associated with one's thumb!

$$\tilde{t}(\boldsymbol{\theta}) = \frac{1}{2}\theta^2 - \Psi(\boldsymbol{\theta}), \tag{7.100}$$

then three useful features of $\tilde{t}$ emerge:

- if there is an image at $\boldsymbol{\theta}$, then the source position is simply $\boldsymbol{\beta} = \nabla_{\theta}\tilde{t}$ [cf. Eq. (7.88)];

- the magnification tensor associated with this image can be calculated by taking the inverse of the Hessian matrix $\mathcal{H} = \nabla_{\theta}\nabla_{\theta}\tilde{t}$ [cf. Eq. (7.86)]; and

- the measured differences in the times of variation observed in multiple images of the same source are just the differences in $d_{eff}\tilde{t}/c$ evaluated at the image positions.[15]

Computing $\tilde{t}$ for a model of a putative lens galaxy allows one to assess whether background sources are being multiply imaged and, if so, to learn about the lens as well as the source.

**EXERCISES**

**Exercise 7.20** *Problem: Catastrophe Optics of an Elliptical Gravitational Lens* T2
Consider an elliptical gravitational lens where the potential $\Psi$ is modeled by

$$\Psi(\boldsymbol{\theta}) = (1 + A\theta_x^2 + 2B\theta_x\theta_y + C\theta_y^2)^q; \quad 0 < q < 1/2. \tag{7.101}$$

Determine the generic form of the caustic surfaces, the types of catastrophe encountered, and the change in the number of images formed when a point source crosses these surfaces. Note that it is in the spirit of catastrophe theory *not* to compute exact expressions but to determine scaling laws and to understand the qualitative features of the images.

**Exercise 7.21** *Challenge: Microlensing in a Galaxy* T2
Our discussion of microlensing assumed a single star and a circularly symmetric potential about it. This is usually a good approximation for stars in our galaxy. However, when the star is in another galaxy and the source is a background quasar (Figs. 7.19, 7.20), it is necessary to include the gravitational effects of the galaxy's other stars and its dark matter. Recast the microlensing analysis (Sec. 7.6.3) in the potential formulation of Eq. (7.98) and add *external magnification* and *external shear* contributions to

---

15. These differences can be used to measure the size and age of the universe. To order of magnitude, the relative delays are the times it takes light to cross the universe (or equivalently, the age of the universe, roughly 10 Gyr) times the square of the scattering angle (roughly 2 arcsec or $\sim 10^{-5}$ radians), which is roughly 1 year. This is very convenient for astronomers (see Fig. 7.20).

FIGURE 7.19   Geometry for gravitational lensing of a quasar source $S$ by a galaxy lens $L$ observed at $O$.



FIGURE 7.20   Gravitational lens in which a distant quasar, Q2237+0305, is quadruply imaged by an intervening galaxy. The four quasar images are denoted A, B, C, and D. The galaxy is much larger than the separation of the images (1–1.5 arcsec) but its bright central core is labeled. There is also a fifth, faint image coincident with this core. There are many examples of gravitational lenses like this, where the source region is very compact and variable so that the delay in the variation seen in the individual images can be used to measure the distance of the source and the size and age of the universe. In addition, microlensing-like variations in the images are induced by individual stars in the lens galaxy moving in front of the quasar. Analyzing these changes can be used to measure the proportion of stars and dark matter in galaxies. Adapted from image by Hubble Space Telescope. (NASA, ESA, STScI.)

$\Psi$ that are proportional to $\theta_x^2 + \theta_y^2$ and $\theta_x^2 - \theta_y^2$, respectively. The latter will break the circular symmetry, and structurally stable caustics will be formed. Explore the behavior of these caustics as you vary the strength and sign of the magnification and shear contributions. Plot a few flux variations that might be observed.

## 7.7 Polarization

In our geometric-optics analyses thus far, we have either dealt with a scalar wave (e.g., a sound wave) or simply supposed that individual components of vector or tensor waves can be treated as scalars. For most purposes, this is indeed the case, and we continue to use this simplification in the following chapters. However, there are some important wave properties that are unique to vector (or tensor) waves. Most of these come under the heading of *polarization* effects. In Secs. 27.3, 27.4, and 27.5, we study polarization effects for (tensorial) gravitational waves. Here and in Secs. 10.5 and 28.6.1, we examine them for electromagnetic waves.

An electromagnetic wave's two polarizations are powerful tools for technology, engineering, and experimental physics. However, we forgo any discussion of this in the present chapter. Instead we focus solely on the geometric-optics propagation law for polarization (Sec. 7.7.1) and an intriguing aspect of it—the geometric phase (Sec. 7.7.2).

### 7.7.1 Polarization Vector and Its Geometric-Optics Propagation Law

A plane electromagnetic wave in a vacuum has its electric and magnetic fields $\mathbf{E}$ and $\mathbf{B}$ perpendicular to its propagation direction $\hat{\mathbf{k}}$ and perpendicular to each other. In a medium, $\mathbf{E}$ and $\mathbf{B}$ may or may not remain perpendicular to $\hat{\mathbf{k}}$, depending on the medium's properties. For example, an Alfvén wave has its vibrating magnetic field perpendicular to the background magnetic field, which can make an arbitrary angle with respect to $\hat{\mathbf{k}}$. By contrast, in the simplest case of an isotropic dielectric medium, where the dispersion relation has our standard dispersion-free form $\Omega = (c/\mathfrak{n})k$, the group and phase velocities are parallel to $\hat{\mathbf{k}}$, and $\mathbf{E}$ and $\mathbf{B}$ turn out to be perpendicular to $\hat{\mathbf{k}}$ and to each other—as in a vacuum. In this section, we confine attention to this simple situation and to linearly polarized waves, for which $\mathbf{E}$ oscillates linearly along a unit polarization vector $\hat{\mathbf{f}}$ that is perpendicular to $\hat{\mathbf{k}}$:

**polarization vector**

$$\mathbf{E} = A e^{i\varphi} \hat{\mathbf{f}}, \qquad \hat{\mathbf{f}} \cdot \hat{\mathbf{k}} \equiv \hat{\mathbf{f}} \cdot \nabla\varphi = 0. \tag{7.102}$$

In the eikonal approximation, $A e^{i\varphi} \equiv \psi$ satisfies the geometric-optics propagation laws of Sec. 7.3, and the polarization vector $\hat{\mathbf{f}}$, like the amplitude $A$, will propagate along the rays. The propagation law for $\hat{\mathbf{f}}$ can be derived by applying the eikonal approximation to Maxwell's equations, but it is easier to infer that law by simple physical reasoning:

1. Since $\hat{\mathbf{f}}$ is orthogonal to $\hat{\mathbf{k}}$ for a plane wave, it must also be orthogonal to $\hat{\mathbf{k}}$ in the eikonal approximation (which, after all, treats the wave as planar on lengthscales long compared to the wavelength).

2. If the ray is straight, then the medium, being isotropic, is unable to distinguish a slow right-handed rotation of $\hat{\mathbf{f}}$ from a slow left-handed rotation, so there will be no rotation at all: $\hat{\mathbf{f}}$ will continue always to point in the same direction (i.e., $\hat{\mathbf{f}}$ will be kept parallel to itself during transport along the ray).

3. If the ray bends, so $d\hat{\mathbf{k}}/ds \neq 0$ (where $s$ is distance along the ray), then $\hat{\mathbf{f}}$ will have to change as well, so as always to remain perpendicular to $\hat{\mathbf{k}}$. The direction of $\hat{\mathbf{f}}$'s change must be $\hat{\mathbf{k}}$, since the medium, being isotropic, cannot provide any other preferred direction for the change. The magnitude of the change is determined by the requirement that $\hat{\mathbf{f}} \cdot \hat{\mathbf{k}}$ remain zero all along the ray and that $\hat{\mathbf{k}} \cdot \hat{\mathbf{k}} = 1$. This immediately implies that the propagation law for $\hat{\mathbf{f}}$ is

**propagation law for polarization vector**

$$\frac{d\hat{\mathbf{f}}}{ds} = -\hat{\mathbf{k}}\left(\hat{\mathbf{f}} \cdot \frac{d\hat{\mathbf{k}}}{ds}\right).$$

(7.103)

This equation states that the vector $\hat{\mathbf{f}}$ is parallel-transported along the ray (cf. Fig. 7.5 in Sec. 24.3.3). Here "parallel transport" means: (i) Carry $\hat{\mathbf{f}}$ a short distance along the ray, keeping it parallel to itself in 3-dimensional space. Because of the bending of the ray and its tangent vector $\hat{\mathbf{k}}$, this will cause $\hat{\mathbf{f}}$ to no longer be perpendicular to $\hat{\mathbf{k}}$. (ii) Project $\hat{\mathbf{f}}$ perpendicular to $\hat{\mathbf{k}}$ by adding onto it the appropriate multiple of $\hat{\mathbf{k}}$. (The techniques of differential geometry for curved lines and surfaces, which we develop in Chaps. 24 and 25 in preparation for studying general relativity, give powerful mathematical tools for analyzing this parallel transport.)

7.7.2

### 7.7.2 Geometric Phase   `T2`

We use the polarization propagation law (7.103) to illustrate a quite general phenomenon known as the *geometric phase,* or sometimes as the *Berry phase,* after Michael Berry who elucidated it. For further details and some history of this concept, see Berry (1990).

As a simple context for the geometric phase, consider a linearly polarized, monochromatic light beam that propagates in an optical fiber. Focus on the evolution of the polarization vector along the fiber's optic axis. We can imagine bending the fiber into any desired shape, thereby controlling the shape of the ray. The ray's shape in turn will control the propagation of the polarization via Eq. (7.103).

If the fiber and ray are straight, then the propagation law (7.103) keeps $\hat{\mathbf{f}}$ constant. If the fiber and ray are circular, then Eq. (7.103) causes $\hat{\mathbf{f}}$ to rotate in such a way as to always point along the generator of a cone, as shown in Fig. 7.21a. This polarization behavior, and that for any other ray shape, can be deduced with the aid of a unit sphere on which we plot the ray direction $\hat{\mathbf{k}}$ (Fig. 7.21b). For example, the ray directions at ray locations $C$ and $H$ of panel a are as shown in panel b of the figure. Notice that the trajectory of $\hat{\mathbf{k}}$ around the unit sphere is a great circle.

On the unit sphere we also plot the polarization vector $\hat{\mathbf{f}}$—one vector at each point corresponding to a ray direction. Because $\hat{\mathbf{f}} \cdot \hat{\mathbf{k}} = 0$, the polarization vectors are always tangent to the unit sphere. Notice that each $\hat{\mathbf{f}}$ on the unit sphere is identical in length and direction to the corresponding one in the physical space of Fig. 7.21a.

The parallel-transport law (7.103) keeps constant the angle $\alpha$ between $\hat{\mathbf{f}}$ and the trajectory of $\hat{\mathbf{k}}$ (i.e., the great circle in panel b of the figure). Translated back to

(a)                    (b)

**FIGURE 7.21** (a) The ray along the optic axis of a circular loop of optical fiber, and the polarization vector $\hat{\mathbf{f}}$ that is transported along the ray by the geometric-optics transport law $d\hat{\mathbf{f}}/ds = -\hat{\mathbf{k}}(\hat{\mathbf{f}} \cdot d\hat{\mathbf{k}}/ds)$. (b) The polarization vector $\hat{\mathbf{f}}$ drawn on the unit sphere. The vector from the center of the sphere to each of the points $A, B, \ldots, H$ is the ray's propagation direction $\hat{\mathbf{k}}$, and the polarization vector (which is orthogonal to $\hat{\mathbf{k}}$ and thus tangent to the sphere) is identical to that in the physical space of the ray (panel a).

panel a, this constancy of $\alpha$ implies that the polarization vector points always along the generators of the cone, whose opening angle is $\pi/2 - \alpha$, as shown.

Next let the fiber and its central axis (the ray) be helical as shown in Fig. 7.22a. In this case, the propagation direction $\hat{\mathbf{k}}$ rotates, always maintaining the same angle $\theta$ to the vertical direction, and correspondingly its trajectory on the unit sphere of Fig. 7.22b is a circle of constant polar angle $\theta$. Therefore (as one can see, e.g., with the aid of a large globe of Earth and a pencil transported around a circle of latitude $90° - \theta$), the parallel-transport law dictates that the angle $\alpha$ between $\hat{\mathbf{f}}$ and the circle *not* remain constant, but instead rotate at the rate

$$d\alpha/d\phi = \cos\theta. \tag{7.104}$$

Here $\phi$ is the angle (longitude on the globe) around the circle. This is the same propagation law as for the direction of swing of a Foucault pendulum as Earth turns (cf. Box 14.5), and for the same reason: the gyroscopic action of the Foucault pendulum is described by parallel transport of its plane along Earth's spherical surface.

In the case where $\theta$ is arbitrarily small (a nearly straight ray), Eq. (7.104) says $d\alpha/d\phi = 1$. This is easily understood: although $\hat{\mathbf{f}}$ remains arbitrarily close to constant, the trajectory of $\hat{\mathbf{k}}$ turns rapidly around a tiny circle about the pole of the unit sphere, so $\alpha$ changes rapidly—by a total amount $\Delta\alpha = 2\pi$ after one trip around the pole, $\Delta\phi = 2\pi$; whence $d\alpha/d\phi = \Delta\alpha/\Delta\phi = 1$. For any other helical pitch angle $\theta$, Eq. (7.104) says that during one round trip, $\alpha$ will change by an amount $2\pi\cos\theta$ that lags behind its change for a tiny circle (nearly straight ray) by the lag angle $\alpha_{\text{lag}} = 2\pi(1 - \cos\theta)$, which is also the solid angle $\Delta\Omega$ enclosed by the path of $\hat{\mathbf{k}}$ on the unit sphere:

$$\alpha_{\text{lag}} = \Delta\Omega. \tag{7.105}$$

7.7 Polarization     **407**

(a)                                        (b)

**FIGURE 7.22** (a) The ray along the optic axis of a helical loop of optical fiber, and the polarization vector $\hat{\mathbf{f}}$ that is transported along this ray by the geometric-optics transport law $d\hat{\mathbf{f}}/ds = -\hat{\mathbf{k}}(\hat{\mathbf{f}} \cdot d\hat{\mathbf{k}}/ds)$. The ray's propagation direction $\hat{\mathbf{k}}$ makes an angle $\theta = 73°$ to the vertical direction. (b) The trajectory of $\hat{\mathbf{k}}$ on the unit sphere (a circle with polar angle $\theta = 73°$), and the polarization vector $\hat{\mathbf{f}}$ that is parallel transported along that trajectory. The polarization vectors in panel a are deduced from the parallel-transport law demonstrated in panel b. The lag angle $\alpha_{\text{lag}} = 2\pi(1 - \cos\theta) = 1.42\pi$ is equal to the solid angle contained inside the trajectory of $\hat{\mathbf{k}}$ (the $\theta = 73°$ circle).

(For the circular ray of Fig. 7.21, the enclosed solid angle is $\Delta\Omega = 2\pi$ steradians, so the lag angle is $2\pi$ radians, which means that $\hat{\mathbf{f}}$ returns to its original value after one trip around the optical fiber, in accord with the drawings in the figure.)

**lag angle for polarization vector in an optical fiber**

Remarkably, Eq. (7.105) is true for light propagation along an optical fiber of any shape: if the light travels from one point on the fiber to another at which the tangent vector $\hat{\mathbf{k}}$ has returned to its original value, then the lag angle is given by the enclosed solid angle on the unit sphere, Eq. (7.105).

By itself, the relationship $\alpha_{\text{lag}} = \Delta\Omega$ is merely a cute phenomenon. However, it turns out to be just one example of a very general property of both classical and quantum mechanical systems when they are forced to make slow, *adiabatic* changes described by circuits in the space of parameters that characterize them. In the more general case, one focuses on a phase lag rather than a direction-angle lag. We can easily translate our example into such a phase lag.

The apparent rotation of $\hat{\mathbf{f}}$ by the lag angle $\alpha_{\text{lag}} = \Delta\Omega$ can be regarded as an advance of the phase of one circularly polarized component of the wave by $\Delta\Omega$ and

**geometric phase change**

a phase retardation of the other circular polarization by the same amount. Thus the phase of a circularly polarized wave will change, after one circuit around the fiber's helix, by an amount equal to the usual phase advance $\Delta\varphi = \int \mathbf{k} \cdot d\mathbf{x}$ (where $d\mathbf{x}$ is displacement along the fiber) plus an extra *geometric* phase change $\pm\Delta\Omega$, where the sign is given by the sense of circular polarization. This type of geometric phase change is found quite generally, when classical vector or tensor waves propagate through

backgrounds that change slowly, either temporally or spatially. The phases of the wave functions of quantum mechanical particles with spin behave similarly.

**Exercise 7.22**  *Derivation: Parallel Transport*  T2

Use the parallel-transport law (7.103) to derive the relation (7.104).

**Exercise 7.23**  *Problem: Martian Rover*  T2

A Martian Rover is equipped with a single gyroscope that is free to pivot about the direction perpendicular to the plane containing its wheels. To climb a steep hill on Mars without straining its motor, it must circle the summit in a decreasing spiral trajectory. Explain why there will be an error in its measurement of North after it has reached the summit. Could it be programmed to navigate correctly? Will a stochastic error build up as it traverses a rocky terrain?

## Bibliographic Note

Modern textbooks on optics deal with the geometric-optics approximation only for electromagnetic waves propagating through a dispersion-free medium. Accordingly, they typically begin with Fermat's principle and then treat in considerable detail the paraxial approximation, applications to optical instruments, and sometimes the human eye. There is rarely any mention of the eikonal approximation or of multiple images and caustics. Examples of texts of this sort that we like are Bennett (2008), Ghatak (2010), and Hecht (2017). For a far more thorough treatise on geometric optics of scalar and electromagnetic waves in isotropic and anisotropic dielectric media, see Kravtsov (2005). A good engineering-oriented text with many contemporary applications is Iizuka (1987).

We do not know of textbooks that treat the eikonal approximation to the degree of generality used in this chapter, though some should, since it has applications to all types of waves (many of which are explored later in this book). For the eikonal approximation specialized to Maxwell's equations, see Kravtsov (2005) and the classic treatise on optics by Born and Wolf (1999), which in this new edition has modern updates by a number of other authors. For the eikonal approximation specialized to the Schrödinger equation and its connection to Hamilton-Jacobi theory, see most any quantum mechanics textbook (e.g., Griffiths, 2004).

Multiple-image formation and caustics are omitted from most standard optics textbooks, except for a nice but out-of-date treatment in Born and Wolf (1999). Much better are the beautiful review by Berry and Upstill (1980) and the much more thorough treatments in Kravtsov (2005) and Nye (1999). For an elementary mathematical treatment of catastrophe theory, we like Saunders (1980). For a pedagogical treatise on gravitational lenses, see Schneider, Ehlers, and Falco (1992). Finally, for some history and details of the geometric phase, see Berry (1990).

# NAME INDEX

Page numbers for entries in boxes are followed by "b," those for epigraphs at the beginning of a chapter by "e," those for figures by "f," and those for notes by "n."

**569**

# SUBJECT INDEX

Second and third level entries are not ordered alphabetically. Instead, the most important or general entries come first, followed by less important or less general ones, with specific applications last.

Page numbers for entries in boxes are followed by "b," those for epigraphs at the beginning of a chapter by "e," those for figures by "f," for notes by "n," and for tables by "t."