

§16.4. THE RODS AND CLOCKS USED TO MEASURE SPACE AND TIME INTERVALS

Turn attention now from the laws of physics in the presence of gravity to the nature of the rods and clocks that must be used for measuring the length and time intervals appearing in those laws.

One need not—and indeed must not!—postulate that proper length s is measured by a certain type of rod (e.g., platinum meter stick), or that proper time τ is measured by a certain type of clock (e.g., hydrogen-maser clock). Rather, one must ask the laws of physics themselves what types of rods and clocks will do the job. Put differently, one *defines an “ideal” rod or clock* to be one which measures proper length as given by $ds = (g_{\alpha\beta} dx^\alpha dx^\beta)^{1/2}$ or proper time as given by $d\tau = (-g_{\alpha\beta} dx^\alpha dx^\beta)^{1/2}$ (the kind of clock to which one was led by physical arguments in §1.5). One must then determine the accuracy to which a given rod or clock is ideal under given circumstances by using the laws of physics to analyze its behavior.

As an obvious example, consider a pendulum clock. If it is placed at rest on the Earth’s surface, if it is tiny enough that redshift effects from one end to the other and time dilation effects due to its swinging velocity are negligible, and if the accuracy one demands is small enough that time variations in the local gravitational acceleration due to Earth tides can be ignored, then the laws of physics report (Box 16.2) that the pendulum clock is “ideal.” However, in any other context (e.g., on a rocket journey to the moon), a pendulum clock should be far from ideal. Wildly changing accelerations, or no acceleration at all, will make it worthless!

Of greater interest are atomic and nuclear clocks of various sorts. Such a clock is analyzed most easily if it is freely falling. One can then study it in its local Lorentz rest frame, using the standard equations of quantum theory; and, of course, one will find that it measures proper time to within the precision ($\Delta t/t \sim 10^{-9}$ to 10^{-14}) of the technology used in its construction. However, one rarely permits his atomic clock to fall freely. (The impact with the Earth’s surface can be expensive!) Nevertheless, even when accelerated at “1 g” = 980 cm/sec² on the Earth’s surface, and even when accelerated at “2 g” in an airliner trying to avoid a midair collision (Box 16.3), an atomic clock—if built solidly—will still measure proper time $d\tau = (-g_{\alpha\beta} dx^\alpha dx^\beta)^{1/2}$ along its world line to nearly the same accuracy as if it were freely falling. To discover this one can perform an experiment. Alternatively, one can analyze the clock in its own “proper reference frame” (§13.6), with Fermi-Walker-transported basis vectors, using the standard local Lorentz laws of quantum mechanics as adapted to accelerated frames (local Lorentz laws plus an “inertial force,” which can be treated as due to a potential with a uniform gradient).

Of course, any clock has a “breaking point,” beyond which it will cease to function properly (Box 16.3). But that breaking point depends entirely on the construction of the clock—and not at all on any “universal influence of acceleration on the march of time.” Velocity produces a universal time dilation; acceleration does not.

The aging of the human body is governed by the same electromagnetic and quantum-mechanical laws as govern the periodicities and level transitions in atoms and molecules. Consequently, aging, like atomic processes, is tied to proper time

“Ideal” rods and clocks defined

How ideal are real clocks?

(1) pendulum clocks

(2) atomic clocks

(3) human clocks

(continued on page 396)

**Box 16.2 PROOF THAT A PENDULUM CLOCK AT REST
ON THE EARTH'S SURFACE IS IDEAL**

That is, a proof that it measures the interval $d\tau = (-g_{\alpha\beta} dx^\alpha dx^\beta)^{1/2}$.

A. Constraint on the Pendulum

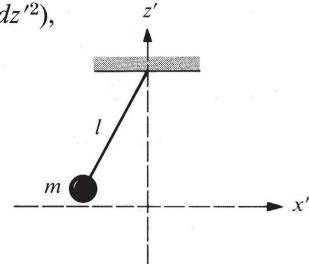
It must be so small that it cannot couple to the spacetime curvature—i.e., so small that the Earth's gravitational field looks uniform in its neighborhood—and that the velocity of its ball is totally negligible compared to the speed of light.

B. Coordinate System and Metric

- (1) General coordinate system: because the Earth's field is nearly Newtonian, one can introduce the coordinates of “linearized theory” (§18.4; one must take this on faith until one reaches that point) in which

$$ds^2 = -(1 + 2\Phi) dt'^2 + (1 - 2\Phi)(dx'^2 + dy'^2 + dz'^2),$$

where Φ is the Newtonian potential.



- (2) Put the origin of coordinates at the pendulum's equilibrium position, and orient the x', z' -plane so the pendulum swings in it.
- (3) Renormalize the coordinates so they measure proper length and proper time at the equilibrium position

$$t = [1 + 2\Phi(0)]^{1/2} t', \quad x^j = [1 - 2\Phi(0)]^{1/2} x'^j.$$

Then near the pendulum (inhomogeneities in the field neglected!)

$$\Phi = \Phi(0) + gz, \quad g = \text{"acceleration of gravity,"} \quad (1)$$

$$ds^2 = -(1 + 2gz) dt^2 + (1 - 2gz)(dx^2 + dy^2 + dz^2). \quad (2)$$

C. Analysis of Pendulum Motion

- (1) Put the total mass m of the pendulum in its ball (negligible mass in its rod). Let its rod have proper length l .

- (2) Calculate the 4-acceleration $\mathbf{a} = \nabla_{\mathbf{u}}\mathbf{u}$ of the pendulum's ball in terms of d^2x^α/dt^2 , using the velocity condition $v \lll 1$ and $dt/d\tau \approx 1$:

$$\begin{aligned} a^x &= d^2x/d\tau^2 + \Gamma^x_{00}(dt/d\tau)^2 = d^2x/dt^2 + \Gamma^x_{00} = d^2x/dt^2, \\ a^z &= d^2z/d\tau^2 + \Gamma^z_{00}(dt/d\tau)^2 = d^2z/dt^2 + \Gamma^z_{00} = d^2z/dt^2 + g. \end{aligned} \quad (3)$$

- (3) This 4-acceleration must be produced by the forces in the rod, and must be directed up the rod so that (for $x \ll l$ so $g \gg d^2z/dt^2$)

$$d^2x/dt^2 = a^x = -(x/l)a^z = -(g/l)x. \quad (4)$$

- (4) Solve this differential equation to obtain

$$x = x_0 \cos(t\sqrt{g/l}). \quad (5)$$

- (5) Thus conclude that the pendulum is periodic in t , which is proper time at the ball's equilibrium position (see equation 2). This means that *the pendulum is an ideal clock when it is at rest on the Earth's surface.*

Note: The above analysis ignores the Earth's rotation; for an alternative analysis including rotation, one can perform a similar calculation at the origin of the pendulum's "proper reference frame" [§13.6; line element (13.71)]. The answer is the same; but now "g" is a superposition of the "gravitational acceleration," and the "centrifugal acceleration produced by Earth's rotation."

Box 16.3 RESPONSE OF CLOCKS TO ACCELERATION AND TO TIDAL GRAVITATIONAL FORCES

Consider an atomic clock with frequency stabilized by some atomic or molecular process—for example, fixed by the “umbrella vibrations” of ammonia molecules [see Feynman *et. al.* (1964)]. When subjected to sufficiently strong accelerations or tidal forces, such a clock will cease to measure proper time with its normal precision. Two types of effects could lead to such departures from “ideality”:

A. Influence of the acceleration or tidal force on the atomic process that provides the frequency stability. Example: If tidal forces are significant over distances of a few angstroms (e.g., near a spacetime “singularity” terminating gravitational collapse), then they can and will deform an ammonia molecule and destroy the regularity of its umbrella vibrations, thereby making useless *any* ammonia atomic clock, no matter how constructed. Similarly, if an ammonia molecule is subjected to accelerations of magnitude comparable to its internal atomic accelerations ($a \sim 10^{12} \text{ "g"} \sim 10^{15} \text{ cm/sec}^2$), which change in times of the order of the “umbrella” vibration period, then it must cease to vibrate regularly, and any clock based on its vibrations must fail. Such limits of principle on the ideality of a clock will vary from one atomic process to another. However, they are far from being a limiting factor on clock construction in 1973. Much more important today is:

B. Influence of the acceleration or tidal force on the macroscopic structure of the clock—a structure dictated by current technology. The crystal oscillator,

which produces the periodic signal output, must be locked to the regulating atomic process in some way. The lock will be disturbed by moderate accelerations. The toughest task for the manufacturer of aircraft clocks is to guarantee that precise locking will be maintained, even when the aircraft is maneuvering desperately to avoid collision with another aircraft or with a missile. In 1972 a solidly built rubidium clock will maintain its lock, with no apparent degradation of stability

[$\Delta t/t \sim 10^{-12}(1 \text{ sec}/t)^{1/2}$ for $1 \text{ sec} \lesssim t \lesssim 10^3 \text{ sec}$] under steady-state accelerations up to 50 “g” or more. But, because of the finite bandwidth of the lock loop (typically $\Delta\nu \sim 20$ to 50 Hz), sudden changes in acceleration will temporarily break the lock, degrading the clock stability to that of the unlocked crystal oscillator—for which an acceleration a produces a change in frequency of about $(a/1 \text{ "g"}) \times 10^{-9}$. But the lock to the rubidium standard is restored quickly ($\delta t \sim 1/\Delta\nu$), bringing the clock back to its normal highly stable performance.*

Tidal forces are so small in the solar system that the clock manufacturer can ignore them. However, a 1973 atomic clock, subjected to the tidal accelerations near a spacetime singularity, should break the “lock” to its atomic process long before the tidal forces become strong enough to influence the atomic process itself.

*For this information on the response of rubidium clocks to acceleration, we thank H. P. Stratemeyer of General Radio Company, Concord, Massachusetts.

as governed by the metric—though, of course, it is also tied to other things, such as cigarette smoking.

In principle, one can build ideal rods and clocks from the geodesic world lines of freely falling test particles and photons. (See Box 16.4.) In other words, spacetime has its own rods and clocks built into itself, even when matter and nongravitational fields are absent!

Box 16.4 IDEAL RODS AND CLOCKS BUILT FROM GEODESIC WORLD LINES*

The Standard Interval. A specific timelike interval—the interval between two particular neighboring events \mathcal{A} and \mathcal{B} —is chosen as the standard interval, and is assigned unit length. It is used to calibrate a huge set of geodesic clocks that pass through \mathcal{A} .

Each *geodesic clock* is constructed and calibrated as follows:

- (1) A timelike geodesic \mathcal{AC} (path of freely falling particle) passes through \mathcal{A} .
- (2) A neighboring world line, everywhere parallel to \mathcal{AC} (and thus not a geodesic), is constructed by the method of Schild's ladder (Box 10.2), which relies only on geodesics.
- (3) Light rays (null geodesics) bounce back and forth between these parallel world lines; each round trip constitutes one “tick.”
- (4) The proper time lapse, τ_0 , between ticks is related to the interval \mathcal{AB} by

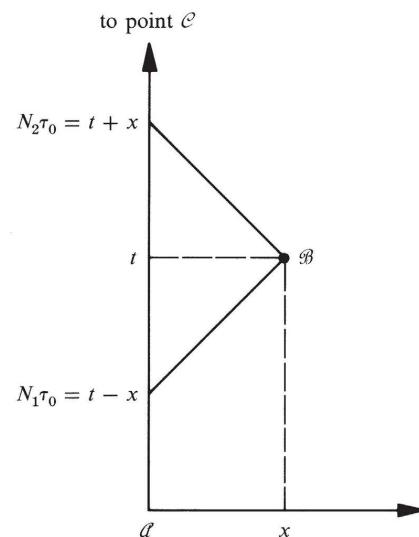
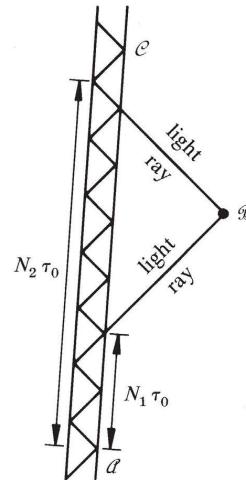
$$-1 \equiv (\mathcal{AB})^2 = -(N_1\tau_0)(N_2\tau_0),$$

where N_1 and N_2 are the number of ticks between the events shown in the diagrams.

[*Proof:* see diagram at right.]

Spacetime is filled with such geodesic clocks. Those that pass through \mathcal{A} are calibrated as above against the standard interval \mathcal{AB} , and are used subsequently to calibrate all other clocks they meet.

* Based on Marzke and Wheeler (1964).



In local Lorentz rest frame of geodesic clock:

$$\begin{aligned} (N_1\tau_0)(N_2\tau_0) &= (t-x)(t+x) \\ &= t^2 - x^2 = -(\mathcal{AB})^2 \end{aligned}$$

Box 16.4 (continued)

Any interval \mathcal{PQ} along the world line of a geodesic clock can be measured by the same method as was used in calibration. The interval \mathcal{PQ} can be timelike, spacelike, or null; its squared length in all three cases will be

$$(\mathcal{PQ})^2 = -(N_3\tau_0)(N_4\tau_0)$$

To achieve a precision of measurement good to one part in N , where N is some large number, take two precautions:

- (1) Demand that the intervals \mathcal{AB} and \mathcal{PQ} be sufficiently small compared to the scale of curvature of spacetime; or specifically,

$$R^{(AB)}(\mathcal{AB})^2 \ll 1/N$$

and

$$R^{(PQ)}(\mathcal{PQ})^2 \ll 1/N,$$

where $R^{(AB)}$ and $R^{(PQ)}$ are the largest relevant components of the curvature tensor in the two regions in question.

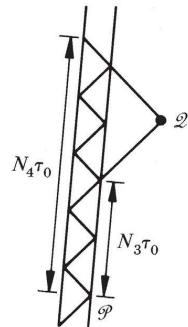
- (2) Demand that the time scale, τ_0 , of the geodesic clocks employed be small compared to \mathcal{AB} and \mathcal{PQ} individually; thus,

$$\tau_0 \ll \mathcal{AB}/N,$$

$$\tau_0 \ll \mathcal{PQ}/N.$$

The Einstein principle that spacetime is described by Riemannian geometry exposes itself to destruction by a “thousand” tests. Thus, from the fiducial interval, \mathcal{AB} , to the interval under measurement, \mathcal{PQ} , there are a “score” of routes of intercomparison, all of which must give the same value for the ratio $\mathcal{PQ}/\mathcal{AB}$. Moreover, one can easily select out “fifty” intervals \mathcal{PQ} to which the same kind of test can be applied. Such tests are not all items for the future.

Some 5×10^9 years ago, electrons arrived by different routes at a common location, a given atom of iron in the core of the earth. This iron atom does not collapse. The Pauli principle of



exclusion keeps the electrons from all falling into the K-orbit. The Pauli principle would not apply if the electrons were not identical or nearly so. From this circumstance it would appear possible to draw an important conclusion (Marzke and Wheeler). With each electron is associated a standard length, its Compton wavelength, \hbar/mc . If these lengths had started different, or changed by different amounts along the different routes, and if the resulting difference in properties were as great as one part in

$$\sim(5 \times 10^9 \text{ yr}) \times (3 \times 10^7 \text{ sec/yr}) \\ \times (5 \times 10^{18} \text{ rev/sec}) \sim 10^{36},$$

by now this difference would have shown up, the varied electrons would have fallen into the K-orbit, and the earth would have collapsed, contrary to observation.

The Marzke-Wheeler construction expresses an arbitrary small interval $\mathcal{P}\mathcal{Q}$, anywhere in spacetime, in terms of the fiducial interval $\mathcal{A}\mathcal{B}$, an interval which itself may be taken for definiteness to be the “geometrodynamic standard centimeter” of §1.5. This construction thus gives a vivid meaning to the idea of Riemannian geometry.

The M-W construction makes no appeal what-

soever to rods and clocks of atomic constitution. This circumstance is significant for the following reasons. The length of the usual platinum meter stick is some multiple, $N_1(\hbar^2/me^2)$, of the Bohr atomic radius. Similarly, the wavelength of the Kr⁸⁶ line is some multiple, $N_2(\hbar c/e^2)(\hbar^2/me^2)$, of a second basic length that depends on the atomic constants in quite a different way. Thus, if there is any change with time in the dimensionless ratio $\hbar c/e^2 = 137.038$, one or the other or both of these atomic standards of length must get out of kilter with the geometrodynamic standard centimeter. In this case, general relativity says, “Stick to the geometrodynamic standard centimeter.”

Hermann Weyl at first thought that one could carry out the comparison of lengths by light rays alone, but H. A. Lorentz pointed out that one can dispense with the geodesics neither of test particles nor of light rays in the measurement process, the construction for which, however, neither Weyl nor Lorentz supplied [literature in Marzke and Wheeler (1964)]. Ehlers, Pirani, and Schild (1972) have given a deeper analysis of the separate parts played in the measurement process by the affine connection, by the conformal part of the metric, and by the full metric.

§16.5. THE MEASUREMENT OF THE GRAVITATIONAL FIELD

“I know how to measure the electromagnetic field using test charges; what is the analogous procedure for measuring the gravitational field?” This question has, at the same time, many answers and none.

It has no answers because nowhere has a precise definition of the term “gravitational field” been given—nor will one be given. Many different mathematical entities are associated with gravitation: the metric, the Riemann curvature tensor, the Ricci curvature tensor, the curvature scalar, the covariant derivative, the connection coefficients, etc. Each of these plays an important role in gravitation theory, and none is so much more central than the others that it deserves the name “gravitational field.” Thus it is that throughout this book the terms “gravitational field” and “gravity” refer in a vague, collective sort of way to all of these entities. Another, equivalent term used for them is the “geometry of spacetime.”

The many faces of gravity,
and how one measures them

To “measure the gravitational field,” then, means to “explore experimentally various properties of the spacetime geometry.” One makes different kinds of measurements, depending on which geometric property of spacetime one is interested in. However, all such measurements must involve a scrutiny of the effects of the spacetime geometry (i.e., of gravity) on particles, on matter, or on nongravitational fields.

For example, to “measure” the metric near a given event, one typically lays out a latticework of rods and clocks (local orthonormal frame, small enough that curvature effects are negligible), and uses it to determine the interval between neighboring events. To measure the Riemann curvature tensor near an event, one typically studies the geodesic deviation (relative accelerations) that curvature produces between the world lines of a variety of neighboring test particles; alternatively, one makes measurements with a “gravity gradiometer” (Box 16.5) if the curvature is static or slowly varying; or with a gravitational wave antenna (Chapter 37) if the curvature fluctuates rapidly. To study the large-scale curvature of spacetime, one examines large-scale effects of gravity, such as the orbits of planets and satellites, or the bending of light by the sun’s gravitational field.

But whatever aspect of gravity one measures, and however one measures it, one is studying the geometry of spacetime.

EXERCISE

Exercise 16.5. GRAVITY GRADIOMETER

The gravity gradiometer of Box 16.5 moves through curved spacetime along an accelerated world line. Calculate the amplitude and phase of oscillation of one arm of the gradiometer relative to the other. [Hint: Perform the calculation in the gradiometer’s “proper reference frame” (§13.6), with Fermi-Walker-transported basis vectors. Use, as the equation for the relative angular acceleration of the two arms,

$$2ml^2(\ddot{\alpha} + \dot{\alpha}/\tau_0 + \omega_0^2\alpha) = \left(\begin{array}{l} \text{Driving torque produced by} \\ \text{Riemann curvature} \end{array} \right),$$

where

$2ml^2$ = (moment of inertia of one arm),

α = (angular displacement of one arm from equilibrium),

$\frac{\pi}{2} + 2\alpha$ = (angular separation of the two arms),

$2ml^2\omega_0^2$ = (torsional spring constant),

ω_0 = (angular frequency of free vibrations),

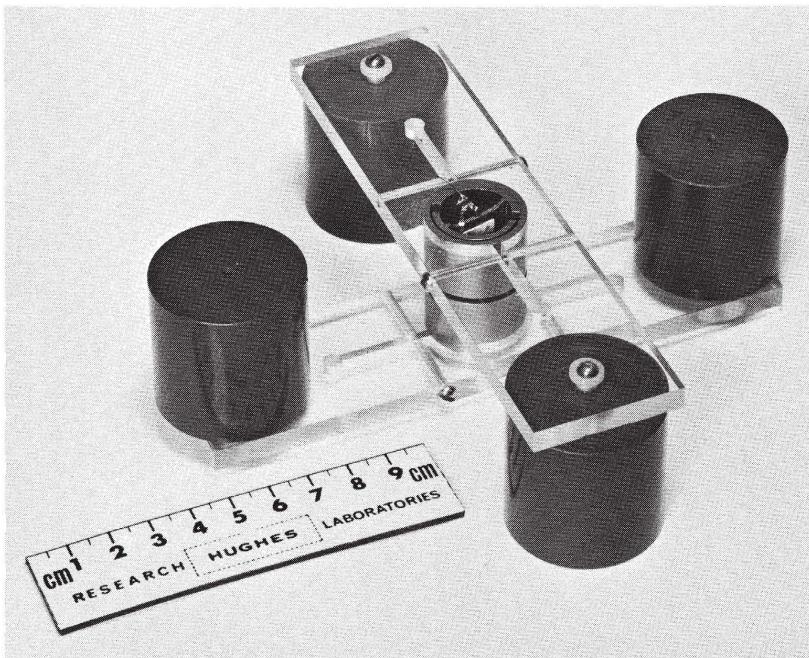
τ_0 = (decay time for free vibrations to damp out due to internal frictional forces).

If ξ is the vector from the center of mass of the gradiometer to mass 1, then one has

$$\left(\begin{array}{l} \text{curvature-produced} \\ \text{acceleration of mass 1} \\ \text{relative to center of} \\ \text{gradiometer} \end{array} \right)_{\hat{k}} = \left(\frac{D^2\xi_k}{d\tau^2} \right)_{\text{geodesic deviation}} = -R_{k\hat{i}\hat{j}\hat{o}}\xi_l;$$

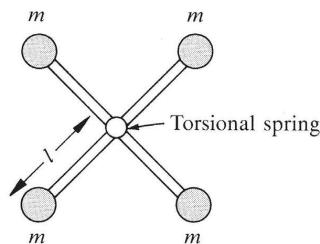
(continued on page 403)

Box 16.5 GRAVITY GRADIOMETER FOR MEASURING THE RIEMANN CURVATURE OF SPACETIME



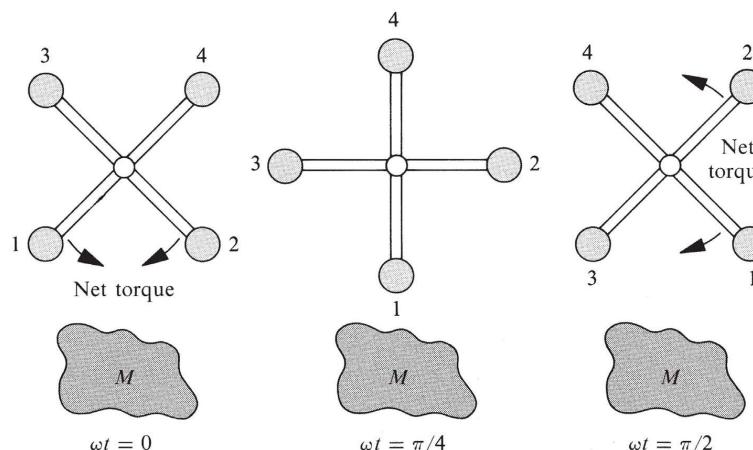
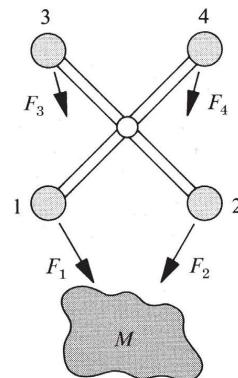
This gravity gradiometer was designed and built by Robert M. Forward and his colleagues at Hughes Research Laboratories, Malibu, California. It measures the Riemann curvature of spacetime produced by nearby masses. By flying a more advanced version of such a gradiometer in an airplane above the Earth's surface, one should be able to measure subsurface mass variations due to varying geological structure. In an Earth-orbiting satellite, such a gradiometer could measure the gravitational multipole moments of the Earth. Technical details of the gradiometer are spelled out in the papers of Forward (1972), and Bell, Forward, and Williams (1970). The principles of its operation are outlined below.

The gradiometer consists of two orthogonal arms with masses m on their ends, connected at their centers by a torsional spring. When the arms are twisted out of orthogonal alignment, they oscillate. A piezoelectric strain transducer is used to measure the oscillation amplitude.



Box 16.5 (continued)

When placed near an external mass, M , the gradiometer experiences a torque: because of the gradient in the gravitational field of M (i.e., because of the spacetime curvature produced by M), the Newtonian forces F_1 and F_2 are greater than F_3 and F_4 ; so a net torque pulls masses 1 and 2 toward each other, and 3 and 4 toward each other. [Note: the forces F_1, F_2, F_3, F_4 depend on whether the gradiometer is in free fall (geodesic motion; $\nabla_u u = 0$) or is moving on an accelerated world line. But the net torque is unaffected by acceleration; acceleration produces equal Newtonian forces on all four masses, with zero net torque.]



When in operation the gradiometer rotates with angular velocity ω about its center. As it rotates, the torques on its arms oscillate:

at $\omega t = 0$ net torque pushes 1 and 2 toward each other;

at $\omega t = \pi/4$ net torque is zero;

at $\omega t = \pi/2$ net torque pushes 1 and 2 away from each other.

The angular frequency of the oscillating torque is 2ω . If 2ω is set equal to $\omega_0 \equiv$ (natural oscillation frequency of the arms), the oscillating torque drives the arms into resonant oscillation. The resulting oscillation amplitude, in the 1970 prototype

of the gradiometer, was easily detectable for gravity gradients (Riemann curvatures) of magnitude

$$\gtrsim 0.0002 \left[\frac{2(\text{mass of earth})}{(\text{radius of earth})^3} \right] \\ \sim 1 \times 10^{-30} \text{ cm}^{-2} \sim .01 \text{ g/cm}^3$$

\sim Riemann curvature produced by a two-kilometer high mountain, idealized as a two-kilometer high cube, at a distance of 15 kilometers. (Neglected in this idealization are isostacy and any lowering of density of Earth's crust in regions of mountain uplift.)

For a mathematical analysis of the gradiometer, see exercise 16.5.

$$\left(\begin{array}{l} \text{torque acting on mass 1} \\ \text{relative to center of} \\ \text{gradiometer} \end{array} \right)_{\hat{i}} = \epsilon_{ijk} \xi_j (-m R_{k\hat{l}\hat{l}} \xi_l).$$

The torque on mass 4 is identical to this (replace ξ by $-\xi$), so the total torque on arm 1-4 is twice this. The components $R_{k\hat{l}\hat{l}}$ of **Riemann** can be regarded as components of a 3×3 symmetric matrix. By appropriate orientation of the reference frame's spatial axes (orientation along "principal axes" of $R_{k\hat{l}\hat{l}}$), one can make $R_{k\hat{l}\hat{l}}$ diagonal at some initial moment of time

$$R_{\hat{x}\hat{0}\hat{x}\hat{0}} \neq 0, R_{\hat{y}\hat{0}\hat{y}\hat{0}} \neq 0, R_{\hat{z}\hat{0}\hat{z}\hat{0}} \neq 0, \text{ all others vanish.}$$

Assume that **Riemann** changes sufficiently slowly along the gradiometer's world line that throughout the experiment $R_{j\hat{l}\hat{l}}$ remains diagonal and constant. For simplicity, place the gradiometer in the \hat{x}, \hat{y} -plane, so it rotates about the \hat{z} axis with angular velocity $\omega \approx \frac{1}{2}\omega_0$:

$$\left(\begin{array}{l} \text{Angle of arm 1-4} \\ \text{relative to } \hat{x} \text{ axis} \end{array} \right) = \omega t.$$

Show that the resultant equation of oscillation is

$$\ddot{\alpha} + \dot{\alpha}/\tau_0 + \omega_0^2 \alpha = \frac{1}{2} (R_{\hat{x}\hat{0}\hat{x}\hat{0}} - R_{\hat{y}\hat{0}\hat{y}\hat{0}}) \sin 2\omega t;$$

and that the steady-state oscillations are

$$\alpha = \text{Im} \left\{ \frac{1}{2} \frac{(R_{\hat{x}\hat{0}\hat{x}\hat{0}} - R_{\hat{y}\hat{0}\hat{y}\hat{0}})}{2\omega_0(\omega_0 - 2\omega + i/2\tau_0)} e^{i2\omega t} \right\}.$$

Thus, for fixed ω (e.g., $2\omega = \omega_0$), by measuring the amplitude and phase of the oscillations, one can learn the magnitude and sign of $R_{\hat{x}\hat{0}\hat{x}\hat{0}} - R_{\hat{y}\hat{0}\hat{y}\hat{0}}$. The other differences, $R_{\hat{y}\hat{0}\hat{y}\hat{0}} - R_{\hat{z}\hat{0}\hat{z}\hat{0}}$ and $R_{\hat{z}\hat{0}\hat{z}\hat{0}} - R_{\hat{x}\hat{0}\hat{x}\hat{0}}$ can be measured by placing the gradiometer's rotation axis along the \hat{x} and \hat{y} axes, respectively.]

CHAPTER 17

HOW MASS-ENERGY GENERATES CURVATURE

The physical world is represented as a four-dimensional continuum. If in this I adopt a Riemannian metric, and look for the simplest laws which such a metric can satisfy, I arrive at the relativistic gravitation theory of empty space. If I adopt in this space a vector field, or the antisymmetrical tensor field derived from it, and if I look for the simplest laws which such a field can satisfy, I arrive at the Maxwell equations for free space. . . . at any given moment, out of all conceivable constructions, a single one has always proved itself absolutely superior to all the rest . . .

ALBERT EINSTEIN (1934, p. 18)

§17.1. AUTOMATIC CONSERVATION OF THE SOURCE AS THE CENTRAL IDEA IN THE FORMULATION OF THE FIELD EQUATION

This section derives the “Einstein field equation”

Turn now from the response of matter to geometry (motion of a neutral test particle on a geodesic; “comma-goes-to-semicolon rule” for the dynamics of matter and fields), and analyze the response of geometry to matter.

Mass is the source of gravity. The density of mass-energy as measured by any observer with 4-velocity \mathbf{u} is

$$\rho = \mathbf{u} \cdot \mathbf{T} \cdot \mathbf{u} = u^\alpha T_{\alpha\beta} u^\beta. \quad (17.1)$$

Therefore the stress-energy tensor \mathbf{T} is the frame-independent “geometric object” that must act as the source of gravity.

This source, this geometric object, is not an arbitrary symmetric tensor. It must have zero divergence

$$\nabla \cdot \mathbf{T} = 0, \quad (17.2)$$

because only so can the law of conservation of momentum-energy be upheld.

Place this source, \mathbf{T} , on the righthand side of the equation for the generation of gravity. On the lefthand side will stand a geometric object that characterizes gravity. That object, like \mathbf{T} , must be a symmetric, divergence-free tensor; and if it is to characterize gravity, it must be built out of the geometry of spacetime and nothing but that geometry. Give this object the name “Einstein tensor” and denote it by \mathbf{G} , so that the equation for the generation of gravity reads

$$\mathbf{G} = \kappa \mathbf{T}. \quad (17.3)$$

↑
[proportionality factor;
to be evaluated later]

(Do not assume that \mathbf{G} is the same Einstein tensor as was encountered in Chapters 8, 13, 14, and 15; that will be proved below!)

The vanishing of the divergence $\nabla \cdot \mathbf{G}$ is not to be regarded as a consequence of $\nabla \cdot \mathbf{T} = 0$. Rather, the obedience of all matter and fields to the conservation law $\nabla \cdot \mathbf{T} = 0$ is to be regarded (1) as a consequence of the way [equation (17.3)] they are wired into the geometry of spacetime, and therefore (2) as required and enforced by an *automatic* conservation law, or *identity*, that holds for any smooth Riemannian spacetime whatsoever, physical or not: $\nabla \cdot \mathbf{G} \equiv 0$. (See Chapter 15 for a fuller discussion and §17.2 below for a fuller justification.) Accordingly, look for a symmetric tensor \mathbf{G} that is an “automatically conserved measure of the curvature of spacetime” in the following sense:

- (1) \mathbf{G} vanishes when spacetime is flat.
- (2) \mathbf{G} is constructed from the Riemann curvature tensor and the metric, and from nothing else.
- (3) \mathbf{G} is distinguished from other tensors which can be built from **Riemann** and \mathbf{g} by the demands (i) that it be linear in **Riemann**, as befits any natural measure of curvature; (ii) that, like \mathbf{T} , it be symmetric and of second rank; and (iii) that it have an automatically vanishing divergence,

$$\nabla \cdot \mathbf{G} \equiv 0. \quad (17.4)$$

Apart from a multiplicative constant, there is only one tensor (exercise 17.1) that satisfies these requirements of being an automatically conserved, second-rank tensor, linear in the curvature, and of vanishing when spacetime is flat. It is the Einstein curvature tensor, \mathbf{G} , expressed in Chapter 8 in terms of the Ricci curvature tensor:

$$\begin{aligned} R_{\mu\nu} &= R^\alpha_{\mu\alpha\nu}, \\ G_{\mu\nu} &= R_{\mu\nu} - \frac{1}{2} g_{\mu\nu} R. \end{aligned} \quad (17.5)$$

Equation describing how matter generates gravity must have form $\mathbf{G} = \kappa \mathbf{T}$, where \mathbf{T} is stress-energy tensor

Properties that the tensor \mathbf{G} must have

Proof that \mathbf{G} must be the Einstein curvature tensor of Chapter 8

This quantity was given vivid meaning in Chapter 15 as the “moment of rotation of the curvature” or, more simply, the “moment of rotation,” constructed by taking the double-dual

$$\mathbf{G} = * \mathbf{Riemann}^* \quad (17.6a)$$

of the Riemann curvature tensor, and then contracting this double dual,

$$G_{\mu\nu} = G^\alpha_{\mu\alpha\nu}. \quad (17.6b)$$

Evaluation of κ (in $\mathbf{G} = \kappa \mathbf{T}$)
by comparing with
Newtonian theory of gravity

In Chapter 15 the vanishing of $\nabla \cdot \mathbf{G}$ was shown to follow as a consequence of the elementary principle of topology that “the boundary of a boundary is zero.”

To evaluate the proportionality constant κ in the “Einstein field equation” $\mathbf{G} = \kappa \mathbf{T}$, one can compare with the well-tested Newtonian theory of gravity. To facilitate the comparison, examine the relative acceleration (geodesic deviation) of particles that fall down a pipe inserted into an idealized Earth of uniform density ρ (Figure 1.12). According to Newton, the relative acceleration is governed by the density; according to Einstein, it is governed by the Riemann curvature of spacetime. Direct comparison of the Newtonian and Einstein predictions using Newtonian coordinates (where $g_{\mu\nu} \approx \eta_{\mu\nu}$) reveals the relation

$$R_{00} \equiv R^\alpha_{0\alpha 0} = 4\pi\rho. \quad (17.7)$$

(See §1.7 for details of the derivation; see Chapter 12 for extensive discussion of Newtonian gravity using this equation.) When applied to the Earth’s interior, the Einstein field equation $\mathbf{G} = \kappa \mathbf{T}$ must thus reduce to $R_{00} = 4\pi\rho$. In component form, the Einstein field equation reads

$$G_{\mu\nu} = R_{\mu\nu} - \frac{1}{2} g_{\mu\nu} R = \kappa T_{\mu\nu}.$$

Its trace reads

$$-R = R - 2R = \kappa T.$$

In consequence, it predicts

$$\begin{aligned} R_{00} &= \frac{1}{2} g_{00} R + \kappa T_{00} = \frac{1}{2} \kappa (2T_{00} - \underbrace{g_{00}}_{-1} T) \\ &= \frac{1}{2} \kappa [2T_{00} + (T^0{}_0 + T^j{}_j)] \\ &= \frac{1}{2} \kappa (T_{00} + T^j{}_j), \end{aligned}$$

which reduces to

$$R_{00} = \frac{1}{2} \kappa \rho \quad (17.8)$$

when one recalls that for the Earth—as for any nearly Newtonian system—the stresses T_{jk} are very small compared to the density of mass-energy $T_{00} = \rho$:

$$\frac{|T_{jk}|}{T_{00}} \sim \frac{\text{pressure}}{\text{density}} \sim \frac{dp}{d\rho} \sim (\text{velocity of sound})^2 \ll 1.$$

The equation $R_{00} = 4\pi\rho$ (derived by comparing relative accelerations in the Newton and Einstein theories) and the equation $R_{00} = \frac{1}{2}\kappa\rho$ (derived directly from the Einstein field equation) can agree only if the proportionality constant κ is 8π .

Thus, the Einstein field equation, describing the generation of curvature by mass-energy, must read

$$\mathbf{G} = 8\pi\mathbf{T}. \quad (17.9)$$

Result: "Einstein field equation" $\mathbf{G} = 8\pi\mathbf{T}$

The lefthand side ("curvature") has units cm^{-2} , since a curvature tensor is a linear machine into which one inserts a displacement (units: cm) and from which one gets a relative acceleration (units: $\text{cm/sec}^2 \sim \text{cm/cm}^2 \sim \text{cm}^{-1}$). The right-hand side also has dimensions cm^{-2} , since it is a linear machine into which one inserts 4-velocity (dimensionless) and from which one gets mass density [units: $\text{g/cm}^3 \sim \text{cm}/\text{cm}^3 \sim \text{cm}^{-2}$; recall from equation (1.12) and Box 1.8 that $1\text{g} = (1\text{g}) \times (G/c^2) = (1\text{g}) \times (0.742 \times 10^{-28} \text{ cm/g}) = 0.742 \times 10^{-28} \text{ cm}]$.

This concludes the simplest derivation of Einstein's field equation that has come to hand, and establishes its correspondence with the Newtonian theory of gravity under Newtonian conditions. That correspondence had to be worked out to determine the factor $\kappa = 8\pi$ on the righthand side of (17.9). Apart from this factor, the central point in the derivation was the demand for, and the existence of, a unique tensorial measure of curvature \mathbf{G} with an identically vanishing divergence.

Exercise 17.1. UNIQUENESS OF THE EINSTEIN TENSOR

EXERCISES

- (a) Show that the most general second-rank, symmetric tensor constructable from **Riemann** and \mathbf{g} , and linear in **Riemann**, is

$$\begin{aligned} & aR_{\alpha\beta} + bRg_{\alpha\beta} + \Lambda g_{\alpha\beta} \\ &= aR^\mu{}_{\alpha\mu\beta} + bR^{\mu\nu}{}_{\mu\nu}g_{\alpha\beta} + \Lambda g_{\alpha\beta}, \end{aligned} \quad (17.10)$$

where a , b , and Λ are constants.

- (b) Show that this tensor has an automatically vanishing divergence if and only if $b = -\frac{1}{2}a$.

- (c) Show that, in addition, this tensor vanishes in flat spacetime, if and only if $\Lambda = 0$ —i.e., if and only if it is a multiple of the Einstein tensor $G_{\alpha\beta} = R_{\alpha\beta} - \frac{1}{2}Rg_{\alpha\beta}$. (Do not bother to prove that $\nabla \cdot \mathbf{G} \equiv 0$; assume it as a result from Chapter 13.)

Exercise 17.2. NO TENSOR CONSTRUCTABLE FROM FIRST DERIVATIVES OF METRIC

Show that there exists *no* tensor with components constructable from the ten metric coefficients $g_{\alpha\beta}$ and their 40 first derivatives $g_{\alpha\beta,\mu}$ —except the metric tensor \mathbf{g} , and products of it with itself; e.g., $\mathbf{g} \otimes \mathbf{g}$. [Hint: Assume there exists some other such tensor, and examine its hypothesized components in a local inertial frame.]

Exercise 17.3. RIEMANN AS THE ONLY TENSOR CONSTRUCTABLE FROM, AND LINEAR IN SECOND DERIVATIVES OF METRIC

Show that (1) **Riemann**, (2) **\mathbf{g}** , and (3) tensors (e.g., **Ricci**) formed from **Riemann** and **\mathbf{g}** but linear in **Riemann**, are the only tensors that (a) are constructable from the ten $g_{\alpha\beta}$, the 40 $g_{\alpha\beta,\mu}$, and the 100 $g_{\alpha\beta,\mu\nu}$, and (b) are linear in the $g_{\alpha\beta,\mu\nu}$. [Hint: Assume there exists some other such tensor, and examine its hypothesized components in an orthonormal, Riemann-normal coordinate system. Use equations (11.30) to (11.32).]

Exercise 17.4. UNIQUENESS OF THE EINSTEIN TENSOR

(a) Show that the Einstein tensor, $G_{\alpha\beta} = R_{\alpha\beta} - \frac{1}{2}Rg_{\alpha\beta}$, is the only second-rank, symmetric tensor that (1) has components constructable solely from $g_{\alpha\beta}, g_{\alpha\beta,\mu}, g_{\alpha\beta,\mu\nu}$; (2) has components linear in $g_{\alpha\beta,\mu\nu}$; (3) has an automatically vanishing divergence, $\nabla \cdot \mathbf{G} = 0$; and (4) vanishes in flat spacetime. This provides added motivation for choosing the Einstein tensor as the left side of the field equation $\mathbf{G} = 8\pi\mathbf{T}$.

(b) Show that, when condition (4) is dropped, the most general tensor is $\mathbf{G} + \Lambda\mathbf{g}$, where Λ is a constant. (See §17.3 for the significance of this.)

§17.2. AUTOMATIC CONSERVATION OF THE SOURCE: A DYNAMIC NECESSITY

The answer $\mathbf{G} = 8\pi\mathbf{T}$ is now on hand; but what is the question? An equation has been derived that connects the Einstein-Cartan “moment of rotation” **\mathbf{G}** with the stress-energy tensor **\mathbf{T}** , but what is the purpose for which one wants this equation in the first place? If geometry tells matter how to move, and matter tells geometry how to curve, does one not have in one’s hands a Gordian knot? And how then can one ever untie it?

The story is no different in character for the dynamics of geometry than it is for other branches of dynamics. To predict the future, one must first specify, on an “initial” hypersurface of “simultaneity,” the position and velocity of every particle, and the amplitude and time-rate of change of every field that obeys a second-order wave equation. One can then evolve the particles and fields forward in time by means of their dynamic equations. Similarly, one must give information about the geometry and its first time-rate of change on the “initial” hypersurface if the Einstein field equation is to be able to predict completely and deterministically the future time-development of the entire system, particles plus fields plus geometry. (See Chapter 21 for details.)

If a prediction is to be made of the geometry, how much information has to be supplied for this purpose? The geometry of spacetime is described by the metric

$$ds^2 = g_{\alpha\beta}(\mathcal{P}) dx^\alpha dx^\beta;$$

Einstein field equation governs the evolution of spacetime geometry

that is, by the ten functions $g_{\alpha\beta}$ of location \mathcal{P} in spacetime. It might then seem that ten functions must be predicted; and, if so, that one would need for the task ten

equations. Not so. Introduce a new set of coordinates $x^{\bar{\mu}}$ by way of the coordinate transformations

$$x^\alpha = x^\alpha(x^{\bar{\mu}}),$$

and find the same spacetime geometry, with all the same bumps, rills, and waves, described by an entirely new set of metric coefficients $g_{\bar{\alpha}\bar{\beta}}(\mathcal{P})$.

It would transgress the power as well as the duty of Einstein's "geometrodynamical law" $\mathbf{G} = 8\pi\mathbf{T}$ if, out of the appropriate data on the "initial-value hypersurface," it were to provide a way to calculate, on out into the future, values for all ten functions $g_{\alpha\beta}(\mathcal{P})$. To predict all ten functions would presuppose a choice of the coordinates; and to make a choice among coordinate systems is exactly what the geometrodynamical law cannot and must not have the power to do. That choice resides of necessity in the man who studies the geometry, not in the Nature that makes the geometry. The geometry in and by itself, like an automobile fender in and by itself, is free of coordinates. The coordinates are the work of man.

It follows that the ten components $G_{\alpha\beta} = 8\pi T_{\alpha\beta}$ of the field equation must not determine completely and uniquely all ten components $g_{\mu\nu}$ of the metric. On the contrary, $G_{\alpha\beta} = 8\pi T_{\alpha\beta}$ must place only six independent constraints on the ten $g_{\mu\nu}(\mathcal{P})$, leaving four arbitrary functions to be adjusted by man's specialization of the four coordinate functions $x^\alpha(\mathcal{P})$.

How can this be so? How can the ten equations $G_{\alpha\beta} = 8\pi T_{\alpha\beta}$ be in reality only six? Answer: by virtue of the "automatic conservation of the source." More specifically, the identity $G^{\alpha\beta}_{;\beta} \equiv 0$ guarantees that the ten equations $G_{\alpha\beta} = 8\pi T_{\alpha\beta}$ contain the four "conservation laws" $T^{\alpha\beta}_{;\beta} = 0$. These four conservation laws—along with other equations—govern the evolution of the source. They do not constrain in any way the evolution of the geometry. The geometry is constrained only by the six remaining, independent equations in $G_{\alpha\beta} = 8\pi T_{\alpha\beta}$.

When viewed in this way, the "automatic conservation of the source" is not merely a philosophically attractive principle. It is, in fact, an absolute dynamic necessity. Without "automatic conservation of the source," the ten $G_{\alpha\beta} = 8\pi T_{\alpha\beta}$ would place ten constraints on the ten $g_{\alpha\beta}$, thus fixing the coordinate system as well as the geometry. With "automatic conservation," the ten $G_{\alpha\beta} = 8\pi T_{\alpha\beta}$ place four constraints (local conservation of energy and momentum) on the source, and six constraints on the ten $g_{\alpha\beta}$, leaving four of the $g_{\alpha\beta}$ to be adjusted by adjustment of the coordinate system.

$\mathbf{G} = 8\pi\mathbf{T}$ must determine only six metric components; the other four are adjustable by changes of coordinates

$\mathbf{G} = 8\pi\mathbf{T}$ leaves four components of metric free because it satisfies four identities
 $0 \equiv \nabla \cdot \mathbf{G} = 8\pi \nabla \cdot \mathbf{T}$
 ("automatic conservation of source")

§17.3. COSMOLOGICAL CONSTANT

In 1915, when Einstein developed his general relativity theory, the permanence of the universe was a fixed item of belief in Western philosophy. "The heavens endure from everlasting to everlasting." Thus, it disturbed Einstein greatly to discover (Chapter 27) that his geometrodynamical law $\mathbf{G} = 8\pi\mathbf{T}$ predicts a *nonpermanent* universe; a dynamic universe; a universe that originated in a "big-bang" explosion,

Einstein's motivation for introducing a cosmological constant

or will be destroyed eventually by contraction to infinite density, or both. Faced with this contradiction between his theory and the firm philosophical beliefs of the day, Einstein weakened; he modified his theory.

The only conceivable modification that does not alter vastly the structure of the theory is to change the lefthand side of the geometrodynamic law $\mathbf{G} = 8\pi\mathbf{T}$. Recall that the lefthand side is *forced* to be the Einstein tensor, $G_{\alpha\beta} = R_{\alpha\beta} - \frac{1}{2}Rg_{\alpha\beta}$, by three assumptions:

- (1) \mathbf{G} vanishes when spacetime is flat.
- (2) \mathbf{G} is constructed from the Riemann curvature tensor and the metric and nothing else.
- (3) \mathbf{G} is distinguished from other tensors that can be built from **Riemann** and **g** by the demands (1) that it be linear in **Riemann**, as befits any natural measure of curvature; (2) that, like **T**, it be symmetric and of second rank; and (3) that it have an automatically vanishing divergence, $\nabla \cdot \mathbf{G} \equiv 0$.

Denote a new, modified lefthand side by “ \mathbf{G} ”, with quotation marks to avoid confusion with the standard Einstein tensor. To abandon $\nabla \cdot \mathbf{G} \equiv 0$ is impossible on dynamic grounds (see §17.2). To change the symmetry or rank of “ \mathbf{G} ” is impossible on mathematical grounds, since “ \mathbf{G} ” must be equated to **T**. To let “ \mathbf{G} ” be nonlinear in **Riemann** would vastly complicate the theory. To construct “ \mathbf{G} ” from anything except **Riemann** and **g** would make “ \mathbf{G} ” no longer a measure of spacetime geometry and would thus violate the spirit of the theory. After much anguish, one concludes that the assumption which one might drop with least damage to the beauty and spirit of the theory is assumption (1), that “ \mathbf{G} ” vanish when spacetime is flat. But even dropping this assumption is painful: (1) although “ \mathbf{G} ” might still be in some sense a measure of geometry, it can no longer be a measure of curvature; and (2) flat, empty spacetime will no longer be compatible with the geometrodynamic law ($\mathbf{G} \neq 0$ in flat, empty space, where $\mathbf{T} = 0$). Nevertheless, these consequences were less painful to Einstein than a dynamic universe.

The only tensor that satisfies conditions (2) and (3) [with (1) abandoned] is the Einstein tensor plus a multiple of the metric:

$$\text{“}G_{\alpha\beta}\text{”} = R_{\alpha\beta} - \frac{1}{2}g_{\alpha\beta}R + \Lambda g_{\alpha\beta} = G_{\alpha\beta} + \Lambda g_{\alpha\beta}$$

(exercise 17.1; see also exercise 17.4). Thus was Einstein (1917) led to his modified field equation

$$\mathbf{G} + \Lambda\mathbf{g} = 8\pi\mathbf{T}. \quad (17.11)$$

The constant Λ he called the “cosmological constant”; it has dimensions cm^{-2} .

The modified field equation, by contrast with the original, admits a static, unchanging universe as one particular solution (see Box 27.5). For this reason, Einstein in 1917 was inclined to place his faith in the modified equation. But thirteen years later Hubble discovered the expansion of the universe. No longer was the cosmological constant necessary. Einstein, calling the cosmological constant “the biggest

Einstein's field equation with the cosmological constant

Why Einstein abandoned the cosmological constant

blunder of my life," abandoned it and returned to his original geometrodynamic law, $\mathbf{G} = 8\pi\mathbf{T}$ [Einstein (1970)].

A great mistake Λ was indeed!—not least because, had Einstein stuck by his original equation, he could have claimed the expansion of the universe as the most triumphant prediction of his theory of gravity.

A mischievous genie, once let out of a bottle, is not easily reconfined. Many workers in cosmology are unwilling to abandon the cosmological constant. They insist that it be abandoned only after cosmological observations reveal it to be negligibly small. As a modern-day motivation for retaining the cosmological constant, one sometimes rewrites the modified field equation in the form

$$\mathbf{G} = 8\pi[\mathbf{T} + \mathbf{T}^{(\text{VAC})}], \quad (17.12\text{a})$$

$$\mathbf{T}^{(\text{VAC})} \equiv -(\Lambda/8\pi)\mathbf{g} \quad (17.12\text{b})$$

A modern-day motivation for the cosmological constant:
vacuum polarization

and interprets $\mathbf{T}^{(\text{VAC})}$ as a stress-energy tensor associated with the vacuum. This viewpoint speculates [Zel'dovich (1967)] that the vacuum polarization of quantum field theory endows the vacuum with the nonzero stress-energy tensor (17.12b), which is completely unobservable except by its gravitational effects. Unfortunately, today's quantum field theory is too primitive to allow a calculation of $\mathbf{T}^{(\text{VAC})}$ from first principles. (See, however, exercise 17.5.)

The mass-energy density that the cosmological constant attributes to the vacuum is

$$\rho^{(\text{VAC})} = T_{\hat{\theta}\hat{\theta}}^{(\text{VAC})} = +\Lambda/8\pi. \quad (17.13)$$

Observational limit on the cosmological constant

If $\Lambda \neq 0$, it must at least be so small that $\rho^{(\text{VAC})}$ has negligible gravitational effects [$|\rho^{(\text{VAC})}| < \rho^{(\text{MATTER})}$] wherever Newton's theory of gravity gives a successful account of observations. The systems of lowest density to which one applies Newtonian theory with some (though not great) success are small clusters of galaxies. Hence, one can place the limit

$$|\rho^{(\text{VAC})}| = |\Lambda|/8\pi \lesssim \rho^{(\text{CLUSTER})} \sim 10^{-29} \text{ g/cm}^3 \sim 10^{-57} \text{ cm}^{-2} \quad (17.14)$$

Why one ignores the cosmological constant everywhere except in cosmology

on the value of the cosmological constant. Evidently, even if $\Lambda \neq 0$, Λ is so small that it is totally unimportant on the scale of a galaxy or a star or a planet or a man or an atom. Consequently it is reasonable to stick with Einstein's original geometrodynamic law ($\mathbf{G} = 8\pi\mathbf{T}$; $\Lambda = 0$) everywhere, except occasionally when discussing cosmology (Chapters 27–30).

Exercise 17.5. MAGNITUDE OF COSMOLOGICAL CONSTANT

EXERCISE

- (a) What is the order of magnitude of the influence of the cosmological constant on the celestial mechanics of the solar system if $\Lambda \sim 10^{-57} \text{ cm}^{-2}$?

(b) Show that the mass-energy density of the vacuum $\rho^{(\text{VAC})} = \Lambda/8\pi \sim 10^{-29} \text{ g/cm}^3$, corresponding to the maximum possible value of Λ , agrees in very rough magnitude with

$$\frac{\text{rest mass of an elementary particle}}{(\text{Compton wavelength of particle})^3} \times (\text{gravitational fine-structure constant}) \\ = \frac{m}{(\hbar/m)^3} \frac{m^2}{\hbar} = \frac{m^6}{\hbar^4}$$

[Zel'dovich (1967, 1968)]. This numerology is suggestive, but has not led to any believable derivation of a stress-energy tensor for the vacuum.

§17.4. THE NEWTONIAN LIMIT

Just as quantum mechanics reduces to classical mechanics in the “correspondence limit” of large actions, $I \gg \hbar$, so general relativity reduces to Newtonian theory in the “correspondence limit” of weak gravity and low velocities. (On “correspondence limits,” see Box 17.1.) This section elucidates, in some mathematical detail, the correspondence between general relativity and Newtonian theory. It begins with “passive” aspects of gravitation (response of matter to gravity) and then turns to “active” aspects (generation of gravity by matter).

Consider an isolated system—e.g., the solar system—in which Newtonian theory is highly accurate. In order that special relativistic effects not be noticeable, all

Box 17.1 CORRESPONDENCE PRINCIPLES

A. General Remarks and Specific Examples

1. As physics develops and expands, its unity is maintained by a network of correspondence principles, through which simpler theories maintain their vitality by links to more sophisticated but more accurate ones.
 - a. Physical optics, with all the new diffraction and interference phenomena for which it accounted, nevertheless also had to account, and did account, for the old, elementary, geometric optics of mirrors and lenses. Geometric optics is recovered from physical optics in the mathematical “correspondence

principle limit” in which the wavelength is made indefinitely small in comparison with all other relevant dimensions of the physical system.

- b. Newtonian mechanics is recovered from the mechanics of special relativity in the mathematical “correspondence principle limit” in which all relevant velocities are negligibly small compared to the speed of light.
- c. Thermodynamics is recovered from its successor theory, statistical mechanics, in the mathematical “correspondence principle limit” in which so many particles are taken into account that fluctuations in pressure,

- particle number, and other physical quantities are negligible compared to the average values of these parameters of the system.
- d. Classical mechanics is recovered from quantum mechanics in the “correspondence principle limit” in which the quantum numbers of the quantum states in question are so large, or the quantities of action that come into play are so great compared to \hbar , that wave and diffraction phenomena make negligible changes in the predictions of standard deterministic classical mechanics. Niels Bohr formulated and took advantage of this correspondence principle even before any proper quantum theory existed. He used it to predict approximate values of atomic energy levels and of intensities of spectral lines. He also expounded it as a guide to all physicists, first in searching for a proper version of the quantum theory, and then in elucidating the content of this theory after it was found.
 2. In all these examples and others, the newer, more sophisticated theory is “better” than its predecessor because it gives a good description of a more extended domain of physics, or a more accurate description of the same domain, or both.
 3. The correspondence between the newer theory and its predecessor (a) gives one the power to recover the older theory from the newer; (b) can be exhibited by straightforward mathematics; and (c), according to the historical record, often guided the development of the newer theory.

B. Correspondence Structure of General Relativity

1. Einstein’s theory of gravity has as distinct limiting cases (a) special relativity; (b) the “linear-

ized theory of gravity”; (c) Newton’s theory of gravity; and (d) the post-Newtonian theory of gravity. Thus, it has a particularly rich correspondence structure.

- a. *Correspondence with special relativity:* General relativity has two distinct kinds of correspondence with special relativity. The first is the limit of vanishing gravitational field everywhere (vanishing curvature); in this limit one can introduce a global inertial frame, set $g_{\mu\nu} = \eta_{\mu\nu}$, and recover completely and precisely the theory of special relativity. The second is local rather than global; it is the demand (“correspondence principle”; “equivalence principle”) that in a local inertial frame all the laws of physics take on their special relativistic forms. As was seen in Chapter 16, this puts no restrictions on the metric (except that $g_{\mu\nu} = \eta_{\mu\nu}$ and $g_{\mu\nu,\alpha} = 0$ in local inertial frames); but it places severe constraints on the behavior of matter and fields in the presence of gravity.
- b. *Correspondence with Newtonian theory:* In the limit of weak gravitational fields, low velocities, and small pressures, general relativity reduces to Newton’s theory of gravity. The correspondence structure is explored mathematically in the text of §17.4.
- c. *Correspondence with post-Newtonian theory:* When Newtonian theory is nearly valid, but “first-order relativistic corrections” might be important, one often uses the “post-Newtonian theory of gravity.” Chapter 39 expounds the post-Newtonian theory and its correspondence with both general relativity and Newtonian theory.
- d. *Correspondence with linearized theory:* In the limit of weak gravitational fields, but possibly large velocities and pressures ($v \sim 1$, $T_{jk} \sim T_{00}$) general relativity reduces to the “linearized theory of gravity”. This correspondence is explored in Chapter 18.

Conditions which a system must satisfy for Newton's theory of gravity to be accurate

velocities in the system, relative to its center of mass and also relative to the Newtonian coordinates, must be small compared to the speed of light

$$v \ll 1. \quad (17.15a)$$

As a particle falls from the outer region of the system to the inner region, gravity accelerates it to a kinetic energy $\frac{1}{2}mv^2 \sim |m\Phi|_{\max}$. [Here $\Phi < 0$ is Newton's gravitational potential, so normalized that $\Phi(\infty) = 0$.] The resulting velocity will be small only if

$$|\Phi| \ll 1. \quad (17.15b)$$

Internal stresses in the system also produce motion—e.g., sound waves. Such waves have characteristic velocities of the order of $|T^{ij}/T^{00}|^{1/2}$ —for example, the speed of sound in a perfect fluid is

$$v = (dp/d\rho)^{1/2} \sim (p/\rho)^{1/2} \sim |T^{ij}/T^{00}|^{1/2}.$$

In order that these velocities be small compared to the speed of light, all stresses must be small compared to the density of mass-energy

$$|T^{ij}/T^{00}| = |T^{ij}|/\rho \ll 1. \quad (17.15c)$$

When, and only when conditions (17.15) hold, one can expect Newtonian theory to describe accurately the system being studied. Correspondence of general relativity with Newtonian theory for gravity in a passive role then demands that the geodesic world lines of freely falling particles reduce to the Newtonian world lines

$$d^2x^i/dt^2 = -\partial\Phi/\partial x^i. \quad (17.16)$$

“Newtonian coordinates” defined

Moreover, they must reduce to this form in *any* relativistic coordinate system where the source and test particles have low velocities $v \ll 1$, and where coordinate lengths and times agree very nearly with the lengths and times of the Newtonian coordinates—which in turn are proper lengths and times as measured by rods and clocks. Thus, the relevant coordinates (called “Galilean” or “Newtonian” coordinates) are ones in which

$$g_{\mu\nu} = \eta_{\mu\nu} + h_{\mu\nu}, \quad |h_{\mu\nu}| \ll 1, \quad |v^j| = |dx^j/dt| \ll 1 \quad (17.17)$$

(weak gravitational field; nearly inertial coordinates; low velocities). In such a coordinate system, the geodesic world lines of test particles have the form

$$\frac{d^2x^i}{dt^2} = \frac{d^2x^i}{d\tau^2} \quad (\text{since } dt/d\tau \approx 1 \text{ when } |h_{\mu\nu}| \ll 1 \text{ and } |v^j| \ll 1)$$

$$= -\Gamma^i_{\alpha\beta} \frac{dx^\alpha}{d\tau} \frac{dx^\beta}{d\tau} \quad (\text{geodesic equation})$$

$$= -\Gamma^i_{00} \quad (\text{since } dt/d\tau \approx 1 \text{ and } |dx^j/d\tau| \ll 1)$$

$$= -\Gamma_{i00} \quad (\text{since } g_{\mu\nu} \approx \eta_{\mu\nu})$$

$$= \frac{1}{2} h_{00,i} - h_{0i,0} \quad (\text{equation for } \Gamma_{\alpha\beta\gamma} \text{ in terms of } g_{\alpha\beta,\gamma})$$

$$= \frac{1}{2} h_{00,i} \quad \begin{cases} \text{all velocities small compared to } c \text{ implies time derivatives small compared to space derivatives} \\ \text{—i.e., } h_{\alpha\beta,0} \sim v h_{\alpha\beta,i} \end{cases}.$$

These geodesic world lines do, indeed, reduce to those of Newtonian theory [equation (17.16)] if one makes the identification

$$\Gamma^i_{00} = -\frac{1}{2} h_{00,i} = \Phi_{,i}. \quad (17.18)$$

Together with the boundary conditions $\Phi(r = \infty) = 0$ and $h_{\mu\nu}(r = \infty) = 0$ (coordinates Lorentz far from the source), this identification implies $h_{00} = -2\Phi$; i.e.,

$$g_{00} = -1 - 2\Phi \text{ for nearly Newtonian systems in Newtonian coordinates.} \quad (17.19)$$

Note that the correspondence tells one the form of h_{00} for nearly Newtonian systems, but not the forms of the other components of the metric perturbation. In fact, the other $h_{\mu\nu}$ could perfectly well be of the same order of magnitude as $h_{00} \sim \Phi$, without influencing the world lines of slowly moving particles, because they always enter the geodesic equation multiplied by the small numbers v or v^2 , or differentiated by t rather than by x^i . The forms of the other $h_{\mu\nu}$ and their small corrections to the Newtonian motion will be explored in Chapters 18, 39, and 40.

The relation $g_{00} = -1 - 2\Phi$ is the mathematical embodiment of the correspondence between general relativity theory and Newtonian theory for passive aspects of gravity. Together with the “validity conditions” (17.15, 17.17), it is a foundation from which one can derive all other aspects of the correspondence for “passive gravity,” including the relation

$$R^i_{0j0} = \partial^2 \Phi / \partial x^i \partial x^j \quad (17.20)$$

(exercise 17.6). Alternatively, all other aspects of this correspondence can be derived by direct comparison of Newton’s predictions with Einstein’s. For example, to derive equation (17.20), examine the relative acceleration of two test particles, one at $x^i + \xi^i$ and the other at x^i . According to Newton

$$\begin{aligned} \frac{d^2\xi^i}{dt^2} &= \frac{d^2(x^i + \xi^i)}{dt^2} - \frac{d^2x^i}{dt^2} \\ &= -\left. \frac{\partial \Phi}{\partial x^i} \right|_{\text{at } x^i + \xi^i} + \left. \frac{\partial \Phi}{\partial x^i} \right|_{\text{at } x^i} = \frac{-\partial^2 \Phi}{\partial x^i \partial x^j} \xi^j. \end{aligned}$$

For comparison, Einstein predicts (equation of geodesic deviation)

$$\frac{D^2\xi^i}{d\tau^2} = \frac{d^2\xi^i}{dt^2} = -R^i_{0j0}\xi^j.$$

↑
[by conditions (17.15) and (17.17)]

Direct comparison gives relation (17.20).

Turn now from correspondence for passive aspects of gravity to correspondence for active aspects. According to Einstein, mass generates gravity (spacetime curvature) by the geometrodynamical law $\mathbf{G} = 8\pi\mathbf{T}$. Apply this law to a nearly Newtonian system, and by the chain of reasoning that precedes equation (17.8) derive the relation

$$R_{00} = 4\pi\rho. \quad (17.21)$$

Einstein gravity reduces to Newton gravity only if, in Newtonian coordinates,
 $g_{00} = -1 - 2\Phi$

The correspondence between Einstein theory and Newton theory for all “passive” aspects of gravity

The Newtonian limit of the Einstein field equation is
 $\nabla^2\Phi = 4\pi\rho$

Combine with the contraction of (17.20),

$$R_{00} = R^i_{0i0} + R^0_{000} = \underset{\substack{\uparrow \\ 0}}{\partial^2 \Phi / \partial x^i \partial x^i} = \nabla^2 \Phi,$$

and thereby obtain Newton's equation for the generation of gravity by mass

$$\nabla^2 \Phi = 4\pi\rho. \quad (17.22)$$

Thus, Einstein's field equation reduces to Newton's field equation in the Newtonian limit.

The correspondence between Newton and Einstein, although clear and straightforward as outlined above, is even more clear and straightforward when Newton's theory of gravity is rewritten in Einstein's language of curved spacetime (Chapter 12; exercise 17.7).

EXERCISES

Exercise 17.6. RAMIFICATIONS OF CORRESPONDENCE FOR GRAVITY IN A PASSIVE ROLE

From the correspondence relation $g_{00} = -1 - 2\Phi$, and from conditions (17.15) and (17.17) for Newtonian physics, derive the correspondence relations

$$\Gamma^i_{00} = \partial\Phi/\partial x^i, \quad R^i_{0j0} = \partial^2\Phi/\partial x^i \partial x^j.$$

Exercise 17.7. CORRESPONDENCE IN THE LANGUAGE OF CURVED SPACETIME [Track 2]

Exhibit the correspondence between the Einstein theory and Cartan's curved-spacetime formulation of Newtonian theory (Chapter 12).

§17.5. AXIOMATIZE EINSTEIN'S THEORY?

Find the most compact and reasonable axiomatic structure one can for general relativity? Then from the axioms derive Einstein's field equation,

$$\mathbf{G} = 8\pi\mathbf{T}?$$

That approach would follow tradition. However, it may be out of date today. More than half a century has gone by since November 25, 1915. For all that time the equation has stood unchanged, if one ignores Einstein's temporary "aberration" of adding the cosmological constant. In contrast the derivations have evolved and become more numerous and more varied. In the beginning axioms told what equation is acceptable. By now the equation tells what axioms are acceptable. Box 17.2 sketches a variety of sets of axioms, and the resulting derivations of Einstein's equation.

There are many ways (Box 17.2) to derive the Einstein field equation

(continued on page 429)

**Box 17.2 SIX ROUTES TO EINSTEIN'S GEOMETRODYNAMIC LAW
OF THE EQUALITY OF CURVATURE AND ENERGY DENSITY
("EINSTEIN'S FIELD EQUATION")**

[Recommended to the attention of Track-1 readers are only route 1 (automatic conservation of the source, plus correspondence with Newtonian theory) and route 2 (Hilbert's variational principle); and even Track-2 readers are advised to finish the rest of this chapter before they study route 3 (physics on a spacelike slice), route 4 (going from superspace to Einstein's equation), route 5 (field of spin 2 in an "unobservable flat spacetime" background), and route 6 (gravitation as an elasticity of space that arises from particle physics).]

1. Model geometrodynamics after electrodynamics and treat "automatic conservation of the source" and correspondence with the Newtonian theory of gravity as the central considerations.
 - a. Particle responds in electrodynamics to field; in general relativity, to geometry.
 - b. The potential for the electromagnetic field is the 4-vector \mathbf{A} (components A_μ).
The potential for the geometry is the metric tensor \mathbf{g} (components $g_{\mu\nu}$).
 - c. The electromagnetic potential satisfies a wave equation with source term (4-current) on the right,

$$\left(\frac{\partial A_\nu}{\partial x^\mu} - \frac{\partial A_\mu}{\partial x^\nu} \right)^{\cdot\nu} = 4\pi j_\mu, \quad (1)$$

so constructed that conservation of the source, $j_\mu^{\cdot\mu} = 0$, is automatic (consequence of an identity fulfilled by the lefthand side). By analogy, the geometrodynamic potential must also satisfy a wave equation with source term (stress-energy tensor) on the right,

$$G_{\mu\nu} = 8\pi T_{\mu\nu}, \quad (2)$$

so constructed that conservation of the source, $T_{\mu\nu}^{\cdot\nu} = 0$ (Chapter 16) is "automatic." This conservation is automatic here because the lefthand side of the equation is a tensor (the Einstein tensor; see Box 8.6 or Chapter 15), built from the metric components and their second derivatives, that fulfills the identity $G_{\mu\nu}^{\cdot\nu} \equiv 0$.

- d. No other tensor which (1) is linear in the second derivatives of the metric components, (2) is free of higher derivatives, and (3) vanishes in flat spacetime, satisfies such an identity.
- e. The constant of proportionality (8π) is fixed by the choice of units [here geometric; see Box 1.8] and by the requirement ("correspondence with Newtonian theory") that a test particle shall oscillate back and forth through a collection of matter of density ρ , or revolve in circular orbit around that collection of matter, at a circular frequency given by $\omega^2 = (4\pi/3)\rho$ (Figure

Box 17.2 (continued)

- 1.12). The foregoing oversimplifies, and omits Einstein's temporary false turns, but otherwise summarizes the reasoning he pursued in arriving at his field equation. This reasoning is spelled out in more detail in the text of Chapter 17.
2. Take variational principle as central.
 - a. Construct out of the metric components the only scalar that exists that (1) is linear in the second derivatives of the metric tensor, (2) contains no higher derivatives, and (3) vanishes in flat spacetime: namely, the Riemann scalar curvature invariant, R .
 - b. Construct the invariant integral,

$$I = \frac{1}{16\pi} \int_{\Omega} R(-g)^{1/2} d^4x. \quad (3)$$

- c. Make small variations, $\delta g^{\mu\nu}$, in the metric coefficients $g^{\mu\nu}$ in the interior of the four-dimensional region Ω , and find that this integral changes by the amount

$$\delta I = \frac{1}{16\pi} \int_{\Omega} G_{\mu\nu} \delta g^{\mu\nu} (-g)^{1/2} d^4x. \quad (4)$$

- d. Demand that I should be an extremum with respect to the choice of geometry in the region interior to Ω ($\delta I = 0$ for arbitrary $\delta g^{\mu\nu}$; "principle of extremal action").
- e. Thus arrive at the Einstein field equation for empty space,

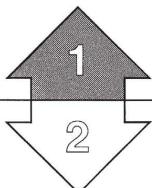
$$G_{\mu\nu} = 0. \quad (5)$$

- f. The continuation of the reasoning leads to the identity

$$G_{\mu\nu}{}^{;v} = 0.$$

Chapter 21, on the variational principle, gives more detail and takes up the additional term that appears on the righthand side of (5) when matter or fields or both are present.

- g. This approach goes back to David Hilbert (1915). No route to the field equations is quicker. Moreover, it connects immediately (see the following section here, 2') with the quantum principle of the "democracy of all histories" [Feynman (1942); Feynman and Hibbs (1965)]. The variational principle is spelled out in more detail in Chapter 21.
- 2'. An aside on the meaning of the classical action integral for the real world of quantum physics.
 - a. A "history of geometry," H , is a spacetime, that is to say, a four-dimensional manifold with four-dimensional $- + + +$ Riemann metric that (1) reduces on one spacelike hypersurface ("hypersurface of simultaneity") to a specified "initial value 3-geometry," A , with positive definite metric and (2) reduces on



another spacelike hypersurface to a specified “final value 3-geometry,” B , also with positive definite metric.

- b. The classical variational principle of Hilbert, as reformulated by Arnowitt, Deser, and Misner, provides a prescription for the dynamical path length, I_H , of any conceivable history H , classically allowed or not, that connects A and B (see Chapter 21 for a fuller statement for what can and must be specified on the initial hypersurface of simultaneity, and on the final one, and for alternative choices of the integrand in the action principle).
- c. Classical physics says that a history H is allowed only if it extremizes the dynamic path length I as compared to all nearby histories. Quantum physics says that all histories occur with equal probability amplitude, in the following sense. The probability amplitude for “the dynamic geometry of space to transit from A to B ” by way of the history H with action integral I_H , and by way of histories that lie within a specified infinitesimal range, $\mathcal{D}H$, of the history H , is given by the expression

$$\left(\begin{array}{l} \text{probability amplitude} \\ \text{to transit from } A \text{ to} \\ B \text{ by way of history } H \\ \text{and histories lying} \\ \text{within the range } \mathcal{D}H \\ \text{about } H \end{array} \right) \sim \exp(iI_H/\hbar)N\mathcal{D}H. \quad (6)$$

Here the normalization factor, N , is the same for all conceivable histories H , allowed or not, that lead from A to B (“principle of democracy of histories”). The quantum of angular momentum, $\hbar = h/2\pi$, expressed in geometric units, has the value

$$\hbar = \hbar_{\text{conv}} G/c^3 = (L^*)^2, \quad (7)$$

where L^* is the Planck length, $L^* = 1.6 \times 10^{-33}$ cm.

- d. The classically allowed history receives “preference without preference.” That history, and histories H that differ from it so little that $\delta I = I_H - I_{\text{class}}$ is only of the order \hbar and less, give contributions to the probability amplitude that interfere constructively. In contrast, destructive interference effectively wipes out the contribution (to the probability amplitude for a transition) that comes from histories that differ more from the classically allowed history. Thus there are quantum fluctuations in the geometry, but they are fluctuations of limited magnitude. The smallness of \hbar ensures that the scale of these fluctuations is unnoticeable at everyday distances (see the further discussion in Chapters 43 and 44). In this sense classical geometrodynamics is a good approximation to the geometrodynamics of the real world of quantum physics.
- 3. “Physics on a spacelike slice or hypersurface of simultaneity,” again with electromagnetism as the model.
 - a. Say over and over “lines of magnetic force never end” and come out with half of Maxwell’s equations. Say over and over “lines of electric force end

The rest of this chapter is Track 2. No previous track-2 material is needed as preparation for it, nor is it necessary preparation for any later chapter, but it will be helpful in Chapter 21 (initial-value equations and variational principle) and in Chapter 39 (other theories of gravity).

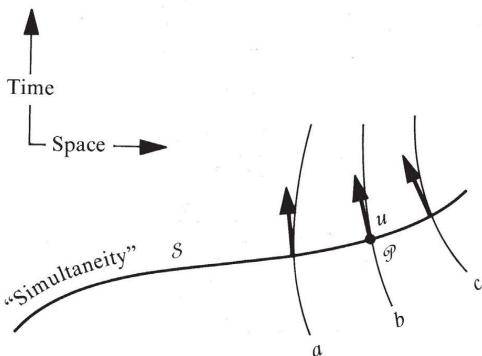
Box 17.2 (continued)

only on charge" and arrive at the other half of Maxwell's equations. Similarly, say over and over

$$\begin{pmatrix} \text{intrinsic} \\ \text{curvature} \\ \text{scalar} \end{pmatrix} + \begin{pmatrix} \text{extrinsic} \\ \text{curvature} \\ \text{scalar} \end{pmatrix} = 16\pi \begin{pmatrix} \text{local density} \\ \text{of mass-} \\ \text{energy} \end{pmatrix} \quad (8)$$

and end up with all ten components of Einstein's equation. To "say over and over" is an abbreviation for demanding that the stated principles hold on every spacelike slice through every event of spacetime.

- b. Spell out explicitly this "spacelike-slice formulation" of the equations of Maxwell and Einstein. Consider an arbitrary point of spacetime, \mathcal{P} ("event"), and an arbitrary "simultaneity" S through \mathcal{P} (hypersurface of simultaneity; spacelike slice through spacetime). Magnetic lines of force run about throughout S , but nowhere is even a single one of them permitted to end. Recall (§3.4) that the demand "lines of magnetic force never end," when imposed on *all* reference frames at \mathcal{P} (for all choices of the "simultaneity" S), guarantees not only $\nabla \cdot \mathbf{B} = 0$, but also $\nabla \times \mathbf{E} + \partial \mathbf{B} / \partial t = 0$. Similarly (§3.4) the demand that "electric lines of force never end except on electric charge," $\nabla \cdot \mathbf{E} = 4\pi J^0$, when imposed on all "simultanities" through \mathcal{P} , guarantees the remaining Maxwell equation $\nabla \times \mathbf{B} = \partial \mathbf{E} / \partial t + 4\pi \mathbf{J}$.
- c. Each simultaneity S through \mathcal{P} has its own slope and curvature. The possibility of different slopes (different local Lorentz frames at \mathcal{P}) is essential for deriving all of Maxwell's equations from the requirements of conservation of flux. Relevant though the slope thus is, the curvature of the hypersurface S never matters for the analysis of electromagnetism. It does matter, however, for any analysis of gravitation modeled on the foregoing treatment of electromagnetism.



"Simultaneity" S (spacelike hypersurface or "slice through spacetime") that cuts through event \mathcal{P} . The "simultaneity" may be considered to be defined by a set of "observers" a, b, c, \dots . Their world lines cross the simultaneity orthogonally, and their clocks all read the same proper time at the instant of crossing. Another simultaneity through \mathcal{P} may have at \mathcal{P} a different curvature or a different slope or both; and it is defined by a different band of observers, with other wrist watches.

- d. “Mass-energy curves space” is the central principle of gravitation. To spell out this principle requires one to examine in succession the terms “space” and “curvature of space” and “density of mass-energy in a given region of space.” “Space” means spacelike hypersurface; or, more specifically, a hypersurface of simultaneity \mathcal{S} that includes the point \mathcal{P} where the physics is under examination.
- e. Denote by \mathbf{u} the 4-vector normal to \mathcal{S} at \mathcal{P} . Then the density of mass-energy in the spacelike hypersurface \mathcal{S} at \mathcal{P} is
- $$\rho = u^\alpha T_{\alpha\beta} u^\beta, \quad (9)$$
- in accordance with the definition of the stress-energy tensor given in Chapter 5.
- f. This density is a single number, dependent on the inclination of the slice one cuts through spacetime, but independent of how curved one cuts this slice. If it is to be equated to “curvature of space,” that curvature must also be independent of how curved one cuts the slice.
- g. Conclude that the geometric quantity, “curvature of space,” must (1) be a single number (a scalar) that (2) depends on the inclination \mathbf{u} of the cut one makes through spacetime at \mathcal{P} in constructing the hypersurface \mathcal{S} , but (3) must be unaffected by how one curves his cut. The demand made here appears paradoxical. One seems to be asking for a measure of curvature that is independent of curvature!
- h. A closer look discloses that three distinct ideas come into consideration here. One is the scalar curvature invariant ${}^{(3)}R$ of the 3-geometry intrinsic to the hypersurface \mathcal{S} at \mathcal{P} : “intrinsic” in the sense that it is defined by, and depends exclusively on, measurements of distance made within the hypersurface. The second is the “extrinsic curvature” of this 3-geometry relative to the 4-geometry of the enveloping spacetime (“how curved one cuts his slice”; see Box 14.1 for more on the distinction between extrinsic and intrinsic curvature). The third is the curvature of the four-dimensional spacetime itself, “normal to \mathbf{u} ,” in some sense yet to be more closely defined. This is the quantity that is independent of how curved one cuts his slice. It is the quantity that is to be identified, up to a factor that depends on the choice of units, with the density of mass-energy.

Box 17.2 (continued)

- i. These three quantities are related in the following way:

$$\begin{aligned}
 & \left(\text{scalar curvature invariant, } {}^{(3)}R, \text{ of the 3-geometry intrinsic to the spacelike hypersurface } S, \text{ a quantity dependent on "how curved one cuts the slice"} \right) + \\
 & \quad \left(\text{a correction term that (a) depends only on the "extrinsic curvature" } K_{\alpha\beta} \text{ (Box 14.1 and Chapter 21) of the hypersurface relative to the four-dimensional geometry in which it is imbedded, and (b) is so calculated (a uniquely determinate calculation) that the sum of this correction term and } {}^{(3)}R \text{ is independent of "how curved one cuts his slice," and (c) has the precise value } (\text{Tr } \mathbf{K})^2 - \text{Tr } \mathbf{K}^2 \equiv (K_\alpha^\alpha)^2 - K_{\alpha\beta} K^{\alpha\beta} \right) \\
 = & \left(\text{a measure of the curvature of spacetime that depends on the 4-geometry of the spacetime and on the inclination } \mathbf{u} \text{ of the spacelike slice } S \text{ cut through spacetime, but is independent, by construction, of "how curved one cuts the slice"} \right) = \left(\text{a scalar quantity that (a) is completely defined by what has just been said and (b) can therefore be calculated in all completeness by standard differential geometry (details in Chapter 21)} \right) \\
 = & \left(2u^\alpha G_{\alpha\beta} u^\beta, \text{ where } G_{\alpha\beta} \text{ is the Einstein curvature tensor of equation 8.49 and Box 8.6} \right) = 2 \left(\text{a quantity interpreted in Track 2, Chapter 15, as the "moment of rotation" associated with a unit element of 3-volume located at } \mathcal{P} \text{ in the hypersurface orthogonal to } \mathbf{u} \right)
 \end{aligned} \tag{10}$$

- j. Conclude that the central principle, "mass-energy curves space," translates to the formula

$${}^{(3)}R + (\text{Tr } \mathbf{K})^2 - \text{Tr } \mathbf{K}^2 = 16\pi\rho, \tag{11}$$

or, in shorthand form,

$$\begin{pmatrix} \text{moment of} \\ \text{rotation} \end{pmatrix} = \begin{pmatrix} \text{intrinsic} \\ \text{curvature} \end{pmatrix} + \begin{pmatrix} \text{extrinsic} \\ \text{curvature} \end{pmatrix} = \begin{pmatrix} \text{density of} \\ \text{mass-energy} \end{pmatrix}, \quad (12)$$

valid for every spacelike slice through spacetime at any arbitrary point \mathcal{P} .

- k. All of Einstein's geometrodynamics is contained in this statement as truly as all of Maxwell's electrodynamics is contained in the statement that the number of lines of force that end in an element of volume is equal to 4π times the amount of charge in that element of volume. The factor 16π is appropriate for the geometric system of units in use in this book (density ρ in cm^{-2} given by $G/c^2 = 0.742 \times 10^{-28} \text{ cm/g}$ multiplied by the density ρ_{conv} expressed in the conventional units of g/cm^3).

- l. Reexpress the principle that "mass-energy curves space" in the form

$$2u^\alpha G_{\alpha\beta} u^\beta = 16\pi u^\alpha T_{\alpha\beta} u^\beta. \quad (13)$$

Demand that this equation should hold for every simultaneity that cuts through \mathcal{P} , whatever its "inclination" \mathbf{u} .

- m. Conclude that the coefficients on the two sides of (13) must agree; thus,

$$G_{\alpha\beta} = 8\pi T_{\alpha\beta}, \quad (14)$$

Einstein's equation in the language of components; or, in the language of abstract geometric quantities,

$$\mathbf{G} = 8\pi \mathbf{T}. \quad (15)$$

- 4. Going from superspace to Einstein's equation rather than from Einstein's equation to superspace.
 - a. A fourth route to Einstein's equation starts with the advanced view of geometrodynamics that is spelled out in Chapter 43. One notes there that the dynamics of geometry unfolds in superspace. Superspace has an infinite number of dimensions. Any one point in superspace describes a complete 3-geometry, ${}^{(3)}\mathcal{G}$, with all its bumps and curvatures. The dynamics of geometry leads from point to point in superspace.
 - b. Like the dynamics of a particle, the dynamics of geometry lends itself to distinct but equivalent mathematical formulations, associated with the names of Lagrange, of Hamilton, and of Hamilton and Jacobi. Of these the most convenient for the present analysis is the last ("H-J").
 - c. In the problem of one particle moving in one dimension under the influence of a potential $V(x)$, the H-J equation reads

$$\underbrace{-\frac{\partial S}{\partial t}}_{\substack{\uparrow \\ \text{total} \\ \text{energy}}} = \underbrace{\frac{1}{2m} \left(\frac{\partial S}{\partial x} \right)^2}_{\substack{\uparrow \\ \text{kinetic} \\ \text{energy}}} + V(x). \quad (16)$$

Box 17.2 (continued)

It has the solution

$$S_E(x,t) = -Et + \int^x [2m(E - V)]^{1/2} dx. \quad (17)$$

Out of this solution one reads the motion by applying the “condition of constructive interference,”

$$\frac{\partial S_E(x,t)}{\partial E} = 0 \quad (18)$$

(one equation connecting the two quantities x and t ; for more on the condition of constructive interference and the H-J method in general, see Boxes 25.3 and 25.4).

- d. In the corresponding equation for the dynamics of geometry, one deals with a function $S = S(^{(3)}\mathcal{G})$ of the 3-geometry. It depends on the 3-geometry itself, and not on the vagaries of one's choice of coordinates or on the corresponding vagaries in the metric coefficients of the 3-geometry,

$$ds^2 = {}^{(3)}g_{mn} dx^m dx^n \quad (19)$$

(${}^{(3)}$ to indicate 3-geometry omitted hereafter for simplicity). This function obeys the H-J equation [the analog of (16)]

$$-(16\pi)^2 \frac{1}{2g} (g_{im}g_{jn} + g_{in}g_{jm} - g_{ij}g_{mn}) \frac{\delta S}{\delta g_{ij}} \frac{\delta S}{\delta g_{mn}} + {}^{(3)}R = 16\pi\rho. \quad (20)$$

- e. Out of this equation for the dynamics of geometry in superspace one can deduce the Einstein field equation by reasoning similar to that employed in going from (17) to (18) (Gerlach 1969).
- f. It would appear that one must break new ground, and establish new foundations, if one is to find out how to regard the “Einstein-Hamilton-Jacobi equation” (20) as more basic than the Einstein field equation that one derives from it. [Since done, by Hojman, Kuchař, and Teitelboim (1973 preprint).]
- 5. Einstein's geometrodynamics viewed as the standard field theory for a field of spin 2 in an “unobservable flat spacetime” background.
 - a. This approach to Einstein's field equation has a long history, references to which will be found in §7.1 and §18.1. (Further discussion of this approach will be found in those two sections and in Box 7.1, exercise 7.3, and Box 18.1).
 - b. The following summary is quoted from Deser (1970): “We wish to give a simple physical derivation of the nonlinearity . . . , using a now familiar argument . . . leading from the linear, massless, spin-2 field to the full Einstein equations

- c. “The Einstein equations may be derived nongeometrically by noting that the free, massless, spin-2 field equations,

$$R^L_{\mu\nu}(\phi) - \frac{1}{2} R^L_{\alpha\alpha}(\phi)\eta_{\mu\nu} \equiv G^L_{\mu\nu}(\phi) \equiv [(\eta_{\mu\alpha}\eta_{\nu\beta} - \eta_{\mu\nu}\eta_{\alpha\beta})\square + \eta_{\mu\nu}\partial_\alpha\partial_\beta + \eta_{\alpha\beta}\partial_\mu\partial_\nu - \eta_{\mu\alpha}\partial_\nu\partial_\beta - \eta_{\nu\beta}\partial_\mu\partial_\alpha]\phi_{\alpha\beta} = 0, \quad (21)$$

whose source is the matter stress-tensor $T_{\mu\nu}$, must actually be coupled to the *total* stress-tensor, including that of the ϕ -field itself. That is, while the free-field equations (21) are of course quite consistent as they stand, [they are not] when there is a dynamic system's $T_{\mu\nu}$ as a source. For then the left side, which is identically divergenceless, is inconsistent with the right, since the coupling implies that $T^{\mu\nu}_{,\nu}$, as computed from the matter equations of motion, is no longer conserved.

- d. “To remedy this [violation of the principle of conservation of momentum and energy], the stress tensor ${}^{(2)}\theta_{\mu\nu}$ arising from the quadratic Lagrangian ${}^{(2)}L$ responsible for equation (21) is then inserted on the right.
- e. “But the Lagrangian ${}^{(3)}L$ leading to these modified equations is then cubic, and itself contributes a cubic ${}^{(3)}\theta_{\mu\nu}$.
- f. “This series continues indefinitely, and sums (if properly derived!) to the full nonlinear Einstein equations, $G_{\mu\nu}$ ([calculated from] $\eta_{\alpha\beta} + \phi_{\alpha\beta} = -\kappa T_{\mu\nu}$ [$+ 8\pi T_{\mu\nu}$ in the geometric units and sign conventions of this book]), which are an infinite series in the deviation $\phi_{\mu\nu}$ of the metric $g_{\mu\nu}$ from its Minkowskian value $\eta_{\mu\nu}$.
- g. Once the iteration is begun (whether or not a $T_{\mu\nu}$ is actually present), it must be continued to all orders, since conservation only holds for the full series

$\sum_{n=2}^{\infty} {}^{(n)}\theta_{\mu\nu}$. Thus, the theory is either left in its (physically irrelevant) free linear

form (21), or it *must* be an infinite series.”

- h. For details, see Deser (1970); the paper goes on (1) to take advantage of a well-chosen formalism (2) to rearrange the calculation, and thus (3) to “derive the full Einstein equations, on the basis of the same self-coupling requirement, but with the advantages that the full theory emerges in closed form with just one added (cubic) term, rather than as an infinite series.”
- i. Deser summarizes the analysis at the end thus: “Consistency has therefore led us to universal coupling, which implies the equivalence principle. It is at this point that the geometric interpretation of general relativity arises, since *all* matter now moves in an effective Riemann space of metric $g^{\mu\nu} \equiv \eta^{\mu\nu} + h^{\mu\nu}$. . . [The] initial flat ‘background’ space is no longer observable.” In other words, this approach to Einstein's field equation can be summarized as “curvature without curvature” or—equally well—as “flat spacetime without flat spacetime”!

Box 17.2 (continued)

6. Sakharov's view of gravitation as an elasticity of space that arises from particle physics.
 - a. The resistance of a homogeneous isotropic solid to deformation is described by two elastic constants, Young's modulus and Poisson's ratio.
 - b. The resistance of space to deformation is described by one elastic constant, the Newtonian constant of gravity. It makes its appearance in the action principle of Hilbert

$$I = \frac{1}{16\pi G} \int {}^{(4)}R(-g)^{1/2} d^4x + \int (L_{\text{matter}} + L_{\text{fields}})(-g)^{1/2} d^4x = \text{extremum.} \quad (22)$$

- c. According to the historical records, it was first learned how many elastic constants it takes to describe a solid from microscopic molecular models of matter (Newton, Laplace, Navier, Cauchy, Poisson, Voigt, Kelvin, Born), not from macroscopic considerations of symmetry and invariance. Thus, count the energy stored up in molecular bonds that are deformed from natural length or natural angle or both. Arrive at an expression for the energy of deformation per unit volume of the elastic material of the form

$$e = A(\text{Tr } \mathbf{s})^2 + B \text{Tr}(\mathbf{s}^2). \quad (23)$$

Here the strain tensor

$$s_{mn} = \frac{1}{2} \left(\frac{\partial \xi_m}{\partial x^n} + \frac{\partial \xi_n}{\partial x^m} \right) \quad (24)$$

measures the strain produced in the elastic medium by motion of the typical point that was at the location x^m to the location $x^m + \xi^m(x)$. The constants A and B are derived out of microscopic physics. They fix the values of the two elastic constants of the macroscopic theory of elasticity.

- d. Andrei Sakharov (1967) (*the Andrei Sakharov*) has proposed a similar microscopic foundation for gravitation or, as he calls it, the "metric elasticity of space." He identifies the action term of Einstein's geometrodynamics [the first term in (22)] "with the change in the action of quantum fluctuations of the vacuum [associated with the physics of particles and fields and brought about] when space is curved."
- e. Sakharov notes that present-day quantum field theory "gets rid by a renormalization process" of an energy density in the vacuum that would formally be infinite if not removed by this renormalization. Thus, in the standard analysis of the degrees of freedom of the electromagnetic field in flat space, one counts the number of modes of vibration per unit volume in the range

of circular wave numbers from k to $k + dk$ as $(2 \cdot 4\pi/8\pi^3)k^2 dk$. Each mode of oscillation, even at the absolute zero of temperature, has an absolute irreducible minimum of “zero-point energy of oscillation,” $\frac{1}{2}\hbar\nu = \frac{1}{2}\hbar ck$ [the fluctuating electric field associated with which is among the most firmly established of all physical effects. It acts on the electron in the hydrogen atom in supplement to the electric field caused by the proton alone, and thereby produces most of the famous Lamb-Rutherford shift in the energy levels of the hydrogen atom, as made especially clear by Welton (1948) and Dyson (1954)]. The totalized density of zero-point energy of the electromagnetic field per unit volume of spacetime (units: cm^4) formally diverges as

$$(\hbar/2\pi^2) \int_0^\infty k^3 dk. \quad (25)$$

Equally formally this divergence is “removed” by “renormalization” [for more on renormalization see, for example, Hepp (1969)].

- f. Similar divergences appear when one counts up formally the energy associated with other fields and with vacuum fluctuations in number of pairs of electrons, μ -mesons, and other particles in the limit of quantum energies large in comparison with the rest mass of any of these particles. Again these divergences in formal calculations are “removed by renormalization.”
- g. Removed by renormalization is a contribution not only to the energy density, and therefore to the stress-energy tensor, but also to the total Lagrange function \mathcal{L} of the variational principle for all these fields and particles,

$$I = \int \mathcal{L} d^4x = \text{extremum.} \quad (26)$$

- h. Curving spacetime alters all these energies, Sakharov points out, extending an argument of Zel'dovich (1967). Therefore the process of “renormalization” or “subtraction” no longer gives zero. Instead, the contribution of zero-point energies to the Lagrangian, expanded as a power series in powers of the curvature, with numerical coefficients A, B, \dots of the order of magnitude of unity, takes a form simplified by Ruzmaikina and Ruzmaikin (1969) to the following:

$$\begin{aligned} \mathcal{L}(R) = & A\hbar \int k^3 dk + B\hbar^{(4)}R \int k dk \\ & + \hbar[C^{(4)}R^2 + DR^{\alpha\beta}R_{\alpha\beta}] \int k^{-1} dk \\ & + (\text{higher-order terms}). \end{aligned} \quad (27)$$

[For the alteration in the number of standing waves per unit frequency in a curved manifold, see also Berger (1966), Sakharov (1967), Hill in De Witt (1967c), Polievktov-Nikoladze (1969), and Berger, Gauduchon, and Mazet (1971).]

- i. Renormalization physics argues that the first term in (27) is to be dropped. The second term, Sakharov notes, is identical in form to the Hilbert action

Box 17.2 (continued)

principle, equation (3) above, with the exception that there the constant that multiplies the Riemann scalar curvature invariant is $-c^3/16\pi G$ (in conventional units), whereas here it is $B\hbar sk dk$ (in the same conventional units). The higher order terms in (27) lead to what Sakharov calls “corrections . . . to Einstein’s equations.”

- j. Overlooking these corrections, one evidently obtains the action principle of Einstein’s theory when one insists on the equality

$$G = \left(\begin{array}{l} \text{Newtonian} \\ \text{constant of gravity} \end{array} \right) = \frac{c^3}{16\pi B\hbar sk dk}. \quad (28)$$

With B a dimensionless numerical factor of the order of unity, it follows, Sakharov argues, that the effective upper limit or “cutoff” in the formally divergent integral in (28) is to be taken to be of the order of magnitude of the reciprocal Planck length [see equation (7)],

$$k_{\text{cutoff}} \sim (c^3/\hbar G)^{1/2} = 1/L^* = 1/1.6 \times 10^{-33} \text{ cm}. \quad (29)$$

In effect Sakharov is saying (1) that field physics suffers a sea change into something new and strange for wavelengths less than the Planck length, and for quantum energies of the order of $\hbar ck_{\text{cutoff}} \sim 10^{28} \text{ eV}$ or 10^{-5} g or more; (2) that in consequence the integral $sk dk$ is cut off; and (3) that the value of this cutoff, arising purely out of the physics of fields and particles, governs the value of the Newtonian constant of gravity, G .

- k. In this sense, Sakharov’s analysis suggests that gravitation is to particle physics as elasticity is to chemical physics: merely a statistical measure of residual energies. In the one case, molecular bindings depend on departures of molecule-molecule bond lengths from standard values. In the other case, particle energies are affected by curvatures of the geometry.
- l. Elasticity, which looks simple, gets its explanation from molecular bindings, which are complicated; but molecular bindings, which are complicated, receive their explanation in terms of Schrödinger’s wave equation and Coulomb’s law of force between charged point-masses, which are even simpler than elasticity.
- m. Einstein’s geometrodynamics, which looks simple, is interpreted by Sakharov as a correction term in particle physics, which is complicated. Is particle physics, which is complicated, destined some day in its turn to unravel into something simple—something far deeper and far simpler than geometry (“pregeometry”; Chapter 44)?

§17.6. "NO PRIOR GEOMETRY": A FEATURE DISTINGUISHING EINSTEIN'S THEORY FROM OTHER THEORIES OF GRAVITY

Whereas Einstein's theory of gravity is exceedingly compelling, one can readily construct less compelling and less elegant alternative theories. The physics literature is replete with examples [see Ni (1972), and Thorne, Ni, and Will (1971) for reviews]. However, when placed among its competitors, Einstein's theory stands out sharp and clear: it agrees with experiment; most of its competitors do not (Chapters 38–40). It describes gravity entirely in terms of geometry; most of its competitors do not. It is free of any "prior geometry"; most of its competitors are not.

Set aside, until Chapter 38, the issue of agreement with experiment. Einstein's theory remains unique. Every other theory either introduces auxiliary gravitational fields [e.g., the scalar field of Brans and Dicke (1961)], or involves "prior geometry," or both. Thus, every other theory is more complicated conceptually than Einstein's theory. Every other theory contains elements of complexity for which there is no experimental motivation.

The concept of "prior geometry" requires elucidation, not least because the rejection of prior geometry played a key role in the reasoning that originally led Einstein to his geometrodynamical equation $\mathbf{G} = 8\pi\mathbf{T}$. By "prior geometry" one means any aspect of the geometry of spacetime that is fixed immutably, i.e., that cannot be changed by changing the distribution of gravitating sources. Thus, prior geometry is not generated by or affected by matter; it is not dynamic. Example: Nordström (1913) formulated a theory in which the physical metric of spacetime \mathbf{g} (the metric that enters into the equivalence principle) is generated by a "background" flat-spacetime metric $\boldsymbol{\eta}$, and by a scalar gravitational field ϕ :

$$\eta^{\alpha\beta}\phi_{,\alpha\beta} = -4\pi\phi\eta^{\alpha\beta}T_{\alpha\beta} \quad \begin{pmatrix} \text{generation of } \phi \text{ by} \\ \text{stress-energy} \end{pmatrix}, \quad (17.23a)$$

$$g_{\alpha\beta} = \phi^2\eta_{\alpha\beta} \quad \begin{pmatrix} \text{construction of } \mathbf{g} \\ \text{from } \phi \text{ and } \boldsymbol{\eta} \end{pmatrix}. \quad (17.23b)$$

In this theory, the physical metric \mathbf{g} (governor of rods and clocks and of test-particle motion) has but one changeable degree of freedom—the freedom in ϕ . The rest of \mathbf{g} is fixed by the flat spacetime metric ("prior geometry") $\boldsymbol{\eta}$. One does not remove the prior geometry by rewriting Nordström's equations (17.23) in a form

$$\begin{array}{c} R = 24\pi T, \\ \left[\begin{array}{l} \text{curvature scalar} \\ \text{constructed from } \mathbf{g} \end{array} \right] \uparrow \\ \left[\begin{array}{l} g^{\alpha\beta}T_{\alpha\beta} \end{array} \right] \uparrow \\ C^{\alpha\beta}_{\mu\nu} = 0 \\ \left[\begin{array}{l} \text{Weyl tensor} \\ \text{constructed from } \mathbf{g} \end{array} \right] \end{array} \quad (17.24)$$

devoid of reference to $\boldsymbol{\eta}$ and ϕ [Einstein and Fokker (1914); exercise 17.8]. Mass can still influence only one degree of freedom in the spacetime geometry. The other degrees of freedom are fixed *à priori*—they are *prior geometry*. And this prior geometry can perfectly well (in principle) be detected by physical experiments that make no reference to any equations (Box 17.3).

Einstein's theory compared with other theories of gravity

All other theories introduce auxiliary gravitational fields or prior geometry

"Prior geometry" defined

Nordström's theory as an illustration of prior geometry

Box 17.3 AN EXPERIMENT TO DETECT OR EXCLUDE CERTAIN TYPES OF PRIOR GEOMETRY

(Based on December 1970 discussions between Alfred Schild and Charles W. Misner)

Choose a momentarily static universe populated with a large supply of suitable pulsars. The pulsars should be absolutely regular, periodically emitting characteristic pulses of both gravitational and electromagnetic waves.

Two fleets of spaceships containing receivers are sent out "on station" to collect the experimental data. Admiral Weber's fleet carries gravitational-wave receivers; Admiral Hertz's fleet, electromagnetic receivers. The captain of each spaceship holds himself "on station" by monitoring three suitably chosen pulsars (of identical frequency) and maneuvering so that their pulses always arrive in coincidence. The experimental data he collects consist of the pulses received from all other pulsars, which he is not using for station keeping, each registered as coincident with or interlaced among the reference (stationary) pulses. [For display purposes, the pattern produced by any single pulsar can be converted to acoustic form. The reference pulses can be played acoustically (by the data-processing computer) on one drum at a fixed rate, and the pulses from other pulsars can be played on a second drum. A pattern of rhythmic beats will result.]

When the data fleet is checked out and tuned up, each captain reports stationary patterns. Now the experiment begins. One or more massive stars are towed in among the fleet. The fleet reacts to stay on station, and reports changes in the data patterns. The spaceships on the outside edges of the fleet verify that no detectable changes occur at their stations; so the incident radiation from the distant pulsars can be regarded as unaffected by the newly placed stars. Data stations nearer the movable stars report the interesting data.

What are the results?

In a universe governed by the laws of special relativity (spacetime always flat), no patterns change. (Weber's fleet was unable to get checked

out in the first place, as no gravitational waves were ever detected from the pulsars). Neither stars, nor anything else, can produce gravitational fields. All aspects of the spacetime geometry are fixed *a priori* (complete prior geometry!). There is no gravity; and no light deflection takes place to make Hertz's captains adjust their positions.

In a universe governed by Nordström's theory of gravity (see text) both fleets get checked out—i.e., both see waves. But neither fleet sees any changes in the rhythmic pattern of beats. The stars being towed about have no influence on either gravitational waves or electromagnetic waves. The prior geometry (η) present in the theory precludes any light deflection or any gravitational-wave deflection.

In a universe governed by Whitehead's (1922) theory of gravity [see Will (1971b) and references cited therein], radio waves propagate along geodesics of the "physical metric" g , and get deflected by the gravitational fields of the stars. But gravitational waves propagate along geodesics of a *flat* background metric η , and are thus unaffected by the stars. Consequently, Hertz's captains must maneuver to keep on station; and they hear a changing beat pattern between the reference pulsars and the other pulsars. But Weber's fleet remains on station and records no changes in the beat pattern. The prior geometry (η) shows itself clearly in the experimental result.

In a universe governed by Einstein's theory, both fleets see effects (no sign of prior geometry because Einstein's theory has no prior geometry). Moreover, if the fleets were originally paired, one Weber ship and one Hertz at each station, they remain paired. No differences exist between the propagation of high-frequency light waves and high-frequency gravitational waves. Both propagate along geodesics of g .

Mathematics was not sufficiently refined in 1917 to cleave apart the demands for “no prior geometry” and for a “geometric, coordinate-independent formulation of physics.” Einstein described both demands by a single phrase, “general covariance.” The “no-prior-geometry” demand actually fathered general relativity, but by doing so anonymously, disguised as “general covariance,” it also fathered half a century of confusion. [See, e.g., Kretschmann (1917).]

A systematic treatment of the distinction between prior geometry (“absolute objects”) and dynamic fields (“dynamic objects”) is a notable feature of Anderson’s (1967) relativity text.

“No prior geometry” as a part of Einstein’s principle of “general covariance”

Exercise 17.8. EINSTEIN-FOKKER REDUCES TO NORDSTRØM

EXERCISE

The vanishing of the Weyl tensor [equation (13.50)] for a spacetime metric \mathbf{g} guarantees that the metric is conformally flat—i.e., that there exists a scalar field ϕ such that $\mathbf{g} = \phi^2 \mathbf{\eta}$, where $\mathbf{\eta}$ is a flat-spacetime metric. [See, e.g., Schouten (1954) for proof.] Thus, the Einstein-Fokker equation (17.24), $C^{\alpha\beta}_{\mu\nu} = 0$, is equivalent to the Nordstrøm equation (17.23b). With this fact in hand, show that the Einstein-Fokker field equation $R = 24\pi T$ reduces to the Nordstrøm field equation (17.23a).

§17.7. A TASTE OF THE HISTORY OF EINSTEIN'S EQUATION

Nothing shows better what an idea is and means today than the battles and changes it has undergone on its way to its present form. A complete history of general relativity would demand a book. Here let a few key quotes from a few of the great papers give a little taste of what a proper history might encompass.

Einstein (1908): “We . . . will therefore in the following assume the complete physical equivalence of a gravitational field and of a corresponding acceleration of the reference system. . . . the clock at a point P for an observer anywhere in space runs $(1 + \Phi/c^2)$ times faster than the clock at the coordinate origin. . . . it follows that light rays are curved by the gravitational field. . . . an amount of energy E has a mass E/c^2 .”

Einstein and Grossmann (1913): “The theory described here originates from the conviction that the proportionality between the inertial and the gravitational mass of a body is an exact law of nature that must be expressed as a foundation principle of theoretical physics. . . . An observer enclosed in an elevator has no way to decide whether the elevator is at rest in a static gravitational field or whether the elevator is located in gravitation-free space in an accelerated motion that is maintained by forces acting on the elevator (equivalence hypothesis). . . . In the decay of radium, for example, that decrease [of mass] amounts to 1/10,000 of the total mass. If those changes in inertial mass did not correspond to changes in gravitational mass, then deviations of inertial from gravitational masses would arise that are far larger than the Eötvös experiments allow. It must therefore be considered as very probable that the identity of gravitational and inertial mass is exact.”

"The sought for generalization will surely be of the form

$$\Gamma_{\mu\nu} = \kappa T_{\mu\nu},$$

where κ is a constant and $\Gamma_{\mu\nu}$ is a contravariant tensor of the second rank that arises out of the fundamental tensor $g_{\mu\nu}$ through differential operations. . . . it proved impossible to find a differential expression for $\Gamma_{\mu\nu}$, that is a generalization of [Poisson's] $\Delta\phi$, and that is a tensor with respect to arbitrary transformations. . . . It seems most natural to demand that the system of equations should be covariant against arbitrary transformations. That stands in conflict with the result that the equations of the gravitational field do not possess this property."

Einstein and Grossman (1914): "In a 1913 treatment . . . we could not show general covariance for these gravitational equations. [Origin of their difficulty: part of the two-index curvature tensor was put on the left, to constitute the second-order part of the field equation, and part was put on the right with $T_{\mu\nu}$ and was called gravitational stress-energy. It was asked that lefthand and righthand sides transform as tensors, which they cannot do under general coordinate transformations.]

Einstein (1915a): "In recent years I had been trying to found a general theory of relativity on the assumption of the relativity even of nonuniform motions. I believed in fact that I had found the only law of gravitation that corresponds to a reasonably formulated postulate of general relativity, and I sought to establish the necessity of exactly this solution in a paper that appeared last year in these proceedings.

"A renewed analysis showed me that that necessity absolutely was not shown in the approach adopted there; that it nevertheless appeared to be shown rested on an error.

"For these reasons, I lost all confidence in the field equations I had set up, and I sought for an approach that would limit the possibilities in a natural way. In this way I was led back to the demand for the general covariance of the field equations, from which I had departed three years ago, while working with my friend Grossmann, only with a heavy heart. In fact we had already at that time come quite near to the solution of the problem that is given in what follows.

"According to what has been said, it is natural to postulate the field equations of gravitation in the form

$$R_{\mu\nu} = -\kappa T_{\mu\nu},$$

since we already know that these equations are covariant with respect to arbitrary transformations of determinant 1. In fact, these equations satisfy all conditions that we have to impose on them. [Here $R_{\mu\nu}$ is a piece of the Ricci tensor that Einstein regarded as covariant.] . . .

"Equations (22a) give in the first approximation

$$\frac{\partial^2 g^{\alpha\beta}}{\partial x^\alpha \partial x^\beta} = 0.$$

By this [condition] the coordinate system is still not determined, in the sense that for this determination four equations are necessary." (Session of Nov. 4, 1915, published Nov. 11.)

Einstein (1915b): "In a recently published investigation, I have shown how a theory of the gravitational field can be founded on Riemann's covariant theory of many-di-

dimensional manifolds. Here it will now be proved that, by introducing a surely bold additional hypothesis on the structure of matter, a still tighter logical structure of the theory can be achieved. . . . it may very well be possible that in the matter to which the given expression refers, gravitational fields play an essential part. Then T^{μ}_{μ} can appear to be positive for the entire structure, although in reality only $T^{\mu}_{\mu} + t^{\mu}_{\mu}$ is positive, and T^{μ}_{μ} vanishes everywhere. We assume in the following that in fact the condition $T^{\mu}_{\mu} = 0$ is fulfilled [quite] generally.

"Whoever does not from the beginning reject the hypothesis that molecular [small-scale] gravitational fields constitute an essential part of matter will see in the following a strong support for this point of view.

"Our hypothesis makes it possible . . . to give the field equations of gravitation in a generally covariant form . . .

$$G_{\mu\nu} = -\kappa T_{\mu\nu}$$

[where $G_{\mu\nu}$ is the Ricci tensor]." (Session of Nov. 11, 1915; published Nov. 18.)

Einstein (1915c): "I have shown that no objection of principle stands in the way of this hypothesis [the field equations], by which space and time are deprived of the last trace of objective reality. In the present work I find an important confirmation of this most radical theory of relativity: it turns out that it explains qualitatively and quantitatively the secular precession of the orbit of Mercury in the direction of the orbital motion, as discovered by Leverrier, which amounts to about 45" per century, without calling on any special hypothesis whatsoever."

Einstein (1915d; session of Nov. 25, 1915; published Dec. 2): "More recently I have found that one can proceed without hypotheses about the energy tensor of matter when one introduces the energy tensor of matter in a somewhat different way than was done in my two earlier communications. The field equations for the motion of the perihelion of Mercury are undisturbed by this modification. . . .

"Let us put

$$G_{im} = -\kappa \left(T_{im} - \frac{1}{2} g_{im} T \right)$$

[where G_{im} is the Ricci tensor]. . . .

. . . these equations, in contrast to (9), contain no new condition, so that no other assumption has to be made about the energy tensor of matter than obedience to the energy-momentum [conservation] laws.

"With this step, general relativity is finally completed as a logical structure. The postulate of relativity in its most general formulation, which makes the spacetime coordinates into physically meaningless parameters, leads compellingly to a completely determinate theory of gravitation that explains the perihelion motion of Mercury. In contrast, the general-relativity postulate is able to open up to us nothing about the nature of the other processes of nature that special relativity has not already taught. The opinion on this point that I recently expressed in these proceedings was erroneous. Every physical theory compatible with special relativity can be aligned into the system of general relativity by means of the absolute differential calculus, without [general relativity] supplying any criterion for the acceptability of that theory."

Hilbert (1915): "Axiom I [notation changed to conform to usage in this book]. The

law of physical events is determined through a world function [Mie's terminology; better known today as "Lagrangian"] L , that contains the following arguments:

$$g_{\mu\nu}, \frac{\partial g_{\mu\nu}}{\partial x^\alpha}, \frac{\partial^2 g_{\mu\nu}}{\partial x^\alpha \partial x^\beta},$$

$$A_\sigma, \frac{\partial A_\sigma}{\partial x^\tau},$$

and specifically the variation of the integral

$$\int L(-g)^{1/2} d^4x$$

must vanish for [changes in] every one of the 14 potentials $g_{\sigma\nu}$, A_σ

"Axiom II (axiom of general invariance). The world function L is invariant with respect to arbitrary transformations of the world parameters [coordinates] x^α

"For the world function L , still further axioms are needed to make its choice unambiguous. If the gravitation equations are to contain only second derivatives of the potentials $g^{\sigma\nu}$, then L must have the form

$$L = R + L_{\text{elec}},$$

where R is the invariant built from the Riemann tensor (curvature of the four-dimensional manifold.) (Session of Nov. 20, 1915.)

Einstein (1916c): "Recently H. A. Lorentz and D. Hilbert have succeeded in giving general relativity an especially transparent form in deriving its equations from a single variation principle. This will be done also in the following treatment. There it is my aim to present the basic relations as transparently as possible and in a way as general as general relativity allows."

Einstein (1916b): "From this it follows, first of all, that gravitational fields spread out with the speed of light. . . [plane] waves transport energy. . . One thus gets . . . the radiation of the system per unit time. . . .

$$\frac{G}{24\pi} \sum_{\alpha,\beta} \left(\frac{\partial^3 J_{\alpha\beta}}{\partial t^3} \right)^2,$$

Hilbert (1917): "As for the principle of causality, the physical quantities and their time-rates of change may be known at the present time in any given coordinate system; a prediction will then have a physical meaning only when it is invariant with respect to all those transformations for which exactly those coordinates used for the present time remain unchanged. I declare that predictions of this kind for the future are all uniquely determined; that is, that the causality principle holds in this formulation:

"From the knowledge of the 14 physical potentials $g_{\mu\nu}$, A_σ , in the present, all predictions about the same quantities in the future follow necessarily and uniquely insofar as they have physical meaning."

CHAPTER 18

WEAK GRAVITATIONAL FIELDS

*The way that can be walked on is not the perfect way.
The word that can be said is not the perfect word.*

LAO-TZU (~3rd century B.C.)

§18.1. THE LINEARIZED THEORY OF GRAVITY

Because of the geometric language and abbreviations used in writing them, Einstein's field equations, $G_{\mu\nu} = 8\pi T_{\mu\nu}$, hardly seem to be differential equations at all, much less ones with many familiar properties. The best way to see that they are is to apply them to weak-field situations

$$g_{\mu\nu} = \eta_{\mu\nu} + h_{\mu\nu}, \quad |h_{\mu\nu}| \ll 1, \quad (18.1)$$

e.g., to the solar system, where $|h_{\mu\nu}| \sim |\Phi| \lesssim M_\odot/R_\odot \sim 10^{-6}$; or to a weak gravitational wave propagating through interstellar space.

In a weak-field situation, one can expand the field equations in powers of $h_{\mu\nu}$, using a coordinate frame where (18.1) holds; and without much loss of accuracy, one can keep only linear terms. The resulting formalism is often called "*the linearized theory of gravity*," because it is an important theory in its own right. In fact, it is precisely this "linearized theory" that one obtains when one asks for the classical field corresponding to quantum-mechanical particles of (1) zero rest mass and (2) spin two in (3) flat spacetime [see Fierz and Pauli (1939)]. Track-2 readers have already explored linearized theory somewhat in §7.1, exercise 7.3, and Box 7.1. There it went under the alternative name, "tensor-field theory of gravity in flat spacetime."

"Linearized theory of gravity":

(1) as weak-field limit of general relativity

(2) as standard "field-theory" description of gravity in "flat spacetime"

(3) as a foundation for
“deriving” general relativity

Details of linearized theory:

(1) connection coefficients

Just as one can “descend” from general relativity to linearized theory by linearizing about flat spacetime (see below), so one can “bootstrap” one’s way back up from linearized theory to general relativity by imposing consistency between the linearized field equations and the equations of motion. or, equivalently, by asking about: (1) the stress-energy carried by the linearized gravitational field $h_{\mu\nu}$; (2) the influence of this stress-energy acting as a source for corrections $h^{(1)}_{\mu\nu}$ to the field; (3) the stress-energy carried by the corrections $h^{(1)}_{\mu\nu}$; (4) the influence of this stress-energy acting as a source for corrections $h^{(2)}_{\mu\nu}$ to the corrections $h^{(1)}_{\mu\nu}$; (5) the stress-energy carried by the corrections to the corrections; and so on. This alternative way to derive general relativity has been developed and explored by Gupta (1954, 1957, 1962), Kraichnan (1955), Thirring (1961), Feynman (1963a), Weinberg (1965), and Deser (1970). But because the outlook is far from geometric (see Box 18.1), the details of the derivation are not presented here. (But see part 5 of Box 17.2.)

Here attention focuses on deriving linearized theory from general relativity. Adopt the form (18.1) for the metric components. The resulting connection coefficients [equations (8.24b)], when linearized in the metric perturbation $h_{\mu\nu}$, read

$$\begin{aligned}\Gamma^\mu_{\alpha\beta} &= \frac{1}{2}\eta^{\mu\nu}(h_{\alpha\nu,\beta} + h_{\beta\nu,\alpha} - h_{\alpha\beta,\nu}) \\ &\equiv \frac{1}{2}(h_\alpha{}^\mu,\beta + h_\beta{}^\mu,\alpha - h_{\alpha\beta},^\mu).\end{aligned}\quad (18.2)$$

The second line here introduces the convention, used routinely whenever one expands in powers of $h_{\mu\nu}$, that indices of $h_{\mu\nu}$ are raised and lowered using $\eta^{\mu\nu}$ and $\eta_{\mu\nu}$, not $g^{\mu\nu}$ and $g_{\mu\nu}$. A similar linearization of the Ricci tensor [equation (8.47)] yields

$$\begin{aligned}R_{\mu\nu} &= \Gamma^\alpha_{\mu\nu,\alpha} - \Gamma^\alpha_{\mu\alpha,\nu} \\ &= \frac{1}{2}(h_\mu{}^\alpha,\nu\alpha + h_\nu{}^\alpha,\mu\alpha - h_{\mu\nu,\alpha}^\alpha - h_{\mu\nu}),\end{aligned}\quad (18.3)$$

where

$$h \equiv h^\alpha_\alpha = \eta^{\alpha\beta}h_{\alpha\beta}. \quad (18.4)$$

After a further contraction to form $R \equiv g^{\mu\nu}R_{\mu\nu} \approx \eta^{\mu\nu}R_{\mu\nu}$, one finds that the Einstein equations, $2G_{\mu\nu} = 16\pi T_{\mu\nu}$, read

$$\begin{aligned}h_{\mu\alpha,\nu}^\alpha + h_{\nu\alpha,\mu}^\alpha - h_{\mu\nu,\alpha}^\alpha - h_{\mu\nu} & \\ - \eta_{\mu\nu}(h_{\alpha\beta},^\alpha\beta - h_{\beta\beta},^\alpha) &= 16\pi T_{\mu\nu}.\end{aligned}\quad (18.5)$$

The number of terms has increased in passing from $R_{\mu\nu}$ (18.3) to $G_{\mu\nu} = R_{\mu\nu} - \frac{1}{2}g_{\mu\nu}R$ (18.5), but this annoyance can be counteracted by defining

(2) “gravitational potentials”
 $\bar{h}_{\mu\nu}$

$$\bar{h}_{\mu\nu} \equiv h_{\mu\nu} - \frac{1}{2}\eta_{\mu\nu}h \quad (18.6)$$

and using a bar to imply a corresponding operation on any other symmetric tensor.

Box 18.1 DERIVATIONS OF GENERAL RELATIVITY FROM GEOMETRIC VIEWPOINT AND FROM SPIN-TWO VIEWPOINT, COMPARED AND CONTRASTED

	<i>Einstein derivation</i>	<i>Spin-2 derivation</i>
Nature of primordial spacetime geometry?	Not primordial; geometry is a dynamic participant in physics	“God-given” flat Lorentz spacetime manifold
Topology (multiple connectedness) of spacetime?	Laws of physics are local; they do not specify the topology	Simply connected Euclidean topology
Vision of physics?	Dynamic geometry is the “master field” of physics	This field, that field, and the other field all execute their dynamics in a flat-spacetime manifold
Starting points for this derivation of general relativity	1. Equivalence principle (world lines of photons and test particles are geodesics of the spacetime geometry) 2. That tensorial conserved quantity which is derived from the curvature (Cartan's moment of rotation) is to be identified with the tensor of stress-momentum-energy (see Chapter 15).	1. Begin with field of spin two and zero rest mass in flat spacetime. 2. Stress-energy tensor built from this field serves as a source for this field.
Resulting equations	Einstein's field equations	Einstein's field equations
Resulting assessment of the spacetime geometry from which derivation started	Fundamental dynamic participant in physics	None. Resulting theory eradicates original flat geometry from all equations, showing it to be unobservable
View about the greatest single crisis of physics to emerge from these equations: complete gravitational collapse	Central to understanding the nature of matter and the universe	Unimportant or at most peripheral

Thus $G_{\mu\nu} = \bar{R}_{\mu\nu}$ to first order in the $h_{\mu\nu}$, and $\bar{h}_{\mu\nu} = h_{\mu\nu}$; i.e., $h_{\mu\nu} = \bar{h}_{\mu\nu} - \frac{1}{2}\eta_{\mu\nu}\bar{h}$. With this notation *the linearized field equations become*

$$-\bar{h}_{\mu\nu,\alpha}{}^\alpha - \eta_{\mu\nu}\bar{h}_{\alpha\beta}{}^{\alpha\beta} + \bar{h}_{\mu\alpha}{}^\alpha{}_\nu + \bar{h}_{\nu\alpha}{}^\alpha{}_\mu = 16\pi T_{\mu\nu}. \quad (18.7) \quad (3) \text{ linearized field equations}$$

The first term in these linearized equations is the usual flat-space d'Alembertian, and the other terms serve merely to keep the equations “gauge-invariant” (see Box

18.2). In Box 18.2 it is shown that, without loss of generality, one can impose the “gauge conditions”

(4) gauge conditions

$$\bar{h}^{\mu\alpha}_{,\alpha} = 0. \quad (18.8a)$$

These gauge conditions are the tensor analog of the Lorentz gauge $A^\alpha_{,\alpha} = 0$ of electromagnetic theory. The field equations (18.7) then become

(5) field equations and metric
in Lorentz gauge

$$-\bar{h}_{\mu\nu,\alpha}^\alpha = 16\pi T_{\mu\nu}. \quad (18.8b)$$

The gauge conditions (18.8a), the field equations (18.8b), and the definition of the metric

$$g_{\mu\nu} = \eta_{\mu\nu} + h_{\mu\nu} = \eta_{\mu\nu} + \bar{h}_{\mu\nu} - \frac{1}{2}\eta_{\mu\nu}\bar{h} \quad (18.8c)$$

are the fundamental equations of the linearized theory of gravity written in Lorentz gauge.

EXERCISES

Exercise 18.1. GAUGE INVARIANCE OF THE RIEMANN CURVATURE

Show that in linearized theory the components of the Riemann tensor are

$$R_{\alpha\mu\beta\nu} = \frac{1}{2}(h_{\alpha\nu,\mu\beta} + h_{\mu\beta,\nu\alpha} - h_{\mu\nu,\alpha\beta} - h_{\alpha\beta,\mu\nu}). \quad (18.9)$$

Then show that these components are left unchanged by a gauge transformation of the form discussed in Box 18.2 [equation (4b)]. Since the Einstein tensor is a contraction of the Riemann tensor, this shows that it is also gauge-invariant.

Exercise 18.2. JUSTIFICATION OF LORENTZ GAUGE

Let a particular solution to the field equations (18.7) of linearized theory be given, in an arbitrary gauge. Show that there necessarily exist four generating functions $\xi_\mu(t, x^j)$ whose gauge transformation [Box 18.2, eq. (4b)] makes

$$\bar{h}^{\text{new}}{}^{\mu\alpha}_{,\alpha} = 0 \quad (\text{Lorentz gauge}).$$

Also show that a subsequent gauge transformation leaves this Lorentz gauge condition unaffected if and only if its generating functions satisfy the sourceless wave equation

$$\xi^{\alpha,\beta}_{,\beta} = 0.$$

Exercise 18.3. EXTERNAL FIELD OF A STATIC, SPHERICAL BODY

Consider the external gravitational field of a static spherical body, as described in the body's (nearly) Lorentz frame—i.e., in a nearly rectangular coordinate system $|h_{\mu\nu}| \ll 1$, in which the body is located at $x = y = z = 0$ for all t . By fiat, adopt Lorentz gauge.

(a) Show that the field equations (18.8b) and gauge conditions (18.8a) imply

$$\begin{aligned} \bar{h}_{00} &= 4M/(x^2 + y^2 + z^2)^{1/2}, & \bar{h}_{0j} &= \bar{h}_{jk} = 0, \\ h_{00} &= h_{xx} = h_{yy} = h_{zz} = 2M/(x^2 + y^2 + z^2)^{1/2}, & h_{\alpha\beta} &= 0 \text{ if } \alpha \neq \beta, \end{aligned}$$

where M is a constant (the mass of the body; see §19.3).

Box 18.2 GAUGE TRANSFORMATIONS AND COORDINATE TRANSFORMATIONS IN LINEARIZED THEORY

A. The Basic Equations of Linearized Theory, written in any coordinate system that is nearly globally Lorentz, are (18.1) and (18.7):

$$g_{\mu\nu} = \eta_{\mu\nu} + h_{\mu\nu}, \quad |h_{\mu\nu}| \ll 1; \quad (1)$$

$$-\bar{h}_{\mu\nu,\alpha}{}^\alpha - \eta_{\mu\nu} \bar{h}_{\alpha\beta}{}^{\alpha\beta} + \bar{h}_{\mu\alpha}{}^{\alpha}{}_\nu + \bar{h}_{\nu\alpha}{}^{\alpha}{}_\mu = 16\pi T_{\mu\nu}. \quad (2)$$

Two different types of coordinate transformations connect nearly globally Lorentz systems to each other: global Lorentz transformations, and infinitesimal coordinate transformations.

1. *Global Lorentz Transformations*:

$$x^\mu = A^\mu{}_{\alpha'} x^{\alpha'}, \quad A^\mu{}_{\alpha'} A^\nu{}_{\beta'} \eta_{\mu\nu} = \eta_{\alpha'\beta'}. \quad (3a)$$

These transform the metric coefficients via

$$\begin{aligned} \eta_{\alpha'\beta'} + h_{\alpha'\beta'} &= g_{\alpha'\beta'} = \frac{\partial x^\mu}{\partial x^{\alpha'}} \frac{\partial x^\nu}{\partial x^{\beta'}} g_{\mu\nu} = A^\mu{}_{\alpha'} A^\nu{}_{\beta'} (\eta_{\mu\nu} + h_{\mu\nu}) \\ &= \eta_{\alpha'\beta'} + A^\mu{}_{\alpha'} A^\nu{}_{\beta'} h_{\mu\nu}. \end{aligned}$$

Thus, $h_{\mu\nu}$ —and likewise $\bar{h}_{\mu\nu}$ —transform like components of a tensor in flat spacetime

$$h_{\alpha'\beta'} = A^\mu{}_{\alpha'} A^\nu{}_{\beta'} h_{\mu\nu}. \quad (3b)$$

2. *Infinitesimal Coordinate Transformations* (creation of “ripples” in the coordinate system):

$$x^{\mu'}(\mathcal{P}) = x^\mu(\mathcal{P}) + \xi^\mu(\mathcal{P}), \quad (4a)$$

where $\xi^\mu(\mathcal{P})$ are four arbitrary functions small enough to leave $|h_{\mu'\nu'}| \ll 1$. Infinitesimal transformations of this sort make tiny changes in the functional forms of all scalar, vector, and tensor fields. *Example:* the temperature T is a unique function of position, $T(\mathcal{P})$; so when written as a function of coordinates it changes

$$\begin{aligned} T(x^{\mu'} = a^\mu) &= T(x^\mu + \xi^\mu = a^\mu) = T(x^\mu = a^\mu - \xi^\mu) \\ &= T(x^\mu = a^\mu) - T_{,\mu} \xi^\mu; \end{aligned}$$

i.e., if $\xi^0 = 0.001 \sin(x^1)$, and if $T = \cos^2(x^0)$, then

$$T = \cos^2(x^{0'}) + 0.002 \sin(x^{1'}) \cos(x^{0'}) \sin(x^{0'}).$$

Box 18.2 (continued)

These tiny changes can be ignored in all quantities except the metric, where tiny deviations from $\eta_{\mu\nu}$ contain all the information about gravity. The usual tensor transformation law for the metric

$$g_{\rho'\sigma'}[x^{\alpha'}(\mathcal{P})] = g_{\mu\nu}[x^\alpha(\mathcal{P})] \frac{\partial x^\mu}{\partial x^{\rho'}} \frac{\partial x^\nu}{\partial x^{\sigma'}},$$

when combined with the transformation law (4a) and with

$$g_{\mu\nu}[x^\alpha(\mathcal{P})] = \eta_{\mu\nu} + h_{\mu\nu}[x^\alpha(\mathcal{P})],$$

reveals that

$$\begin{aligned} g_{\rho'\sigma'}(x^{\alpha'} = a^\alpha) &= \eta_{\rho\sigma}(x^\alpha = a^\alpha) - \xi_{\rho,\sigma} - \xi_{\sigma,\rho} \\ &\quad + \text{negligible corrections } \sim h_{\rho\sigma,\alpha}\xi^\alpha \text{ and } \sim h_{\rho\alpha}\xi^\alpha_{,\sigma}. \end{aligned}$$

Hence, *the metric perturbation functions in the new ($x^{\mu'}$) and old (x^μ) coordinate systems are related by*

$$h_{\mu\nu}^{\text{new}} = h_{\mu\nu}^{\text{old}} - \xi_{\mu,\nu} - \xi_{\nu,\mu}, \quad (4b)$$

whereas the functional forms of all other scalars, vectors, and tensors are unaltered, to within the precision of linearized theory.

B. Gauge Transformations and Gauge Invariance. In linearized theory one usually regards equation (4b) as gauge transformations, analogous to those

$$A_\mu^{\text{new}} = A_\mu^{\text{old}} + \Psi_{,\mu}, \quad (5a)$$

of electromagnetic theory. The fact that gravitational gauge transformations do not affect the functional forms of scalars, vectors, or tensors (i.e., observables) is called “gauge invariance.” Just as a straightforward calculation reveals the gauge invariance of the electromagnetic field,

$$F_{\mu\nu}^{\text{new}} = A_{\nu,\mu}^{\text{new}} - A_{\mu,\nu}^{\text{new}} = A_{\nu,\mu}^{\text{old}} + \Psi_{,\nu\mu} - A_{\mu,\nu}^{\text{old}} - \Psi_{,\mu\nu} = F_{\mu\nu}^{\text{old}}, \quad (5b)$$

so a straightforward calculation (exercise 18.1) reveals the gauge invariance of the Riemann tensor

$$R_{\mu\nu\alpha\beta}^{\text{new}} = R_{\mu\nu\alpha\beta}^{\text{old}}. \quad (6)$$

Such gauge invariance was already guaranteed by the fact that $R_{\mu\nu\alpha\beta}$ are the components of a tensor, and are thus essentially the same whether calculated in an orthonormal frame $g_{\hat{\mu}\hat{\nu}} = \eta_{\mu\nu}$, in the old coordinates where $g_{\mu\nu} = \eta_{\mu\nu} + h_{\mu\nu}^{\text{old}}$, or in the new coordinates where $g_{\mu\nu} = \eta_{\mu\nu} + h_{\mu\nu}^{\text{new}}$.

Like the Riemann tensor, the Einstein tensor and the stress-energy tensor are unaffected by gauge transformations. Hence, if one knows a specific solution $\bar{h}_{\mu\nu}$ to the linearized field equations (2) for a given $T^{\mu\nu}$, one can obtain another solution that describes precisely the same physical situation (all observables unchanged) by the change of gauge (4), in which ξ_μ are four arbitrary but small functions.

C. Lorentz Gauge. One can show (exercise 18.2) that for any physical situation, one can specialize the gauge (i.e., the coordinates) so that $\bar{h}^{\mu\alpha}_{,\alpha} = 0$. This is the Lorentz gauge introduced in §18.1. The Lorentz gauge is not fixed uniquely. The gauge condition $\bar{h}^{\mu\alpha}_{,\alpha} = 0$ is left unaffected by any gauge transformation for which

$$\xi^{\alpha,\beta}_{\beta} = 0.$$

(See exercise 18.2.)

D. Curvilinear Coordinate Systems. Once the gauge has been fixed by fiat for a given system (e.g., the solar system), one can regard $h_{\mu\nu}$ and $\bar{h}_{\mu\nu}$ as components of tensors in flat spacetime; and one can regard the field equations (2) and the chosen gauge conditions as geometric, coordinate-independent equations in flat spacetime. This viewpoint allows one to use curvilinear coordinates (e.g., spherical coordinates centered on the sun), if one wishes. But in doing so, one must everywhere replace the Lorentz components of the metric, $\eta_{\mu\nu}$, by the metric's components $g_{\mu\nu}$ flat in the flat-spacetime curvilinear coordinate system; and one must replace all ordinary derivatives ("commas") in the field equations and gauge conditions by covariant derivatives whose connection coefficients come from $g_{\mu\nu}$ flat. See exercise 18.3 for an example.

(b) Adopt spherical polar coordinates,

$$x = r \sin \theta \cos \phi, \quad y = r \sin \theta \sin \phi, \quad z = r \cos \theta.$$

By regarding $h_{\mu\nu}$ and $\bar{h}_{\mu\nu}$ as components of tensors in flat spacetime (see end of Box 18.2), and by using the usual tensor transformation laws, put the solution found in (a) into the form

$$\begin{aligned} \bar{h}_{00} &= 4M/r, & \bar{h}_{0j} &= \bar{h}_{jk} = 0, \\ h_{00} &= \frac{2M}{r}, & h_{0j} &= 0, & h_{jk} &= \frac{2M}{r} g_{jk} \text{ flat} \end{aligned}$$

where $g_{\alpha\beta}$ flat are the components of the flat-spacetime metric in the spherical coordinate system

$$\begin{aligned} g_{00} \text{ flat} &= -1, & g_{rr} \text{ flat} &= 1, & g_{\theta\theta} \text{ flat} &= r^2, \\ g_{\phi\phi} \text{ flat} &= r^2 \sin^2 \theta, & g_{\alpha\beta} \text{ flat} &= 0 \text{ when } \alpha \neq \beta. \end{aligned}$$

Thereby conclude that the general relativistic line element, accurate to linearized order, is

$$ds^2 = -(1 - 2M/r) dt^2 + (1 + 2M/r)(dr^2 + r^2 d\theta^2 + r^2 \sin^2\theta d\phi^2).$$

- (c) Derive this general, static, spherically symmetric, Lorentz-gauge, vacuum solution to the linearized field equations from scratch, working entirely in spherical coordinates. [Hint: As discussed at the end of Box 18.2, $\eta_{\mu\nu}$ in equation (18.8c) must be replaced by $g_{\mu\nu}^{\text{flat}}$; and in the field equations and gauge conditions (18.8a, b), all commas (partial derivatives) must be replaced by covariant derivatives, whose connection coefficients come from $g_{\mu\nu}^{\text{flat}}$]
 (d) Calculate the Riemann curvature tensor for this gravitational field. The answer should agree with equation (1.14).
-

§18.2. GRAVITATIONAL WAVES

Linearized theory and electromagnetic theory compared

The gauge conditions and field equations (18.8a, b) of linearized theory bear a close resemblance to the equations of electromagnetic theory in Lorentz gauge and flat spacetime,

$$A^\alpha_{,\alpha} = 0, \quad (18.10a)$$

$$-A^\mu_{,\alpha}{}^\alpha = 4\pi J^\mu. \quad (18.10b)$$

They differ only in the added index ($h^{\mu\nu}$ versus A^μ , $T^{\mu\nu}$ versus J^μ). Consequently, from past experience with electromagnetic theory, one can infer much about linearized gravitation theory.

For example, the field equations (18.8b) must have gravitational-wave solutions. The analog of the electromagnetic plane wave

$$A^x = A^x(t - z), \quad A^y = A^y(t - z), \quad A^z = 0, \quad A^0 = 0,$$

Plane gravitational waves

will be the gravitational plane wave

$$\begin{aligned} \bar{h}^{xx} &= \bar{h}^{xx}(t - z), & \bar{h}^{xy} &= \bar{h}^{xy}(t - z), & \bar{h}^{yy} &= \bar{h}^{yy}(t - z), \\ \bar{h}^{\mu 0} &= \bar{h}^{\mu z} = 0 \text{ for all } \mu. \end{aligned} \quad (18.11)$$

Although a detailed study of such waves will be delayed until Chapters 35–37, some properties of these waves are explored in the exercises at the end of the next section.

§18.3. EFFECT OF GRAVITY ON MATTER

How to analyze effects of weak gravity on matter

The effects of weak gravitational fields on matter can be computed by using the linearized metric (18.1) and Christoffel symbols (18.2) in the appropriate equations of motion—i.e., in the geodesic equation (for the motion of particles or light rays), in the hydrodynamic equations (for fluid matter), in Maxwell's equations (for electromagnetic waves), or in the equation $\nabla \cdot \mathbf{T} = 0$ for the total stress-energy tensor

of whatever fields and matter may be present. Exercises 18.5, 18.6 and 18.7 provide examples, as do the Newtonian-limit calculations in exercises 16.1 and 16.4, and in §17.4. If, however, the lowest-order (linearized) gravitational “forces” (Christoffel-symbol terms) have a significant influence on the motion of the sources of the gravitational field, one finds that the linearized field equation (18.7) is inadequate, and better approximations to Einstein’s equations must be considered. [Thus emission of gravitational waves by a mechanically or electrically driven oscillator falls within the scope of linearized theory, but emission by a double-star system, or by stellar oscillations that gravitational forces maintain, will require discussion of nonlinear terms (gravitational “stress-energy”) in the Einstein equations; see §§36.9 to 36.11.]

The above conclusions follow from a consideration of conservation laws associated with the linearized field equation. Just as the electromagnetic equations (18.10a, b) guarantee charge conservation

$$J^{\mu}_{,\mu} = 0, \quad \int_{\text{all space}} J^0(t, x) dx dy dz \equiv Q = \text{const},$$

so the gravitational equations (18.8a, b) guarantee conservation of the total 4-momentum and angular momentum of any body bounded by vacuum:

Conservation of 4-momentum
and angular momentum in
linearized theory

$$T^{\mu\nu}_{,\nu} = 0, \quad (18.12a)$$

$$\int_{\text{body}} T^{\mu 0}(t, x) dx dy dz \equiv P^{\mu} = \text{const}; \quad (18.12b)$$

$$(x^{\alpha} T^{\beta\mu} - x^{\beta} T^{\alpha\mu})_{,\mu} = 0, \quad (18.13a)$$

$$\int_{\text{body}} (x^{\alpha} T^{\beta 0} - x^{\beta} T^{\alpha 0}) dx dy dz \equiv J^{\alpha\beta} = \text{const}. \quad (18.13b)$$

(See §5.11 for the basic properties of angular momentum in special relativity. The angular momentum here is calculated relative to the origin of the coordinate system.) Now it is important that the stress-energy components $T^{\mu\nu}$, which appear in the linearized field equations (18.7) and in these conservation laws, are precisely the components one would calculate using special relativity (with $g_{\mu\nu} = \eta_{\mu\nu}$). As a result, the energy-momentum conservation formulated here contains no contributions or effects of gravity! From this one sees that linearized theory assumes that gravitational forces do no significant work. For example, energy losses due to gravitational radiation-damping forces are neglected by linearized theory. Similarly, conservation of 4-momentum P^{μ} for each of the bodies acting as sources of $h_{\mu\nu}$ means that each body moves along a geodesic of $\eta_{\mu\nu}$ (straight lines in the nearly Lorentz coordinate system) rather than along a geodesic of $g_{\mu\nu} = \eta_{\mu\nu} + h_{\mu\nu}$. Thus, linearized theory can be used to calculate the motion of test particles and fields, using $g_{\mu\nu} = \eta_{\mu\nu} + h_{\mu\nu}$; but to include gravitational corrections to the motion of the sources themselves—to allow them to satisfy $T^{\mu\nu}_{,\nu} = 0$ rather than $T^{\mu\nu}_{,\nu} = 0$ —one must reinsert into the field equations the nonlinear terms that linearized theory discards. (See, e.g., Chapter 20 on conservation laws; §§36.9–36.11 on the generation of gravitational waves and radiation reaction; and Chapter 39 on the post-Newtonian approximation.)

Limit on validity of linearized
theory: gravity must not
affect motions of sources
significantly

The energy, momentum, and angular momentum radiated by gravitational waves in linearized theory can be calculated by special-relativistic methods analogous to those used in electromagnetic theory for electromagnetic waves [Fierz and Pauli (1939)], but it will be more informative and powerful to use a fully gravitational approach (Chapters 35 and 36).

EXERCISES

Exercise 18.4. SPACETIME CURVATURE FOR A PLANE GRAVITATIONAL WAVE

Calculate the components of the Riemann curvature tensor [equations (18.9)] for the gravitational plane wave (18.11). [Answer:

$$\begin{aligned} R_{x0x0} &= -R_{y0y0} = -R_{x0xz} = +R_{y0yz} = +R_{xzxz} = -R_{yzyz} = -\frac{1}{4}(\bar{h}_{xx} - \bar{h}_{yy})_{tt}; \\ R_{x0y0} &= -R_{x0yz} = +R_{xzyz} = -R_{xzy0} = -\frac{1}{2}\bar{h}_{xy,tt}; \end{aligned}$$

all other components vanish except those obtainable from the above by the symmetries $R_{\alpha\beta\gamma\delta} = R_{[\alpha\beta][\gamma\delta]} = R_{\gamma\delta\alpha\beta}$.

Exercise 18.5. A PRIMITIVE GRAVITATIONAL-WAVE DETECTOR (see Figure 18.1)

Two beads slide almost freely on a smooth stick; only slight friction impedes their sliding. The stick falls freely through spacetime, with its center moving along a geodesic and its ends attached to gyroscopes, so they do not rotate. The beads are positioned equidistant (distance $\frac{1}{2}\ell$) from the stick's center. Plane gravitational waves [equation (18.11) and exercise 18.4], impinging on the stick, push the beads back and forth ("geodesic deviation"; "tidal gravitational forces"). The resultant friction of beads on stick heats the stick; and the passage of the waves is detected by measuring the rise in stick temperature.* (Of course, this is not the best of all conceivable designs!) Neglecting the effect of friction on the beads' motion, calculate the proper distance separating them as a function of time. [Hints: Let ξ be the separation between the beads; and let $\mathbf{n} = \xi/|\xi|$ be a unit vector that points along the stick in the stick's own rest frame. Then their separation has magnitude $\ell = \xi \cdot \mathbf{n}$. The fact that the stick is nonrotating is embodied in a parallel-transport law for \mathbf{n} , $\nabla_u \mathbf{n} = 0$. ("Fermi-Walker transport" of §§6.5, 6.6, and 13.6 reduces to parallel transport, because the stick moves along a geodesic with $\mathbf{a} = \nabla_u \mathbf{u} = 0$.) Thus,

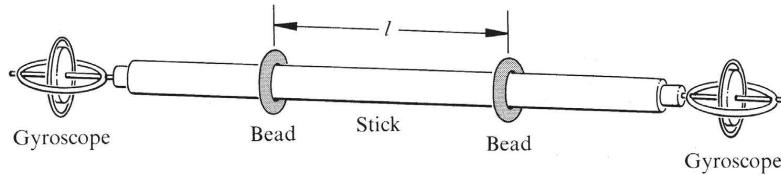
$$\begin{aligned} d\ell/d\tau &= \nabla_u(\xi \cdot \mathbf{n}) = (\nabla_u \xi) \cdot \mathbf{n}, \\ d^2\ell/d\tau^2 &= \nabla_u \nabla_u(\xi \cdot \mathbf{n}) = (\nabla_u \nabla_u \xi) \cdot \mathbf{n}, \end{aligned}$$

where τ is the stick's proper time. But $\nabla_u \nabla_u \xi$ is produced by the Riemann curvature of the wave (geodesic deviation):

$$\nabla_u \nabla_u \xi = \text{projection along } \mathbf{n} \text{ of } [-\mathbf{Riemann}(\dots, \mathbf{u}, \xi, \mathbf{u})].$$

(The geodesic-deviation forces perpendicular to the stick, i.e., perpendicular to \mathbf{n} , are coun-

*This thought experiment was devised by Bondi [1957, 1965; Bondi and McCrea (1960)] as a means for convincing skeptics of the reality of gravitational waves.

**Figure 18.1.**

A primitive detector for gravitational waves, consisting of a beaded stick with gyroscopes on its ends [Bondi (1957)]. See exercise 18.5 for discussion.

terbalanced by the stick's pushing back on the beads to stop them from passing through it—no penetration of matter by matter! Thus,

$$d^2\ell/d\tau^2 = -\mathbf{Riemann}(\dots, \mathbf{u}, \boldsymbol{\xi}, \mathbf{u}) \cdot \mathbf{n} = -\mathbf{Riemann}(\mathbf{n}, \mathbf{u}, \boldsymbol{\xi}, \mathbf{u}).$$

Evaluate this acceleration in the stick's local Lorentz frame. Orient the coordinates so the waves propagate in the z -direction and the stick's direction has components $n^z = \cos \theta$, $n^x = \sin \theta \cos \phi$, $n^y = \sin \theta \sin \phi$. Solve the resulting differential equation for $\ell(\tau)$. [Answer:

$$\ell = \ell_0 \left[1 + \frac{1}{4} (\bar{h}_{xx} - \bar{h}_{yy}) \sin^2 \theta \cos 2\phi + \frac{1}{2} \bar{h}_{xy} \sin^2 \theta \sin 2\phi \right],$$

where \bar{h}_{jk} are evaluated on the stick's world line ($x = y = z = 0$). Notice that, if the stick is oriented along the direction of wave propagation (if $\theta = 0$), the beads do not move. In this sense, the effect of the waves (geodesic deviation) is purely transverse. For further discussion, see §§35.4 to 35.6.]

§18.4. NEARLY NEWTONIAN GRAVITATIONAL FIELDS

The general solution to the linearized field equations in Lorentz gauge [equations (18.8a, b)] lends itself to expression as a retarded integral of the form familiar from electromagnetic theory:

$$\bar{h}_{\mu\nu}(t, \mathbf{x}) = \int \frac{4T_{\mu\nu}(t - |\mathbf{x} - \mathbf{x}'|, \mathbf{x}')}{|\mathbf{x} - \mathbf{x}'|} d^3x'. \quad (18.14)$$

The gravitational-wave aspects of this solution will be studied in Chapter 36. Here focus attention on a nearly Newtonian source: $T_{00} \gg |T_{0j}|$, $T_{00} \gg |T_{jk}|$, and velocities slow enough that retardation is negligible. In this case, (18.14) reduces to

$$\bar{h}_{00} = -4\Phi, \quad \bar{h}_{0j} = \bar{h}_{jk} = 0, \quad (18.15a)$$

$$\Phi(t, \mathbf{x}) = -\int \frac{T_{00}(t, \mathbf{x}')}{|\mathbf{x} - \mathbf{x}'|} d^3x' = \text{Newtonian potential}. \quad (18.15b)$$

The corresponding metric (18.8c) is

$$ds^2 = -(1 + 2\Phi) dt^2 + (1 - 2\Phi)(dx^2 + dy^2 + dz^2) \quad (18.15c)$$

$\approx -(1 - 2M/r) dt^2 + (1 + 2M/r)(dx^2 + dy^2 + dz^2)$ far from source.

Retarded-integral solution of linearized field equation

Newtonian gravity as a limit of linearized theory

Bending of light and gravitational redshift predicted by linearized theory

The errors in this metric are: (1) missing corrections of order Φ^2 due to nonlinearities of which linearized theory is oblivious; (2) missing corrections due to setting $\bar{h}_{0j} = 0$ (these are of order $\bar{h}_{0j} \sim \Phi v$, where $v \sim |T_{0j}|/T_{00}$ is a typical velocity in the source); (3) missing corrections due to setting $\bar{h}_{jk} = 0$ [these are of order $\bar{h}_{jk} \sim \Phi(|T_{jk}|/T_{00})$]. In the solar system all these errors are $\sim 10^{-12}$, whereas $\Phi \sim 10^{-6}$.

Passive correspondence with Newtonian theory demanded only that $g_{00} = -(1 + 2\Phi)$; see equation (17.19). However, linearized theory determines all the metric coefficients, up to errors of $\sim \Phi v$, $\sim \Phi^2$, and $\sim \Phi(|T_{jk}|/T_{00})$. This is sufficient accuracy to predict correctly (fractional errors $\sim 10^{-6}$) the bending of light and the gravitational redshift in the solar system, but not perihelion shifts.

EXERCISES

Exercise 18.6. BENDING OF LIGHT BY THE SUN

To high precision, the sun is static and spherical, so its external line element is (18.15c) with $\Phi = -M/r$; i.e.,

$$ds^2 = -(1 - 2M/r) dt^2 + (1 + 2M/r)(dx^2 + dy^2 + dz^2) \text{ everywhere outside sun.} \quad (18.16)$$

A photon moving in the equatorial plane ($z = 0$) of this curved spacetime gets deflected very slightly from the world line

$$x = t, \quad y = b \equiv \text{"impact parameter,"} \quad z = 0. \quad (18.17)$$

Calculate the amount of deflection as follows.

(a) Write down the geodesic equation (16.4a) for the photon's world line,

$$\frac{dp^\alpha}{d\lambda^*} + \Gamma^\alpha_{\beta\gamma} p^\beta p^\gamma = 0. \quad (18.18)$$

[Here $\mathbf{p} = d/d\lambda^* = (4\text{-momentum of photon}) = (\text{tangent vector to photon's null geodesic}).]$

(b) By evaluating the connection coefficients in the equatorial plane, and by using the approximate values, $|p^y| \ll p^0 \approx p^x$, of the 4-momentum components corresponding to the approximate world line (18.17), show that

$$\frac{dp^y}{d\lambda^*} = \frac{-2Mb}{(x^2 + b^2)^{3/2}} p^x \frac{dx}{d\lambda^*}, \quad p^x = p^0 \left[1 + O\left(\frac{M}{b}\right) \right] = \text{const} \left[1 + O\left(\frac{M}{b}\right) \right].$$

(c) Integrate this equation for p^y , assuming $p^y = 0$ at $x = -\infty$ (photon moving precisely in x -direction initially); thereby obtain

$$p^y(x = +\infty) = -\frac{4M}{b} p^x.$$

(d) Show that this corresponds to deflection of light through the angle

$$\Delta\phi = 4M/b = 1''.75 (R_\odot/b), \quad (18.19)$$

where R_\odot is the radius of the sun. For a comparison of this prediction with experiment, see Box 40.1.

Exercise 18.7. GRAVITATIONAL REDSHIFT

(a) Use the geodesic equation for a photon, written in the form

$$dp_\mu/d\lambda^* - \Gamma^\alpha_{\mu\beta} p_\alpha p^\beta = 0,$$

to prove that any photon moving freely in the sun's gravitational field [line element (18.16)] has $dp_0/d\lambda^* = 0$; i.e.,

$$p_0 = \text{constant along photon's world line.} \quad (18.20)$$

(b) An atom at rest on the sun's surface emits a photon of wavelength λ_e , as seen in its orthonormal frame. [Note:

$$h\nu_e = h/\lambda_e = (\text{energy atom measures}) = -\mathbf{p} \cdot \mathbf{u}_e, \quad (18.21)$$

where \mathbf{p} is the photon's 4-momentum and \mathbf{u}_e is the emitter's 4-velocity.] An atom at rest far from the sun receives the photon, and measures its wavelength to be λ_r . [Note: $h/\lambda_r = -\mathbf{p} \cdot \mathbf{u}_r$.] Show that the photon is redshifted by the amount

$$z \equiv \frac{\lambda_r - \lambda_e}{\lambda_e} = \frac{M_\odot}{R_\odot} = 2 \times 10^{-6}. \quad (18.22)$$

[Hint: $\mathbf{u}_r = \partial/\partial t$; $\mathbf{u}_e = (1 - 2M/r)^{-1/2} \partial/\partial t$. Why?] For further discussion of the gravitational redshift and experimental results, see §§7.4 and 38.5; also Figures 38.1 and 38.2.

CHAPTER 19

MASS AND ANGULAR MOMENTUM OF A GRAVITATING SYSTEM

§19.1. EXTERNAL FIELD OF A WEAKLY GRAVITATING SOURCE

Metric far from a weakly gravitating system, as a power series in $1/r$:

Consider an isolated system with gravity so weak that in calculating its structure and motion one can completely ignore self-gravitational effects. (This is true of an asteroid, and of a nebula with high-energy electrons and protons spiraling in a magnetic field; it is not true of the Earth or the sun.) Assume nothing else about the system—for example, by contrast with Newtonian theory, allow velocities to be arbitrarily close to the speed of light, and allow stresses T^{jk} and momentum densities T^{0j} to be comparable to the mass-energy density T^{00} .

(1) derivation

Calculate the weak gravitational field,

$$g_{\mu\nu} = \eta_{\mu\nu} + h_{\mu\nu}, \quad (19.1)$$

$$\bar{h}_{\mu\nu} \equiv h_{\mu\nu} = \int \frac{4\bar{T}_{\mu\nu}(t - |\mathbf{x} - \mathbf{x}'|, \mathbf{x}')}{|\mathbf{x} - \mathbf{x}'|} d^3x', \quad (19.2)$$

produced by such a system [see “barred” version of equation (18.14)]. Restrict attention to the spacetime region far outside the system, and expand $h_{\mu\nu}$ in powers of $\mathbf{x}'/r \equiv \mathbf{x}'/|\mathbf{x}|$, using the relations

$$\bar{T}_{\mu\nu}(t - |\mathbf{x} - \mathbf{x}'|, \mathbf{x}') = \sum_{n=0}^{\infty} \frac{1}{n!} \left[\frac{\partial^n}{\partial t^n} \bar{T}_{\mu\nu}(t - r, \mathbf{x}') \right] (r - |\mathbf{x} - \mathbf{x}'|)^n, \quad (19.3a)$$

$$r - |\mathbf{x} - \mathbf{x}'| = x^j \left(\frac{x'^j}{r} \right) + \frac{1}{2} \frac{x^j x^k}{r} \left(\frac{x'^j x'^k - r'^2 \delta_{jk}}{r^2} \right) + \dots, \quad (19.3b)$$

$$\frac{1}{|\mathbf{x} - \mathbf{x}'|} = \frac{1}{r} + \frac{x^j}{r^2} \frac{x'^j}{r} + \frac{1}{2} \frac{x^j x^k}{r^3} \frac{(3x'^j x'^k - r'^2 \delta_{jk})}{r^2} + \dots \quad (19.3c)$$

Perform the calculation in the system's rest frame, where

$$P^j \equiv \int T^{0j} d^3x = 0, \quad (19.4a)$$

with origin of coordinates at the system's center of mass

$$\int x^j T^{00} d^3x = 0. \quad (19.4b)$$

The result, after a change of gauge to simplify h_{00} and h_{0j} , is

$$ds^2 = -\left[1 - \frac{2M}{r} + 0\left(\frac{1}{r^3}\right)\right]dt^2 - \left[4\epsilon_{jk\ell}S^k \frac{x^\ell}{r^3} + 0\left(\frac{1}{r^3}\right)\right]dt dx^j + \left[\left(1 + \frac{2M}{r}\right)\delta_{jk} + \left(\begin{array}{l} \text{gravitational radiation terms} \\ \text{that die out as } O(1/r) \end{array}\right)\right]dx^j dx^k. \quad (19.5)$$

(see exercise 19.1 for derivation.) Here M and S^k are the body's mass and intrinsic angular momentum.

$$M = \int T^{00} d^3x, \quad (19.6a)$$

$$S_k = \int \epsilon_{k\ell m} x^\ell T^{m0} d^3x. \quad (19.6b)$$

The corresponding Newtonian potential is

$$\Phi = -\frac{1}{2}(g_{00} - \eta_{00}) = -\frac{M}{r} + 0\left(\frac{1}{r^3}\right). \quad (19.6c)$$

Conclusion: With an appropriate choice of gauge, Φ and g_{00} far from any weak source are time-independent and are determined uniquely by the source's mass M ; g_{0j} is time-independent and is fixed by the source's intrinsic angular momentum S^j ; but g_{jk} has time-dependent terms (gravitational waves!) of $O(1/r)$.

How metric depends on
system's mass M and
angular momentum S

The rest of this chapter focuses on the “imprints” of the mass and angular momentum in the gravitational field; the gravitational waves will be ignored, or almost so, until Chapter 35.

Exercise 19.1. DERIVATION OF METRIC FAR OUTSIDE A WEAKLY GRAVITATING BODY

EXERCISE

- (a) Derive equation (19.5). [*Hints:* (1) Follow the procedure outlined in the text. (2) When calculating h_{00} , write out explicitly the $n = 0$ and $n = 1$ terms of (19.2), to precision $O(1/r^2)$, and simplify the $n = 0$ term using the identities

$$T^{jk} = \frac{1}{2}(T^{00}x^j x^k)_{,00} + (T^{ij}x^k + T^{ik}x^j)_{,i} - \frac{1}{2}(T^{lm}x^j x^k)_{,lm}, \quad (19.7a)$$

$$T^{ll}x^m = \left(T^{0l}x^l x^m - \frac{1}{2}T^{0m}r^2\right)_{,0} + \left(T^{lk}x^k x^m - \frac{1}{2}T^{lm}r^2\right)_{,l}. \quad (19.7b)$$

(Verify that these identities follow from $T^{\alpha\beta}_{,\beta} = 0$.) (3) When calculating h_{0m} , write out explicitly the $n = 0$ term of (19.2), to precision $O(1/r^2)$, and simplify it using the identity

$$T^{0k}x^j + T^{0j}x^k = (T^{00}x^jx^k)_{,0} + (T^{0l}x^jx^k)_{,l}. \quad (19.7c)$$

(Verify that this follows from $T^{\alpha\beta}_{,\beta} = 0$.) (4) Simplify h_{00} and h_{0m} by the gauge transformation generated by

$$\begin{aligned} \xi_0 &= \frac{1}{2r}\frac{\partial}{\partial t}\int T^{00'}r'^2 d^3x' + \frac{x^j}{r^3}\int \left(T^{0k'}x^{k'}x^{j'} - \frac{1}{2}T^{0j'}r'^2\right) d^3x' \\ &\quad + \int (T_{00'} + T_{ll'})\left[\frac{x^jx^{j'}}{r^2} + \frac{(3x^{j'}x^{k'} - r'^2\delta_{jk})x^jx^k}{2r^4}\right] d^3x' \\ &\quad + \sum_{n=2}^{\infty} \frac{1}{n!}\frac{\partial^{n-1}}{\partial t^{n-1}}\int (T_{00'} + T_{kk'})\frac{(r - |\mathbf{x} - \mathbf{x}'|)^n}{|\mathbf{x} - \mathbf{x}'|} d^3x', \\ \xi_m &= -\frac{2x^j}{r^3}\int T_{00'}x^{j'}x^{m'} d^3x' + 4\sum_{n=1}^{\infty} \frac{1}{n!}\frac{\partial^{n-1}}{\partial t^{n-1}}\int T_{0m'}\frac{(r - |\mathbf{x} - \mathbf{x}'|)^n}{|\mathbf{x} - \mathbf{x}'|} d^3x' \\ &\quad + \frac{x^m}{r}\xi_0 - \frac{1}{2}\left(\frac{1}{r}\right)_{,m}\int T_{00'}r'^2 d^3x' - \left(\frac{x^k}{r^2}\right)_{,m}\int \left(T^{0j'}x^{j'}x^{k'} - \frac{1}{2}T^{0k'}r'^2\right) d^3x' \\ &\quad - \sum_{n=2}^{\infty} \frac{1}{n!}\frac{\partial^{n-2}}{\partial t^{n-2}}\int (T_{00'} + T_{kk'})\left[\frac{(r - |\mathbf{x} - \mathbf{x}'|)^n}{|\mathbf{x} - \mathbf{x}'|}\right]_{,m} d^3x' \end{aligned}$$

Here $T_{\mu\nu}'$ denotes $T_{\mu\nu}(t - r, \mathbf{x}')$.]

(b) Prove that the system's mass and angular momentum are conserved. [Note: Because $T^{\alpha\beta}_{,\beta} = 0$ (self-gravity has negligible influence), the proof is no different here than in flat spacetime (Chapter 5).]

§19.2. MEASUREMENT OF THE MASS AND ANGULAR MOMENTUM

For a weakly gravitating system:

(1) total mass M can be measured by applying Kepler's "1-2-3" law to orbiting particles

The values of a system's mass and angular momentum can be measured by probing the imprint they leave in its external gravitational field. Of all tools one might use to probe, the simplest is a test particle in a gravitationally bound orbit. If the particle is sufficiently far from the source, its motion is affected hardly at all by the source's angular momentum or by the gravitational waves; only the spherical, Newtonian part of the gravitational field has a significant influence. Hence, the particle moves in an elliptical Keplerian orbit. To determine the source's mass M , one need only apply Kepler's third law (perhaps better called "Kepler's 1-2-3 law"):

$$M = \left(\frac{2\pi}{\text{orbital period}}\right)^2 \left(\text{Semi-major axis}\right)^3; \quad \text{i.e., } M^1 = \omega^2 a^3. \quad (19.8)$$

The source's angular momentum is not measured quite so easily. One must use a probe that is insensitive to Newtonian gravitational effects, but "feels" the off-diagonal term,

$$g_{0j} = -2\epsilon^{jk\ell} S^k x^\ell / r^3, \quad (19.9)$$

in the metric (19.5). One such probe is the precession of the perihelion of a corevolving satellite, relative to the precession for a counterrevolving satellite. A gyroscope is another such probe. Place a gyroscope at rest in the source's gravitational field. By a force applied to its center of mass, prevent it from falling. As time passes, the g_{0j} term in the metric will force the gyroscope to precess relative to the basis vectors $\partial/\partial x^j$; and since these basis vectors are "tied" to the coordinate system, which in turn is tied to the Lorentz frames at infinity, which in turn are tied to the "fixed stars" (cf. §39.12), the precession is relative to the "fixed stars." The angular velocity of precession, as derived in exercise 19.2, is

$$\Omega = \frac{1}{r^3} \left[-S + \frac{3(S \cdot x)x}{r^2} \right]. \quad (19.10)$$

One sometimes says that the source's rotation "drags the inertial frames near the source," thereby forcing the gyroscope to precess. For further discussion, see §§21.12, 40.7, and 33.4.

(2) total angular momentum S can be measured by examining the precession of gyroscopes

Exercise 19.2. GYROSCOPE PRECESSION

Derive equation (19.10) for the angular velocity of gyroscope precession. [Hints: Place an orthonormal tetrad at the gyroscope's center of mass. Tie the tetrad rigidly to the coordinate system, and hence to the "fixed stars"; more particularly, choose the tetrad to be that basis $\{e_{\hat{\alpha}}\}$ which is dual to the following 1-form basis:

$$\omega^i = [1 - (2M/r)]^{1/2} dt + 2\epsilon_{jk\ell} S^k(x^\ell/r^3) dx^j, \quad \omega^j = [1 + (2M/r)]^{1/2} dx^j. \quad (19.11)$$

The spatial legs of the tetrad, e_j , rotate relative to the gyroscope with an angular velocity ω given by [see equation (13.69)]

$$-\epsilon_{ijk}\omega^k = \Gamma_{ij0}.$$

Consequently, the gyroscope's angular momentum vector L precesses relative to the tetrad with angular velocity $\Omega = -\omega$:

$$\frac{dL^j}{dt} = \epsilon_{jk\ell} \Omega^k L^\ell, \quad \epsilon_{ijk} \Omega^k = \Gamma_{ij0}. \quad (19.12)$$

Calculate Γ_{ij0} for the given orthonormal frame, and thereby obtain equation (19.10) for Ω .]

EXERCISE

§19.3. MASS AND ANGULAR MOMENTUM OF FULLY RELATIVISTIC SOURCES

Abandon, now, the restriction to weakly gravitating sources. Consider an isolated, gravitating system inside which spacetime may or may not be highly curved—a black hole, a neutron star, the Sun, . . . But refuse, for now, to analyze the system's interior or the "strong-field region" near the system. Instead, restrict attention to the weak

gravitational field far from the source, and analyze it using linearized theory in vacuum. Expand $h_{\mu\nu}$ in multipole moments and powers of $1/r$; and adjust the gauge, the Lorentz frame, and the origin of coordinates to simplify the resulting metric. The outcome of such a calculation is a gravitational field identical to that for a weak source [equation (19.5)!] (Details of the calculation are not spelled out here because of their length; but see exercise 19.3.)

But before accepting this as the distant field of an arbitrary source, one should examine the nonlinear effects in the vacuum field equations. Two types of nonlinearities turn out to be important far from the source: (1) nonlinearities in the static, Newtonian part of the metric, which generate metric corrections

$$\delta g_{00} = -2M^2/r^2, \quad \delta g_{jk} = \frac{3}{2}(M^2/r^2)\delta_{jk},$$

(see exercise 19.3 and §39.8), thereby putting the metric into the form

$$ds^2 = - \left[1 - \frac{2M}{r} + \frac{2M^2}{r^2} + 0\left(\frac{1}{r^3}\right) \right] dt^2 - \left[4\epsilon_{jkl}S^k \frac{x^l}{r^3} + 0\left(\frac{1}{r^3}\right) \right] dt dx^j + \left[\left(1 + \frac{2M}{r} + \frac{3M^2}{2r^2}\right) \delta_{jk} + \left(\text{gravitational radiation terms}\right) \right] dx^j dx^k; \quad (19.13)$$

Metric far from any gravitating system, as a power series in $1/r$

(2) a gradual decrease in the source's mass, gradual changes in its angular momentum, and gradual changes in its "rest frame" to compensate for the mass, angular momentum, and linear momentum carried off by gravitational waves (see Box 19.1, which is best read only after finishing this section).

By measuring the distant spacetime geometry (19.13) of a given source, one cannot discover whether that source has strong internal gravity, or weak. But when one expresses the constants M and S_j , which determine g_{00} and g_{0j} , as integrals over the interior of the source, one discovers a crucial difference: if the internal gravity is weak, then linearized theory is valid throughout the source, and

$$M = \int T_{00} d^3x, \quad S_j = \int \epsilon_{jkl} x^k T^{l0} d^3x; \quad (19.14)$$

Failure of volume integrals for M and S when source has strong internal gravity

but if the gravity is strong, these formulas fail. Does this failure prevent one, for strong gravity, from identifying the constants M and S_j of the metric (19.13) as the source's mass and angular momentum? Not at all, according to the following argument.

Consider, first, the mass of the sun. For the sun one expects Newtonian theory to be highly accurate (fractional errors $\sim M_\odot/R_\odot \sim 10^{-6}$); so one can assert that the constant M appearing in the line element (19.13) is, indeed

$$M = \int \rho d^3x = \int T_{00} d^3x = \text{total mass.}$$

But might this assertion be wrong? To gain greater confidence and insight, adopt the viewpoint of "*controlled ignorance*"; i.e., do not pretend to know more than what is needed. (This style of physical argument goes back to Newton's famous "*Hypotheses non fingo*," i.e. "I do not feign hypotheses.") In evaluating the volume integral of T_{00} (usual Newtonian definition of M), one needs a theory of the internal structure

of the sun. For example, one must know that the visible surface layers of the sun do not hide a massive central core, so dense and large that relativistic gravitational fields $|\Phi| \sim 1$ exist there. If one makes use in the analysis of a fluid-type stress-energy tensor $T^{\mu\nu}$, one needs to know equations of state, opacities, and theories of energy generation and transport. One needs to justify the fluid description as an adequate approximation to the atomic constitution of matter. One needs to assume that an ultimate theory of matter explaining the rest masses of protons and electrons will not assign an important fraction of this mass to strong (nonlinear) gravitational fields on a submicroscopic scale. It is plausible that one could do all this, but it is also obvious that this is not the way the mass of the sun is, in fact, determined by astronomers! Theories of stellar structure are adjusted to give the observed mass; they are not constructed to let one deduce the mass from nongravitational observations. The mass of the sun is measured in practice by studying the orbits of planets in its external gravitational field, a procedure equivalent to reading the mass M off the line element (19.13), rather than evaluating the volume integral $\int T^{00} d^3x$.

To avoid all the above uncertainties, and to make theory correspond as closely as possible to experiment, *one defines the “total mass-energy” M of the sun or any other body to be the constant that appears in the line element (19.13) for its distant external spacetime geometry. Similarly, one defines the body’s intrinsic angular momentum as the constant 3-vector \mathbf{S} appearing in its line element (19.13).* Operationally, the total mass-energy M is measured via Kepler’s third law; the angular momentum \mathbf{S} is measured via its influence on the precession of a gyroscope or a planetary orbit. This is as true when the body is a black hole or a neutron star as when it is the sun.

What kind of a geometric object is the intrinsic angular momentum \mathbf{S} ? It is defined by measurements made far from the source, where, with receding distance, spacetime is becoming flatter and flatter (asymptotically flat). Thus, it can be regarded as a 3-vector in the “asymptotically flat spacetime” that surrounds the source. But in what Lorentz frame is \mathbf{S} a 3-vector? Clearly, in the asymptotic Lorentz frame where the line element (19.13) is valid; i.e., in the asymptotic Lorentz frame where the source’s distant “coulomb” (“ M/r ”) field is static; i.e., in the “asymptotic rest frame” of the source. Alternatively, one can regard \mathbf{S} as a 4-vector, \mathbf{S} , which is purely spatial ($S^0 = 0$) in the asymptotic rest frame. If one denotes the 4-velocity of the asymptotic rest frame by \mathbf{U} , then the fact that \mathbf{S} is purely spatial can be restated geometrically as $\mathbf{S} \cdot \mathbf{U} = 0$, or

$$\mathbf{S} \cdot \mathbf{P} = 0, \quad (19.15)$$

where

$$\mathbf{P} \equiv M\mathbf{U} \equiv \text{“total 4-momentum of source”} \quad (19.16)$$

is still another vector residing in the asymptotically flat region of spacetime.

The total 4-momentum \mathbf{P} and intrinsic angular momentum \mathbf{S} satisfy conservation laws that are summarized in Box 19.1. These conservation laws are valuable tools in gravitation theory and relativistic astrophysics, but the derivation of these laws (Chapter 20) does not compare in priority to topics such as neutron stars and basic cosmology; so most readers will wish to skip it on a first reading of this book.

Definition of “total mass-energy” M and “angular momentum” \mathbf{S} in terms of external gravitational field

\mathbf{S} as a geometric object in an asymptotically flat region far outside source

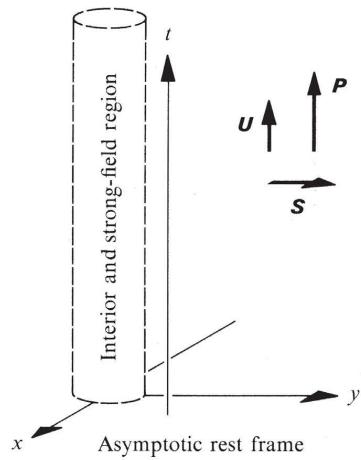
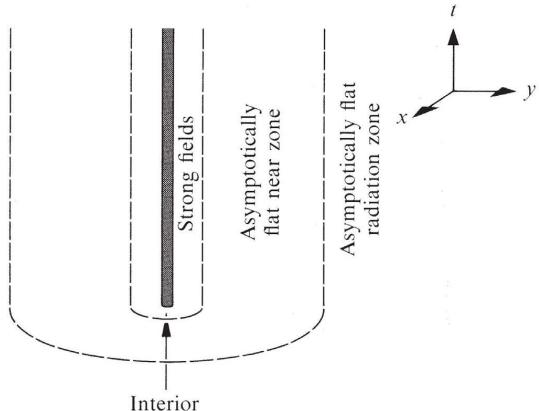
“Asymptotic rest frame” and “total 4-momentum”

Conservation laws for total 4-momentum and angular momentum

(continued on page 456)

Box 19.1 TOTAL MASS-ENERGY, 4-MOMENTUM, AND ANGULAR MOMENTUM OF AN ISOLATED SYSTEM

A. Spacetime is divided into (1) the source's interior; which is surrounded by (2) a strong-field vacuum region; which in turn is surrounded by (3) a weak-field, asymptotically flat, near-zone region; which in turn is surrounded by (4) a weak-field, asymptotically flat, radiation-zone region. This box and this chapter treat only the asymptotically flat regions. The interior and strong-field regions are treated in the next chapter.



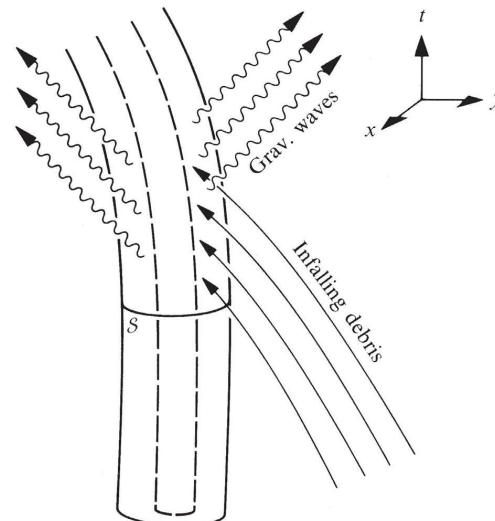
B. The asymptotic rest frame of the source is that global, asymptotically Lorentz frame (coordinates t , x , y , z) in which the distant, "coulomb" part of the source's field is at rest (see diagram). The asymptotic rest frame does not extend into the strong-field region; any such extension of it would necessarily be forced by the curvature into a highly non-Lorentz, curvilinear form. The spatial origin of the asymptotic rest frame is so adjusted that the source is centered on it—i.e., that the distant Newtonian potential is $\Phi = -M/(x^2 + y^2 + z^2)^{1/2} + 0(1/r^3)$; i.e., that Φ has no dipole term, $\mathbf{D} \cdot \mathbf{x}/r^3$, such as would originate from an offset of the coordinates.

C. To the source one can attribute a total mass-energy M , a 4-velocity \mathbf{U} , a total 4-momentum \mathbf{P} , and an intrinsic angular momentum vector, \mathbf{S} . The 4-vectors \mathbf{U} , \mathbf{P} , and \mathbf{S} reside in the asymptotically flat region of spacetime and can be moved about freely there (negligible curvature \Rightarrow parallel transport around closed curves does not change \mathbf{U} , \mathbf{P} , or \mathbf{S}). The source's 4-velocity \mathbf{U} is defined to equal the 4-velocity of the asymptotic rest frame ($U^0 = 1$, $\mathbf{U} = 0$ in rest frame). The total mass-energy M is measured via Kepler's third ("1-2-3") law [equation (19.8)]. The total 4-momentum is defined by $\mathbf{P} \equiv M\mathbf{U}$. The intrinsic angular momentum \mathbf{S} is orthogonal to the 4-velocity \mathbf{U} , $\mathbf{S} \cdot \mathbf{U} = 0$ (so $S^0 = 0$; $\mathbf{S} \neq 0$ in general in asymptotic rest frame); \mathbf{S} is measured via gyroscope precession or differential perihelion precession (§19.2).

In the asymptotic rest frame, with an appropriate choice of gauge (i.e., of ripples in the coordinates), *the slight deviations from flat-spacetime geometry are described by the line element*

$$ds^2 = - \left[1 - \frac{2M}{r} + \frac{2M^2}{r^2} + 0\left(\frac{1}{r^3}\right) \right] dt^2 - \left[4\epsilon_{jk\ell} S^k \frac{x^\ell}{r^3} + 0\left(\frac{1}{r^3}\right) \right] dt dx^j + \left[\left(1 + \frac{2M}{r} + \frac{3M^2}{2r^2}\right) \delta_{jk} + (\text{gravitational radiation terms}) \right] dx^j dx^k. \quad (1)$$

D. Conservation of 4-momentum and angular momentum: Suppose that particles fall into a source or are ejected from it; suppose that electromagnetic waves flow in and out; suppose the source emits gravitational waves. All such processes break the source's isolation and can change its total 4-momentum \mathbf{P} , its intrinsic angular momentum \mathbf{S} , and its asymptotic rest frame. Surround the source with a spherical shell S , which is far enough out to be in the asymptotically flat region. Keep this shell always at rest in the source's momentary asymptotic rest frame. By probing the source's gravitational field near S , measure its 4-momentum \mathbf{P} and intrinsic angular momentum \mathbf{S} as functions of the shell's proper time τ . An analysis given in the next chapter reveals that the 4-momentum is conserved, in the sense that



Interstellar debris falls into a black hole, and gravitational waves emerge.

$$\frac{dP^\alpha}{d\tau} = - \int_S T^{\alpha j} n_j d(\text{area}) = \begin{pmatrix} \text{rate at which 4-momentum} \\ \text{flows inward through shell} \end{pmatrix}, \quad (2)$$

where \mathbf{n} is the unit outward normal to S and the integral is performed in the shell's momentary rest frame. In words: *the rate at which 4-momentum flows through the shell, as measured in the standard special relativistic manner, equals the rate of change of the source's gravitationally measured 4-momentum. Similarly, the angular momentum is conserved in the sense that*

$$\frac{dS_i}{d\tau} = - \int_S (\epsilon_{ijk} x^j T^{kl}) n_l d(\text{area}) = \begin{pmatrix} \text{rate at which angular} \\ \text{momentum flows inward} \\ \text{through the shell} \end{pmatrix}, \quad (3a)$$

$$\frac{dS_0}{d\tau} = - \frac{dU^\alpha}{d\tau} S_\alpha = \begin{pmatrix} \text{change required to keep } \mathbf{S} \text{ orthogonal to } \mathbf{U}; \\ \text{"Fermi-Walker-transport law"; cf. §§6.5, 13.6.} \end{pmatrix}. \quad (3b)$$

In these conservation laws $T^{\alpha\beta}$ is the total stress-energy tensor at the shell, including contributions from matter, electromagnetic fields, and gravitational waves. The gravitational-wave contribution, called $T^{(GW)\alpha\beta}$, is treated in Chapter 35.

Note: The conservation laws in the form stated above contain fractional errors of order M/r (contributions from "gravitational potential energy" of infalling material), but such errors go to zero in the limit of a very large shell ($r \rightarrow \infty$).

Note: The formulation of these conservation laws given in the next chapter is more precise and more rigorous, but less physically enlightening than the one here.

EXERCISE**Exercise 19.3. GRAVITATIONAL FIELD FAR FROM A STATIONARY, FULLY RELATIVISTIC SOURCE**

Derive the line element (19.13) for the special case of a source that is time-independent ($g_{\mu\nu,t} = 0$). This can be a difficult problem, if one does not proceed carefully along the following outlined route. (1) Initially ignore all nonlinearities in the Einstein field equations. The field is weak far from the source. These nonlinearities will be absent from the dominant terms. (2) Calculate the dominant terms using linearized theory in the Lorentz gauge [equations (18.8)]. (3) In particular, write the general solution to the vacuum, time-independent wave equation (18.8b) in the following form involving $n^j \equiv x^j/r \equiv$ (unit vector in radial direction):

$$\begin{aligned}\bar{h}_{00} &= \frac{A^0}{r} + \frac{B^j n^j}{r^2} + O\left(\frac{1}{r^3}\right), \\ \bar{h}_{0j} &= \frac{A^j}{r} + \frac{B^{jk} n^k}{r^2} + O\left(\frac{1}{r^3}\right), \\ \bar{h}_{jk} &= \frac{A^{jk}}{r} + \frac{B^{jk\ell} n^\ell}{r^2} + O\left(\frac{1}{r^3}\right), \\ A^{jk} &= A^{(jk)}, \quad B^{jk\ell} = B^{(jk)\ell}.\end{aligned}\tag{19.17}$$

(Round brackets denote symmetrization.) (4) Then impose the Lorentz gauge conditions $\bar{h}_{\alpha\beta,\beta} = 0$ on this general solution, thereby learning

$$\begin{aligned}A^j &= 0, \quad A^{jk} = 0, \\ B^{jk}(\delta^{jk} - 3n^j n^k) &= 0, \\ B^{jk\ell}(\delta^{jk\ell} - 3n^k n^\ell) &= 0.\end{aligned}\tag{19.18}$$

(5) Write B^{jk} as the sum of its trace $3B$, its traceless symmetric part S^{jk} , and its traceless antisymmetric part (these are its “irreducible parts”):

$$B^{jk} = B\delta^{jk} + S^{jk} + \epsilon^{jk\ell} F^\ell, \quad S^{jj} = 0.\tag{19.19}$$

Show that any tensor B^{jk} can be put into such a form. Then show that the gauge conditions (19.18) imply $S^{jk} = 0$. (6) Similarly show that any tensor $B^{jk\ell}$ that is symmetric on its first two indices can be put into the form

$$\begin{aligned}B^{jk\ell} &= \delta^{jk} A^\ell + C^{(j} \delta^{k)\ell} + \epsilon^{m\ell(j} E^{k)m} + S^{jk\ell}, \\ E^{km} \text{ symmetric and traceless, i.e., } E^{km} &= E^{(km)}, \quad E^{kk} = 0, \\ S^{jk\ell} \text{ symmetric and traceless, i.e., } S^{jk\ell} &= S^{(jk)\ell}, \\ S^{jj\ell} &= S^{ikk} = S^{ijk} = 0.\end{aligned}\tag{19.20}$$

Then show that the gauge conditions (19.18) imply $C^j = -2A^j$ and $E^{km} = S^{jk\ell} = 0$. (7) Combining all these results, conclude that

$$\begin{aligned}\bar{h}_{00} &= \frac{A^0}{r} + \frac{B^j n^j}{r^2} + O\left(\frac{1}{r^3}\right), \\ \bar{h}_{0j} &= \frac{\epsilon^{jk\ell} n^k F^\ell}{r^2} + \frac{B n^j}{r^2} + O\left(\frac{1}{r^3}\right), \\ \bar{h}_{jk} &= \frac{\delta^{jk} A^\ell n^\ell - A^j n^k - A^k n^j}{r^2} + O\left(\frac{1}{r^3}\right).\end{aligned}\tag{19.21}$$

Then use gauge transformations, which stay within Lorentz gauge, to eliminate B and A^j from \bar{h}_{0j} and \bar{h}_{jk} ; so

$$\begin{aligned}\bar{h}_{00} &= \frac{A^0}{r} + \frac{(B^j + A^j)n^j}{r^2} + O\left(\frac{1}{r^3}\right), \\ \bar{h}_{0j} &= \frac{\epsilon^{jkl}n^kF^l}{r^2} + O\left(\frac{1}{r^3}\right), \\ \bar{h}_{jk} &= O\left(\frac{1}{r^3}\right).\end{aligned}\quad (19.22)$$

- (8) Translate the origin of coordinates so $x^j_{\text{new}} = x^j_{\text{old}} - (B^j + A^j)/A^0$; in the new coordinate system $\bar{h}_{\alpha\beta}$ has the same form as (19.22), but with $B^j + A^j$ removed. From the resultant $\bar{h}_{\alpha\beta}$, construct the metric and redefine the constants A^0 and F^l to agree with equation (19.13).
(9) All linear terms in the metric are now accounted for. The dominant nonlinear terms must be proportional to the square, $(M/r)^2$, of the dominant linear term. The easiest way to get the proportionality constant is to take the Schwarzschild geometry for a fully relativistic, static, spherical source [equation (31.1)], by a change of coordinates put it in the form

$$ds^2 = -\left(\frac{1-M/2r}{1+M/2r}\right)^2 dt^2 + \left(1+\frac{M}{2r}\right)^4 (dx^2 + dy^2 + dz^2) \quad (19.23)$$

(exercise 25.8), and expand it in powers of M/r .

§19.4. MASS AND ANGULAR MOMENTUM OF A CLOSED UNIVERSE

“There are no snakes in Ireland.”

Statement of St. Patrick
after driving the snakes
out of Ireland (legend*)

There is no such thing as “the energy (or angular momentum, or charge) of a closed universe,” according to general relativity, and this for a simple reason. To weigh something one needs a platform on which to stand to do the weighing.

To weigh the sun, one measures the periods and semimajor axes of planetary orbits, and applies Kepler’s “1-2-3” law, $M = \omega^2 a^3$. To measure the angular momentum, S , of the sun (a task for space technology in the 1970’s or 1980’s!), one measures the precession of a gyroscope in a near orbit about the sun, or one examines some other aspect of the “dragging of inertial frames.” To determine the electric charge

For a closed universe the total mass-energy M and angular momentum S are undefined and undefinable

* Stokes (1887) and other standard references deny this legend. In part I of Stokes the basic manuscript references are listed, including especially codex manuscript Rawlinson B.512 in 154 folios, in double columns, written by various hands in the fourteenth and fifteenth centuries (*cf. Catalogi codicum manuscriptorum Bibliothecae Bodleianaee Partis Quintae Fasciculus Primus*, Oxford, 1862, col. 728–732). In this manuscript, folio 97b.1, line 14, reads in the translation of Stokes, Part I, p. xxx: “as Paradise is without beasts, without a snake, without a lion, without a dragon, without a scorpion, without a mouse, without a frog, so is Ireland in the same manner without any harmful animal, save only the wolf . . .”

of a body, one surrounds it by a large sphere, evaluates the electric field normal to the surface at each point on this sphere, integrates over the sphere, and applies the theorem of Gauss. But within any closed model universe with the topology of a 3-sphere, a Gaussian 2-sphere that is expanded widely enough from one point finds itself collapsing to nothingness at the antipodal point. Also collapsed to nothingness is the attempt to acquire useful information about the “charge of the universe”: the charge is trivially zero. By the same token, every “surface integral” (see details in Chapter 20) to determine mass-energy or angular momentum collapses to nothingness. To make the same point in another way: around a closed universe there is no place to put a test object or gyroscope into Keplerian orbit to determine either any so-called “total mass” or “rest frame” or “4-momentum” or “angular momentum” of the system. These terms are undefined and undefinable. Words, yes; meaning, no.

Not having a defined 4-momentum for a closed universe may seem at first sight disturbing; but it would be far more disturbing to be given four numbers and to be told authoritatively that they represent the components of some purported “total energy-momentum 4-vector of the universe.” Components with respect to what local Lorentz frame? At what point? And what about the change in this vector on parallel transport around a closed path leading back to that strangely preferred point? It is a happy salvation from these embarrassments that the issue does not and cannot arise!

Imagine a fantastically precise measurement of the energy of a γ -ray. The experimenter wishes to know how much this γ -ray contributes to the total mass-energy of the universe. Having measured its energy in the laboratory, he then corrects it for the negative gravitational energy by which it is bound to the Earth. The result,

$$E_{\text{corrected}} = h\nu(1 - M_\oplus/R_\oplus),$$

is the energy the photon will have after it climbs out of the Earth’s gravitational field. But this is only the first in a long chain of corrections for energy losses (redshifts) as the photon climbs out of the gravitational fields of the solar system, the galaxy, the local cluster of galaxies, the supercluster, and then what? These corrections show no sign of converging, unless to $E_{\text{corrected}} = 0$.

Asymptotic flatness as the key to the definability of M and \mathbf{S}

Quite in contrast to the charge-energy-angular-momentum facelessness of a closed universe are the attractive possibilities of defining and measuring all three quantities in any space that is asymptotically flat. One does not have to revolutionize present-day views of cosmology to talk of asymptotically flat space. It is enough to note how small is the departure from flatness, as measured by the departure of $(-g_{00})^{1/2}$ from unity, in cases of astronomical or astrophysical interest (Box 19.2). Surrounding a region where any dynamics, however complicated, is going on, whenever the geometry is asymptotically flat to some specified degree of precision, then to that degree of precision it makes sense to speak of the total energy-momentum 4-vector of the dynamic region, \mathbf{P} , and its total intrinsic angular momentum, \mathbf{S} . Parallel transport of either around any closed curve in the flat region brings it back to its

Box 19.2 METRIC CORRECTION TERM NEAR SELECTED HEAVENLY BODIES

	m	m	r	$\frac{m}{r} = 1 - (-g_{00})^{1/2}$
At shoulder of Venus de Milo	$2 \times 10^5 \text{ g}$	$= 1.5 \times 10^{-23} \text{ cm}$	30 cm	5×10^{-25}
At surface of Earth	$6 \times 10^{27} \text{ g}$	$= 4 \times 10^{-1} \text{ cm}$	$6.4 \times 10^8 \text{ cm}$	6×10^{-10}
At Earth's distance from sun	$2 \times 10^{33} \text{ g}$	$= 1.5 \times 10^5 \text{ cm}$	$1.5 \times 10^{13} \text{ cm}$	1×10^{-8}
At sun's distance from center of galaxy	$2 \times 10^{44} \text{ g}$	$= 1.5 \times 10^{16} \text{ cm}$	$2.5 \times 10^{22} \text{ cm}$	6×10^{-7}
At distance of galaxy from center of Virgo cluster of galaxies	$6 \times 10^{47} \text{ g}$	$= 4 \times 10^{19} \text{ cm}$	$3 \times 10^{25} \text{ cm}$	1×10^{-6}

starting point unchanged. Moreover, it makes no difference how enormous are the departures from flatness in the dynamic region (black holes, collapsing stars, intense gravitational waves, etc.); far away the curvature will be weak, and the 4-momentum and angular momentum will reveal themselves by their imprints on the spacetime geometry.

CHAPTER 20

CONSERVATION LAWS FOR 4-MOMENTUM AND ANGULAR MOMENTUM

We denote as energy of a material system in a certain state the contribution of all effects (measured in mechanical units of work) produced outside the system when it passes in an arbitrary manner from its state to a reference state which has been defined ad hoc.

WILLIAM THOMPSON (later Lord Kelvin),
as quoted by Max von Laue in Schilpp (1949), p. 514.

All forms of energy possess inertia.

ALBERT EINSTEIN, conclusion
from his paper of September 26, 1905,
as summarized by von Laue in Schilpp (1949), p. 523.

§20.1. OVERVIEW

Chapter 5 (stress-energy tensor) is needed as preparation for this chapter, which in turn is needed as preparation for the Track-2 portion of Chapter 36 (generation of gravitational waves) and will be useful in understanding Chapter 35 (propagation of gravitational waves).

Chapter 19 expounded the key features of total 4-momentum \mathbf{P} and total angular momentum \mathbf{S} for an arbitrary, gravitating system. But one crucial feature was left unproved: the conservation laws for \mathbf{P} and \mathbf{S} (Box 19.1). To prove those conservation laws is the chief purpose of this chapter. But other interesting, if less important, aspects of \mathbf{P} and \mathbf{S} will be encountered along the route to the proof—Gaussian flux integrals for 4-momentum and angular momentum; a stress-energy “pseudotensor” for the gravitational field, which is a tool in constructing volume integrals for \mathbf{P} and \mathbf{S} ; and the nonlocalizability of the energy of the gravitational field.

§20.2. GAUSSIAN FLUX INTEGRALS FOR 4-MOMENTUM AND ANGULAR MOMENTUM

In electromagnetic theory one can determine the conserved total charge of a source by adding up the number of electric field lines emanating from it—i.e., by performing a Gaussian flux integral over a closed 2-surface surrounding it:

$$Q = \frac{1}{4\pi} \oint E^j d^2 S_j = \frac{1}{4\pi} \oint F^{0j} d^2 S_j. \quad (20.1)$$

Gaussian flux integrals for charge and Newtonian mass

Similarly, in Newtonian theory one can determine the mass of a source by evaluating the Gaussian flux integral

$$M = \frac{1}{4\pi} \oint \Phi_{,j} d^2 S_j. \quad (20.2)$$

These flux integrals work because the charge and mass of a source place indelible imprints on the electromagnetic and gravitational fields that envelop it.

The external gravitational field (spacetime geometry) in general relativity possesses similar imprints, imprints not only of the source's total mass-energy M , but also of its total 4-momentum \mathbf{P} and its intrinsic angular momentum \mathbf{S} (see Box 19.1). Hence, it is reasonable to search for Gaussian flux integrals that represent the 4-momentum and angular momentum of the source.

To simplify the search, carry it out initially in linearized theory, and use Maxwell electrodynamics as a guide. In electrodynamics the Gaussian flux integral for charge follows from Maxwell's equations $F^{\mu\nu}_{,\nu} = 4\pi J^\mu$, plus the crucial fact that $F^{\mu\nu}$ is antisymmetric, so that $F^{0\mu}_{,\mu} = F^{0j}_{,j}$:

$$Q = \int J^0 d^3 x = \frac{1}{4\pi} \int F^{0\nu}_{,\nu} d^3 x = \frac{1}{4\pi} \int F^{0j}_{,j} d^3 x = \frac{1}{4\pi} \oint F^{0j} d^2 S_j.$$

↑
[Gauss's theorem]

To find analogous flux integrals in linearized theory, rewrite the linearized field equations (18.7) in an analogous form involving an entity with analogous crucial symmetries. The entity needed turns out to be

$$H^{\mu\alpha\nu\beta} \equiv -(\bar{h}^{\mu\nu}\eta^{\alpha\beta} + \eta^{\mu\nu}\bar{h}^{\alpha\beta} - \bar{h}^{\alpha\nu}\eta^{\mu\beta} - \bar{h}^{\mu\beta}\eta^{\alpha\nu}). \quad (20.3) \quad H^{\mu\alpha\nu\beta} \text{ defined}$$

As one readily verifies from this expression, it has the same symmetries as the Riemann tensor

$$\begin{aligned} H^{\mu\alpha\nu\beta} &= H^{\nu\beta\mu\alpha} = H^{[\mu\alpha][\nu\beta]}, \\ H^{\mu[\alpha\nu\beta]} &= 0. \end{aligned} \quad (20.4)$$

This entity, like $\bar{h}^{\mu\nu}$, transforms as a tensor under the Lorentz transformations of linearized theory; but it is not gauge-invariant, so it is not a tensor in the general relativistic sense.

Linearized field equations in terms of $H^{\mu\alpha\nu\beta}$

Gaussian flux integrals in linearized theory: (1) for 4-momentum

(2) for angular momentum

Generalization of Gaussian flux integrals to full general relativity

In terms of $H^{\mu\alpha\nu\beta}$, the linearized field equations (18.7) take on the much simplified form

$$2G^{\mu\nu} = H^{\mu\alpha\nu\beta}_{,\alpha\beta} = 16\pi T^{\mu\nu}; \quad (20.5)$$

and from these, by antisymmetry of $H^{\mu\alpha\nu\beta}$ in ν and β , follow the source conservation laws of linearized theory,

$$T^{\mu\nu}_{,\nu} = \frac{1}{16\pi} H^{\mu\alpha\nu\beta}_{,\alpha\beta\nu} = 0,$$

which were discussed back in §18.3. The same antisymmetry also produces a Gaussian flux integral for the source's total 4-momentum:

$$\begin{aligned} P^\mu &= \int T^{\mu 0} d^3x = \frac{1}{16\pi} \int H^{\mu\alpha 0\beta}_{,\alpha\beta} d^3x = \frac{1}{16\pi} \int H^{\mu\alpha 0j}_{,\alpha j} d^3x \\ &= \frac{1}{16\pi} \oint_S H^{\mu\alpha 0j}_{,\alpha} d^2S_j. \end{aligned} \quad [Gauss's \text{ theorem}] \quad (20.6)$$

Here the closed 2-surface of integration S must completely surround the source and must lie in a 3-surface of constant time x^0 . The integral (20.6) for the source's energy P^0 , which is used more frequently than the integrals for P^j , reduces to an especially simple form in terms of $g_{\alpha\beta} = \eta_{\alpha\beta} + h_{\alpha\beta}$:

$$P^0 = \frac{1}{16\pi} \int_S (g_{jk,k} - g_{kk,j}) d^2S_j \quad (20.7)$$

(see exercise 20.1).

A calculation similar to (20.6), but more lengthy (exercise 20.2), yields a flux integral for total angular momentum about the origin of coordinates:

$$\begin{aligned} J^{\mu\nu} &= \int (x^\mu T^{\nu 0} - x^\nu T^{\mu 0}) d^3x \\ &= \frac{1}{16\pi} \oint_S (x^\mu H^{\nu\alpha 0j}_{,\alpha} - x^\nu H^{\mu\alpha 0j}_{,\alpha} + H^{\mu j 0\nu} - H^{\nu j 0\mu}) d^2S_j. \end{aligned} \quad (20.8)$$

To evaluate the flux integrals in (20.6) to (20.8) (by contrast with the volume integrals), one need utilize only the gravitational field far outside the source. Since that gravitational field has the same form in full general relativity for strong sources as in linearized theory for weak sources, the flux integrals can be used to calculate P^μ and $J^{\mu\nu}$ for any isolated source whatsoever, weak or strong:

$$\left. \begin{aligned} P^\mu &= \frac{1}{16\pi} \oint_S H^{\mu\alpha 0j}_{,\alpha} d^2S_j, \\ P^0 &= \frac{1}{16\pi} \oint_S (g_{jk,k} - g_{kk,j}) d^2S_j, \\ J^{\mu\nu} &= \frac{1}{16\pi} \oint_S (x^\mu H^{\nu\alpha 0j}_{,\alpha} - x^\nu H^{\mu\alpha 0j}_{,\alpha} \\ &\quad + H^{\mu j 0\nu} - H^{\nu j 0\mu}) d^2S_j. \end{aligned} \right\} \begin{array}{l} \text{in full general relativity} \\ \text{theory, for any isolated} \\ \text{source, when the closed} \\ \text{surface of integration } S \\ \text{is in the asymptotically} \\ \text{flat region surrounding} \\ \text{the source, and when} \\ \text{asymptotically Minkows-} \\ \text{kian coordinates are used.} \end{array} \quad (20.9)$$

Knowing P^μ and $J^{\mu\nu}$, one can calculate the source's total mass-energy M and intrinsic angular momentum S^μ by the standard procedure of Box 5.6:

$$M = (-P^\mu P_\mu)^{1/2}, \quad (20.10)$$

$$Y^\mu = -J^{\mu\nu}P_\nu/M^2 = \begin{cases} \text{vector by which the source's asymptotic,} \\ \text{"M/r", spherical field is displaced from} \\ \text{being centered on the origin of coordinates} \end{cases} \quad (20.11)$$

$$S_\rho = \frac{1}{2}\epsilon_{\mu\nu\sigma\rho}(J^{\mu\nu} - Y^\mu P^\nu + Y^\nu P^\mu)P^\sigma/M. \quad (20.12)$$

Note especially that *the integrands of the flux integrals (20.9) are not gauge-invariant*. In any local inertial frame at an event $\mathcal{P}_0 [g_{\mu\nu}(\mathcal{P}_0) = \eta_{\mu\nu}, g_{\mu\nu,\alpha}(\mathcal{P}_0) = 0]$ they vanish, since

$$g_{\mu\nu,\alpha} = h_{\mu\nu,\alpha} = 0 \Rightarrow H^{\mu\nu\alpha\beta}_{,\alpha} = 0; \quad g_{\mu\nu} = \eta_{\mu\nu} \Rightarrow H^{\mu\nu\alpha\beta} = 0.$$

This is reasonable behavior; their Newtonian analog, the integrand Φ_j = (gravitational acceleration) of the Newtonian flux integral (20.2), similarly vanishes in local inertial frames.

Although the integrands of the flux integrals are not gauge-invariant, the total integrals P^μ (4-momentum) and $J^{\mu\nu}$ (angular momentum) most assuredly are! They have meaning and significance independent of any coordinate system and gauge. They are tensors in the asymptotically flat region surrounding the source.

The spacetime must be asymptotically flat if there is to be any possibility of defining energy and angular momentum. Only then can linearized theory be applied; and only on the principle that linearized theory applies far away can one justify using the flux integrals (20.9) in the full nonlinear theory. Nobody can compel a physicist to move in close to define energy and angular momentum. He has no need to move in close; and he may have compelling motives not to: the internal structure of the sources may be inaccessible, incomprehensible, uninteresting, dangerous, expensively distant, or frightening. This requirement for far-away flatness is a remarkable feature of the flux integrals (20.9); it is also a decisive feature. Even the coordinates must be asymptotically Minkowskian; otherwise most formulas in this chapter fail or require modification. In particular, *when evaluating the 4-momentum and angular momentum of a localized system, one must apply the flux integrals (20.9) only in asymptotically Minkowskian coordinates. If such coordinates do not exist (spacetime not flat at infinity), one must completely abandon the flux integrals, and the quantities that rely on them for definition: the total mass, momentum, and angular momentum of the gravitating source.* In this connection, recall the discussion of §19.4. It described, in physical terms, why "total mass-energy" is a limited concept, useful only when one adopts a limited viewpoint that ignores cosmology. (Compare "light ray" or "particle," concepts of enormous value, but concepts that break down when wave optics or wave mechanics enter significantly.)

Summary: Attempts to use formulas (20.9) in ways that lose sight of the Minkowski boundary conditions (and especially simply adopting them unmodified in curvilinear coordinates) easily and unavoidably produce nonsense.

Total mass-energy, center of mass, and intrinsic angular momentum

Gaussian flux integrals valid only in asymptotically flat region of spacetime and in asymptotically Minkowskian coordinates

EXERCISES**Exercise 20.1. FLUX INTEGRAL FOR TOTAL MASS-ENERGY IN LINEARIZED THEORY**

Show that the flux integral (20.6) for P^0 reduces to (20.7). Then show that, when applied to a nearly Newtonian source [line element (18.15c)], it reduces further to the familiar Newtonian flux integral (20.2).

Exercise 20.2. FLUX INTEGRAL FOR ANGULAR MOMENTUM IN LINEARIZED THEORY

Derive the Gaussian flux integral (20.8) for $J^{\mu\nu}$. [*Hint:* use the field equations (20.5) to show

$$16\pi x^\mu T^{r0} = (x^\mu H^{\nu\alpha 0k})_{,\alpha} - H^{\nu j 0\mu}_{,\nu} - H^{\nu 00\mu}_{,\nu}; \quad (20.13)$$

and then use Gauss's theorem to evaluate the volume integral of equation (20.8)].

Exercise 20.3. FLUX INTEGRALS FOR AN ARBITRARY STATIONARY SOURCE

(a) Use the flux integrals (20.9) to calculate P^μ and $J^{\mu\nu}$ for an arbitrary stationary source. For the asymptotically flat metric around the source, use (19.13), with the gravitational radiation terms set to zero.

(b) Verify that the “auxiliary equations” (20.10) to (20.12) give the correct answer for this source’s total mass-energy M and intrinsic angular momentum S^μ .

§20.3. VOLUME INTEGRALS FOR 4-MOMENTUM AND ANGULAR MOMENTUM

The full Einstein field equations in terms of $H^{\mu\alpha\nu\beta}$

Volume integrals for 4-momentum and angular momentum in full general relativity

It is easy, in linearized theory, to convert the surface integrals for P^μ and $J^{\mu\nu}$ into volume integrals over the source; one can simply trace backward the steps that led to the surface integrals in the first place [equation (20.6); exercise 20.2]. How, in full general relativity, can one similarly convert from the surface integrals to volume integrals? The answer is rather easy, if one thinks in the right direction. One need only put the full Einstein field equations into the form

$$H^{\mu\alpha\nu\beta}_{,\alpha\beta} = 16\pi T_{\text{eff}}^{\mu\nu} \quad (20.14)$$

analogous to equations (20.5) of linearized theory. Here $H^{\mu\alpha\nu\beta}$ is to be defined in terms of $h_{\mu\nu} \equiv g_{\mu\nu} - \eta_{\mu\nu}$ by equation (20.3), even deep inside the source where $|h_{\mu\nu}|$ might be $\gtrsim 1$. This form of the Einstein equations then permits a conversion of the Gaussian flux integrals into volume integrals, just as in linearized theory:

$$\begin{aligned} P^\mu &= \frac{1}{16\pi} \oint H^{\mu\alpha 0j}_{,\alpha} d^2 S_j = \frac{1}{16\pi} \int H^{\mu\alpha 0j}_{,\alpha j} d^3 x = \frac{1}{16\pi} \int H^{\mu\alpha 0\beta}_{,\alpha\beta} d^3 x \\ &= \int T_{\text{eff}}^{\mu 0} d^3 x. \end{aligned} \quad (20.15)$$

Similarly,

$$J^{\mu\nu} = \int (x^\mu T_{\text{eff}}^{\nu 0} - x^\nu T_{\text{eff}}^{\mu 0}) d^3 x. \quad (20.16)$$

[Crucial to the conversion is the use of partial derivatives rather than covariant derivatives in equations (20.14).] In these volume integrals, as throughout the preceding discussion, the coordinates must become asymptotically Lorentz ($g_{\mu\nu} \rightarrow \eta_{\mu\nu}$) far from the source.

The form of $T_{\text{eff}}^{\mu\nu}$ can be calculated by recalling that $H^{\mu\alpha\nu\beta}_{,\alpha\beta}$ is a linearized approximation to the Einstein curvature tensor (20.5). Define the nonlinear corrections by

$$16\pi t^{\mu\nu} \equiv H^{\mu\alpha\nu\beta}_{,\alpha\beta} - 2G^{\mu\nu}. \quad (20.17) \quad t^{\mu\nu} \text{ (''stress-energy pseudotensor'' defined}$$

(To calculate them in terms of $g_{\mu\nu}$ or $h_{\mu\nu} = g_{\mu\nu} - \eta_{\mu\nu}$ is straightforward but lengthy. The precise form of these corrections will never be needed in this book.) Then Einstein's equations read

$$H^{\mu\alpha\nu\beta}_{,\alpha\beta} = 16\pi t^{\mu\nu} + 2G^{\mu\nu} = 16\pi(t^{\mu\nu} + T^{\mu\nu}),$$

so that

$$T_{\text{eff}}^{\mu\nu} = T^{\mu\nu} + t^{\mu\nu}. \quad (20.18) \quad T_{\text{eff}}^{\mu\nu} \text{ defined}$$

The quantity $t^{\mu\nu}$ is sometimes called a “stress-energy pseudotensor for the gravitational field.” The Einstein field equations (20.14) imply, because $H^{\mu\alpha\nu\beta}_{,\alpha\beta}$ is anti-symmetric in ν and β , that

$$T_{\text{eff},\nu}^{\mu\nu} = (T^{\mu\nu} + t^{\mu\nu})_{,\nu} = 0. \quad (20.19) \quad \text{Conservation law for } T_{\text{eff}}^{\mu\nu}$$

These equations are equivalent to $T^{\mu\nu}_{;\nu} = 0$, but they are written with partial derivatives rather than covariant derivatives—a fact that permits conversions back and forth between volume integrals and surface integrals.

All the quantities $H^{\mu\alpha\nu\beta}$, $T_{\text{eff}}^{\mu\nu}$, and $t^{\mu\nu}$ depend for their definition and existence on the choice of coordinates; they have no existence independent of coordinates; they are not components of tensors or of any other geometric object. Correspondingly, the equations (20.14) to (20.19) involving $T_{\text{eff}}^{\mu\nu}$ and $t^{\mu\nu}$ have no geometric, coordinate-free significance; they are not “covariant tensor equations.” There is, nevertheless, adequate invariance under general coordinate transformations to give the values P^μ and $J^{\mu\nu}$ of the volume integrals (20.15) and (20.16) geometric, coordinate-free significance in the asymptotically flat region far outside the source. Although this invariance is hard to see in the volume integrals themselves, it is clear from the surface-integral forms (20.9) that no coordinate transformation which changes the coordinates only inside some spatially bounded region can influence the values of the integrals. For coordinate changes in the distant, asymptotically flat regions, linearized theory guarantees that under Lorentz transformations the integrals for P^μ and $J^{\mu\nu}$ will transform like special relativistic tensors, and that under infinitesimal coordinate transformations (gauge changes) they will be invariant.

$H^{\mu\alpha\nu\beta}$, $t^{\mu\nu}$, and $T_{\text{eff}}^{\mu\nu}$ are coordinate-dependent objects

Because $t^{\mu\nu}$ are not tensor components, they can vanish at a point in one coordinate system but not in another. The resultant ambiguity in defining a localized energy density t^{00} for the gravitational field has a counterpart in ambiguities that exist in

Other, equally good versions of $H_{\text{eff}}^{\mu\alpha\nu\beta}$, $t^{\mu\nu}$, $T_{\text{eff}}^{\mu\nu}$:

the formal definition of $t^{\mu\nu}$. It is clear that any quantities $H_{\text{new}}^{\mu\alpha\nu\beta}$ which agree with the original $H^{\mu\alpha\nu\beta}$ in the asymptotic weak-field region will give the same values as $H^{\mu\alpha\nu\beta}$ does for the P^μ and $J^{\mu\nu}$ surface integrals (20.9). One especially convenient choice has been given by Landau and Lifshitz (1962; §100), who define

$$(1) \quad H_{\text{L-L}}^{\mu\alpha\nu\beta}$$

$$H_{\text{L-L}}^{\mu\alpha\nu\beta} = g^{\mu\nu}g^{\alpha\beta} - g^{\alpha\nu}g^{\mu\beta}, \quad (20.20)$$

where $g^{\mu\nu} \equiv (-g)^{1/2}g^{\mu\nu}$. Landau and Lifshitz show that Einstein's equations can be written in the form

$$H_{\text{L-L},\alpha\beta}^{\mu\alpha\nu\beta} = 16\pi(-g)(T^{\mu\nu} + t_{\text{L-L}}^{\mu\nu}), \quad (20.21)$$

$$(2) \quad t_{\text{L-L}}^{\alpha\beta}$$

where the Landau-Lifshitz pseudotensor components

$$\begin{aligned} (-g)t_{\text{L-L}}^{\alpha\beta} = & \frac{1}{16\pi} \left\{ g^{\alpha\beta}_{,\lambda}g^{\lambda\mu}_{,\mu} - g^{\alpha\lambda}_{,\lambda}g^{\beta\mu}_{,\mu} + \frac{1}{2}g^{\alpha\beta}g_{\lambda\mu}g^{\lambda\nu}_{,\rho}g^{\mu\lambda}_{,\nu} \right. \\ & - (g^{\alpha\lambda}g_{\mu\nu}g^{\beta\nu}_{,\rho}g^{\mu\rho}_{,\lambda} + g^{\beta\lambda}g_{\mu\nu}g^{\alpha\nu}_{,\rho}g^{\mu\rho}_{,\lambda}) + g_{\lambda\mu}g^{\nu\rho}g^{\alpha\lambda}_{,\nu}g^{\beta\mu}_{,\rho} \\ & \left. + \frac{1}{8}(2g^{\alpha\lambda}g^{\beta\mu} - g^{\alpha\beta}g^{\lambda\mu})(2g_{\nu\rho}g_{\sigma\tau} - g_{\rho\sigma}g_{\nu\tau})g^{\nu\tau}_{,\lambda}g^{\rho\sigma}_{,\mu} \right\} \end{aligned} \quad (20.22)$$

$$(3) \quad T_{\text{L-L eff}}^{\mu\nu}$$

are precisely quadratic in the first derivatives of the metric. (Einstein also gave a pseudotensor $t_E^\mu_\nu$ with this property, but it was not symmetric and so did not lead to an integral for $J^{\mu\nu}$.) Because $H_{\text{L-L}}^{\mu\alpha\nu\beta}$ has the same symmetries as $H^{\mu\alpha\nu\beta}$ and equals $H^{\mu\alpha\nu\beta}$ far from the source (exercise 20.4), and because the field equations (20.21) in terms of $H_{\text{L-L}}^{\mu\alpha\nu\beta}$ have the same form as in terms of $H^{\mu\alpha\nu\beta}$, it follows that

$$T_{\text{L-L eff}}^{\mu\nu} \equiv (-g)(T^{\mu\nu} + t_{\text{L-L}}^{\mu\nu}) \quad (20.23a)$$

has all the properties of the $T_{\text{eff}}^{\mu\nu}$ introduced earlier in this section:

$$T_{\text{L-L eff},\nu}^{\mu\nu} = 0, \quad (20.23b)$$

$$P^\mu = \int T_{\text{L-L eff}}^{\mu 0} d^3x, \quad (20.23c)$$

$$J^{\mu\nu} = \int (x^\mu T_{\text{L-L eff}}^{\nu 0} - x^\nu T_{\text{L-L eff}}^{\mu 0}) d^3x. \quad (20.23d)$$

EXERCISE

Exercise 20.4. FORM OF $H_{\text{L-L}}^{\mu\alpha\nu\beta}$ FAR FROM SOURCE

Show that the entities $H_{\text{L-L}}^{\mu\alpha\nu\beta}$ of equations (20.20) reduce to $H^{\mu\alpha\nu\beta}$ (20.3) in the weak-field region far outside the source.

§20.4. WHY THE ENERGY OF THE GRAVITATIONAL FIELD CANNOT BE LOCALIZED

Consider an element of 3-volume $d\Sigma_\nu$ and evaluate the contribution of the “gravitational field” in that element of 3-volume to the energy-momentum 4-vector, using

in the calculation either the pseudotensor $t^{\mu\nu}$ or the pseudotensor $t_{\text{L-L}}^{\mu\nu}$ discussed in the last section. Thereby obtain

$$\mathbf{p} = \mathbf{e}_\mu t^{\mu\nu} d\Sigma_\nu$$

or

$$\mathbf{p} = \mathbf{e}_\mu t_{\text{L-L}}^{\mu\nu} d\Sigma_\nu.$$

Right? No, the question is wrong. The motivation is wrong. The result is wrong. The idea is wrong.

To ask for the amount of electromagnetic energy and momentum in an element of 3-volume makes sense. First, there is one and only one formula for this quantity. Second, and more important, this energy-momentum in principle “has weight.” It curves space. It serves as a source term on the righthand side of Einstein’s field equations. It produces a relative geodesic deviation of two nearby world lines that pass through the region of space in question. It is observable. Not one of these properties does “local gravitational energy-momentum” possess. There is no unique formula for it, but a multitude of quite distinct formulas. The two cited are only two among an infinity. Moreover, “local gravitational energy-momentum” has no weight. It does not curve space. It does not serve as a source term on the righthand side of Einstein’s field equations. It does not produce any relative geodesic deviation of two nearby world lines that pass through the region of space in question. It is not observable.

Why one cannot define a localized energy-momentum for the gravitational field

Anybody who looks for a magic formula for “local gravitational energy-momentum” is looking for the right answer to the wrong question. Unhappily, enormous time and effort were devoted in the past to trying to “answer this question” before investigators realized the futility of the enterprise. Toward the end, above all mathematical arguments, one came to appreciate the quiet but rock-like strength of Einstein’s equivalence principle. One can always find in any given locality a frame of reference in which all local “gravitational fields” (all Christoffel symbols; all $\Gamma^\alpha_{\mu\nu}$) disappear. No Γ ’s means no “gravitational field” and no local gravitational field means no “local gravitational energy-momentum.”

Nobody can deny or wants to deny that gravitational forces make a contribution to the mass-energy of a gravitationally interacting system. The mass-energy of the Earth-moon system is less than the mass-energy that the system would have if the two objects were at infinite separation. The mass-energy of a neutron star is less than the mass-energy of the same number of baryons at infinite separation. Surrounding a region of empty space where there is a concentration of gravitational waves, there is a net attraction, betokening a positive net mass-energy in that region of space (see Chapter 35). At issue is not the existence of gravitational energy, but the localizability of gravitational energy. It is not localizable. The equivalence principle forbids.

Look at an old-fashioned potato, replete with warts and bumps. With an orange marking pen, mark on it a “North Pole” and an “equator”. The length of the equator is very far from being equal to 2π times the distance from the North Pole to the

equator. The explanation, “curvature,” is simple, just as the explanation, “gravitation”, for the deficit in mass of the earth-moon system (or deficit for the neutron star, or surplus for the region of space occupied by the gravitational waves) is simple. Yet it is not possible to ascribe the deficit in the length of the equator in the one case, or in mass in the other case, in any uniquely right way to different elements of the manifold (2-dimensional in the one case, 3-dimensional in the other). Look at a small region on the surface of the potato. The geometry there is locally flat. Look at any small region of space in any of the three gravitating systems. In an appropriate coordinate system it is free of gravitational field. The over-all effect one is looking at is a global effect, not a local effect. That is what the mathematics cries out. That is the lesson of the nonuniqueness of the $t^{\mu\nu}$!

§20.5. CONSERVATION LAWS FOR TOTAL 4-MOMENTUM AND ANGULAR MOMENTUM

Consider a system such as our galaxy or the solar system, which is made up of many gravitating bodies. Some of the bodies may be highly relativistic (black holes; neutron stars), while others are not. However, insist that in the regions between the bodies spacetime be nearly flat (gravity be weak)—so flat, in fact, that one can cover the entire system with coordinates which are (almost) globally inertial, except in a small neighborhood of each body where gravity may be strong. Such coordinates can exist only if the Newtonian gravitational potential, $\Phi \approx \frac{1}{2}(\eta_{00} - g_{00})$, in the interbody region is small:

$$\Phi_{\text{interbody}} \sim (\text{Mass of system})/(\text{radius of system}) \ll 1.$$

The solar system certainly satisfies this condition ($\Phi_{\text{interbody}} \sim 10^{-7}$), as does the Galaxy ($\Phi_{\text{interbody}} \sim 10^{-6}$), as do clusters of galaxies ($\Phi_{\text{interbody}} \sim 10^{-6}$); *but the universe as a whole does not* ($\Phi_{\text{interbody}} \sim 1$)!

In evaluating volume integrals for the system’s total 4-momentum, split its volume into a region containing each body (denoted “ A ”) plus an interbody region; and neglect the pseudotensor contribution from the almost-flat interbody region:

$$\begin{aligned} P_{\text{system}}^\mu &= \sum_A \int_A T_{\text{eff}}^{\mu 0} d^3x + \int_{\text{interbody region}} T_{\text{eff}}^{\mu 0} d^3x \\ &= \sum_A P_A^\mu + \int_{\text{interbody region}} T^{\mu 0} d^3x. \end{aligned} \tag{20.24a}$$

Total 4-momentum and angular momentum for a system of gravitating bodies

Because spacetime is asymptotically flat around each body, P_A^μ is the 4-momentum of body A as measured gravitationally by an experimenter near it. The integral of $T^{\mu 0}$ over the interbody region is the contribution of any gas, particles, or magnetic

fields out there to the total 4-momentum. A similar breakup of the angular momentum reads

$$J_{\text{system}}^{\mu\nu} = \sum_A J_A^{\mu\nu} + \int_{\substack{\text{interbody} \\ \text{region}}} (x^\mu T^{\nu 0} - x^\nu T^{\mu 0}) d^3x. \quad (20.24b)$$

In operational terms, these breakups show that *the total 4-momentum and angular momentum of the system, as measured gravitationally by an experimenter outside it, are sums of P^μ and $J^{\mu\nu}$ for each individual body, as measured gravitationally by an experimenter near it, plus contributions of the usual special-relativistic type from the interbody matter and fields.* This is true even if some of the bodies are hurtling through the system with speeds near that of light; their gravitationally measured P^μ and $J^{\mu\nu}$ contribute, on an equal footing with anyone else's, to the system's total P^μ and $J^{\mu\nu}$!

Surround this asymptotically flat system by a two-dimensional surface S that is at rest in some asymptotic Lorentz frame. Then the 4-momentum and angular momentum inside S change at a rate (as measured in S 's rest frame) given by

$$\begin{aligned} \frac{dP^\mu}{dt} &= \frac{d}{dt} \int T_{\text{eff}}^{\mu 0} d^3x = \int T_{\text{eff},0}^{\mu 0} d^3x = - \int T_{\text{eff},j}^{\mu j} d^3x \\ &= - \oint T_{\text{eff}}^{\mu j} d^2S_j, \end{aligned} \quad (20.25)$$

and similarly

$$\frac{dJ^{\mu\nu}}{dt} = - \oint_{S_2} (x^\mu T_{\text{eff}}^{\nu j} - x^\nu T_{\text{eff}}^{\mu j}) d^2S_j. \quad (20.26)$$

Although the pseudotensor $t^{\mu\nu}$, in the interbody region and outside the system, contributes negligibly to the total 4-momentum and angular momentum (by assumption), its contribution via gravitational waves to the time derivatives dP^μ/dt and $dJ^{\mu\nu}/dt$ can be important when added up over astronomical periods of time. Thus, one must not ignore it in the flux integrals (20.25), (20.26).

In evaluating these flux integrals, it is especially convenient to use the Landau-Lifshitz form of $T_{\text{eff}}^{\mu\nu}$, since that form contains no second derivatives of the metric. Thus set

$$T_{\text{eff}}^{\mu\nu} = (-g)(T^{\mu\nu} + t_{\text{L-L}}^{\mu\nu}) \approx (T^{\mu\nu} + t_{\text{L-L}}^{\mu\nu}),$$

where $t_{\text{L-L}}^{\mu\nu}$ are given by equations (20.22). Only those portions of $t_{\text{L-L}}^{\mu\nu}$ that die out as $1/r^2$ or $1/r^3$ at large r can contribute to the flux integrals (20.25), (20.26). For static solutions [$g_{\mu\nu} \sim \text{const.} + O(1/r)$], $t_{\text{L-L}}^{\mu\nu}$ dies out as $1/r^4$. Hence, the only contributions come from dynamic parts of the metric, which, at these large distances, are entirely in the form of gravitational waves. The study of gravitational waves in Chapter 35 will reveal that when $t_{\text{L-L}}^{\mu\nu}$ is averaged over several wavelengths, it becomes a stress-energy tensor $T^{(\text{GW})\mu\nu}$ for the waves, which has all the properties one ever requires of any stress-energy tensor. (For example, via Einstein's equations

Rates of change of total 4-momentum and angular momentum:

(1) expressed as flux integrals of $T_{\text{eff}}^{\mu\nu}$

$G^{(B)\mu\nu} = 8\pi T^{(GW)\mu\nu}$, it contributes to the “background” curvature of the spacetime through which the waves propagate.) Moreover, averaging $t_{L-L}^{\mu\nu}$ over several wavelengths before evaluating the flux integrals (20.25), (20.26) cannot affect the values of the integrals. Therefore, one can freely make in these integrals the replacement

$$T_{\text{eff}}^{\mu\nu} = T^{\mu\nu} + T^{(GW)\mu\nu},$$

thereby obtaining

(2) expressed as flux integrals
of $T^{\mu\nu} + T^{(GW)\mu\nu}$

$$-\frac{dP^\mu}{dt} = \oint_S (T^{\mu j} + T^{(GW)\mu j}) d^2S_j, \quad (20.27)$$

$$-\frac{dJ^{\mu\nu}}{dt} = \oint_S [x^\mu(T^{\nu j} + T^{(GW)\nu j}) - x^\nu(T^{\mu j} + T^{(GW)\mu j})] d^2S_j. \quad (20.28)$$

These are tensor equations in the asymptotically flat spacetime surrounding the system. All reference to pseudotensors and other nontensorial entities has disappeared.

Equations (20.27) and (20.28) say that *the rate of loss of 4-momentum and angular momentum from the system, as measured gravitationally, is precisely equal to the rate at which matter, fields, and gravitational waves carry off 4-momentum and angular momentum.*

This theorem is extremely useful in thought experiments where one imagines changing the 4-momentum or angular momentum of a highly relativistic body (e.g., a rotating neutron star) by throwing particles onto it from far away [see, e.g., Hartle (1970)].

EXERCISE

Exercise 20.5. TOTAL MASS-ENERGY IN NEWTONIAN LIMIT

(a) Calculate $t_{L-L}^{\alpha\beta}$ for the nearly Newtonian metric

$$ds^2 = -(1 + 2\Phi) dt^2 + (1 - 2\Phi) \delta_{jk} dx^j dx^k$$

(see §18.4). Assume the source is slowly changing, so that time derivatives of Φ can be neglected compared to space derivatives. [Answer:

$$\begin{aligned} t_{L-L}^{00} &= -\frac{7}{8\pi} \Phi_{,j} \Phi_{,j}, \\ t_{L-L}^{0j} &= 0, \\ t_{L-L}^{jk} &= \frac{1}{4\pi} (\Phi_{,j} \Phi_{,k} - \frac{1}{2} \delta_{jk} \Phi_{,l} \Phi_{,l}). \end{aligned} \quad (20.29)$$

(Note: t_{L-L}^{jk} as given here is the “stress tensor for a Newtonian gravitational field”; cf. exercises 39.5 and 39.6.)

(b) Let the source of the gravitational field be a perfect fluid with

$$T^{\mu\nu} = (\rho + p) u^\mu u^\nu + pg^{\mu\nu}, \quad p/\rho \sim r^2 \equiv (dx/dt)^2 \sim |\Phi|.$$

Let the Newtonian potential satisfy the source equation

$$\Phi_{,jj} = 4\pi\rho.$$

Show that the energy of the source is

$$\begin{aligned} P^0 &= \int (T^{00} + t^{00})(-g) d^3x \\ &= \int [\underbrace{\rho/(1-v^2)^{1/2}}_{\text{Lorentz contraction factor}} + \frac{1}{2}\rho v^2 + \frac{1}{2}\rho\Phi] \underbrace{(g_{xx}g_{yy}g_{zz})^{1/2}}_{\text{proper volume}} dx dy dz \quad (20.30) \\ &\quad + \text{higher-order corrections.} \end{aligned}$$

- (c) Show that the “equations of motion” $T_{\text{L-Eff},\nu}^{\mu\nu} = 0$ reduce to the standard equations (16.3) of Newtonian hydrodynamics.
-

§20.6. EQUATIONS OF MOTION DERIVED FROM THE FIELD EQUATION

Consider the Einstein field equation

$$\mathbf{G} = 8\pi\mathbf{T} \quad (20.31)$$

under conditions where space is empty of everything except a source-free electromagnetic field:

$$T^{\mu\nu} = \frac{1}{4\pi} \left(F^{\mu\alpha} g_{\alpha\beta} F^{\nu\beta} - \frac{1}{4} g^{\mu\nu} F_{\sigma\tau} F^{\sigma\tau} \right) \quad (20.32)$$

(cf. the expression for stress-energy tensor of the electromagnetic field in §5.6). To predict from (20.31) how the geometry changes with time, one has to know how the electromagnetic field changes with time. The field is expressed as the “exterior derivative” of the 4-potential,

$$\mathbf{F} = \mathbf{dA} \text{ (language of forms)}$$

or

$$F_{\mu\nu} = \frac{\partial A_\nu}{\partial x^\mu} - \frac{\partial A_\mu}{\partial x^\nu} \text{ (language of components),} \quad (20.33)$$

and the time rate of change of the field is governed by the Maxwell equation

$$\mathbf{d}^* \mathbf{F} = 0$$

or

$$F^{\mu\nu}_{;\nu} = 0. \quad (20.34)$$

Vacuum Maxwell equations
derived from Einstein field
equation

If it seems a fair division of labor for the Maxwell equation to predict the development in time of the Maxwell field and the Einstein equation to do the same for the Einstein field, then it may come as a fresh surprise to discover that the Einstein equation (20.31), plus expression (20.32) for the Maxwell stress-energy, can do both jobs. One does not have to be given the Maxwell “equation of motion” (20.34). One can derive it fresh from (20.31) plus (20.32). The proof proceeds in five steps (see also exercise 3.18 and §5.10). Step one: The Bianchi identity $\nabla \cdot \mathbf{G} \equiv 0$ implies conservation of energy-momentum $\nabla \cdot \mathbf{T} = 0$. Step two: Conservation expresses itself in the language of components in the form

$$0 = 8\pi T^{\mu\nu}_{;\nu} = 2F^{\mu\alpha}_{;\nu}g_{\alpha\beta}F^{\nu\beta} + 2F^{\mu\alpha}g_{\alpha\beta}F^{\nu\beta}_{;\nu} - g^{\mu\nu}F_{\sigma\tau;\nu}F^{\sigma\tau}. \quad (20.35)$$

Step three: Leaving the middle term unchanged, rearrange the first term so that, like the last term, it carries a factor $F^{\sigma\tau}$. Thus in that first term let the indices $\nu\beta$ of $F^{\nu\beta}$ be replaced in turn by $\sigma\tau$ and by $\tau\sigma$, to subdivide that term into

$$\begin{aligned} & F^{\mu\alpha}_{;\sigma}g_{\alpha\tau}F^{\sigma\tau} + F^{\mu\alpha}_{;\tau}g_{\alpha\sigma}F^{\tau\sigma} \\ &= (F^\mu_{\tau;\sigma} - F^\mu_{\sigma;\tau})F^{\sigma\tau} \\ &= g^{\mu\nu}(F_{\nu\tau;\sigma} + F_{\sigma\nu;\tau})F^{\sigma\tau}. \end{aligned} \quad (20.36)$$

Step four: Combine the first and the last terms in (20.35) to give

$$g^{\mu\nu}(F_{\nu\tau;\sigma} + F_{\sigma\nu;\tau} + F_{\tau\sigma;\nu})F^{\sigma\tau}. \quad (20.37)$$

The indices on the derivatives of the field quantities stand in cyclic order. This circumstance annuls all the terms in the connection coefficients $\Gamma^\alpha_{\beta\gamma}$ when one writes out the covariant derivatives explicitly. Thus one can replace the covariant derivatives by ordinary derivatives. Moreover, these three derivatives annul one another identically when one substitutes for the fields their expressions (20.33) in terms of the potentials. Consequently, nothing remains in the conservation law (20.35) except the middle term, giving rise to four statements ($\mu = 0, 1, 2, 3$)

$$F^\mu_\beta F^{\beta\nu}_{;\nu} = 0 \quad (20.38)$$

about the four quantities ($\beta = 0, 1, 2, 3$)

$$F^{\beta\nu}_{;\nu}. \quad (20.39)$$

Step five: The determinant of the coefficients in the four equations (20.38) for the four unknowns (20.39) has the value

$$\begin{vmatrix} F^0_0 F^0_1 F^0_2 F^0_3 \\ \dots\dots\dots\dots \\ \dots\dots\dots\dots \\ F^3_0 F^3_1 F^3_2 F^3_3 \end{vmatrix} = -(\mathbf{E} \cdot \mathbf{B})^2 \quad (20.40)$$

(see exercise 20.6, part i). In the generic case, this one function of the four variables (t, x, y, z) vanishes on one or more hypersurfaces; but off any such hypersurface (i.e., at “normal points” in spacetime) it differs from zero. At all normal points, the solution of the four linear equations (20.38) with their nonvanishing determinant gives identically zero for the four unknowns (20.39); that is to say, Maxwell’s “equations of motion”

$$F^{\beta\nu}_{;\nu} = 0$$

are fulfilled and must be fulfilled as a straight consequence of Einstein’s field equation (20.31)—plus expression 20.32 for the stress-energy tensor. Special cases admit counterexamples (see exercise 20.8); but in the generic case one need not invoke Maxwell’s equations of motion; one can deduce them from the Einstein field equation.

Turn from the dynamics of the Maxwell field itself to the dynamics of a charged particle moving under the influence of the Maxwell field. Make no more appeal to outside providence for the Lorentz equation of motion than for the Maxwell equation of motion. Instead, to generate the Lorentz equation call once more on the Einstein field equation or, more directly, on its consequence, the principle of the local conservation of energy-momentum.

Keep track of the world line of the particle from $t = t$ to $t = t + \Delta t$ (Figure 20.1). Generate a “world tube” around this world line. Thus, at each value of the time coordinate t , take the location of the particle as center; construct a sphere of radius ϵ around this center; and note how the successive spheres sweep out the desired world tube. Construct “caps” on this tube at times t and $t + \Delta t$. The two caps, together with the world tube proper, bound a region of spacetime in which energy and momentum can be neither created nor destroyed (“no creation of moment of rotation,” in the language of the Bianchi identities, Chapter 15). Therefore the energy-momentum emerging out of the “top” cap has to equal the energy-momentum entering the “bottom” cap, supplemented by the amount of energy-momentum carried in across the world-tube by the Maxwell field. Out of such an analysis, as performed in flat spacetime, one ends up with the Lorentz equation of motion in its elementary form (see Chapters 3 and 4),

Lorentz force equation
derived from the Einstein
field equation

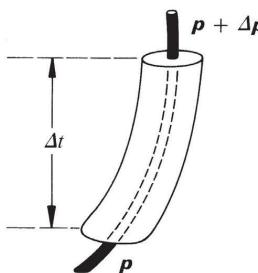


Figure 20.1.

“World tube.” The change in the 4-momentum of the particle is governed by the flow of 4-momentum across the boundary of the world tube.

$$d\mathbf{p}/d\tau = e\langle \mathbf{F}, \mathbf{u} \rangle \quad (\text{language of forms})$$

or in curved spacetime, the Lorentz equation of motion in covariant form,

$$\nabla_{\mathbf{u}}\mathbf{p} = m \nabla_{\mathbf{u}}\mathbf{u} = e\langle \mathbf{F}, \mathbf{u} \rangle \quad (\text{form language})$$

or

$$m \left[\frac{d^2x^\alpha}{d\tau^2} + \Gamma^\alpha_{\mu\nu} \frac{dx^\mu}{d\tau} \frac{dx^\nu}{d\tau} \right] = eF^\alpha_\beta \frac{dx^\beta}{d\tau} \quad (\text{component language}). \quad (20.41)$$

“One ends up with the Lorentz equation of motion”—but only after hurdling problems of principle along the way. One would understand what a particle is if one understood how to do the calculation of balance of energy-momentum with all rigor! Few calculations in all of physics have been done in so many ways by so many leading investigators, from Lorentz and his predecessors to Dirac and Rohrlich [see Teitelboim (1970, 1971) for still further insights]. Among the issues that develop are two that never cease to compel attention. (1) The particle responds according to the Lorentz force law (20.41) to a field. This field is the sum of a contribution from external sources and from the particle itself. How is the field exerted by the particle on itself to be calculated? Insofar as it is not already included in its effects in the “experimental mass” m in (20.41), this force is to be calculated as half the difference between the retarded field and the advanced field caused by that particle (see §36.11 for a more detailed discussion of the corresponding point for an emitter of gravitational radiation). This difference is singularity-free. On the world line, it has the following simple value [valid in general for point particles; valid for finite-sized particles when and only when the particle changes its velocity negligibly compared to the speed of light during the light-travel time across itself—see, e.g., Burke (1970)]

$$\frac{1}{2}(F_{\text{ret}} - F_{\text{adv}})^{\mu\nu} = \frac{2e}{3} \left(\frac{dx^\mu}{d\tau} \frac{d^3x^\nu}{d\tau^3} - \frac{d^3x^\mu}{d\tau^3} \frac{dx^\nu}{d\tau} \right). \quad (20.42)$$

A particle acted on by its own electromagnetic field (“radiation damping”)

Infinite self-energy of a point particle

Every acceptable line of reasoning has always led to expression (20.42). It also represents the field required to reproduce the long-known and thoroughly tested law of radiation damping. (2) “Infinite self-energy.” Around a particle at rest, or close to a particle in an arbitrary state of motion, the field is e/r^2 and the field energy is

$$(1/8\pi) \int_{r_{\min}}^{\epsilon} (e/r^2)^2 4\pi r^2 dr = (e^2/2)(r_{\min}^{-1} - \epsilon^{-1}). \quad (20.43)$$

This expression diverges as r_{\min} is allowed to go to zero. To hurdle this difficulty, one arranges the calculation of energy balance in such a way that there always appears the sum of this “self-energy” and the “bare mass.” The two terms individually are envisaged as “going to infinity” as r_{\min} goes to zero; but the sum is identified with the “experimental mass” and is required to remain finite. Of course, no particle is a classical object. A proper calculation of the energy has to be conducted at the quantum level. There it is easier to hide from sight the separate infinities—but they

are still present, and promise to remain until the structure of a particle is understood.

Before one turns from the Maxwell and Lorentz equations of motion to a final example (deriving the geodesic equations of motion for an uncharged particle), is it not time to object to the whole program of “deriving an equation of motion from Einstein’s field equation”? First, is it not a pretentious parade of pomposity to say it comes “from Einstein’s field equation” (and even more, “from Einstein’s field equations”) when it really comes from a principle so elementary and long established as the law of conservation of 4-momentum? It cannot be contested that this conservation principle, in historical fact, came before geometrodynamics, just as it came before electrodynamics and before the theories of all other established fields. However, in no theory but Einstein’s is this principle incorporated as an identity. Only here does the conservation of energy-momentum appear as a fully automatic consequence of the inner working of the machinery of the world (energy density tied to moment of rotation, and moment of rotation automatically conserved; see Chapter 17). Out of Einstein’s theory one can derive the equation of motion of a particle. Out of Maxwell’s one cannot. Thus, nothing prevents one from acting on a charge with an “external” force, over and above the Lorentz force, nor from tailoring this force in such a way that the charge follows some prescribed world line (“engine-driven source”). It makes no difficulties whatsoever for Maxwell’s equations that one has shifted attention from a world line that follows the Lorentz equation of motion to one that does not. Quite the contrary is true in general relativity. To shift from right world line (geodesic) to wrong world line makes the difference between satisfying Einstein’s field equation in the vicinity of that world line and being unable to satisfy Einstein’s field equation.

The Maxwell field equations are so constructed that they automatically fulfill and demand the conservation of charge; but not everything has charge. The Einstein field equation is so constructed that it automatically fulfills and demands the conservation of momentum-energy; and everything does have energy. The Maxwell field equations are indifferent to the interposition of an “external” force, because that force in no way threatens the principle of conservation of charge. The Einstein field equation cares about every force, because every force is a medium for the exchange of energy.

Electromagnetism has the motto, “I count all the electric charge that’s here.” All that bears no charge escapes its gaze.

“I weigh all that’s here” is the motto of spacetime curvature. No physical entity escapes this surveillance.

Why, then, is the derivation of the geodesic equation of motion of an object said to be based on “Einstein’s geometrodynamic field equation” rather than on “the principle of conservation of 4-momentum”? Because geometry responds by its curvature to mass-energy in every form. Most of all, because geometry outside tells about mass-energy inside, free of all concern about issues of internal structure (violent motions, unknown forces, tortuously curved and even multiply-connected geometry).

If one objection to the plan to derive the equation of motion of a particle “from the field equation” has been disposed of, then the moment has come to deal with

Why one is justified to regard equations of motion as consequences of the Einstein field equation

How one can avoid complexities of particle structure when deriving equations of motion: the “external viewpoint”

Derivation of geodesic motion from Einstein field equation:

(1) derivation in brief

(2) derivation with care

Coupling of curvature to particle moments produces deviations from geodesic motion

the other natural objection: Is there not an inner contradiction in trying to apply to a “particle” (implying idealization to a point) a field equation that deals with the continuum? Answer: There *is* a contradiction in dealing with a point. Therefore do not deal with a point. Do not deal with internal structure at all. Analyze the motion by looking at the geometry outside the object. That geometry provides all the handle one needs to follow the motion.

Already here one sees the difference from the derivation of the Lorentz equation of motion as sketched out above. There (1) no advantage was taken of geometry outside as indicator of motion inside; (2) a detailed bookkeeping was envisaged of the localization in space of the electromagnetic energy; and (3) this bookkeeping brought up the issue of the internal structure of the particle, which could not be satisfactorily resolved.

Now begin the analysis in the new geometrodynamical spirit. Surrounding “the Schwarzschild zone of influence” of the object, mark out a “buffer zone” (Figure 20.2) that extends out to the region where the “background geometry” begins to depart substantially from flatness. Idealize the geometry in the buffer zone as that of an unchanging source merging asymptotically (“boundary \mathcal{B} of buffer zone”) into flat space. It suffices to recall the properties of the spacetime geometry far outside an unchanging (i.e., nonradiating) source (exercise 19.3) to draw the key conclusion: relative to this flat spacetime and regardless of its internal structure, the object remains at rest, or continues to move in a straight line at uniform velocity (conservation of total 4-momentum; §20.5). In other words, it obeys the geodesic equation of motion. If this is the result in a flash, then it is appropriate to go back a step to review it, to find out what it means and what it demands.

When the object is absent and the background geometry alone has to be considered, then the geodesic is a well-defined mathematical construct. Moreover, Fermi-Walker transport along this geodesic gives a well-defined way to construct a comoving local inertial frame (see §13.6). Relative to this frame, the representative point of the geodesic remains for all time at rest at the origin.

In what way does the presence of the object change this picture? The object possesses an angular momentum, mass quadrupole moments, and higher multipole moments. They interact with the tide-producing accelerations (Riemann curvature) of the background geometry. Depending on the orientation in space of these moments, the interactions drive the object off its geodesic course in one direction or another (see §40.9). These anomalies in the motion go hand in hand with anomalies in the geometry. On and near the ideal mathematical geodesic the metric is Minkowskian. At a point removed from this geodesic by a displacement with Riemann normal coordinates ξ^1, ξ^2, ξ^3 (see §11.6), the metric components differ from their canonical values $(-1, 1, 1, 1)$ by amounts proportional (1) to the squares and products of the ξ^m and (2) to the components of the Riemann curvature tensor (tide-producing acceleration) of the background geometry. These second-order terms produce departures from ideality in the buffer zone, departures that may be described symbolically as of order

$$\delta(\text{metric}) \sim r^2 \cdot R \cdot (\text{spherical harmonic of order two}). \quad (20.44)$$

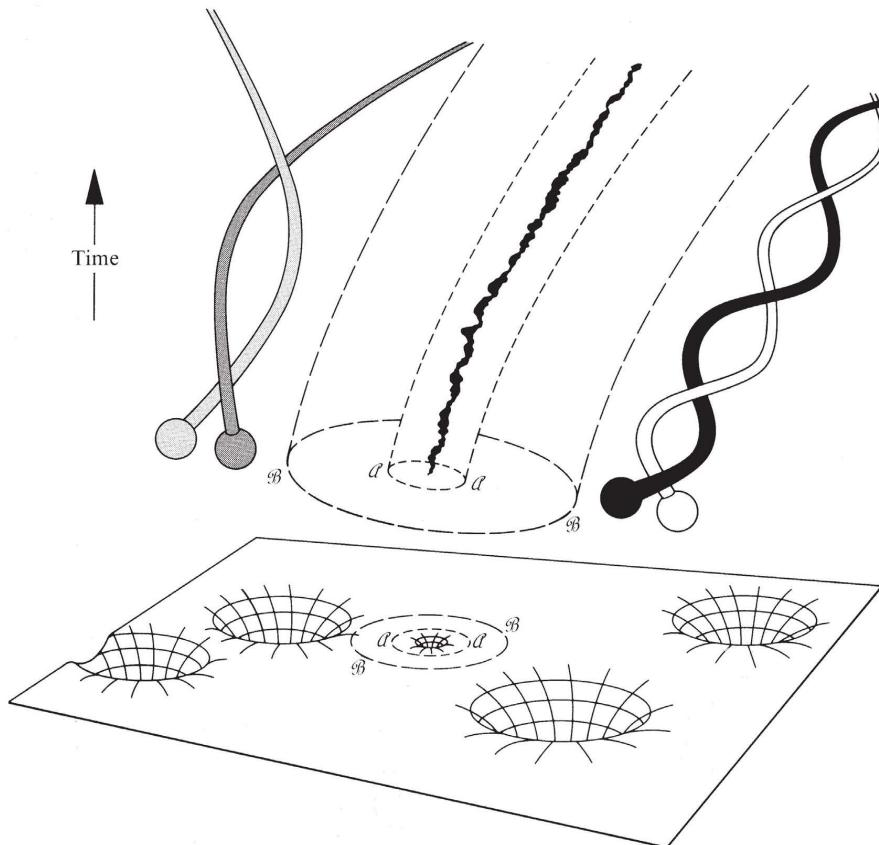


Figure 20.2.

“Buffer zone”: the shell of space between α and β , where the geometry is appropriately idealized as the spherically symmetric “Schwarzschild geometry” of a localized center of attraction (the object under study) in an asymptotically flat space. Inside α : the “zone of influence” of the object. In the general case where this object lacks all symmetry, the metric is found to depart more and more from ideal “Schwarzschild character” as the exploration of the geometry is carried inward from α (effect of angular momentum of the object on the metric; effect of quadrupole moment; effect of higher moments). Outside β : the “background geometry.” As this geometry is explored at greater and greater distances outside β , it is found to depart more and more from flatness (effect of concentrations of mass, gravitational waves, and other geometrodynamics).

Here r is the distance from the geodesic and R is the magnitude of the significant components of the curvature tensor. The object produces not only the standard “Schwarzschild” departure from flatness,

$$\delta(\text{metric}) \sim m/r, \quad (20.45)$$

which by itself (in a flat background) would bring about no departure from geodesic motion, but also correction terms which may be symbolized as

$$\delta(\text{metric}) \sim (S/r^2) \text{ (spherical harmonic of order one)} \quad (20.46)$$

and

$$\delta(\text{metric}) \sim (\mathcal{I}/r^3) \text{ (spherical harmonic of order two)} \quad (20.47)$$

and higher-order terms. Here $S(\text{cm}^2)$ is a typical component of the angular momentum vector or “spin”; $\mathcal{I}(\text{cm}^3)$ is a representative component of the moment of inertia or quadrupole tensor (see Chapter 36 for details), and higher terms have higher-order coefficients.

Coupling of spin to curvature

The tide-producing acceleration generated by the surroundings of the object (“background geometry”) acts on the spin of the object with a force of order RS and pulls it away from geodesic motion with an acceleration of the order

$$\text{acceleration (cm}^{-1}\text{)} \sim \frac{R(\text{cm}^{-2})S(\text{cm}^2)}{m(\text{cm})} \quad (20.48)$$

(see exercise 40.8). Otherwise stated, the surrounding and the spin both put warps in the geometry, and these warps conspire together to push the object off track.

The sum of the relevant two perturbations in the metric is qualitatively of the form

$$\delta g \sim r^2 R + S/r^2. \quad (20.49)$$

The sum is least where r has a value of the order

$$r \sim (S/R)^{1/4}, \quad (20.50)$$

and there it has the magnitude

$$\delta g \sim (SR)^{1/2}. \quad (20.51)$$

To “derive the geodesic equation of motion with some preassigned accuracy ϵ ” may be defined to mean that the metric in the buffer zone is Minkowskian within the latitude ϵ . In the illustrative example, this means that $(SR)^{1/2}$ is required to be of the order of ϵ or less. Nothing can be done about the value of R because the background curvature R is a feature of the background geometry. One can meet the requirement only by imposing limits on the mass and moments of the object. In the example, where the dominating moment is the angular momentum, one must require that this parameter of the object be less in order of magnitude than the limit

$$S \sim \epsilon^2/R. \quad (20.52)$$

Evidently this and similar conditions on the higher moments are most easily satisfied by demanding that the object have spherical symmetry ($S = 0$, $\mathcal{I} = 0$, higher

moments = 0). Then the perturbation in the metric, again disregarding angle factors and indices, is qualitatively of the form

$$\delta g \sim r^2 R + m/r, \quad (20.53)$$

and the buffer zone is best designed to bracket the minimizing value of r ,

$$r_{\alpha} \leq [r \sim (m/R)^{1/3}] \leq r_{\beta}. \quad (20.54)$$

The departure of the metric from Minkowskian perfection in the buffer zone is of the order

$$\delta g \sim (m^2 R)^{1/3}. \quad (20.55)$$

To achieve any preassigned accuracy ϵ for δg , one must demand that the mass be less than a limit of the order

$$m \sim \epsilon^{3/2}/R^{1/2}. \quad (20.56)$$

No object of finite mass moving under the influence of a complex background will admit a buffer zone where the geometry approaches Minkowskian values with arbitrary precision. Therefore it is incorrect to say that such an object follows a geodesic world line. It is meaningless to say that an object of finite rest mass follows a geodesic world line. World line of what? If the object is a black hole, there is no point inside its “horizon” (capture surface; one-way membrane; see Chapters 33 and 34) that is relevant to the physics going on outside. Geodesic world line within what background geometry? It has no sense to speak of a geometry that “lies behind” or is “background to” a black hole.

The sense in which no body can move on a geodesic of spacetime

Turn from one motion of one object in one spacetime to a continuous one-parameter family of spacetimes, with the mass m of the object being the parameter that distinguishes one of these solutions of Einstein’s field equation from another. Go to the limit $m = 0$. Then the size of the buffer zone shrinks to zero and the departure of the metric from Minkowskian perfection in the buffer zone also goes to zero. In this limit (“test particle”), it makes sense to say that the object moves in a straight line with uniform velocity in the local inertial frame or, otherwise stated, it pursues a geodesic in the background geometry. Moreover, this background geometry is well-defined: it is the limit of the spacetime geometry as the parameter m goes to zero [see Infeld and Schild (1949)]. In this sense, the geodesic equation of motion follows as an inescapable consequence of Einstein’s field equation.

The sense in which test particles do move on geodesics of a background geometry

The concept of “background” as limit of a one-parameter family of spacetimes extends itself to the case where the object bears charge as well as mass, and where the surrounding space is endowed with an electromagnetic field. This time the one-parameter family consists of solutions of the combined Einstein-Maxwell equations. The charge-to-mass ratio e/m is fixed. The mass m is again the adjustable parameter. In the limit when m goes to zero, one is left with (1) a background geometry, (2) a background electromagnetic field, and (3) a world line that obeys

Motion of a charged test particle in curved spacetime

the general-relativity version of the Lorentz equation of motion in this background as a consequence of the field equations [Chase (1954)]. In contrast, a so-called “unified field theory of gravitation and electromagnetism” that Einstein tentatively put forward at one stage of his thinking, as a conceivable alternative to the combination of his standard 1915 geometrodynamics with Maxwell’s standard electrodynamics, has been shown [Callaway (1953)] to lead to the wrong equation of motion for a charged particle. It moves as if uncharged no matter how much charge is piled on its back. If that theory were correct, no cyclotron could operate, no atom could exist, and life itself would be impossible.

Thus the ability to yield the correct equation of motion of a particle has today become an added ace in the hand of general relativity. The idea for such a treatment dates back to Einstein and Grommer (1927). Corrections to the geodesic equation of motion arising from interaction between the spin of the object (when it has finite dimensions) and the curvature of the background geometry are treated by Papapetrou (1951) and more completely by Pirani (1956) (see exercise 40.8). A book on the subject exists [Infeld and Plebanski (1960)]. Section 40.9 describes how corrections to geodesic motion affect lunar and planetary orbits. Some of the problems that arise when the object under study fragments or emits a directional stream of radiation, and unresolved issues of principle, are discussed by Wheeler (1961).

When one turns from the limit of infinitesimal mass to an object of finite mass, no simpler situation presents itself than a system of uncharged black holes (Chapter 33). Everything about the motion of these objects follows from an application of the source-free Einstein equation $\mathbf{G} = 0$ to the region of spacetime outside the horizons (see Chapter 34) of the several objects. The theory of motion is then geometrodynamics and nothing but geometrodynamics.

It has to be emphasized that all the considerations on motion in this section are carried out in the context of classical theory. In the real world of quantum physics, the geometry everywhere experiences unavoidable, natural, zero-point fluctuations (Chapter 43). The calculated local curvatures associated with these fluctuations at the Planck scale of distances [$L = (\hbar G/c^3)^{1/2} = 1.6 \times 10^{-33}$ cm] are enormous [$R \sim 1/L^2 \sim 0.4 \times 10^{66}$ cm $^{-2}$] compared to the curvature produced on much larger scales by any familiar object (electron or star). No detailed analysis of the interaction of these two curvatures has ever been made. Such an analysis would define a smoothed-out average of the geometry over regions larger than the local quantum fluctuations. With respect to this average geometry, the object will follow geodesic motion: this is the expectation that no one has ever seen any reason to question—but that no one has proved.

References on derivation of equations of motion from Einstein field equation

Quantum mechanical limitations on the derivation

EXERCISES

Exercise 20.6. SIMPLE FEATURES OF THE ELECTROMAGNETIC FIELD AND ITS STRESS-ENERGY TENSOR

- (a) Show that the “scalar” $-1/2F_{\alpha\beta}F^{\alpha\beta}$ (invariant with respect to coordinate transformations) and the “pseudoscalar” $1/4F_{\alpha\beta}{}^*F^{\alpha\beta}$ (reproduces itself under a coordinate transformation up to a \pm sign, according as the sign of the Jacobian of the transformation is positive

or negative) have in any local inertial frame the values $\mathbf{E}^2 - \mathbf{B}^2$ and $\mathbf{E} \cdot \mathbf{B}$, respectively ("the two Lorentz invariants" of the electromagnetic field).

(b) Show that the Poynting flux $(\mathbf{E} \times \mathbf{B})/4\pi$ is less in magnitude than the energy density $(\mathbf{E}^2 + \mathbf{B}^2)/8\pi$, save for the exceptional case where both Lorentz invariants of the field vanish (case where the field is locally "null").

(c) A charged pith ball is located a small distance from the North Pole of a bar magnet. Draw the pattern of electric and magnetic lines of force, indicating where the electromagnetic field is "null" in character. Is it legitimate to say that a "null field" is a "radiation field"?

(d) A plane wave is traveling in the z -direction. Show that the corresponding electromagnetic field is everywhere null.

(e) Show that the superposition of two monochromatic plane waves traveling in different directions is null on at most a set of points of measure zero.

(f) In the "generic case" where the field (\mathbf{E}, \mathbf{B}) at the point of interest is not null, show that the Poynting flux is reduced to zero by viewing the field from a local inertial frame that is traveling in the direction of $\mathbf{E} \times \mathbf{B}$ with a velocity

$$v = \tanh \alpha, \quad (20.57)$$

where the velocity parameter α is given by the formula

$$\tanh 2\alpha = \frac{(\text{Poynting flux})}{(\text{energy density})} = \frac{2|\mathbf{E} \times \mathbf{B}|}{\mathbf{E}^2 + \mathbf{B}^2}. \quad (20.58)$$

(g) Show that all components of the electric and magnetic field in this new frame can be taken to be zero except E_x and B_x .

(h) Show that the 4×4 determinant built out of the components of the field in mixed representation, F_α^β , is invariant with respect to general coordinate transformations. (Hint: Use the theorem that the determinant of the product of three matrices is equal to the product of the determinants of those three matrices.)

(i) Show that this determinant has the value $-(\mathbf{E} \cdot \mathbf{B})^2$ by evaluating it in the special local inertial frame of (f).

(j) Show that in this special frame the Maxwell stress-energy tensor has the form

$$\|T^\mu_\nu\| = \frac{E_x^2 + B_x^2}{8\pi} \begin{vmatrix} -1 & 0 & 0 & 0 \\ 0 & -1 & 0 & 0 \\ 0 & 0 & +1 & 0 \\ 0 & 0 & 0 & +1 \end{vmatrix} \quad (20.59)$$

(Faraday tension along the lines of force; Faraday pressure at right angles to the lines of force).

(k) In the other case, where the field is locally null, show that one can always find a local inertial frame in which the field has the form $\mathbf{E} = (0, F, 0)$, $\mathbf{B} = (0, 0, F)$ and the stress-energy tensor has the value

$$\|T^\mu_\nu\| = \frac{F^2}{4\pi} \begin{vmatrix} -1 & 1 & 0 & 0 \\ -1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{vmatrix} \quad (\mu \text{ for row, } \nu \text{ for column}). \quad (20.60)$$

(l) Regardless of whether the electromagnetic field is or is not null, show that the Maxwell stress-energy tensor has zero trace, $T^\mu_\mu = 0$, and that its square is a multiple of the unit tensor,

$$\begin{aligned} T^\mu_\alpha T^\alpha_\nu &= \frac{\delta^\mu_\nu}{(8\pi)^2} [(\mathbf{E}^2 - \mathbf{B}^2)^2 + (2\mathbf{E} \cdot \mathbf{B})^2] \\ &= \frac{\delta^\mu_\nu}{(8\pi)^2} [(\mathbf{E}^2 + \mathbf{B}^2)^2 - (2\mathbf{E} \times \mathbf{B})^2]. \end{aligned} \quad (20.61)$$

Exercise 20.7. THE STRESS-ENERGY TENSOR DETERMINES THE ELECTROMAGNETIC FIELD EXCEPT FOR ITS COMPLEXION

(a) Given a non-zero symmetric 4×4 tensor $T^{\mu\nu}$ which has zero trace $T^\mu_{\mu} = 0$ and whose square is a multiple, $M^4/(8\pi)^2$, of the unit matrix, show that, according as this multiple is zero (“null case”) or positive, the tensor can be transformed to the form (20.60) or (20.59) by a suitable rotation in 3-space or by a suitable choice of local inertial frame, respectively.

(b) In the generic (non-null) case in the frame in question, show that $T^{\mu\nu}$ is the Maxwell tensor of the “extremal electromagnetic field” $\xi_{\mu\nu}$ with components

$$\begin{aligned}\mathbf{E}^{(\text{extremal})} &= (M, 0, 0), \\ \mathbf{B}^{(\text{extremal})} &= (0, 0, 0).\end{aligned}\quad (20.62)$$

Show that it is also the Maxwell tensor of the “dual extremal field” $*\xi_{\mu\nu}$ with components

$$\begin{aligned}* \mathbf{E}^{(\text{extremal})} &= (0, 0, 0), \\ * \mathbf{B}^{(\text{extremal})} &= (M, 0, 0).\end{aligned}\quad (20.63)$$

(c) Recalling that the duality operation $*$ applied twice to an antisymmetric second-rank tensor (2-form) in four-dimensional space leads back to the negative of that tensor, show that the operator $e^{*\alpha}$ (“duality rotation”) has the value

$$e^{*\alpha} = (\cos \alpha) + (\sin \alpha)*. \quad (20.64)$$

(d) Show that the most general electromagnetic field which will reproduce the non-null tensor $T^{\mu\nu}$ in the frame in question, and therefore in any coordinate system, is

$$F_{\mu\nu} = e^{*\alpha} \xi_{\mu\nu}. \quad (20.65)$$

(e) Derive a corresponding result for the null case. [The field $F_{\mu\nu}$ defined in the one frame and therefore in every coordinate system by (d) and (e) is known as the “Maxwell square root” of $T^{\mu\nu}$; $\xi_{\mu\nu}$ is known as the “extremal Maxwell square root” of $T^{\mu\nu}$; and the angle α is called the “complexion of the electromagnetic field.” See Misner and Wheeler (1957); see also Boxes 20.1 and 20.2, adapted from that paper.]

Box 20.1 CONTRAST BETWEEN PROPER LORENTZ TRANSFORMATION AND DUALITY ROTATION

Quantity	General proper Lorentz transformation	Duality rotation
Components of the Maxwell stress-energy tensor or the “Maxwell square” of the field \mathbf{F}	Transformed	Unchanged
The invariants $\mathbf{E}^2 - \mathbf{B}^2$ and $(\mathbf{E} \cdot \mathbf{B})^2$	Unchanged	Transformed
The combination $[(\mathbf{E}^2 - \mathbf{B}^2)^2 + (2\mathbf{E} \cdot \mathbf{B})^2] = [(\mathbf{E}^2 + \mathbf{B}^2)^2 - (2\mathbf{E} \times \mathbf{B})^2]$	Unchanged	Unchanged

Box 20.2 TRANSFORMATION OF THE GENERIC (NON-NULL) ELECTROMAGNETIC FIELD TENSOR $F = (E, B)$ IN A LOCAL INERTIAL FRAME

Field values	At start	After simplifying duality rotation
At start	E, B	E and B perpendicular, and E greater than B
After simplifying Lorentz transformation	E and B parallel to each other and parallel to x -axis	E parallel to x -axis and $B = 0$

Exercise 20.8. THE MAXWELL EQUATIONS CANNOT BE DERIVED FROM THE LAW OF CONSERVATION OF STRESS-ENERGY WHEN $(E \cdot B) = 0$ OVER AN EXTENDED REGION

Supply a counter-example to the idea that the Maxwell equations,

$$F^{\mu\nu}_{;\nu} = 0,$$

follow from the Einstein equation; or, more precisely, show that (1) the condition that the Maxwell stress-energy tensor should have a vanishing divergence plus (2) the condition that this Maxwell field is the curl of a 4-potential A_μ can both be satisfied, while yet the stated Maxwell equations are violated. [Hint: It simplifies the analysis without obscuring the main point to consider the problem in the context of flat spacetime. Refer to the paper of Teitelboim (1970) for the decomposition of the retarded field of an arbitrarily accelerated charge into two parts, of which the second, there called $F^{\mu\nu}_{II}$, meets the stated requirements, and has everywhere off the worldline $(E \cdot B) = 0$, but does not satisfy the stated Maxwell equations.]

Exercise 20.9. EQUATION OF MOTION OF A SCALAR FIELD AS CONSEQUENCE OF THE EINSTEIN FIELD EQUATION

The stress-energy tensor of a massless scalar field is taken to be

$$T_{\mu\nu} = (1/4\pi)(\phi_{,\mu}\phi_{,\nu} - 1/2g_{\mu\nu}\phi_{,\alpha}\phi^{,\alpha}). \quad (20.66)$$

Derive the equation of motion of this scalar field from Einstein's field equation.

CHAPTER 21

VARIATIONAL PRINCIPLE AND INITIAL-VALUE DATA

Whenever any action occurs in nature, the quantity of action employed by this change is the least possible.

PIERRE MOREAU DE MAUPERTUIS (1746)

In the theory of gravitation, as in all other branches of theoretical physics, a mathematically correct statement of a problem must be determinate to the extent allowed by the nature of the problem; if possible, it must ensure the uniqueness of its solution.

VLADIMIR ALEXANDROVITCH FOCK (1959)

Things are as they are because they were as they were.

THOMAS GOLD (1972)

Calcalemus

G. W. LEIBNIZ

§21.1. DYNAMICS REQUIRES INITIAL-VALUE DATA

This chapter is entirely Track 2. No earlier Track-2 material is needed as preparation for it, but Chapters 9–11 and 13–15 will be helpful. It is needed as preparation for Box 30.1 (mixmaster universe) and for Chapters 42 and 43.

No plan for predicting the dynamics of geometry could be at the same time more mistaken and more right than this: “Give the distribution of mass-energy; then solve Einstein’s second-order equation,

$$\mathbf{G} = 8\pi\mathbf{T}, \quad (21.1)$$

for the geometry.” Give the distribution of mass-energy in spacetime and solve for the spacetime geometry? No. Give the fields that generate mass-energy, and their

To Karel Kuchař, Claudio Teitelboim, and James York go warm thanks for their collaboration in the preparation of this chapter, and for permission to draw on the lecture notes of K. K. and to quote results of K. K. [especially exercise 21.10] and of J. Y. [especially equations (21.87), (21.88), and (21.152)] prior to publication elsewhere.

time-rates of change, and give 3-geometry of space and its time-rate of change, all at one time, and solve for the 4-geometry of spacetime at that one time? Yes. And only then let one's equations for geometrodynamics and field dynamics go on to predict for all time, in and by themselves, needing no further prescriptions from outside (needing only work!), both the spacetime geometry and the flow of mass-energy throughout this spacetime. This, in brief, is the built-in "plan" of geometrodynamics, the plan spelled out in more detail in this chapter.

Contest the plan. Point out that the art of solving any coupled set of equations lies in separating the unknowns from what is known or to be prescribed. Insist that this separation is already made in (21.1). On the right already stands the source of curvature. On the left already stands the receptacle of curvature in the form of what one wants to know, the metric coefficients, twice differentiated. Claim therefore that one has nothing to do except to go ahead and solve these equations for the metric coefficients. However, in analyzing the structure of the equations to greater depth [see Cartan (1922a) for the rationale of analyzing a coupled set of partial differential equations], one discovers that one can only make the split between "the source and the receptacle" in the right way when one has first recognized the still more important split between "the initial-value data and the future." Thus—to summarize the results before doing the analysis—four of the ten components of Einstein's law connect the curvature of space here and now with the distribution of mass-energy here and now, and the other six equations tell how the geometry as thus determined then proceeds to evolve.

In determining what are appropriate initial-value data to give, one discovers no guide more useful than the Hilbert variational principle,

$$I = \int \mathcal{L} d^4x = \int L(-g)^{1/2} d^4x = \underset{\substack{\uparrow \\ [\text{exercise 8.16}]}}{\int L d(\text{proper 4-volume})} = \text{extremum} \quad (21.2)$$

or the Arnowitt-Deser-Misner ("ADM") variant of it (§21.6) and generalizations thereof by Kuchař (§21.9). Out of this principle one can recognize most directly what one must hold fixed at the limits (on an initial spacelike hypersurface and on a final spacelike hypersurface) as one varies the geometry (§21.2) throughout the spacetime "filling of this sandwich," if one is to have a well-defined extremum problem.

The Lagrange function L (scalar function) or the Lagrangian density $\mathcal{L} = (-g)^{1/2}L$ (quantity to be integrated over coordinate volume) is built of geometry alone, when one deals with curved empty space, but normally fields are present as well, and contribute also to the Lagrangian; thus,

$$\begin{aligned} \mathcal{L} &= \mathcal{L}_{\text{geom}} + \mathcal{L}_{\text{field}} = (-g)^{1/2}L; \\ L &= L_{\text{geom}} + L_{\text{field}}. \end{aligned} \quad (21.3)$$

The variation of the field Lagrangian with respect to the typical metric coefficient proves to be, of all ways, the one most convenient for generating (that is, for calculating) the corresponding component of the symmetric stress-energy tensor of the field (§21.3).

Give initial data, predict geometry

Four of ten components of Einstein equation are conditions on initial-value data

New view of stress-energy tensor

A computer, allowing for the effect of this field on the geometry and computing ahead from instant to instant the evolution of the metric with time, imposes its own ordering on the events of spacetime. In effect, it slices spacetime into a great number of spacelike slices. It finds it most convenient (§21.4) to do separate bookkeeping on (1) the 3-geometry of the individual slices and (2) the relation between one such slice and the next, as expressed in a “lapse function” N and a 3-vector “shift function” N_i .

The 3-geometry internal to the individual slice or “simultaneity” defines in and by itself the three-dimensional Riemannian curvature intrinsic to this hypersurface; but for a complete account of this hypersurface one must know also the extrinsic curvature (§21.5) telling how this hypersurface is curved with respect to the enveloping four-dimensional spacetime manifold.

In terms of the space-plus-time split of the 4-geometry, the action principle of Hilbert takes a simple and useful form (§21.6).

In the most elementary example of the application of an action principle in mechanics, where one writes

$$I = \int_{x', t'}^{x, t} L(dx/dt, x, t) dt \quad (21.4)$$

and extremizes the integral, one already knows that the resultant “dynamic path length” or “dynamic phase” or “action,”

$$S(x, t) = I_{\text{extremum}}, \quad (21.5)$$

is an important quantity, not least because it gives (up to a factor \hbar) the phase of the quantum-mechanical wave function. Moreover, the rate of change of this action function with position is what one calls momentum,

$$p = \partial S(x, t)/\partial x; \quad (21.6)$$

and the (negative of the) rate of change with time gives energy (Figure 21.1),

$$E = -\partial S(x, t)/\partial t; \quad (21.7)$$

and the relation between these two features of a system of wave crests,

$$E = H(p, x), \quad (21.8)$$

Hamiltonian as a dispersion relation

call it “dispersion relation” or call it what one will, is the central topic of mechanics.

When dealing with the dynamics of geometry in the Arnowitt-Deser-Misner formulation,* one finds it convenient to think of the specified quantities as being

* *Historical remark.* No one knew until recently what coordinate-free geometric-physical quantity *really* is fixed at limits in the Hilbert-Palatini variational principle. In his pioneering work on the Hamiltonian formulation of general relativity, Dirac paid no particular attention to any variational principle. He had to generalize the Hamiltonian formalism to accommodate it to general relativity, introducing “first- and second-class constraints” and generalizations of the Poisson brackets of classical mechanics. The work of Arnowitt, Deser, and Misner, by contrast, took the variational principle as the foundation for the whole treatment, even though they too did not ask what it is that is fixed at limits in the sense of

Figure 21.1.

Momentum and (the negative of the) energy viewed as rate of change of “dynamic phase” or “action,”

$$S(x, t) = I_{\text{extremum}}(x, t) = \left(\begin{array}{l} \text{extremum} \\ \text{value of} \end{array} \right) \int_{x', t'}^{x, t} L(x, \dot{x}, t) dt, \quad (1)$$

with respect to position and time; thus,

$$\delta S = p \delta x - E \delta t. \quad (2)$$

The variation of the integral I with respect to changes of the history along the way, $\delta x(t)$, is already zero by reason of the optimization of the history; so the only change that takes place is

$$\begin{aligned} \delta S &= \delta I_{\text{extremum}} = L(x, \dot{x}, t) \delta t + \int_{x', t'}^{x + \Delta x, t} \delta L dt \\ &= L \delta t + \int_{x', t'}^{x + \Delta x, t} \left(\frac{\partial L}{\partial \dot{x}} \delta \dot{x} + \frac{\partial L}{\partial x} \delta x \right) dt \\ &= L \delta t + \frac{\partial L}{\partial \dot{x}} \Delta x + \int_{x', t'}^{x + \Delta x, t} \left(\frac{\partial L}{\partial x} - \frac{d}{dt} \frac{\partial L}{\partial \dot{x}} \right) \delta x dt. \quad (3) \end{aligned}$$

[zero by reason
of extremization]

When one contemplates only a change δx in the coordinates (x, t) of the end point (change of history from $\mathcal{O}\mathcal{P}$ to $\mathcal{O}\mathcal{Q}$), one has $\Delta x = \delta x$. When one makes only a change δt in the end point (change of history from $\mathcal{O}\mathcal{P}$ to $\mathcal{O}\mathcal{S}$), one has $\Delta t = (\text{indicator of change from } \mathcal{P} \text{ to } \mathcal{R}) = -\dot{x} \delta t$. For the general variation of the final point, one thus has $\Delta x = \delta x - \dot{x} \delta t$ and

$$\delta S = \frac{\partial L}{\partial \dot{x}} \delta x - \left(\dot{x} \frac{\partial L}{\partial \dot{x}} - L \right) \delta t. \quad (4)$$

One concludes that the “dispersion relation” is obtained by taking the relations [compare (2) and (4)]

$$\left(\begin{array}{l} \text{rate of change of} \\ \text{dynamic phase} \\ \text{with position} \end{array} \right) = (\text{momentum}) = p = \frac{\partial L(x, \dot{x}, t)}{\partial \dot{x}} \quad (5)$$

and

$$-\left(\begin{array}{l} \text{rate of change of} \\ \text{dynamic phase} \\ \text{with time} \end{array} \right) = (\text{energy}) = E = \dot{x} \frac{\partial L}{\partial \dot{x}} - L, \quad (6)$$

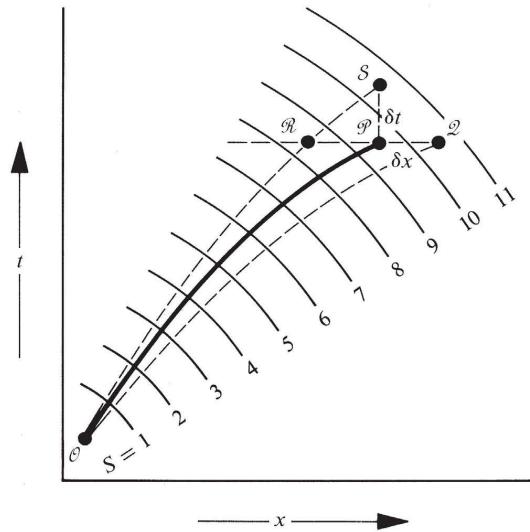
and eliminating \dot{x} from them [solve (5) for \dot{x} and substitute that value of \dot{x} into (6)]; thus

$$E = H(p, x, t) \quad (7)$$

or

$$-\frac{\partial S}{\partial t} = H\left(\frac{\partial S}{\partial x}, x, t\right). \quad (8)$$

Every feature of this elementary analysis has its analog in geometrodynamics.



a coordinate-free geometric-physical quantity. The great payoff of this work was recognition of the lapse and shift functions of equation (21.40) as Lagrange multipliers, the coefficients of which gave directly and simply Dirac's constraints. They did not succeed in arriving at a natural and simple time-coordinate, but that goal has in the meantime been achieved in the “extrinsic time” of Kuchar and York (§21.11). However, the Arnowitt-Deser Misner approach opened the door to the “intrinsic time” of Sharp, Baierlein, and Wheeler, where 3-geometry is fixed at limits, and 3-geometry is the carrier of information about time; and this led directly to Wheeler's “superspace version” of the treatment of Arnowitt, Deser, and Misner.

Action viewed as dependent
on 3-geometry

the 3-geometry $(^{(3)}\mathcal{G})$ of the initial spacelike hypersurface and the 3-geometry $(^{(3)}\mathcal{G})$ of the final spacelike hypersurface. One envisages the action integral as extremized with respect to the choice of the spacetime that fills the “sandwich” between these two faces. If one has thus determined the spacetime, one has automatically by that very act determined the separation in proper time of the two hypersurfaces. There is no additional time-variable to be brought in or considered. The one concept $(^{(3)}\mathcal{G})$ thus takes the place in geometrodynamics of the two quantities x, t of particle dynamics. The action S that there depended on x and t here depends on the 3-geometry of the face of the sandwich; thus,

$$S = S(^{(3)}\mathcal{G}). \quad (21.9)$$

A change in the 3-geometry changes the action. The amount of the change in action per elementary change in 3-geometry defines the “field momentum” π_{true}^{ij} conjugate to the geometrodynamic field coordinate g_{ij} , according to the formula

$$\delta S = \int \pi_{\text{true}}^{ij} \delta g_{ij} d^3x. \quad (21.10)$$

Comparing this equation out of the Arnowitt, Deser, and Misner (ADM) canonical formulation of geometrodynamics (§21.7) with the expression for change of action with change of endpoint in elementary mechanics,

$$\delta S = p \delta x - E \delta t, \quad (21.11)$$

one might at first think that something is awry, there being no obvious reference to time in (21.10). However, the 3-geometry is itself automatically the carrier of information about time; and (21.10) is complete. Moreover, with no “time” variable other than the information that $(^{(3)}\mathcal{G})$ itself already carries about time, there is also no “energy.” Thus the “dispersion relation” that connects the rates of change of action with respect to the several changes that one can make in the “field coordinates” or 3-geometry takes the form

$$\mathcal{K}(\pi^{ij}, g_{mn}) = 0, \quad (21.12)$$

Hamiltonian versus
super-Hamiltonian

with the E-term of (21.8) equal to zero (details in §21.7). All the content of Einstein’s general relativity can be extracted from this one Hamiltonian, or “super-Hamiltonian,” to give it a more appropriate name [see DeWitt (1967a), pp. 1113–1118, for an account of the contributions of Dirac, of Arnowitt, Deser, and Misner, and of others to the Hamiltonian formulation of geometrodynamics; and see §21.7 and subsequent sections of this chapter for the meaning and payoffs of this formulation].

The difference between a Hamiltonian and a super-Hamiltonian [see, for example, Kramers (1957)] shows nowhere more clearly than in the problem of a charged particle moving in flat space under the influence of the field derived from the electromagnetic 4-potential, $A_\mu(x^\alpha)$. The Hamiltonian treatment derives the equation of motion from the action principle,

$$0 = \delta I = \delta \int \left[p_i \frac{dx^i}{dt} - H(p_j, x^k, t) \right] dt$$

with

$$H = -\frac{e}{c}\phi + \left[m^2 + \eta^{ij} \left(p_i + \frac{e}{c}A_i\right) \left(p_j + \frac{e}{c}A_j\right)\right]^{1/2}.$$

The super-Hamiltonian analysis gets the equations of motion from the action principle

$$0 = \delta I' = \delta \int \left[p_\mu \frac{dx^\mu}{d\lambda} - \mathcal{K}(p_\alpha, x^\beta) \right] d\lambda.$$

Here the super-Hamiltonian is given by the expression

$$\mathcal{K}(p_\alpha, x^\beta) = \frac{1}{2} \left[m^2 + \eta^{\mu\nu} \left(p_\mu + \frac{e}{c}A_\mu\right) \left(p_\nu + \frac{e}{c}A_\nu\right) \right].$$

The variational principle gives Hamilton's equations for the rates of change

$$dx^\alpha/d\lambda = \partial \mathcal{K}/\partial p_\alpha$$

and

$$dp_\beta/d\lambda = -\partial \mathcal{K}/\partial x^\beta.$$

From these equations, one discovers that \mathcal{K} itself must be a constant, independent of the time-like parameter λ . The value of this constant has to be imposed as an initial condition, $\mathcal{K} = 0$ ("specification of particle mass"), thereafter maintained by the Hamiltonian equations themselves. This vanishing of \mathcal{K} in no way kills the partial derivatives,

$$\partial \mathcal{K}/\partial p_\alpha \quad \text{and} \quad -\partial \mathcal{K}/\partial x^\beta,$$

that enter Hamilton's equations for the rates of change,

$$dx^\alpha/d\lambda \quad \text{and} \quad dp_\beta/d\lambda.$$

Whether derived in the one formalism or the other, the equations of motion are equivalent, but the covariance shows more clearly in the formalism of the super-Hamiltonian, and similarly in general relativity.

Granted values of the "field coordinates" $g_{ij}(x, y, z)$ ^(3.2) and field momenta $\pi_{\text{true}}^{ij}(x, y, z) = \delta S/\delta g_{ij}$ compatible with (21.12), one has what are called "compatible initial-value data on an initial spacelike hypersurface." One can proceed as described in §21.8 to integrate ahead in time step by step from one spacelike hypersurface to another and another, and construct the whole 4-geometry. Here one is dealing with what in mathematical terminology are hyperbolic differential equations that have the character of a wave equation.

In contrast, one deals with elliptic differential equations that have the character of a Poisson potential equation when one undertakes in the first place to construct the needed initial-value data (§21.9). In the analysis of these elliptic equations, it

Another choice of what to fix at boundary hypersurface:
conformal part of 3-geometry
plus extrinsic time

Mach updated: mass-energy
there governs inertia here

proves helpful to distinguish in the 3-geometry between (1) the part of the metric that determines relative lengths at a point, which is to say angles (“the conformal part of the metric”) and (2) the common multiplicative factor that enters all the components of the g_{ij} at a point to determine the absolute scale of lengths at that point. This breakdown of the 3-geometry into two parts provides a particularly simple way to deal with two special initial-value problems known as the time-symmetric and time-antisymmetric initial-value problems (§21.10).

The ADM formalism is today in course of development as summarized in §21.11. In Wheeler’s (1968a) “superspace” form, the ADM treatment takes the 3-geometry to be fixed on each of the bounding spacelike hypersurfaces. In contrast, York (§21.11) goes back to the original Hilbert action principle, and discovers what it takes to be fixed on each of the bounding spacelike hypersurfaces. The appropriate data turn out to be the “conformal part of the 3-geometry” plus something closely related to what Kuchař (1971a and 1972) calls the “extrinsic time.” The contrast between Wheeler’s approach and the Kuchař-York approach shows particularly clearly when one (1) deals with a flat spacetime manifold, (2) takes a flat spacelike section through this spacetime, and then (3) introduces a slight bump on this slice, of height ϵ . The 3-geometry intrinsic to this deformed slice differs from Euclidean geometry only to the second order in ϵ . Therefore to read back from the full 3-geometry to the time (“the forward advance of the bump”) requires in this case an operation something like extracting a square root. In contrast, the Kuchař-York treatment deals with the “extrinsic curvature” of the slice, something proportional to the first power of ϵ , and therefore provides what is in some ways a more convenient measure of time [see especially Kuchař (1971) for the construction of “extrinsic time” for arbitrarily strong cylindrical gravitational waves; see also Box 30.1 on “time” as variously defined in “mixmaster cosmology”]. York shows that the time-variable is most conveniently identified with the variable “dynamically conjugate to the conformal factor in the 3-geometry.”

The initial-value problem of geometrodynamics can be formulated either in the language of Wheeler or in the language of Kuchař and York. In either formulation (§21.9 or §21.11) it throws light on what one ought properly today to understand by Mach’s principle (§21.12). That principle meant to Mach that the “acceleration” dealt with in Newtonian mechanics could have a meaning only if it was acceleration with respect to the fixed stars or to something equally well-defined. It guided Einstein to general relativity. Today it is summarized in the principle that “mass-energy there governs inertia here,” and is given mathematical expression in the initial-value equations.

The analysis of the initial-value problem connected past and future across a spacelike hypersurface. In contrast, one encounters a hypersurface that accommodates a timelike vector when one deals (§21.13) with the junction conditions between one solution of Einstein’s field equation (say, the Friedmann geometry interior to a spherical cloud of dust of uniform density) and another (say, the Schwarzschild geometry exterior to this cloud of dust). Section 21.13, and the chapter, terminate with notes on gravitational shock waves and the characteristic initial-value problem (the statement of initial-value data on a light cone, for example).

§21.2. THE HILBERT ACTION PRINCIPLE AND THE PALATINI METHOD OF VARIATION

Five days before Einstein presented his geometrodynamic law in its final and now standard form, Hilbert, animated by Einstein's earlier work, independently discovered (1915a) how to formulate this law as the consequence of the simplest action principle of the form (21.2–21.3) that one can imagine:

$$L_{\text{geom}} = (1/16\pi)^{(4)}R. \quad (21.13)$$

(Replace $1/16\pi$ by $c^3/16\pi G$ when going from the present geometric units to conventional units; or divide by $\hbar \sim L^{*2}$ to convert from dynamic phase, with the units of action, to actual phase of a wave function, with the units of radians). Here ${}^{(4)}R$ is the four-dimensional scalar curvature invariant, as spelled out in Box 8.4.

This action principle contains second derivatives of the metric coefficients. In contrast, the action principle for mechanics contains only first derivatives of the dynamic variables; and similarly only derivatives of the type $\partial A_\alpha / \partial x^\beta$ appear in the action principle for electrodynamics. Therefore one might also have expected only first derivatives, of the form $\partial g_{\mu\nu} / \partial x^\gamma$, in the action principle here. However, no scalar invariant lets itself be constructed out of these first derivatives. Thus, to be an invariant, L_{geom} has to have a value independent of the choice of coordinate system. But in the neighborhood of a point, one can always so choose a coordinate system that all first derivatives of the $g_{\mu\nu}$ vanish. Apart from a constant, there is no scalar invariant that can be built homogeneously out of the metric coefficients and their first derivatives.

When one turns from first derivatives to second derivatives, one has all twenty distinct components of the curvature tensor to work with. Expressed in a local inertial frame, these twenty components are arbitrary to the extent of the six parameters of a local Lorentz transformation. There are thus $20 - 6 = 14$ independent local features of the curvature ("curvature invariants") that are coordinate-independent, any one of which one could imagine employing in the action principle. However, ${}^{(4)}R$ is the only one of these 14 quantities that is linear in the second derivatives of the metric coefficients. Any choice of invariant other than Hilbert's complicates the geometrodynamic law, and destroys the simple correspondence with the Newtonian theory of gravity (Chapter 17).

Hilbert originally conceived of the independently adjustable functions of x, y, z, t in the variational principle as being the ten distinct components of the metric tensor in contravariant representation, $g^{\mu\nu}$. Later Palatini (1919) discovered a simpler and more instructive listing of the independently adjustable functions: not the ten $g^{\mu\nu}$ alone, but the ten $g^{\mu\nu}$ plus the forty $\Gamma_{\mu\nu}^\alpha$ of the affine connection.

To give up the standard formula for the connection Γ in terms of the metric g and let Γ "flap in the breeze" is not a new kind of enterprise in mathematical physics. Even in the simplest problem of mechanics, one can give up the standard formula for the momentum p in terms of a time-derivative of the coordinate x and also let

Variational principle the simplest route to Einstein's equation

Scalar curvature invariant the only natural choice

Idea of varying coordinate and momentum independently

p “flap in the breeze.” Then $x(t)$ and $p(t)$ become two independently adjustable functions in a new variational principle,

$$I = \int_{x', t'}^{x, t} \left[p(t) \frac{dx(t)}{dt} - H(p(t), x(t), t) \right] dt = \text{extremum.} \quad (21.14)$$

Happily, out of the extremization with respect to choice of the function $p(t)$, one recovers the standard formula for the momentum in terms of the velocity. The extremization with respect to choice of the other function, $x(t)$, gives the equation of motion just as does the more elementary variational analysis of Euler and Lagrange, where $x(t)$ is the sole adjustable function. A further analysis of this equivalence between the two kinds of variational principles in particle mechanics appears in Box 21.1. In that box, one also sees the two kinds of variational principle as applied to electrodynamics.

To express the Hilbert variational principle in terms of the $\Gamma_{\mu\nu}^\lambda$ and $g^{\alpha\beta}$ regarded as the primordial functions of t, x, y, z , note that the Lagrangian density is

$$L_{\text{geom}}(-g)^{1/2} = (1/16\pi)^{(4)}R(-g)^{1/2} = (1/16\pi)g^{\alpha\beta}R_{\alpha\beta}(-g)^{1/2}. \quad (21.15)$$

Here, as in any spacetime manifold with an affine connection, one has (Chapter 14)

$$R_{\alpha\beta} = R^\lambda_{\alpha\lambda\beta}, \quad (21.16)$$

where

$$R^\lambda_{\alpha\mu\beta} = \partial\Gamma^\lambda_{\alpha\beta}/\partial x^\mu - \partial\Gamma^\lambda_{\alpha\mu}/\partial x^\beta + \Gamma^\lambda_{\sigma\mu}\Gamma^\sigma_{\alpha\beta} - \Gamma^\lambda_{\sigma\beta}\Gamma^\sigma_{\alpha\mu}, \quad (21.17)$$

and every Γ is given in advance (in a coordinate frame) as symmetric in its two lower indices. In order that the integral I of (21.2–21.3) should be an extremum, one requires that the variation in I caused by changes both in the $g^{\mu\nu}$ and in the Γ 's should vanish; thus,

$$0 = \delta I = (1/16\pi) \int \delta[g^{\alpha\beta}R_{\alpha\beta}(-g)^{1/2}] d^4x + \int \delta[L_{\text{field}}(-g)^{1/2}] d^4x. \quad (21.18)$$

Consider now the variations of the individual factors in the first and second integrals in (21.18). The variation of the first factor is trivial, $\delta g^{\alpha\beta}$. In the variation of the second factor, $R_{\alpha\beta}$, changes in the $g^{\alpha\beta}$ play no part; only changes in the Γ 's appear. Moreover, the variation $\delta\Gamma_{\alpha\beta}^\lambda$ is a tensor even though $\Gamma_{\alpha\beta}^\lambda$ itself is not. Thus in the transformation formula

$$\Gamma^{\bar{\gamma}}_{\bar{\alpha}\bar{\beta}} = \left[\Gamma_{\sigma\tau}^\lambda \frac{\partial x^\sigma}{\partial x^{\bar{\alpha}}} \frac{\partial x^\tau}{\partial x^{\bar{\beta}}} + \frac{\partial^2 x^\lambda}{\partial x^{\bar{\alpha}} \partial x^{\bar{\beta}}} \right] \frac{\partial x^{\bar{\gamma}}}{\partial x^\lambda}, \quad (21.19)$$

Variation of connection is a tensor

the last term destroys the tensor character of any set of $\Gamma_{\sigma\tau}^\lambda$ individually, but subtracts out in the difference $\delta\Gamma_{\sigma\tau}^\lambda$ between two alternative sets of Γ 's. Note that the variation $\delta R^\lambda_{\alpha\mu\beta}$ of the typical component of the curvature tensor consists of two terms of

(continued on page 500)

Box 21.1 RATE OF CHANGE OF ACTION WITH DYNAMIC COORDINATE (= "MOMENTUM") AND WITH TIME, AND THE DISPERSION RELATION (= "HAMILTONIAN") THAT CONNECTS THEM IN PARTICLE MECHANICS AND IN ELECTRODYNAMICS

A. PROLOG ON THE PARTICLE-MECHANICS ANALOG OF THE PALATINI METHOD

In particle mechanics, one considers the history $x = x(t)$ to be adjustable between the end points (x', t') and (x, t) and varies it to extremize the integral $I = \int L(x, \dot{x}, t) dt$ taken between these two limits.

Expressed in terms of coordinates and momenta (see Figure 21.1), the integral has the form

$$I = \int [p\dot{x} - H(p, x, t)] dt, \quad (1)$$

where $x(t)$ is again the function to be varied and p is only an abbreviation for a certain function of x and \dot{x} ; thus, $p = \partial L(x, \dot{x}, t) / \partial \dot{x}$. Viewed in this way, the variation, $\delta p(t)$, of the momentum is governed by, and is only a reflection of, the variation $\delta x(t)$.

1. Momentum Treated as Independently Variable

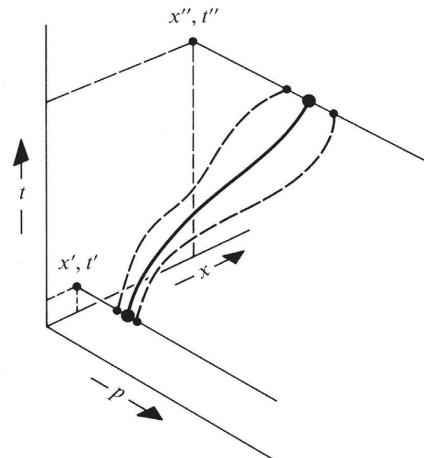
There miraculously exists, however, quite another way to view the problem (see inset). One can regard $x(t)$ and $p(t)$ as two quite uncorrelated and independently adjustable functions. One abandons the formula $p = \partial L(x, \dot{x}, t) / \partial \dot{x}$, only to recover it, or the equivalent of it, from the new "independent-coordinate-and-momentum version" of the variation principle.

The variation of (1), as defined and calculated in this new way, becomes

$$\delta I = p \delta x \Big|_{x', t'}^{x'', t''} + \int_{x', t'}^{x'', t''} \left[\left(\dot{x} - \frac{\partial H}{\partial p} \right) \delta p + \left(-\dot{p} - \frac{\partial H}{\partial x} \right) \delta x \right] dt. \quad (2)$$

Demand that the coefficient of δp vanish and have the sought-for new version,

$$\dot{x} = \frac{\partial H(p, x, t)}{\partial p}$$



Box 21.1 (continued)

of the old relation, $p = \partial L(x, \dot{x}, t) / \partial \dot{x}$, between momentum and velocity. The vanishing of the coefficient of δx gives the other Hamilton equation,

$$\dot{p} = - \frac{\partial H(p, x, t)}{\partial x}, \quad (3)$$

equivalent in content to the original Lagrange equation of motion,

$$\frac{d}{dt} \frac{\partial L}{\partial \dot{x}} - \frac{\partial L}{\partial x} = 0. \quad (4)$$

That $p(t)$ in this double variable conception is—before the extremization!—a function of time quite separate from and independent of the function $x(t)$ shows nowhere more clearly than in the circumstance that $p(t)$ has no end point conditions imposed on it, whereas x' and x'' are specified. Thus not only is the shape of the history subject to adjustment in x, p, t space in the course of achieving the extremum, but even the end points are subject to being slid along the two indicated lines in the inset, like beads on a wire.

2. Action as Tool for Finding Dispersion Relation

Denote by $S(x, t)$ the “action,” or extremal value of I , for the classical history that starts with (x', t') and ends at (x, t) ($= \hbar$ times phase of de Broglie wave). To change the end points to $(x + \delta x, t)$ makes the change in action

$$\delta S = p \delta x. \quad (5)$$

Thus momentum is “rate of change of action with dynamic coordinate.”

To change the end point to

$$(x + \delta x, t + \delta t) = ([x + \dot{x} \delta t] + [\delta x - \dot{x} \delta t], t + \delta t) \quad (6)$$

makes the change in action

$$\delta S = p[\delta x - \dot{x} \delta t] + L \delta t = p \delta x - H \delta t. \quad (7)$$

Thus the Hamiltonian is the negative of “the rate of change of action with time.”

In terms of the Hamiltonian $H = H(p, x)$, the “dispersion relation” for de Broglie waves becomes

$$-\frac{\partial S}{\partial t} = H\left(\frac{\partial S}{\partial x}, x\right). \quad (8)$$

In the derivation of this dispersion relation, one can profitably short-cut all talk of $p(t)$ and $x(t)$ as independently variable quantities, and derive the result in hardly

more than one step from the definition $I = \int L(x, \dot{x}, t) dt$. Similarly in electrodynamics.

The remainder of this box best follows a first perusal of Chapter 21.

B. ANALOG OF THE PALATINI METHOD IN ELECTRODYNAMICS

In source-free electrodynamics, one considers as given two spacelike hypersurfaces S' and S'' , and the magnetic fields-as-a-function-of-position in each, B' and B'' (this second field will later be written without the " superscript to simplify the notation). To be varied is an integral extended over the region of spacetime between the two hypersurfaces,

$$I_{\text{Maxwell}} \equiv \int \mathcal{L}_{\text{Maxwell}} d^4x = -\frac{1}{16\pi} \int F^{\mu\nu} F_{\mu\nu} (-g)^{1/2} d^4x. \quad (9)$$

1. Variation of Field on Hypersurface and Variation of Location of Hypersurface are Cleanly Separated Concepts in Electromagnetism

The electromagnetic field \mathbf{F} is the physically relevant quantity in electromagnetism (compare the 3-geometry in geometrodynamics). By contrast, the 4-potential \mathbf{A} has no direct physical significance. A change of gauge in the potentials,

$$A_\mu = A_{\mu_{\text{new}}} + \partial\lambda/\partial x^\mu$$

leaves unchanged the field components

$$F_{\mu\nu} = \partial A_\nu/\partial x^\mu - \partial A_\mu/\partial x^\nu$$

(compare the coordinate transformation that changes the $g_{\mu\nu}$ while leaving unchanged the $(^3)\mathcal{S}$). The variation of the fields within the body of the sandwich is nevertheless expressed most conveniently in terms of the effect of changes δA_μ in the potentials.

One also wants to see how the action integral is influenced by changes in the location of the upper spacelike hypersurface ("many-fingered time"). Think of the point of the hypersurface that is presently endowed with coordinates $x, y, z, t(x, y, z)$ as being displaced to $x, y, z, t + \delta t(x, y, z)$. Now renounce this use of a privileged coordinate system. Describe the displacement of the simultaneity in terms of a 4-vector δn (not a unit 4-vector) normal to the hypersurface Σ . The element of 4-volume $\delta\Omega$ included between the original upper face of the sandwich and the new upper face, that had in the privileged coordinate system the form $(-g)^{1/2} \delta t(x, y, z) d^3x$, in the notation of Chapter 20 becomes

$$\delta\Omega = \delta n^\mu d^3\Sigma_\mu = (\delta n \cdot d^3\Sigma), \quad (10)$$

where the element of surface $d^3\Sigma_\mu$ already includes the previously listed factor $(-g)^{1/2}$.

Box 21.1 (continued)

Counting together the influence of changes in the field values on the upper hypersurface and changes in the location of that hypersurface, one has

$$\begin{aligned} \delta S = \delta I_{\text{extremal}} &= -(1/16\pi) \int_{\text{upper } \Sigma} F^{\mu\nu} F_{\mu\nu} (\delta \mathbf{n} \cdot d^3 \boldsymbol{\Sigma}) \\ &\quad + (1/4\pi) \int_{\text{upper } \Sigma} F^{\mu\nu} \underbrace{\Delta A_\mu}_{\substack{\text{replace by} \\ \text{its equivalent} \\ (\delta A_\mu - \delta n^\alpha A_{\mu;\alpha})}} d^3 \boldsymbol{\Sigma}_\nu \quad (11) \\ &\quad + (1/4\pi) \int_{4\text{-volume}} \underbrace{F^{\mu\nu}_{;\nu} \delta A_\mu}_{\substack{\text{has to vanish} \\ \text{because integral has} \\ \text{been extremized}}} (-g)^{1/2} d^4x. \end{aligned}$$

Simplify this expression by arranging the coordinates so that the hypersurface shall be a hypersurface of constant t , and so that lines of constant x, y, z shall be normal to this hypersurface. Then it follows that the element of volume on that hypersurface contains a single nonvanishing component, $d^3 \boldsymbol{\Sigma}_0 = (-g)^{1/2} d^3x$. The antisymmetry of the field quantity $F^{0\nu}$ in its two indices requires that ν be a spacelike label, $i = 1, 2, 3$. The variation of the action becomes

$$\delta S = \int \left[\frac{(-g)^{1/2} F^{i0}}{4\pi} \delta A_i - \underbrace{\left\{ \frac{(-g)^{1/2} F^{i0}}{4\pi} A_{i;0} - \mathcal{L}_{\text{Maxwell}} \right\}}_{\substack{\text{add and subtract} \\ \left\{ \frac{(-g)^{1/2} F^{i0}}{4\pi} A_0 \delta t \right\}_{,i}}} \delta t \right] d^3x. \quad (12)$$

2. Meaning of Field "Momentum" in Electrodynamics

Identify this expression with the quantity

$$\delta S = \int \pi_{EM}^i \delta A_i d^3x - \int \mathcal{K} \delta \Omega, \quad (13)$$

where

$$\pi_{EM}^i = \frac{\delta S}{\delta A_i} = \left(\begin{array}{l} \text{"density of electromagnetic} \\ \text{momentum dynamically canon-} \\ \text{ically conjugate to } A_i \end{array} \right) = \frac{(-g)^{1/2} F^{i0}}{4\pi} = -\frac{\mathcal{E}^i}{4\pi} \quad (14)$$

is a simple multiple of the electric field and where

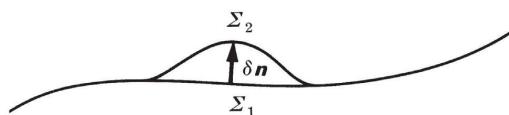
$$\mathcal{K} = -\frac{\delta S}{\delta \Omega} = \begin{pmatrix} \text{"density of} \\ \text{electromagnetic} \\ \text{Hamiltonian"} \end{pmatrix} = (1/16\pi)[F^{\mu\nu}F_{\mu\nu} + 4F^{i0}(A_{i;0} - A_{0;i})] \quad (15)$$

$$= (1/8\pi)(\mathbf{E}^2 + \mathbf{B}^2).$$

The concept of dynamic Hamiltonian density agrees with the usual concept of density of electromagnetic energy, despite the very different context in which the two quantities are derived and used. However, the canonical momentum π_{EM}^i has nothing directly whatsoever to do with the density of electromagnetic momentum as defined, for example, by the Poynting vector, despite the confusing similarity in the standard names for the two quantities. Note that there is no term δA_0 in (13); that is, $\pi_{EM}^0 \equiv 0$.

3. Bubble Differentiation

The “bubble differentiation” with respect to “many-fingered time” that appears in (15) was first introduced by Tomonaga (1946). One thinks of a spacelike hypersurface Σ_1 , a magnetic field \mathbf{B} defined as a function of position on this hypersurface (by an observer on a world line normal to this hypersurface), and a prescription S that carries one from this information to a single number, the action. (Divided by \hbar , this action gives the phase of the “wave function” or “probability amplitude” for the occurrence of this particular distribution of field values over this particular hypersurface.) One goes to a second hypersurface Σ_2 (see inset), which is identical with Σ_1 , except in the immediate vicinity of a given point. Take a distribution of field values over Σ_2 that is identical with the original distribution over Σ_1 , “identity of location” being defined by means of the normal. Evaluate the difference, δS , in the value of the dynamic phase or action in the two cases. Divide this difference by the amount of proper 4-volume $\delta\Omega = \int(\delta\mathbf{n} \cdot d^3\boldsymbol{\Sigma})$ contained in the “bubble” between the two hypersurfaces. Take the quotient, evaluate it in the limit in which the size of the bubble goes to zero, and in this way get the “bubble-time derivative,” $\delta S/\delta\Omega$, of the action.



Box 21.1 (continued)

What does it mean to say that the action, S , besides depending on the hypersurface, Σ , depends also on the distribution of the magnetic field, B , over that hypersurface? The action depends on the physical quantity, $\mathbf{B} = \nabla \times \mathbf{A}$, not on the prephysical quantity, \mathbf{A} . Thus a change in gauge $\delta A_i = \partial\lambda/\partial x^i$, cannot make any change in S . On the other hand, the calculated value of the change in S for this alteration in \mathbf{A} is

$$\begin{aligned}\delta(\text{action}) &= \delta S = \int \frac{\delta S}{\delta A_i} \delta A_i d^3x \\ &= \int \frac{\delta S}{\delta A_i} \frac{\partial \lambda}{\partial x^i} d^3x = - \int \left(\frac{\delta S}{\delta A_i} \right)_{,i} \lambda(x, y, z) d^3x.\end{aligned}\quad (16)$$

In order that there shall be no dependence of action on gauge, it follows that this expression must vanish for arbitrary $\lambda(x, y, z)$, a result only possible if $S(\Sigma, \mathbf{B}) = S(\text{hypersurface, field on hypersurface})$ satisfies the identity

$$\left(\frac{\delta S}{\delta A_i} \right)_{,i} = \pi_{EM,i}^i = -(1/4\pi)\mathcal{E}^i_{,i} = 0. \quad (17)$$

4. Hamilton-Jacobi "Propagation Law" for Electrodynamics

The "dispersion relation" or "Hamilton-Jacobi equation" for electromagnetism relates (1) the changes of the "dynamic phase" or "action" brought about by alterations in the dynamic variables A_i (the generalization of the x of particle dynamics) with (2) the changes brought about by alterations in many-fingered time (the generalization of the single time t of particle dynamics); thus (15) translates into

$$-\frac{\delta S}{\delta \Omega} = \frac{(4\pi)^2}{8\pi} \left(\frac{\delta S}{\delta \mathbf{A}} \right)^2 + \frac{1}{(8\pi)} (\nabla \times \mathbf{A})^2 \quad (18)$$

C. DISPERSION RELATIONS FOR GEOMETRODYNAMICS AND ELECTRODYNAMICS COMPARED AND CONTRASTED

Geometrodynamics possesses a direct analog of equation (17) ("action depends on no information carried by the vector potential \mathbf{A} except the magnetic field $\mathbf{B} = \nabla \times \mathbf{A}$ "), in an equation that says the action depends on no information carried by the metric g_{ij} on the "upper face of the sandwich" except the 3-geometry there, ⁽³⁾ \mathcal{G} . It also possesses a direct analog of equation (18) ("dynamic equation for the propagation of the action") with this one difference: in electrodynamics the field variable \mathbf{B} and the many-fingered time are distinct in character, whereas in geometrodynamics the "field" and the "many-fingered time" can be regarded as two aspects of one and the same ⁽³⁾ \mathcal{G} :

D. ACTION PRINCIPLE AND DISPERSION RELATION ARE ROOTED IN THE QUANTUM PRINCIPLE; FEYNMAN'S PRINCIPLE OF THE DEMOCRATIC EQUALITY OF ALL HISTORIES

For more on action principles in physics, see for example Mercier (1953), Lanczos (1970), and Yourgrau and Mandelstam (1968).

Newton (1687) in the first page of the preface to the first edition of his *Principia* notes that “The description of right lines . . . , upon which geometry is founded, belongs to mechanics. Geometry does not teach us to draw these lines, but requires them to be drawn.”

Newton’s remark is also a question. Mechanics moves a particle along a straight line, but what is the machinery by which mechanics accomplishes this miracle? The quantum principle gives the answer. The particle moves along the straight line only by not moving along the straight line. In effect it “feels out” every conceivable world line that leads from the start, (x', t') , to the point of detection, (x'', t'') , “compares” one with another, and takes the extremal world line. How does it accomplish this miracle?

The particle is governed by a “probability amplitude to transit from (x', t') to (x'', t'') .” This amplitude or “propagator,” $\langle x'', t'' | x', t' \rangle$, is the democratic sum with equal weight of contributions from every world line that leads from start to finish; thus,

$$\langle x'', t'' | x', t' \rangle = N \int e^{iI_H/\hbar} \mathcal{D}x. \quad (15)$$

Here N is a normalization factor, the same for all histories.

$\mathcal{D}x$ is the “volume element” for the sum over histories. For a “skeleton history” defined by giving x_n at $t_n = t_0 + n \Delta t$, one has $\mathcal{D}x$ equal, up to a multiplicative constant, to $dx_1 dx_2 \dots dx_N$. When the history is defined by the Fourier coefficients in such an expression as

$$x(t) = \frac{x'(t'' - t) + x''(t - t')}{(t'' - t')} + \sum_n a_n \sin n\pi \frac{(t - t')}{(t'' - t')}, \quad (16)$$

the volume element, again up to a multiplicative factor, is $da_1 da_2 \dots$

Destructive interference in effect wipes out the contribution to the transition probability from histories that differ significantly from the “extremal history” or “classical history.” Histories that are near that extremal history, on the other hand, contribute constructively, and for a simple reason: a small departure of the first order from the classical history brings about a change in phase which is only of the second order in the departure.

In this elementary example, one sees illustrated why it is that extremal principles play such a large part in classical dynamics. They remind one that all classical physics rests on a foundation of quantum physics. The central ideas are (1) the principle

Box 21.1 (continued)

of superposition of probability amplitudes, (2) constructive and destructive interference, (3) the “democracy of all histories,” and (4) the probability amplitude associated with a history H is $e^{iI_H/\hbar}$, apart from a normalizing factor that is a multiplicative constant.

For more on the democracy of histories and the sum over histories see Feynman (1942, 1948, 1949, 1951, and 1955), and the book of Feynman and Hibbs (1965); also Hibbs (1951), Morette (1951), Choquard (1955), Polkinghorne (1955), Fujiwara (1962), and the survey and literature references in Kursunoglu (1962); also reports of Dempster (1963) and Symanzik (1963). This outlook has been applied by many workers to discuss the quantum formulation of geometrodynamics, the first being Misner (1957) and one of the latest being Faddeev (1971).

the form $\delta\Gamma_{\alpha\beta,\mu}^\lambda$ and four terms of the form $\Gamma \delta\Gamma$ (indices being dropped for simplicity). One coordinate system is as good as another in dealing with a tensor. Therefore pick a coordinate system in which all the Γ 's vanish at the point under study. The terms $\Gamma \delta\Gamma$ drop out. In this coordinate system, the variation of the curvature is expressed in terms of first derivatives of quantities like $\delta\Gamma_{\alpha\beta}^\lambda$. One then need only replace the ordinary derivatives by covariant derivatives to obtain a formula correct in any coordinate system,

$$\delta R_{\alpha\mu\beta}^\lambda = \delta\Gamma_{\alpha\beta;\mu}^\lambda - \delta\Gamma_{\alpha\mu;\beta}^\lambda, \quad (21.20)$$

along with its contraction,

$$\delta R_{\alpha\beta} = \delta\Gamma_{\alpha\beta;\lambda}^\lambda - \delta\Gamma_{\alpha\lambda;\beta}^\lambda. \quad (21.21)$$

The third factor that appears in the variation principle is $(-g)^{1/2}$. Its variation (exercise 21.1) is

$$\delta(-g)^{1/2} = -\frac{1}{2} (-g)^{1/2} g_{\mu\nu} \delta g^{\mu\nu}. \quad (21.22)$$

The other integrand, the Lagrange density L_{field} , will depend on the fields present and their derivatives, but will be assumed to contain the metric only as $g^{\mu\nu}$ itself, never in the form of any derivatives of $g^{\mu\nu}$.

In order for an extremum to exist, the following expression has to vanish:

$$(1/16\pi) \int \left[\left(R_{\alpha\beta} - \frac{1}{2} g_{\alpha\beta} R \right) \delta g^{\alpha\beta} + g^{\alpha\beta} (\delta\Gamma_{\alpha\beta;\lambda}^\lambda - \delta\Gamma_{\alpha\lambda;\beta}^\lambda) \right] (-g)^{1/2} d^4x \\ + \int \left(\frac{\delta L_{\text{field}}}{\delta g^{\alpha\beta}} - \frac{1}{2} g_{\alpha\beta} L_{\text{field}} \right) \delta g^{\alpha\beta} (-g)^{1/2} d^4x = 0 \quad (21.23)$$

Focus attention on the term in (21.23) that contains the variations of Γ ,

$$(1/16\pi) \int g^{\alpha\beta} (\delta\Gamma_{\alpha\beta;\lambda}^\lambda - \delta\Gamma_{\alpha\lambda;\beta}^\lambda) (-g)^{1/2} d^4x,$$

and integrate by parts to eliminate the derivatives of the $\delta\Gamma$. To prepare the way for this integration, introduce the concept of *tensor density*, a notational device widely applied in general relativity. The concept of tensor density aims at economy. Without this concept, one will treat the tensor

$$\epsilon_{\mu\alpha\beta\gamma} = (-g)^{1/2} [\mu\alpha\beta\gamma]$$

(see exercise 3.13) as having $4^4 = 256$ components, and its covariant derivative as having $4^5 = 1,024$ components, of which one is

$$\begin{aligned}\epsilon_{0123;\rho} &= \partial(-g)^{1/2}/\partial x^\rho \epsilon_{[0123]} - \Gamma_{0\rho}^\sigma \epsilon_{\sigma 123} - \Gamma_{1\rho}^\sigma \epsilon_{0\sigma 23} \\ &\quad - \Gamma_{2\rho}^\sigma \epsilon_{01\sigma 3} - \Gamma_{3\rho}^\sigma \epsilon_{012\sigma} \\ &= [(-g)^{1/2}_{,\rho} - \Gamma_{\sigma\rho}^\sigma (-g)^{1/2}] [0123].\end{aligned}$$

Concept of tensor density

The symbol $[\alpha\beta\gamma\delta]$, with values $(0, -1, +1)$, introduces what is largely excess baggage, doing mere bookkeeping on alternating indices. Drop this unhandiness. Introduce instead the non-tensor $(-g)^{1/2}$ and *define* for it the law of covariant differentiation,

$$(-g)^{1/2}_{,\rho} = (-g)^{1/2}_{,\rho} - \Gamma_{\sigma\rho}^\sigma (-g)^{1/2}. \quad (21.24)$$

These four components take the place of the 1,024 components and communicate all the important information that was in them.

Associated with the vector j_μ is the vector density

$$j_\mu = (-g)^{1/2} j_\mu;$$

with the tensor $T_{\mu\nu}$, the tensor density

$$\mathfrak{T}_{\mu\nu} = (-g)^{1/2} T_{\mu\nu};$$

and so on; the German gothic letter is a standard indicator for the presence of the factor $(-g)^{1/2}$. On some occasions (see, for example, §21.11) it is convenient to multiply the components of a tensor with a power of $(-g)^{1/2}$ other than 1. According to the value of the exponent, the resulting assemblage of components is then called a tensor density of this or that *weight*.

The law of differentiation of an ordinary or standard tensor density formed from a tensor of arbitrary order,

$$\mathfrak{A}^{\cdot\cdot\cdot} = (-g)^{1/2} A^{\cdot\cdot\cdot},$$

is

$$(\mathfrak{A}^{\cdot\cdot\cdot})_{;\rho} = (\mathfrak{A}^{\cdot\cdot\cdot})_{,\rho} + (\text{standard } \Gamma^{\cdot\cdot\cdot} \text{ terms of a standard covariant derivative multiplied into } \mathfrak{A}^{\cdot\cdot\cdot}) - (\mathfrak{A}^{\cdot\cdot\cdot}) \Gamma_{\sigma\rho}^\sigma.$$

The covariant derivative of a product is the sum of two terms: the covariant deriva-

tive of the first, times the second, plus the first times the covariant derivative of the second.

Now return to the integral to be evaluated. Combine the factors $g^{\alpha\beta}$ and $(-g)^{1/2}$ into the tensor density $\mathbf{g}^{\alpha\beta}$. Integrate covariantly by parts, as justified by the rule for the covariant derivative of a product. Get a “term at limits,” plus the integral

$$-(1/16\pi) \int (\mathbf{g}^{\alpha\beta}_{;\lambda} - \delta_\lambda^\beta \mathbf{g}^{\alpha\gamma}_{;\gamma}) \delta\Gamma_{\alpha\beta}^\lambda d^4x.$$

This integral is the only term in the action integral that contains the variations of the Γ 's at the “interior points” of interest here. For the integral to be an extremum, the symmetrized coefficient of $\delta\Gamma_{\alpha\beta}^\lambda$ must vanish,

$$\mathbf{g}^{\alpha\beta}_{;\lambda} - \frac{1}{2} \delta_\lambda^\alpha \mathbf{g}^{\beta\gamma}_{;\gamma} - \frac{1}{2} \delta_\lambda^\beta \mathbf{g}^{\alpha\gamma}_{;\gamma} = 0.$$

This set of forty equations for the forty covariant derivative $\mathbf{g}^{\alpha\beta}_{;\lambda}$ has only the zero solution,

$$\mathbf{g}^{\alpha\beta}_{;\lambda} = 0. \quad (21.25)$$

Thus the “density formed from the reciprocal metric tensor” is covariantly constant.

This simple result (1) brings many simple results in its train: the covariant constancy of (2) $(-g)^{1/2}$, (3) $\mathbf{g}^{\alpha\beta}$, (4) $g_{\alpha\beta}$, and (5) $\mathbf{g}_{\alpha\beta}$. Of these, (4) is of special interest here, and (2) is needed in proving it, as follows. Take definition (21.24) for the covariant derivative of $(-g)^{1/2}$, and calculate the ordinary derivative that appears in the first term from exercise 21.1. One encounters in this calculation terms of the form $\partial\mathbf{g}^{\alpha\beta}/\partial x^\lambda$. Use (21.25) to evaluate them, and end up with the result

$$(-g)^{1/2}_{;\lambda} = 0.$$

From this result it follows that the covariant derivative of the $(\frac{1}{2})$ -tensor density $(-g)^{1/2} \delta_\gamma^\alpha$ is also zero. But this tensor density is the product of the tensor density $\mathbf{g}^{\alpha\beta}$ by the ordinary metric tensor $g_{\beta\gamma}$. In the covariant derivative of this product by x^λ , one already knows that the derivative of the first factor is zero. Therefore the first factor times the derivative of the second must be zero,

$$\mathbf{g}^{\alpha\beta} g_{\beta\gamma;\lambda} = 0,$$

and from this it follows that

$$g_{\beta\gamma;\lambda} = 0, \quad (21.26)$$

as was to be proven; or, explicitly,

$$\frac{\partial g_{\beta\gamma}}{\partial x^\lambda} - g_{\gamma\sigma} \Gamma_{\beta\lambda}^\sigma - g_{\beta\sigma} \Gamma_{\gamma\lambda}^\sigma = 0.$$

Solve these equations for the Γ 's, which up to now have been independent of the $g_{\beta\gamma}$, and end up with the standard equation for the connection coefficients,

$$\Gamma_{\mu\nu}^\rho = \frac{1}{2} g^{\rho\sigma} (g_{\mu\sigma,\nu} + g_{\sigma\nu,\mu} - g_{\mu\nu,\sigma}), \quad (21.27)$$

as required for Riemannian geometry.

Similarly, equate to zero the coefficient of $\delta g^{\alpha\beta}$ in the variation (21.23), and find all ten components of Einstein's field equation, in the form

$$G_{\alpha\beta} = 8\pi \underbrace{\left(g_{\alpha\beta} L_{\text{field}} - 2 \frac{\delta L_{\text{field}}}{\delta g^{\alpha\beta}} \right)}_{\substack{\uparrow \text{identified in §21.3 with} \\ \text{the stress-energy tensor } T_{\alpha\beta}}} \quad (21.28)$$

Among variations of the metric, one of the simplest is the change

$$g_{\text{new}\mu\nu} = g_{\mu\nu} + \delta g_{\mu\nu} = g_{\mu\nu} + \xi_{\mu;\nu} + \xi_{\nu;\mu} \quad (21.29)$$

brought about by the infinitesimal coordinate transformation

$$x_{\text{new}}^\mu = x^\mu - \xi^\mu. \quad (21.30)$$

Although the metric changes, the 3-geometry does not. It does not matter whether the spacetime geometry that one is dealing with extremizes the action principle or not, whether it is a solution of Einstein's equations or not; the action integral I is a scalar invariant, a number, the value of which depends on the physics but not at all on the system of coordinates in which that physics is expressed. This invariance even obtains for both parts of the action principle individually (I_{geom} and I_{fields}). Therefore neither part will be affected in value by the variation (21.29). In other words, the quantity

$$\delta I_{\text{geom}} = (1/16\pi) \int G_{\alpha\beta} (\xi^{\alpha;\beta} + \xi^{\beta;\alpha}) (-g)^{1/2} d^4x \underset{\substack{\uparrow \\ \text{“covariant integration by parts”}}}{=} -(1/8\pi) \int G_{\alpha\beta} ;^\beta \xi^\alpha (-g)^{1/2} d^4x \quad (21.31)$$

Action unaffected by mere change in coordinatization

must vanish whatever the 4-geometry and whatever the change ξ^α . In this way, one sees from a new angle the contracted Bianchi identities of Chapter 15,

$$G_{\alpha\beta} ;^\beta = 0. \quad (21.32)$$

The “neutrality” of the action principle with respect to a mere coordinate transformation such as (21.29) shows once again that the variational principle—and with it Einstein's equation—cannot determine the coordinates or the metric, but only the 4-geometry itself.

Exercise 21.1. VARIATION OF THE DETERMINANT OF THE METRIC TENSOR

EXERCISE

Recalling that the change in the value of any determinant is given by multiplying the change in each element of that determinant by its cofactor and adding the resulting products (exercise 5.5) prove that

$$\delta(-g)^{1/2} = \frac{1}{2} (-g)^{1/2} g^{\mu\nu} \delta g_{\mu\nu} \quad \text{and} \quad \delta(-g)^{1/2} = -\frac{1}{2} (-g)^{1/2} g_{\mu\nu} \delta g^{\mu\nu}.$$

Also show that

$$g = \det \|g^{\mu\nu}\| \quad \text{and} \quad \delta(-g)^{1/2} = +\frac{1}{2} g_{\mu\nu} \delta g^{\mu\nu}.$$

§21.3. MATTER LAGRANGIAN AND STRESS-ENERGY TENSOR

The derivation of Einstein's geometrodynamical law from Hilbert's action principle puts on the righthand side a source term that is derived from the field Lagrangian. In contrast, the derivation of Chapter 17 identified the source term with the stress-energy tensor of the field. For the two derivations to be compatible, the stress-energy tensor must be given by the expression

Lagrangian generates
stress-energy tensor

$$T_{\alpha\beta} = -2 \frac{\delta L_{\text{field}}}{\delta g^{\alpha\beta}} + g_{\alpha\beta} L_{\text{field}}, \quad (21.33a)$$

or

$$(-g)^{1/2} T^{\alpha\beta} \equiv \mathcal{T}^{\alpha\beta} = 2 \frac{\delta \mathcal{L}_{\text{field}}}{\delta g_{\alpha\beta}}. \quad (21.33b)$$

What are the consequences of this identification?

By the term "Lagrange function of the field" as employed here, one means the Lagrange function of the classical theory as formulated in flat spacetime, with the flat-spacetime metric replaced wherever it appears by the actual metric, and with the "comma-goes-to-semicolon rule" of Chapter 16 applied to all derivatives.

Were one dealing with a general tensorial field, the comma-goes-to-semicolon rule would introduce, in addition to the derivative of the tensorial field with all its indices, a number of Γ 's equal to the number of indices. The presence of these Γ 's in the field Lagrangian would have unhappy consequences for the Palatini variational procedure described in §21.2. No longer would the Γ 's end up given in terms of the metric coefficients by the standard formula (21.27). No longer would the geometry, as derived from the Hilbert-Palatini variation principle, be Riemannian. Then what?

These troublesome issues do not arise in two well-known simple cases, a scalar field and an electromagnetic field. In the one case, the field Lagrangian becomes

$$L_{\text{field}} = (1/8\pi)[-g^{\alpha\beta}(\partial\phi/\partial x^\alpha)(\partial\phi/\partial x^\beta) - m^2\phi^2]. \quad (21.34)$$

Electromagnetism as an example

No connection coefficient comes in; the quantity being differentiated is a scalar. In the other case, the field Lagrangian is built on first derivatives of the 4-potential A_μ . Therefore Γ 's should appear, according to the standard rules for covariant differentiation (Box 8.4). However, the derivatives of the A 's appear, never alone, but always in an antisymmetric combination where the Γ 's cancel, making covariant derivatives equivalent to ordinary derivatives:

$$F_{\mu\nu} = A_{\nu;\mu} - A_{\mu;\nu} = A_{\nu,\mu} - A_{\mu,\nu}. \quad (21.35)$$

Contrast to stress-energy tensor of "canonical field theory"

In both cases, the differentiations of (21.33) to generate the stress-energy tensor are easily carried out (exercises 21.2 and 21.3) and give the standard expressions already seen [(5.22) and (5.23)] for $T_{\mu\nu}$ in one of these two cases in an earlier chapter.

Field theory provides a quite other method to generate a so-called canonical expression for the stress-energy tensor of a field [see, for example, Wentzel (1949)].

By the very manner of construction, such an expression is guaranteed also to satisfy the law of conservation of momentum and energy, and by this circumstance it too becomes useful in certain contexts. However, the canonical tensor is often not symmetric in its two indices, and in such cases violates the law of conservation of angular momentum (see discussion in §5.7). Even when symmetric, it may give a quite different localization of stress and energy than that given by (21.33). Field theory in and by itself is unable to decide between these different pictures of where the field energy is localized. However, direct measurements of the pull of gravitation provide in principle [see, for example, Feynman (1964)] a means to distinguish between alternative prescriptions for the localization of stress-energy, because gravitation responds directly to density of mass-energy and momentum. It is therefore a happy circumstance that the theory of gravity in the variational formulation gives a unique prescription for fixing the stress-energy tensor, a prescription that, besides being symmetric, also automatically satisfies the laws of conservation of momentum and energy (exercises 21.2 and 21.3). [For an early discussion of the symmetrization of the stress-energy tensor, see Rosenfeld (1940) and Belinfante (1940). A more extensive discussion is given by Corson (1953) and Davis (1970), along with extensive references to the literature.]

When one deals with a spinor field, one finds it convenient to take as the quantities to be varied, not the metric coefficients themselves, but the components of a tetrad of orthonormal vectors defined as a tetrad field over all space [see Davis (1970) for discussion and references].

Exercise 21.2. STRESS-ENERGY TENSOR FOR A SCALAR FIELD

Given the Lagrange function (21.34) of a scalar field, derive the stress-energy tensor for this field. Also write down the field equation for the scalar field that one derives from this Lagrange function (in the general case where the field executes its dynamics within the arena of a curved spacetime). Show that as a consequence of this field equation, the stress-energy tensor satisfies the conservation law, $T_{\alpha\beta}^{;\beta} = 0$.

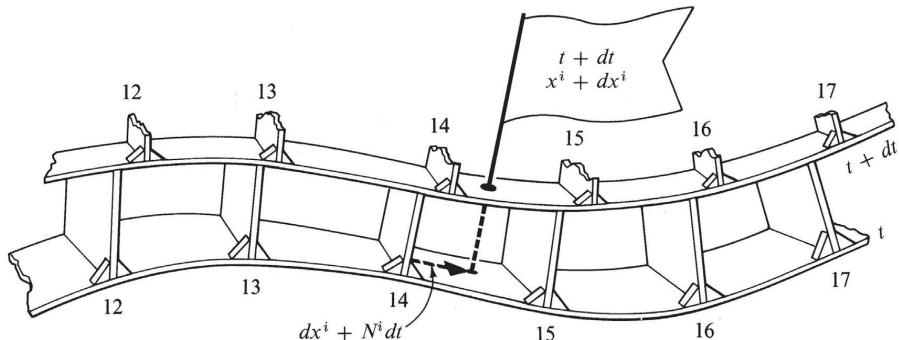
Exercise 21.3. FARADAY-MAXWELL STRESS-ENERGY TENSOR

Given the Lagrangian density $-F_{\mu\nu}F^{\mu\nu}/16\pi$, reexpress it in terms of the variables A_μ and $g^{\mu\nu}$, and by use of (21.33) derive the stress-energy tensor as discussed in §5.6. Also derive from the Lagrange variation principle the field equation $F_{\alpha\beta}^{;\beta} = 0$ (curved spacetime, but—for simplicity—a charge-free region of space). As a consequence of this field equation, show that the Faraday-Maxwell stress-energy tensor satisfies the conservation law, $T_{\alpha\beta}^{;\beta} = 0$. For a more ambitious project, show that any stress-energy tensor derived from a field Lagrangian by the prescription of equation (21.33) will automatically satisfy the conservation law $T_{\alpha\beta}^{;\beta} = 0$.

§21.4. SPLITTING SPACETIME INTO SPACE AND TIME

There are many ways to “push forward” many-fingered time and explore spacetime faster here and slower there, or faster there and slower here. However, a computer is most efficiently programmed only when it follows one definite prescription. The

EXERCISES

**Figure 21.2.**

Building two 3-geometries into a thin sandwich 4-geometry, by interposing perpendicular connectors between the two, with preassigned lengths and shifts. What would otherwise be flexible thereupon becomes rigid. The flagged point illustrates equation (21.40).

Slice spacetime to compute spacetime

successive hypersurfaces on which it gives the geometry are most conveniently described by successive values of a time-parameter t . One treats on a different footing the 3-geometries of these hypersurfaces and the 4-geometry that fills in between these laminations.

The slicing of spacetime into a one-parameter family of spacelike hypersurfaces is called for, not only by the analysis of the dynamics along the way, but also by the boundary conditions as they pose themselves in any action principle of the form, “Give the 3-geometries on the two faces of a sandwich of spacetime, and adjust the 4-geometry in between to extremize the action.”

Thin sandwich 4-geometry

There is no simpler sandwich to consider than one of infinitesimal thickness (Figure 21.2). Choosing coordinates adapted to the $(3 + 1)$ -space-time split, designate the “lower” (earlier) hypersurface in the diagram as $t = \text{constant}$ and the “upper” (later) one as $t + dt = \text{constant}$ (names, only names; no direct measure whatsoever of proper time). Compare the two hypersurfaces with two ribbons of steel out of which one wants to construct a rigid structure. To give the geometry on the two ribbons by no means fixes this structure; for that purpose, one needs cross-connectors between the one ribbon and the other. It is not even enough (1) to specify that these connectors are to be welded on perpendicular to the lower ribbon; (2) to specify where each is to be welded; and (3) to give its length. One must in addition tell where each connector joins the upper surface. If the proper distances between tops of the connectors are everywhere shorter than the distances between the bases of the connectors, the double ribbon will have the curve of the cable of a suspension bridge; if everywhere longer, the curve of the arch of a masonry bridge. The data necessary for the construction of the sandwich are thus (1) the metric of the 3-geometry of the lower hypersurface,

$$g_{ij}(t, x, y, z) dx^i dx^j, \quad (21.36)$$

telling the $(\text{distance})^2$ between one point in that hypersurface and another; (2) the metric on the upper hypersurface,

$$g_{ij}(t + dt, x, y, z) dx^i dx^j; \quad (21.37)$$

(3) a formula for the proper length,

$$\begin{pmatrix} \text{lapse of} \\ \text{proper time} \\ \text{between lower} \\ \text{and upper} \\ \text{hypersurface} \end{pmatrix} = \begin{pmatrix} \text{"lapse} \\ \text{function"} \end{pmatrix} dt = N(t, x, y, z) dt, \quad (21.38)$$

of the connector that is based on the point (x, y, z) of the lower hypersurface; and
(4) a formula for the place on the upper hypersurface,

$$x_{\text{upper}}^i(x^m) = x^i - N^i(t, x, y, z) dt, \quad (21.39)$$

where this connector is to be welded. Omit part of this information, and find the structure deprived of rigidity.

The rigidity of the structure of the thin sandwich is most immediately revealed in the definiteness of the 4-geometry of the spacetime filling of the sandwich. Ask for the proper interval ds or $d\tau$ between $x^\alpha = (t, x^i)$ and $x^\alpha + dx^\alpha = (t + dt, x^i + dx^i)$. The Pythagorean theorem in its 4-dimensional form

$$ds^2 = \left(\begin{array}{c} \text{proper distance} \\ \text{in base 3-geometry} \end{array} \right)^2 - \left(\begin{array}{c} \text{proper time from} \\ \text{lower to upper 3-geometry} \end{array} \right)^2$$

Metric of 4-geometry
depends on lapse and shift of
connectors of the two
3-geometries

yields the result (see Figure 21.2).

$$ds^2 = g_{ij}(dx^i + N^i dt)(dx^j + N^j dt) - (N dt)^2 \quad (21.40)$$

Here as in (21.36) the g_{ij} are the metric coefficients of the 3-geometry, distinguished by their Latin labels from the Greek-indexed components of the 4-metric,

$$ds^2 = {}^4g_{\alpha\beta} dx^\alpha dx^\beta, \quad (21.41)$$

labeled here with a suffix ⁽⁴⁾ to reduce the possibility of confusion. Comparing (21.41) and (21.40), one arrives at the following construction of the 4-metric out of the 3-metric and the lapse and shift functions [Arnowitt, Deser, and Misner (1962)]:

$$\begin{vmatrix} {}^4g_{00} & {}^4g_{0k} \\ {}^4g_{i0} & {}^4g_{ik} \end{vmatrix} = \begin{vmatrix} (N_s N^s - N^2) & N_k \\ N_i & g_{ik} \end{vmatrix}. \quad (21.42)$$

Details of the 4-geometry

The welded connectors do the job!

In (21.42), the quantities N^m are the components of the shift in its original primordial contravariant form, whereas the $N_i = g_{im} N^m$ are the covariant components, as calculated within the 3-geometry with the 3-metric. To invert this relation,

$$N^m = g^{ms} N_s \quad (21.43)$$

is to deal with the reciprocal 3-metric, a quantity that has to be distinguished sharply from the reciprocal 4-metric. Thus, the reciprocal 4-metric is

$$\begin{vmatrix} {}^4g^{00} & {}^4g^{0m} \\ {}^4g^{k0} & {}^4g^{km} \end{vmatrix} = \begin{vmatrix} -(1/N^2) & (N^m/N^2) \\ (N^k/N^2) & (g^{km} - N^k N^m / N^2) \end{vmatrix}, \quad (21.44)$$

a result that one checks by calculating out the product

$${}^4g_{\alpha\beta} {}^4g^{\beta\gamma} = {}^4\delta_\alpha^\gamma$$

according to the standard rules for matrix multiplication.

The volume element has the form

$$(-{}^4g)^{1/2} dx^0 dx^1 dx^2 dx^3 = Ng^{1/2} dt dx^1 dx^2 dx^3. \quad (21.45)$$

Welding the connectors to the two steel ribbons, or adding the lapse and shift functions to the 3-metric, by rigidifying the 4-metric, also automatically determines the components of the unit timelike normal vector \mathbf{n} . The condition of normalization on this 4-vector is most easily formulated by saying that there exists a 1-form, also called \mathbf{n} for the sake of convenience, dual to \mathbf{n} , and such that the product of this vector by this 1-form has the value

$$\langle \mathbf{n}, \mathbf{n} \rangle = -1. \quad (21.46)$$

This 1-form has the value

$$\mathbf{n} = n_\beta dx^\beta = -N dt + 0 + 0 + 0. \quad (21.47)$$

Only so can this 1-form, this structure of layered surfaces, automatically yield a value of unity, one bong of the bell, when pierced as in Figure 2.4 by a vector that represents an advance of one unit in proper time, regardless of what x , y , and z displacements it also has. Thus the unit timelike normal vector in covariant 1-form representation necessarily has the components

$$n_\beta = (-N, 0, 0, 0) \quad (21.48)$$

The components of the unit normal

Raise the indices via (21.44) to obtain the contravariant components of the same normal, represented as a tangent vector; thus,

$$n^\alpha = [(1/N), -(N^m/N)]. \quad (21.49)$$

This result receives a simple interpretation on inspection of Figure 21.2. Thus the typical “perpendicular connector” in the diagram can be said to have the components

$$(dt, -N^m dt)$$

and to have the proper length $d\tau = N dt$; so, ratioed down to a vector \mathbf{n} of unit proper length, the components are precisely those given by (21.49).

§21.5. INTRINSIC AND EXTRINSIC CURVATURE

The central concept in Einstein's account of gravity is curvature, so it is appropriate to analyze curvature in the language of the $(3 + 1)$ -space-time split. The curvature intrinsic to the 3-geometry of a spacelike hypersurface may be defined and calculated by the same methods described and employed in the calculation of four-dimensional curvature in Chapter 14. Of all measures of the intrinsic curvature, one of the simplest is the Riemann scalar curvature invariant 3R (written for simplicity of notation in what follows without the prefix, as R); and of all ways to define this invariant (see Chapter 14), one of the most compact uses the limit (see exercise 21.4)

$$R \left(\begin{array}{l} \text{at point} \\ \text{under study} \end{array} \right) = \lim_{\epsilon \rightarrow 0} 18 \frac{4\pi\epsilon^2 - \left(\begin{array}{l} \text{proper area of a surface (approximately)} \\ \text{a 2-sphere) defined as the locus of the} \\ \text{points at a proper distance } \epsilon \end{array} \right)}{4\pi \epsilon^4}$$

Scalar curvature as measure
of area deficit

(21.50)

For a more detailed description of the curvature intrinsic to the 3-geometry, capitalize on differential geometry as already developed in Chapters 8 through 14, amending it only as required to distinguish what is three-dimensional from what is four-dimensional. Begin by considering a displacement

$$d\mathcal{P} = \mathbf{e}_i dx^i \quad (21.51)$$

within the hypersurface. Here the \mathbf{e}_i are the basis tangent vectors $\mathbf{e}_i = \partial/\partial x^i$ (in one notation) or $\mathbf{e}_i = \partial\mathcal{P}/\partial x^i$ (in another notation) dual to the three coordinate 1-forms dx^i . Any field of tangent vectors \mathbf{A} that happens to lie in the hypersurface lets itself be expressed in terms of the same basis vectors:

$$\mathbf{A} = \mathbf{e}_i A^i. \quad (21.52)$$

The scalar product of this vector with the base vector \mathbf{e}_j is

$$(\mathbf{A} \cdot \mathbf{e}_j) = A^i (\mathbf{e}_i \cdot \mathbf{e}_j) = A^i g_{ij} = A_j. \quad (21.53)$$

Now turn attention from a vector at one point to the parallel transport of the vector to a nearby point.

A vector lying on the equator of the Earth and pointing toward the North Star, transported parallel to itself along a meridian to a point still on the Earth's surface, but 1,000 km to the north, will no longer lie in the 2-geometry of the surface of the Earth. A telescope located in the northern hemisphere has to raise its tube to see the North Star! The generalization to a three-dimensional hypersurface imbedded in a 4-geometry is immediate. Take vector \mathbf{A} , lying in the hypersurface, and transport it along an elementary route lying in the hypersurface, and in the course of this transport displace it at each stage parallel to itself, where "parallel" means parallel with respect to the geometry of the enveloping 4-manifold. Then \mathbf{A} will ordinarily

end up no longer lying in the hypersurface. Thus the “covariant derivative” of \mathbf{A} in the direction of the i -th coordinate direction in the geometry of the enveloping spacetime (that is, the \mathbf{A} at the new point diminished by the transported \mathbf{A}) has the form (see §10.4)

$${}^{(4)}\nabla_{\mathbf{e}_i}\mathbf{A} = {}^{(4)}\nabla_i\mathbf{A} = {}^{(4)}\nabla_i(\mathbf{e}_j A^j) = \mathbf{e}_j \frac{\partial A^j}{\partial x^i} + ({}^{(4)}\Gamma_{ji}^\mu \mathbf{e}_\mu) A^j. \quad (21.54)$$

A special instance of this formula is the equation for the covariantly measured change of the base vector \mathbf{e}_m itself,

$${}^{(4)}\nabla_i \mathbf{e}_m = {}^{(4)}\Gamma_{mi}^\mu \mathbf{e}_\mu. \quad (21.55)$$

In both (21.54) and (21.55) the presence of the “out-of-the-hypersurface component”

$$(A^j {}^{(4)}\Gamma_{ji}^0)(\mathbf{e}_0 \cdot \mathbf{n}) \quad (21.56)$$

From parallel transport in 4-geometry to parallel transport in 3-geometry

is quite evident. Now kill this component. Project ${}^{(4)}\nabla\mathbf{A}$ orthogonally onto the hypersurface. In this way arrive at a parallel transport and a covariant derivative that are intrinsic to the 3-geometry of the hypersurface. By rights this covariant derivative should be written ${}^{(3)}\nabla$; but for simplicity of notation it will be written as ∇ in the rest of this chapter, except where ambiguity might arise. To get the value of the new covariant derivative, one has only to rewrite (21.54) with the suffix ${}^{(4)}$ replaced everywhere by a ${}^{(3)}$, or, better, dropped altogether and with the “dummy index” of summation $\mu = (0, 1, 2, 3)$ replaced by $m = (1, 2, 3)$. However, it is more convenient, following Israel (1966), to turn from an expression dealing with contravariant components A^i of \mathbf{A} to one dealing with covariant components $A_i = (\mathbf{A} \cdot \mathbf{e}_i)$. Thus the covariant derivative of \mathbf{A} in the direction of the i -th coordinate direction in the hypersurface, calculated with respect to the 3-geometry intrinsic to the hypersurface itself, has for its h -th covariant component the quantity [see equation (10.18)]

$$A_{h|i} = \mathbf{e}_h \cdot {}^{(3)}\nabla_{\mathbf{e}_i}\mathbf{A} \equiv \mathbf{e}_h \cdot \nabla_i\mathbf{A} = \frac{\partial A_h}{\partial x^i} - A^m \Gamma_{mhi} (= A_{h;i} \text{ for } \mathbf{A} \text{ in } \Sigma). \quad (21.57)$$

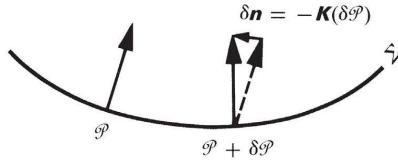
Here the notation of the vertical stroke distinguishes this covariant derivative from the covariant derivative taken with respect to the 4-geometry, as, for example, in equations (10.17ff). The connection coefficients here for three dimensions, like those dealt with earlier for four dimensions [see the equations leading from (14.14) through (14.15)], allow themselves to be expressed in terms of the metric coefficients and their first derivatives, and have the interpretation

$${}^{(3)}\Gamma_{mhi} \equiv \Gamma_{mhi} = \mathbf{e}_m \cdot \nabla_i \mathbf{e}_h. \quad (21.58)$$

From the connection coefficients in turn, one calculates as in Chapter 14 the full Riemann curvature tensor ${}^{(3)}R_{jmn}^i$ of the 3-geometry intrinsic to the hypersurface.

Over and above the curvature intrinsic to the simultaneity, one now encounters a concept not covered in previous chapters (except fleetingly in Box 14.1), the *extrinsic curvature* of the 3-geometry. This idea has no meaning for a 3-geometry

A new covariant derivative, taken with respect to the 3-geometry

**Figure 21.3.**

Extrinsic curvature measures the fractional shrinkage and deformation of a figure lying in the spacelike hypersurface Σ that takes place when each point in the figure is carried forward a unit interval of proper time “normal” to the hypersurface out into the enveloping spacetime. (No enveloping spacetime? No extrinsic curvature!) The extrinsic curvature tensor is a positive multiple of the unit tensor when elementary displacements $\delta\mathcal{P}$, in whatever direction within the surface they point, all experience the same fractional shrinkage. Thus the extrinsic curvature of the hypersurface illustrated in the figure is positive. The dashed arrow represents the normal vector n at the fiducial point \mathcal{P} after parallel transport to the nearby point $\mathcal{P} + \delta\mathcal{P}$.

conceived in and by itself. It depends for its existence on this 3-geometry’s being imbedded as a well-defined slice in a well-defined enveloping spacetime. It measures the curvature of this slice relative to that enveloping 4-geometry (Figure 21.3).

Take the normal that now stands at the point \mathcal{P} and, “keeping its base in the hypersurface” Σ , transport it parallel to itself as a “fiducial vector” to the point $\mathcal{P} + \delta\mathcal{P}$, and there subtract it from the normal vector that already stands at that point. The difference, δn , may be regarded in the appropriate approximation as a “vector,” the value of which is governed by and depends linearly on the “vector” of displacement $\delta\mathcal{P}$.

To obviate any appeal to the notion of approximation, go from the finite displacement $\delta\mathcal{P}$ to the limiting concept of the vector-valued “displacement 1-form” $d\mathcal{P}$ [see equation 15.13]. Also replace the finite but not rigorously defined vector δn by the limiting concept of a vector-valued 1-form dn . This quantity, regarded as a vector, being the change in a vector n that does not change in length, must represent a change in direction and thus stand perpendicular to n . Therefore it can be regarded as lying in the hypersurface Σ . Depending linearly on $d\mathcal{P}$, it can be represented in the form

$$dn = -K(d\mathcal{P}). \quad (21.59)$$

Extrinsic curvature as an operator

Here the linear operator K is the extrinsic curvature presented as an abstract coordinate-independent geometric object. The sign of K as defined here is positive when the tips of the normals in Figure 21.3 are closer than their bases, as they are, for example, during the recontraction of a model universe, in agreement with the conventions employed by Eisenhart (1926), Schouten (1954), and Arnowitt, Deser and Misner (1962), but opposite to the convention of Israel (1966).

Into the slots in the 1-forms that appear on the lefthand and righthand sides of (21.59), insert in place of the general tangent vector [which is to describe the general

local displacement, so far left open, as in the discussion following (2.12a)] a very special tangent vector, the basis vector \mathbf{e}_i , for a displacement in the i -th coordinate direction. Thus find (21.59) reading

$${}^{(4)}\nabla_i \mathbf{n} = -\mathbf{K}(\mathbf{e}_i) = -K_i^j \mathbf{e}_j, \quad (21.60)$$

where the K_i^j are the components of the linear operator \mathbf{K} in a coordinate representation. Take the scalar product of both sides of (21.60) with the basis vector \mathbf{e}_m . Recall $(\mathbf{e}_m \cdot \mathbf{n}) = 0$. Thus establish the symmetry of the tensor K_{im} , covariantly presented, in its two indices:

$$\begin{aligned} K_{im} &= K_i^j g_{jm} = K_i^j (\mathbf{e}_j \cdot \mathbf{e}_m) = -\mathbf{e}_m \cdot {}^{(4)}\nabla_i \mathbf{n} = \mathbf{n} \cdot {}^{(4)}\nabla_i \mathbf{e}_m \\ &= (\mathbf{n} \cdot \mathbf{e}_0) {}^{(4)}\Gamma_{mi}^0 = \mathbf{n} \cdot {}^{(4)}\nabla_m \mathbf{e}_i = K_{mi}. \end{aligned} \quad (21.61)$$

↑
[see (21.55)]

A knowledge of the tensor K_{ij} of extrinsic curvature assists in revealing the changes of the four vectors $\mathbf{n}, \mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3$ under parallel transport. Equation (21.60) already tells how \mathbf{n} changes under parallel transport. The change of \mathbf{e}_m is to be read off from (21.55) as a vector. It is adequate identification of this vector to know its scalar product with each of four independent vectors: with the basis vectors $\mathbf{e}_1, \mathbf{e}_2$, and \mathbf{e}_3 , or, more briefly, with \mathbf{e}_s , in (21.58); and with the normal vector \mathbf{n} in (21.61). Thus one arrives, following Israel (1966), at what are known as the equations of Gauss and Weingarten, in happy oversight of all change of notation in the intervening century:

Gauss-Weingarten equation
for 4-transport in terms of
extrinsic curvature

$${}^{(4)}\nabla_i \mathbf{e}_j = K_{ij} \frac{\mathbf{n}}{\mathbf{n} \cdot \mathbf{n}} + {}^{(3)}\Gamma_{ji}^h \mathbf{e}_h. \quad (21.62)$$

Knowing from this equation how each basis vector in Σ changes, one also knows how to rewrite (21.54) for the change in any vector field \mathbf{A} that lies in Σ . The change in both cases is expressed relative to a fiducial vector transported from a fiducial nearby point. By the term “parallel transport” one now means “parallel with respect to the geometry of the enveloping spacetime”:

$${}^{(4)}\nabla_i \mathbf{A} = A^j {}_{|i} \mathbf{e}_j + K_{ij} A^j \frac{\mathbf{n}}{(\mathbf{n} \cdot \mathbf{n})}. \quad (21.63)$$

Of special importance is the evaluation of extrinsic curvature when spacetime is sliced up into spacelike slices according to the plan of Arnowitt, Deser, and Misner as described in §21.4. The 4-geometry of the thin sandwich illustrated in Figure 21.2, rudimentary though it is, is fully defined by the 3-metric on the two faces of the sandwich and by the lapse and shift functions N and N^i . The normal in covariant representation according to (21.47) has the components

$$(n_0, n_1, n_2, n_3) = (-N, 0, 0, 0). \quad (21.64)$$

The change in \mathbf{n} relative to “ \mathbf{n} transported parallel to itself in the enveloping 4-geometry,” according to the definition of parallel transport, is

$$\begin{aligned}
 (\mathbf{d}n)_i &= n_{i;k} \mathbf{dx}^k \\
 &= \left[\frac{\partial n_i}{\partial x^k} - {}^{(4)}\Gamma_{ik}^\sigma n_\sigma \right] \mathbf{dx}^k \\
 &= N {}^{(4)}\Gamma_{ik}^0 \mathbf{dx}^k
 \end{aligned} \tag{21.65}$$

Compare to the same change as expressed in terms of the extrinsic curvature tensor,

$$(\mathbf{d}n)_i = -K_{ik} \mathbf{dx}^k. \tag{21.66}$$

Conclude that this tensor has the value

$$K_{ik} = -n_{i;k} = -N {}^{(4)}\Gamma_{ik}^0 = -N[{}^{(4)}g^{00} {}^{(4)}\Gamma_{0ik} + {}^{(4)}g^{0p} {}^{(4)}\Gamma_{pik}],$$

or, with the help of equations (21.42) and (21.44),

$$\begin{aligned}
 K_{ik} &= (1/N)[{}^{(4)}\Gamma_{0ik} - N^p {}^{(3)}\Gamma_{pik}] \\
 &= \frac{1}{2N} \left[\frac{\partial N_i}{\partial x^k} + \frac{\partial N_k}{\partial x^i} - \frac{\partial g_{ik}}{\partial t} - 2\Gamma_{pik} N^p \right] \\
 &= \frac{1}{2N} \left[N_{i|k} + N_{k|i} - \frac{\partial g_{ik}}{\partial t} \right].
 \end{aligned} \tag{21.67}$$

Extrinsic curvature in terms
of shift and change of
3-metric

This is the extrinsic curvature expressed in terms of the ADM lapse and shift functions [Arnowitt, Deser, and Misner (1962)].

As an example, let Σ have the geometry of a 3-sphere

$$ds^2 = a^2[dx^2 + \sin^2\chi(d\theta^2 + \sin^2\theta d\phi^2)]. \tag{21.68}$$

Extrinsic curvature of
expanding 3-sphere

Let the nearby spacelike slice in the one-parameter family of slices, the slice with the label $t + dt$ (only a label!) have a 3-metric given by the same formula with the radius a replaced by $a + da$. The 4-geometry of the thin sandwich between these two slices is completely undetermined until one gives the lapse and shift functions. For simplicity, take the shift vector N^i (see Figure 21.2) to be everywhere zero and the lapse function at every point on Σ to have the same value N . The separation in proper time between the two spheres is thus $d\tau = N dt$. Any geometric figure located in Σ expands with time. The fractional increase of any length in this figure per unit of proper time is the same in whatever direction that length is oriented, and has the value

$$\left(\begin{array}{l} \text{fractional increase} \\ \text{of length per unit} \\ \text{of proper time} \end{array} \right) = \frac{1}{a} \frac{da}{d\tau} = \frac{1}{2N} \frac{1}{a^2} \frac{d(a^2)}{dt}. \tag{21.69}$$

The negative of this quantity, multiplied by the (1_1) unit tensor, $\mathbf{1} = \mathbf{d}\mathcal{P}$, gives the extrinsic curvature tensor in (1_1) representation,

$$\mathbf{K} = -\frac{1}{2N} \frac{1}{a^2} \frac{d(a^2)}{dt} \mathbf{1}. \tag{21.70}$$

One confirms this result (exercise 21.5) by direct calculation of the components K_i^j using the ADM formula (21.67) as the starting point.

The Riemann curvature $R^a_{bcd} = {}^{(3)}R^a_{bcd}$ intrinsic to the hypersurface Σ , together with the extrinsic curvature K_{ij} , give one information on the Riemann and Einstein curvatures of the 4-geometry. In the calculation, it is not convenient to use the coordinate basis,

$$\begin{array}{ll} \text{basis vectors,} & \text{basis 1-forms} \\ \mathbf{e}_0 = \partial_t, & dt, \\ \mathbf{e}_i = \partial_i, & dx^i, \end{array}$$

because ordinarily the basis vector \mathbf{e}_0 does not stand perpendicular to the hypersurface (see Figure 21.2). Adopt a different basis but one that is still self-dual,

Basic forms for calculating 4-curvature

$$\begin{array}{ll} \text{basis vectors,} & \text{basis 1-forms,} \\ \mathbf{e}_n \equiv \mathbf{n} = N^{-1}(\partial_t - N^m \partial_m), & \omega^n = N dt = (\mathbf{n} \cdot \mathbf{n}) \mathbf{n} \\ \mathbf{e}_i = \partial_i, & \omega^i \equiv dx^i + N^i dt. \end{array} \quad (21.71)$$

Also use Greek labels $\bar{\alpha} = n, 1, 2, 3$, instead of Greek labels $\alpha = 0, 1, 2, 3$, to list components.

Recall that curvature is measured by the change in a vector on transport around a closed route; or, from equation (14.23),

$$\mathcal{R}(\mathbf{u}, \mathbf{v})\mathbf{w} = \nabla_{\mathbf{u}} \nabla_{\mathbf{v}} \mathbf{w} - \nabla_{\mathbf{v}} \nabla_{\mathbf{u}} \mathbf{w} - \nabla_{[\mathbf{u}, \mathbf{v}]} \mathbf{w}. \quad (21.72)$$

Let the vector transported be \mathbf{e}_i and let the route be defined by \mathbf{e}_j and \mathbf{e}_k . The latter two vectors belong to a coordinate basis. Therefore the “route closes automatically”, $[\mathbf{e}_j, \mathbf{e}_k] = 0$, and the final term in (21.72) drops out of consideration. Call on (21.62) and (21.60) to find

$$\begin{aligned} {}^{(4)}\nabla_{\mathbf{e}_j} {}^{(4)}\nabla_{\mathbf{e}_k} \mathbf{e}_i &= {}^{(4)}\nabla_{\mathbf{e}_j} \left[K_{ik} \frac{\mathbf{n}}{(\mathbf{n} \cdot \mathbf{n})} + {}^{(3)}\Gamma_{ik}^m \mathbf{e}_m \right] \\ &= K_{ik,j} \frac{\mathbf{n}}{(\mathbf{n} \cdot \mathbf{n})} - K_{ik} K_j^m \mathbf{e}_m \frac{1}{(\mathbf{n} \cdot \mathbf{n})} + {}^{(3)}\Gamma_{ik,j}^m \mathbf{e}_m \\ &\quad + {}^{(3)}\Gamma_{ik}^m \left[K_{mj} \frac{\mathbf{n}}{(\mathbf{n} \cdot \mathbf{n})} + {}^{(3)}\Gamma_{mj}^s \mathbf{e}_s \right]. \end{aligned} \quad (21.73)$$

Evaluate similarly the term with indices j and k reversed, subtract from (21.73), simplify, and find

$$\begin{aligned} \mathcal{R}(\mathbf{e}_j, \mathbf{e}_k) \mathbf{e}_i &= (K_{ik|j} - K_{ij|k}) \frac{\mathbf{n}}{(\mathbf{n} \cdot \mathbf{n})} \\ &\quad + [(\mathbf{n} \cdot \mathbf{n})^{-1}(K_{ij} K_k^m - K_{ik} K_j^m) + {}^{(3)}R^m_{ijk}] \mathbf{e}_m. \end{aligned} \quad (21.74)$$

Gauss-Codazzi: 4-curvature in terms of intrinsic 3-geometry and extrinsic curvature

The coefficients give directly the desired components of the curvature tensor

$${}^{(4)}R^m_{ijk} = {}^{(3)}R^m_{ijk} + (\mathbf{n} \cdot \mathbf{n})^{-1}(K_{ij} K_k^m - K_{ik} K_j^m) \quad (21.75)$$

and

$${}^{(4)}R^n_{ijk} = (\mathbf{n} \cdot \mathbf{n})^{-1} {}^{(4)}R_{nijk} = -(\mathbf{n} \cdot \mathbf{n})^{-1}(K_{ij|k} - K_{ik|j}). \quad (21.76)$$

Equations (21.75) and (21.76) are known as the equations of Gauss and Codazzi [for literature, see Eisenhart (1926)]. It follows from (21.75) that the components of the curvature of the 3-geometry will normally only then agree with the corresponding components of the curvature of the 4-geometry when the imbedding happens to be accomplished at the point under study with a hypersurface free of extrinsic curvature. The directly opposite situation is illustrated by the familiar example of a 2-sphere imbedded in a flat 3-space, where the lefthand side of (21.75) (with dimensions lowered by one unit throughout!) is zero, and the extrinsic and intrinsic curvature on the right exactly cancel.

Important components of the Einstein curvature let themselves be evaluated from the Gauss-Codazzi results. In doing the calculation, it is simplest to think of \mathbf{e}_i , \mathbf{e}_j and \mathbf{e}_k as being an orthonormal tetrad, \mathbf{n} being itself already normalized and orthogonal to every vector in the hypersurface. Then, employing (14.7) and (21.75), one finds

$$\begin{aligned} -G_0^0 &= {}^{(4)}R^{12}_{12} + {}^{(4)}R^{23}_{23} + {}^{(4)}R^{31}_{31} \\ &= {}^{(3)}R^{12}_{12} + {}^{(3)}R^{23}_{23} + {}^{(3)}R^{31}_{31} \\ &\quad + (\mathbf{n} \cdot \mathbf{n})^{-1}[(K_1^2 K_2^1 - K_2^2 K_1^1) + (K_2^3 K_3^1 - K_3^3 K_2^1) \\ &\quad + (K_3^1 K_1^3 - K_1^1 K_3^3)] \\ &= \frac{1}{2} R - \frac{1}{2} (\mathbf{n} \cdot \mathbf{n})^{-1}[(\text{Tr } \mathbf{K})^2 - \text{Tr } (\mathbf{K}^2)]. \end{aligned} \tag{21.77}$$

Einstein curvature in terms of
extrinsic curvature

Here R is the 3-dimensional scalar curvature invariant and Tr stands for “trace of”; thus,

$$\text{Tr } \mathbf{K} = g^{ij} K_{ij} = g_{ij} K^{ij} = K_j^j \tag{21.78}$$

and

$$\text{Tr } \mathbf{K}^2 = (K^2)_j^j = K_j^m K_m^j = g_{js} K^{sm} g_{mi} K^{ij}. \tag{21.79}$$

The result, though obtained in an orthonormal tetrad, plainly is covariant with respect to general coordinate transformations within the spacelike hypersurface; and it makes no explicit reference whatever to any time coordinate, in this respect providing a coordinate-free description of the Einstein curvature.

The Einstein field equation equates (21.77) to $8\pi\rho$, where ρ is the density of mass-energy. Expression (21.77) is the “measure of curvature that is independent of how curved one cuts a spacelike slice.” This measure of curvature is central to the derivation of Einstein’s field equation that is sketched in Box 17.2, item 3, “Physics on a Spacelike Slice.”

The other component of the Einstein curvature tensor that is easily evaluated by (14.7) from the results at hand has the form

$$\begin{aligned} G_1^n &= {}^{(4)}R^{n2}_{12} + {}^{(4)}R^{n3}_{13} \\ &= -(\mathbf{n} \cdot \mathbf{n})^{-1}(K_{1|2}^2 - K_{2|1}^2 + K_{1|3}^3 - K_{3|1}^3), \end{aligned} \tag{21.80}$$

when referred to an orthonormal frame. One immediately translates to a form valid for any frame \mathbf{e}_1 , \mathbf{e}_2 , \mathbf{e}_3 in the hypersurface, orthonormal or not,

$$G_i^n = -(\mathbf{n} \cdot \mathbf{n})^{-1}[K_{i|m}^m - (\text{Tr } \mathbf{K})_{|i}]. \tag{21.81}$$

The other initial-value
equation

Equation (21.77) is the
central Einstein equation,
“mass-energy fixes
curvature”

The Einstein field equation equates this quantity to 8π times the i -th covariant component of the density of momentum carried by matter and fields other than gravity.

The four components of the Einstein field equation so far written down will have a central place in what follows as “initial-value equations” of general relativity. The other six components will not be written out: (1) the dynamics lets itself be analyzed more simply by Hamiltonian methods; and (2) the calculation takes work. It demands that one evaluate the remaining type of object, $\mathcal{R}(\mathbf{e}_j, \mathbf{n})\mathbf{e}_i$. One step towards that calculation will be found in exercise 21.7. Sachs does the calculation (1964, equation 10) but only after specializing to Gaussian normal coordinates. These coordinates presuppose a very special slicing of spacetime: (1) geodesics issuing normally from the spacelike hypersurface $n = 0$ cut all subsequent simultaneities $n = \text{constant}$ normally; and (2) the n coordinate directly measures lapse of proper time, or proper length, whichever is appropriate,* along these geodesics. In coordinates so special it is not surprising that the answer looks simple:

$${}^{(4)}R^n_{ink} = (\mathbf{n} \cdot \mathbf{n})^{-1} \left(\frac{\partial K_{ik}}{\partial n} + K_{im} K^m{}_k \right). \quad \begin{matrix} \text{(Gaussian normal)} \\ \text{coordinates} \end{matrix} \quad (21.82)$$

Additional terms come into (21.82) when one uses, instead of the Gaussian normal coordinate system, the coordinate system of Arnowitt, Deser, and Misner. The ADM coordinates are employed here because they allow one to analyze the dynamics as one *wants* to analyze the dynamics, with freedom to push the spacelike hypersurface ahead in time at different rates in different places (“many-fingered time”). Fischer (1971) shows how to evaluate and understand the geometric content of such formulas in a coordinate-free way by using the concept of Lie derivative of a tensor field, an introduction to which is provided by exercise 21.8.

*Here Sachs’ equation (10) is generalized to the case where the unit normal \mathbf{n} is not necessarily timelike. Sachs used $\mathbf{n} = \partial/\partial t$.

EXERCISES

Exercise 21.4. SCALAR CURVATURE INVARIANT IN TERMS OF AREA DEFICIT

It being 10,000 km from North Pole to equator, one would have 62,832 km for the length of the “equator” if the earth were flat, as contrasted to the actual $\sim 40,000$ km, a difference reflecting the fact that the surface is curved up into closure. Turn from this “pre-problem” to the actual problem, a 3-sphere

$$ds^2 = a^2[d\chi^2 + \sin^2\chi(d\theta^2 + \sin^2\theta d\phi^2)].$$

Measure off from $\chi = 0$ a 2-sphere of proper radius $\epsilon = a\chi$. Determine the proper area of this 2-sphere as a function of χ . Verify that relation (21.50) on the area deficit gives in the limit $\epsilon \rightarrow 0$ the correct result $R = 6/a^2$. For a more ambitious exercise: (1) take a general (smooth) 3-geometry; (2) express the metric near any chosen point in terms of Riemann’s normal coordinates as given in §11.6; (3) determine the locus of the set of points at the proper distance ϵ to the lowest interesting power of ϵ in terms of the spherical polar angles θ and ϕ (direction of start of geodesic of length ϵ); (4) determine to the lowest interesting power of ϵ the proper area of the figure defined by these points; and thereby establish (21.50) [for more on this topic see, for example, Cartan (1946), pp. 252–256].

Exercise 21.5. EXTRINSIC CURVATURE TENSOR FOR SLICE OF FRIEDMANN GEOMETRY

Confirm the result (21.70) for the extrinsic curvature by direct calculation from formula (21.67).

Exercise 21.6. EVALUATION OF $\mathcal{R}(\mathbf{e}_j, \mathbf{e}_k)\mathbf{n}$

Evaluate this quantity along the model of (21.74) or otherwise. How can it be foreseen that the coefficient of \mathbf{n} in the result must vanish identically? Comparing coefficients of \mathbf{e}_m , find ${}^{(4)}R^m_{\mu jk}$ and test for equivalence to equation (21.76).

Exercise 21.7. EVALUATION OF THE COMMUTATOR $[\mathbf{e}_j, \mathbf{n}]$

The evaluation of this commutator is a first step toward the calculation of a quantity like $\mathcal{R}(\mathbf{e}_j, \mathbf{n})\mathbf{e}_i$. Expressing \mathbf{e}_j as the differential operator $\partial/\partial x^j$, use (21.49) to represent \mathbf{n} also as a differential operator. In this way, show that the commutator in question has the value $-(N_j/N)\mathbf{n} - (N^m_{,j}/N)\mathbf{e}_m$.

Exercise 21.8. LIE DERIVATIVE OF A TENSOR (exercise provided by J. W. York, Jr.)

Define the Lie derivative of a tensor field and explore some of its properties. The Lie derivative along a vector field \mathbf{n} is a differential operator that operates on tensor fields \mathbf{T} of type (ℓ, s) , converting them into tensors $\mathcal{L}_{\mathbf{n}}\mathbf{T}$, also of type (ℓ, s) . The Lie differentiation process obeys the usual chain rule and has additivity properties [compare equations (10.2b, 10.2c, 10.2d) for the covariant derivative]. For scalar functions f , one has $\mathcal{L}_{\mathbf{n}}f \equiv \mathbf{n}[f] = f_{,\mu}n^{\mu}$. The Lie derivative of a vector field \mathbf{u} along a vector field \mathbf{v} was defined in exercise 9.11 by

$$\mathcal{L}_{\mathbf{u}}\mathbf{v} \equiv [\mathbf{u}, \mathbf{v}].$$

If the action of $\mathcal{L}_{\mathbf{n}}$ on 1-forms is defined, the extension to tensors of general type will be simple, because the latter can always be decomposed into a sum of tensor products of vectors and 1-forms. If σ is a 1-form and \mathbf{v} is a vector, then one defines $\mathcal{L}_{\mathbf{n}}\sigma$ to be that 1-form satisfying

$$\langle \mathcal{L}_{\mathbf{n}}\sigma, \mathbf{v} \rangle = \mathbf{n}[\langle \sigma, \mathbf{v} \rangle] - \langle \sigma, [\mathbf{n}, \mathbf{v}] \rangle$$

for arbitrary \mathbf{v} .

(a) Show that in a coordinate basis

$$\mathcal{L}_{\mathbf{n}}\sigma = (\sigma_{\alpha,\beta} n^{\beta} + \sigma_{\beta} n^{\beta}_{,\alpha}) \mathbf{dx}^{\alpha}.$$

(b) Show that in a coordinate basis

$$\mathcal{L}_{\mathbf{n}}\mathbf{T} = (T_{\alpha\beta,\mu} n^{\mu} + T_{\mu\beta} n^{\mu}_{,\alpha} + T_{\alpha\mu} n^{\mu}_{,\beta}) \mathbf{dx}^{\alpha} \otimes \mathbf{dx}^{\beta}$$

where \mathbf{T} is of type (ℓ, s) .

(c) Show that in (a) and (b), all partial derivatives can be replaced by covariant derivatives. [Observe that Lie differentiation is defined independently of the existence of an affine connection. For more information, see, for example, Bishop and Goldberg (1968) and Schouten (1954).]

Exercise 21.9. EXPRESSION FOR DYNAMIC COMPONENTS OF THE CURVATURE TENSOR (exercise provided by J. W. York, Jr.)

The Gauss-Codazzi equations can be viewed as giving 14 of the 20 algebraically independent components of the spacetime curvature tensor in terms of the intrinsic and extrinsic geometry of three-dimensional (non-null) hypersurfaces. In order to accomplish a space-plus-time splitting of the Hilbert Lagrangian $\sqrt{-g}{}^{(4)}R$, one must express, in addition, the remaining

6 components of the curvature tensor in an analogous manner. It is convenient for this purpose to express all tensors as spacetime tensors, and to use Lie derivation in the direction of the timelike unit normal field of the spacelike hypersurfaces as a generalized notion of time differentiation. A number of preliminary results must be proven:

$$(a) \quad \mathcal{L}_u g_{\mu\nu} = u_{\mu;\nu} + u_{\nu;\mu},$$

$$(b) \quad \mathcal{L}_u(g_{\mu\nu} + u_\mu u_\nu) \equiv \mathcal{L}_u(\gamma_{\mu\nu}) \\ = u_{\mu;\nu} + u_{\nu;\mu} + u_\mu a_\nu + a_\mu u_\nu,$$

where $\gamma_{\mu\nu}$ is the metric of the spacelike hypersurface, expressed in the spacetime coordinate basis, and $a^\mu \equiv u^\lambda \nabla_\lambda u^\mu$ is the curvature vector (4-acceleration) of the timelike normal curves whose tangent field is u^μ . (Recall that $u_\mu a^\mu = 0$.)

(c) Prove that the extrinsic curvature tensor is given by

$$K_{\mu\nu} = -\frac{1}{2} \mathcal{L}_u \gamma_{\mu\nu}.$$

(d) The unit tensor of projection into the hypersurface is defined by

$$\perp_\nu^\mu \equiv \delta_\nu^\mu + u^\mu u_\nu.$$

In terms of \perp show that one can write

$$u_{\alpha;\beta} \equiv -K_{\alpha\beta} - \omega_{\alpha\beta} - a_\alpha u_\beta,$$

where

$$K_{\alpha\beta} = -\perp_\alpha^\mu \perp_\beta^\nu u_{(\mu;\nu)}$$

and

$$\omega_{\alpha\beta} = -\perp_\alpha^\mu \perp_\beta^\nu u_{[\mu;\nu]}.$$

(e) From the fact that u^μ is the unit normal field for a family of spacelike hypersurfaces, show that $\omega_{\alpha\beta} = 0$.

(f) The needed tools are now on hand. To obtain the result:

- (i) Write down $\mathcal{L}_u K_{\mu\nu}$ (see exercise 21.8);
- (ii) Insert this expression into the Ricci identity in the form

$$u^\sigma \nabla_\sigma \nabla_\mu u_\nu = u^\sigma \nabla_\mu \nabla_\sigma u_\nu + {}^{(4)}R_{\rho\nu\mu\sigma} u^\sigma u^\rho;$$

(iii) Project the two remaining free indices into the hypersurface using \perp , and show that one obtains

$$\perp_\alpha^\mu \perp_\beta^\rho {}^{(4)}R_{\mu\nu\rho\sigma} u^\nu u^\sigma = \mathcal{L}_u K_{\alpha\beta} + K_{\alpha\gamma} K_\beta^\gamma \\ + {}^{(3)}\nabla_{(\alpha} a_{\beta)} + a_\alpha a_\beta,$$

where ${}^{(3)}\nabla_\alpha a_\beta \equiv \perp_\alpha^\mu \perp_\beta^\nu \nabla_\mu a_\nu$ can be shown to be the three-dimensional covariant derivative of a_β . In Gaussian normal coordinates, show that one obtains from this result

$$R_{0i0j} = \frac{\partial}{\partial t} K_{ij} + K_{ik} K_j^k.$$

(g) Finally, in the construction of ${}^{(4)}R$, one needs to show that

$$\gamma^{\mu\nu} [{}^{(3)}\nabla_{(\mu} a_{\nu)} + a_\mu a_\nu] = g^{\mu\nu} [{}^{(3)}\nabla_{(\mu} a_{\nu)} + a_\mu a_\nu] = a_{;\lambda}^\lambda.$$

**Exercise 21.10. EXPRESSION OF ${}^{(4)}R^i_{n\mu}$ IN TERMS OF EXTRINSIC CURVATURE, PLUS A COVARIANT DIVERGENCE
(exercise provided by K. Kuchař)**

Let α' be an arbitrary smooth set of four coordinates, not necessarily coordinated in any way with the choice of the 1-parameter family of hypersurfaces.

(a) Show that

$${}^{(4)}R^i_{n\mu} = g^{\alpha'\gamma'} n^{\beta'} (n_{\alpha';\beta'\gamma'} - n_{\alpha';\gamma'\beta'}).$$

(b) Show that the covariant divergences

$$(n^{\beta'} n^{\gamma'};_{\beta'})_{;\gamma'}$$

and

$$-(n^{\beta'} n^{\gamma'};_{\gamma'})_{;\beta'}$$

can be removed from this expression in such a way that what is left behind contains only first derivatives of the unit normal vector \mathbf{n} .

(c) Noting that the basis vectors \mathbf{e}_i and \mathbf{n} form a complete set, justify the formula

$$g^{\beta'\mu'} = e_i^{\beta'} \omega^{i\mu'} + (\mathbf{n} \cdot \mathbf{n})^{-1} n^{\beta'} n^{\mu'},$$

where ω^i is the 1-form dual to \mathbf{e}_i .

(d) Noting that $n_{\alpha';\beta'} n^{\alpha'} = 0$ and

$$K_{ij} = -e_{i\alpha'} n^{\alpha'};_{\beta'} e_j^{\beta'},$$

show that

$${}^{(4)}R^i_{n\mu} = (\text{Tr } \mathbf{K})^2 - \text{Tr } \mathbf{K}^2 \text{ plus a covariant divergence.}$$

§21.6. THE HILBERT ACTION PRINCIPLE AND THE ARNOWITT-DESER-MISNER MODIFICATION THEREOF IN THE SPACE-PLUS-TIME SPLIT

For analyzing the dynamics, it happily proves unnecessary to possess the missing formula for ${}^{(4)}R^n_{ink}$. It is essential, however, to have the Lagrangian density,

$$16\pi\mathcal{L}_{\text{geom}} = (-{}^{(4)}g)^{1/2} {}^{(4)}R, \quad (21.83)$$

in the Hilbert action principle as the heart of all the dynamic analysis. In the present ADM (1962) notation, this density has the form

$$\begin{aligned} (-{}^{(4)}g)^{1/2} {}^{(4)}R &= (-{}^{(4)}g)^{1/2} [{}^{(4)}R^{ij}_{ij} + 2 {}^{(4)}R^{in}_{in}] \\ &= (-{}^{(4)}g)^{1/2} [R + (\mathbf{n} \cdot \mathbf{n})(\text{Tr } \mathbf{K}^2 - (\text{Tr } \mathbf{K})^2) + 2(\mathbf{n} \cdot \mathbf{n}) {}^{(4)}R^i_{n\mu}]. \end{aligned} \quad (21.84)$$

Kuchař (1971b; see also exercise 21.10) shows how to calculate a sufficient part of this quantity without calculating all of it. The difference between the “sufficient part” and the “whole” is a time derivative plus a divergence, a quantity of the form

$$[(-{}^{(4)}g)^{1/2} A^\alpha]_\alpha = (-{}^{(4)}g)^{1/2} A^\alpha_{;\alpha}. \quad (21.85)$$

Drop a complete derivative from the Hilbert action principle to get the ADM principle

When one multiplies (21.83) by $dt dx^1 dx^2 dx^3$ and integrates to obtain the action integral, the term (21.85) integrates out to a surface term. Variations of the geometry interior to this surface make no difference in the value of this surface term. Therefore it has no influence on the equations of motion to drop the term (21.85). The result of the calculation (exercise 21.10) is simple: what is left over after dropping the divergence merely changes the sign of the terms in $\text{Tr } \mathbf{K}^2$ and $(\text{Tr } \mathbf{K})^2$ in (21.84). Thus the variation principle becomes

$$\begin{aligned} (\text{extremum}) &= I_{\text{modified}} = \int \mathcal{L}_{\text{modified}} d^4x \\ &= (1/16\pi) \int [R + (\mathbf{n} \cdot \mathbf{n})(\text{Tr } \mathbf{K})^2 - \text{Tr } \mathbf{K}^2] Ng^{1/2} dt d^3x + \int \mathcal{L}_{\text{fields}} d^4x. \end{aligned} \quad (21.86)$$

This expression, rephrased, is the starting point for Arnowitt, Deser, and Misner's analysis of the dynamics of geometry.

Two supplements from a paper of York (1972b; see also exercise 21.9) enlarge one's geometric insight into what is going on in the foregoing analysis. First, the tensor of extrinsic curvature lets itself be defined [see also Fischer (1971)] most naturally in the form

$$\mathbf{K} = -\frac{1}{2} \mathcal{L}_n \mathbf{g}, \quad (21.87)$$

where \mathbf{g} is the metric tensor of the 3-geometry, \mathbf{n} is the timelike unit normal field, and \mathcal{L} is the Lie derivative as defined in exercise 21.8. Second, the divergence (21.85), which has to be added to the Lagrangian of (21.86) to obtain the full Hilbert Lagrangian, is

$$-2[(-{}^4g)^{1/2}(n^{\alpha'} \text{Tr } \mathbf{K} + a^{\alpha'})]_{,\alpha'}, \quad (21.88)$$

where the coordinates are general (see exercise 21.10), and

$$a^{\alpha'} = n^{\alpha'}_{;\beta} n^\beta \quad (21.89)$$

is the 4-acceleration of an observer traveling along the timelike normal \mathbf{n} to the successive slices.

§21.7. THE ARNOWITT, DESER, AND MISNER FORMULATION OF THE DYNAMICS OF GEOMETRY

Dirac (1959, 1964, and earlier references cited therein) formulated the dynamics of geometry in a (3 + 1)-dimensional form, using generalizations of Poisson brackets and of Hamilton equations. Arnowitt, Deser, and Misner instead made the Hilbert-Palatini variational principle the foundation for this dynamics. Because of its simplicity, this ADM (1962) approach is followed here. The gravitational part of the integrand in the Hilbert-Palatini action principle is rewritten in the condensed but standard form (after inserting a 16π that ADM avoid by other units) as

$$\begin{aligned} 16\pi \mathcal{L}_{\text{geom true}} &= \mathcal{L}_{\text{geom ADM}} = -g_{ij} \partial \pi^{ij} / \partial t - N \mathcal{K} - N_i \mathcal{K}^i \\ &\quad - 2 \left[\pi^{ij} N_j - \frac{1}{2} N^i \text{Tr } \mathbf{n} + N^{ij} (g)^{1/2} \right]_{,i}. \end{aligned} \quad (21.90)$$

Here each item of abbreviation has its special meaning and will play its special part, a part foreshadowed by the name now given it:

$$\pi_{\text{true}}^{ij} = \frac{\delta(\text{action})}{\delta g_{ij}} = \begin{pmatrix} \text{"geometrodynamic} \\ \text{field momentum"} \text{ dyn-} \\ \text{amically conjugate to} \\ \text{the "geometrodynamic} \\ \text{field coordinate" } g_{ij} \end{pmatrix} = \frac{\pi^{ij}}{16\pi}; \pi^{ij} = g^{1/2}(g^{ij}\text{Tr } \mathbf{K} - K^{ij}) \quad (21.91)$$

Momenta conjugate to the dynamic g_{ij}

(here the π^{ij} of ADM is usually more convenient than π_{true}^{ij}); and

$$\begin{aligned} \mathcal{H}_{\text{true}} &= \mathcal{H}(\pi_{\text{true}}^{ij}, g_{ij}) = (\text{"super-Hamiltonian"}) = \mathcal{H}/16\pi; \\ \mathcal{H}(\pi^{ij}, g_{ij}) &= g^{-1/2} \left(\text{Tr } \mathbf{n}^2 - \frac{1}{2} (\text{Tr } \mathbf{n})^2 \right) - g^{1/2}R; \end{aligned} \quad (21.92)$$

and

$$16\pi\mathcal{H}_{\text{true}}^i = \mathcal{H}^i = \mathcal{H}^i(\pi^{ij}, g_{ij}) = (\text{"supermomentum"}) = -2\pi^{ik}{}_{|k}. \quad (21.93)$$

Here the covariant derivative is formed treating π^{ik} as a tensor density, as its definition in (21.91) shows it to be (see §21.2). The quantities to be varied to extremize the action are the coefficients in the metric of the 4-geometry, as follows: the six g_{ij} and the lapse function N and shift function N_i ; and also the six “geometrodynamic momenta,” π^{ij} . To vary these momenta as well as the metric is (1) to follow the pattern of elementary Hamiltonian dynamics (Box 21.1), where, by taking the momentum p to be as independently variable as the coordinate x , one arrives at two Hamilton equations of the first order instead of one Lagrange equation of the second order, and (2) to follow in some measure the lead of the Palatini variation principle of §21.2. There, however, one had 40 connection coefficients to vary, whereas here one has come down to only six π^{ij} . To know these momenta and the 3-metric is to know the extrinsic curvature. Before carrying out the variation, drop the divergence $-2[\]_i$ from (21.90), since it gives rise only to surface integrals and therefore in no way affects the equations of motion that will come out of the variational principle. Also rewrite the first term in (21.90) in the form

$$-(\partial/\partial t)(g_{ij}\pi^{ij}) + \pi^{ij}\partial g_{ij}/\partial t, \quad (21.94)$$

and drop the complete time-derivative from the variation principle, again because it is irrelevant to the resulting equations of motion. The action principle now takes the form

$$\begin{aligned} \text{extremum} &= I_{\text{true}} = I_{\text{ADM}}/16\pi \\ &= (1/16\pi) \int [\pi^{ij}\partial g_{ij}/\partial t - N\mathcal{H}(\pi^{ij}, g_{ij}) - N_i\mathcal{H}^i(\pi^{ij}, g_{ij})] d^4x \\ &\quad + \int \mathcal{L}_{\text{field}} d^4x. \end{aligned} \quad (21.95)$$

The action principle itself, here as always, tells one what must be fixed to make the action take on a well-defined value (if and when the action possesses an extremum). Apart from appropriate potentials having to do with fields other than geom-

Action principle says, fix 3-geometry on each face of sandwich

What a 3-geometry is

Electromagnetism gives example of momentum conjugate to "field coordinate"

etry, the only quantities that have to be fixed appear at first sight to be the values of the six g_{ij} on the initial and final spacelike hypersurfaces. However, the ADM action principle is invariant with respect to any change of coordinates $x^1, x^2, x^3 \rightarrow x^{\bar{1}}, x^{\bar{2}}, x^{\bar{3}}$ within the successive spacelike slices. Therefore the quantities that really have to be fixed on the two faces of the sandwich are the 3-geometries ${}^{(3)}\mathcal{G}'$ (on the initial hypersurface) and ${}^{(3)}\mathcal{G}$ (on the final hypersurface) and nothing more.

In mathematical terms, a 3-geometry ${}^{(3)}\mathcal{G}$ is the "equivalence class" of a set of differentiable manifolds that are isometrically equivalent to each other under diffeomorphisms. In the terms of the everyday physicist, a 3-geometry is the equivalence class of 3-metrics $g_{ij}(x, y, z)$ that are equivalent to one another under coordinate transformations. In more homely terms, two automobile fenders have one and the same 2-geometry if they have the same shape, regardless of how much the coordinate rulings painted on the one may differ from the coordinate rulings painted on the other.

To have in equation (21.95) an example of a field Lagrangian that is at the same time physically relevant and free of avoidable complications, take the case of a source-free electromagnetic field. It would be possible to take the field Lagrangian to have the standard Maxwell value,

$$(1/8\pi)(\mathbf{E}^2 - \mathbf{B}^2) \rightarrow -(1/16\pi)F_{\mu\nu}F^{\mu\nu}, \quad (21.96)$$

with

$$F_{\mu\nu} = \partial A_\nu / \partial x^\mu - \partial A_\mu / \partial x^\nu. \quad (21.97)$$

The variation of the Lagrangian with respect to the independent dynamic variables of the field, the four potentials A_α , would then immediately give the four second-order partial differential wave equations for these four potentials. However, to have instead a larger number of first-order equations is as convenient for electrodynamics as it is for geometrodynamics. One seeks for the analog of the Hamiltonian equations of particle dynamics,

$$\begin{aligned} dx/dt &= \partial H(x, p)/\partial p, \\ dp/dt &= -\partial H(x, p)/\partial x. \end{aligned} \quad (21.98)$$

One gets those equations by replacing the Lagrange integral $\int L(x, \dot{x}) dt$ by the Hamilton integral $\int [p\dot{x} - H(x, p)] dt$. Likewise, here one replaces the action integrand of (21.96) by what in flat spacetime would be

$$(1/4\pi) \left[A_{\mu,\nu} F^{\mu\nu} + \frac{1}{4} F_{\mu\nu} F^{\mu\nu} \right]. \quad (21.99)$$

In actuality, spacetime is to be regarded as not only curved but also sliced up into spacelike hypersurfaces. This $(3+1)$ split of the geometry made it desirable to split the ten geometrodynamical potentials into the six g_{ij} and the four lapse and shift functions. Here one similarly splits the four A_μ into the three components A_i of the vector potential and the scalar potential $A_0 = -\phi$ (with the sign so chosen that, in flat spacetime in a Minkowski coordinate system, $\phi = A^0$). In this notation, the

Lagrange density function, including the standard density factor $(-{}^{(4)}g)^{1/2}$ but dropping a complete time integral $(\partial/\partial t)(A_i \mathcal{E}^i)$ that has no influence on the equations of motion, is given by the formula

$$4\pi\mathcal{L}_{\text{field}} = -\mathcal{E}^i \partial A_i / \partial t + \phi \mathcal{E}^i_{,i} - \frac{1}{2} Ng^{-1/2} g_{ij} (\mathcal{E}^i \mathcal{E}^j + \mathcal{B}^i \mathcal{B}^j) + N^i [ijk] \mathcal{E}^j \mathcal{B}^k. \quad (21.100)$$

Lagrange density for electromagnetism

Here use is made of the alternating symbol $[ijk]$, defined as changing sign on the interchange of any two labels, and normalized so that $[123] = 1$. Note that the 3-tensor ϵ^{ijk} and the alternating symbol $[ijk]$ are related much as are the corresponding four-dimensional objects in equation (8.10), so that one can write

$$\mathcal{B}^i = \frac{1}{2} [ijk] (A_{k,i} - A_{j,k}). \quad (21.101)$$

The quantities \mathcal{B}^i are the components of the magnetic field in the spacelike slice. They are not regarded as independently variable. They are treated as fully fixed by the choice of the three potentials A_i . The converse is the case for the components \mathcal{E}^i of the electric field: they are treated like momenta, and as independently variable.

Extremizing the action with respect to the \mathcal{E}^i (exercise 21.11) gives the analog of the equation $dx/dt = p/m$ in particle mechanics, and the analog of the equation

$$E_i = -\partial A_i / \partial t - \partial \phi / \partial x^i \quad (21.102)$$

of flat-spacetime electrodynamics; namely,

$$-\partial A_i / \partial t - \phi_{,i} - Ng^{-1/2} g_{ij} \mathcal{E}^j - [ijk] N^j \mathcal{B}^k = 0. \quad (21.103)$$

The initial-value equation of electromagnetism

Here the last term containing the shift functions N^j , arises from the obliquity of the coordinate system. ADM give the following additional but equivalent ways to state the result (21.103):

$$\begin{aligned} \mathcal{E}^i &= \frac{1}{2} [ijk] * F_{jk} \\ &= \frac{1}{2} [ijk] \left\{ \frac{1}{2} [jkl] \mu^l (-{}^{(4)}g)^{1/2} {}^{(4)}g^{\mu\alpha} {}^{(4)}g^{\nu\beta} F_{\alpha\beta} \right\}. \end{aligned} \quad (21.104)$$

They note that \mathcal{E}^j and \mathcal{B}^j are not directly the contravariant components of the fields in the simultaneity Σ ,

$$\mathbf{E} = E^j \mathbf{e}_j, \quad \mathbf{B} = B^j \mathbf{e}_j, \quad (21.105)$$

but the contravariant densities,

$$\mathcal{E}^j = g^{1/2} E^j, \quad \mathcal{B}^j = g^{1/2} B^j. \quad (21.106)$$

Extremizing the action with respect to the three A_i (exercise 21.12) gives the curved-spacetime analog of the Maxwell equations,

$$\partial \mathbf{E} / \partial t = \nabla \times \mathbf{B}. \quad (21.107)$$

Divergence relation by extremization with respect to ϕ

Action principle tells what to fix at limits

At limits, fix not potentials but magnetic field itself

The remaining potential, ϕ , enters the action principle at only one point. Extremizing with respect to it gives immediately the divergence relation of source-free electromagnetism,

$$\mathcal{E}^i_{,i} = 0. \quad (21.108)$$

If an action principle tells in and by itself what quantities are to be fixed at the limits, what lessons does (21.100) give on this score? One can go back to the example of particle mechanics in Hamiltonian form, as in Box 21.1, and note that there the momentum p could “flap in the breeze.” Only the coordinate x had to be fixed at the limits. Thus the variation of the action was

$$\begin{aligned} \delta I &= \delta \int [p\dot{x} - H(x, p)] dt \\ &= \int \{[\dot{x} - \partial H/\partial p] \delta p + (d/dt)(p \delta x) + [-\dot{p} - \partial H/\partial x] \delta x\} dt. \end{aligned} \quad (21.109)$$

To arrive at a well-defined extremum of the action integral I , it was not enough to annul the coefficients, in square brackets, of δp and δx ; that is, to impose Hamilton’s equations of motion. It was necessary in addition to annul the quantities at limits, $p \delta x$; that is, to specify x at the start and at the end of the motion. Similarly here. The quantities ϕ and \mathcal{E}^i flap in the breeze, but the magnetic field has to be specified on the two faces of the sandwich to allow one to speak of a well-defined extremum of the action principle. Why the magnetic field, or the three quantities

$$\partial A_j / \partial x^i - \partial A_i / \partial x^j; \quad (21.110)$$

why not the three A_i themselves? When one varies (21.100) with respect to the A_i , and integrates the variation of the first term by parts, as one must to arrive at the dynamic equations, one obtains a term at limits

$$\int_{\Sigma_{\text{initial}}} \mathcal{E}^i \delta A_i d^3x - \int_{\Sigma_{\text{final}}} \mathcal{E}^i \delta A_i d^3x. \quad (21.111)$$

One demands that both these terms at limits must vanish in order to have a well-defined variational problem. Go from the given vector potential to another vector potential, $A_{i_{\text{new}}}$, by the gauge transformation

$$A_{i_{\text{new}}} = A_i + \delta A_i = A_i + \partial \lambda / \partial x^i. \quad (21.112)$$

The magnetic-field components given by the three $A_{i_{\text{new}}}$ differ in no way from those listed in (21.110). Moreover the “variation at limits,”

$$\int \mathcal{E}^i \delta A_i d^3x = \int \mathcal{E}^i \partial \lambda / \partial x^i d^3x = - \int \lambda \mathcal{E}^i_{,i} d^3x, \quad (21.113)$$

is automatically zero by virtue of the divergence condition (21.108), for any arbitrary choice of λ . Therefore the quantities fixed at limits are not the three A_i themselves (mere potentials) but the physically significant quantities (21.110), the components of the magnetic field. Moreover, the divergence condition $\mathcal{E}^i_{,i} = 0$ now becomes the initial-value equation for the determination of the potential ϕ .

In the preceding paragraph one need only replace “the three A_i ” by “the six g_{ij} ” and “the components of the magnetic field” by “the 3-geometry ${}^{(3)}\mathcal{G}$ ” and “the potential ϕ ” by “the lapse and shift functions N and N^i ” to pass from electrodynamics to geometrodynamics.

With this parallelism in view, turn back to the variational principle (21.95) of general relativity in the ADM formulation. With the 3-geometry fixed on the two faces of the sandwich, vary conditions in between to extremize the action, varying in turn the π^{ij} , the g_{ij} , and the lapse and shift functions. The geometrodynamic momenta appear everywhere only algebraically in the action principle, except in the term $-2N_i\pi^{ij}|_j$. Variation and integration by parts gives $2N_{i|j}\delta\pi^{ij}$. Collecting coefficients of $\delta\pi^{ij}$ and annuling the sum of these coefficients, one arrives at one of the several conditions required for an extremum,

$$\partial g_{ij}/\partial t = 2Ng^{-1/2} \left(\pi_{ij} - \frac{1}{2} g_{ij} \text{Tr } \boldsymbol{\pi} \right) + N_{i|j} + N_{j|i}. \quad (21.114)$$

This result agrees with what one gets from equations (21.91) defining geometrodynamic momentum in terms of extrinsic curvature, together with expression (21.67) for extrinsic curvature in terms of lapse and shift. The result (21.114) here is no less useful than the result

$$dx/dt = \partial H(x, p)/\partial p = p/m$$

in the most elementary problem in mechanics: it marks the first step in splitting a second-order equation or equations into twice as many first-order equations.

Now vary the action with respect to the g_{ij} and again, after appropriate integration by parts and rearrangement, find the remaining first-order dynamic equations of general relativity [simplified by use of equations (21.116) and (21.117)],

$$\begin{aligned} \partial\pi^{ij}/\partial t &= -Ng^{1/2} \left(R^{ij} - \frac{1}{2} g^{ij}R \right) + \frac{1}{2} Ng^{-1/2}g^{ij} \left(\text{Tr } \boldsymbol{\pi}^2 - \frac{1}{2} (\text{Tr } \boldsymbol{\pi})^2 \right) \\ &\quad - 2Ng^{-1/2} \left(\pi^{im}\pi_m^j - \frac{1}{2} \pi^{ij}\text{Tr } \boldsymbol{\pi} \right) \\ &\quad + g^{1/2}(N^{ij} - g^{ij}N^m|_m) + (\pi^{ij}N^m)|_m \\ &\quad - N^i|_m\pi^{mj} - N^j|_m\pi^{mi} + \left[\begin{array}{l} \text{source terms arising from fields} \\ \text{other than geometry, omitted here for} \\ \text{simplicity, but discussed by ADM (1962)} \end{array} \right]^{ij} \end{aligned} \quad (21.115)$$

Finally extremize the action (21.95) with respect to the lapse function N and the shift functions N_i , and find the four so-called initial-value equations of general relativity, equivalent to (21.77) and (21.81) or to $G_n^\alpha = 8\pi T_n^\alpha$; thus,

$$-(1/16\pi)\mathcal{H}(\pi^{ij}, g_{ij}) = (1/8\pi)Ng^{-1/2}g_{ij}(\mathcal{E}^i\mathcal{E}^j + \mathcal{B}^i\mathcal{B}^j), \quad (21.116)$$

$$-(1/16\pi)\mathcal{H}^i(\pi^{ij}, g_{ij}) = -(1/4\pi)[ijk]\mathcal{E}^j\mathcal{B}^k. \quad (21.117)$$

ADM principle reproduces formula for geometrodynamic momentum

Dynamic and initial-value equations out of ADM formalism

EXERCISES**Exercise 21.11. FIRST EXPLOITATION OF THE ADM VARIATIONAL PRINCIPLE FOR THE ELECTROMAGNETIC FIELD**

Extremize the action principle (21.100) with respect to the \mathcal{E}^i and derive the result (21.103).

Exercise 21.12. SECOND EXPLOITATION OF THE ADM VARIATIONAL PRINCIPLE FOR THE ELECTROMAGNETIC FIELD

Extremize (21.100) with respect to the A_i , and verify that the resulting equations in any Minkowski-flat region are equivalent to (21.107).

Exercise 21.13. FARADAY-MAXWELL SOURCE TERM IN THE DYNAMIC EQUATIONS OF GENERAL RELATIVITY

Evaluate the final indicated source terms in (21.115) from the Lagrangian (21.100) of Maxwell electrodynamics, regarded as a function of the A_i and the g_{ij} .

Exercise 21.14. THE CHOICE OF ϕ DOESN'T MATTER

Prove the statement in the text that the dynamic development of the electric and magnetic fields themselves is independent of the choice made for the scalar potential $\phi(t, x, y, z)$ in the analysis (a) in flat spacetime in Minkowski coordinates and (b) in general relativity, according to equations (21.103), and (21.107) as generalized in exercise 21.12.

Exercise 21.15. THE CHOICE OF SLICING OF SPACETIME DOESN'T MATTER

Given a metric ${}^{(3)}g_{ij}(x, y, z)$ and an extrinsic curvature $K^{ij}(x, y, z)$ on a spacelike hypersurface Σ , and given that these quantities satisfy the initial-value equations (21.116) and (21.117), and given two alternative choices for the lapse and shift functions (N, N_i) and $(N + \delta N, N_i + \delta N_i)$, show that the curvature itself (as distinguished from its components in these two distinct coordinate systems), as calculated at a point \mathcal{P} a “little way” (first order of small quantities) off the hypersurface, by way of the dynamic equations (21.114) and (21.115), is independent of this choice of lapse and shift.

§21.8. INTEGRATING FORWARD IN TIME

In the Hamiltonian formalism of Arnowitt, Deser, and Misner [see also the many papers by many workers on the quantization of general relativity—primarily putting Einstein's theory into Hamiltonian form—cited, for example, in references 1 and 2 of Wheeler (1968)], the dynamics of geometry takes a form quite similar to the Hamiltonian dynamics of geometry. There one gives x and p at a starting time and integrates two first-order equations for dx/dt and dp/dt ahead in time to find these dynamically conjugate variables at all future times. Here one gives appropriate values of g_{ij} and π^{ij} over an initial spacelike hypersurface and integrates the two first-order equations (21.114) and (21.115) ahead in time to find the geometry at future times. For example, one can rewrite the differential equations as difference equations according to the practice by now familiar in modern hydrodynamics, and then carry out the integration on an electronic digital computer of substantial memory capacity.

Time in general relativity has a many-fingered quality very different from the one-parameter nature of time in nonrelativistic particle mechanics [see, however, Dirac, Fock, and Podolsky (1932) for a many-time formalism for treating the relativistic dynamics of a system of many interacting particles]. He who is studying the geometry is free to push ahead the spacelike hypersurface faster at one place than another, so long as he keeps it spacelike. This freedom expresses itself in the lapse function $N(t, x, y, z)$ at each stage, t , of the integration. Equations (21.114) and (21.115) are not a conduit to feed out information on N to the analyst. They are a conduit for the analyst to feed in information on N . The choice of N is to be made, not by nature, but by man. The dynamic equations cannot begin to fulfill their purpose until this choice is made. The “time parameter” t is only a label to distinguish one spacelike hypersurface from another in a one-parameter family of hypersurfaces; but N thus tells the spacing in proper time, as it varies from place to place, between the successive slices on which one chooses to record the time-evolution of the geometry. A cinema camera can record what happens only one frame at a time, but the operator can make a great difference in what that camera sees by his choice of angle for the filming of the scene. So here, with the choice of slicing.

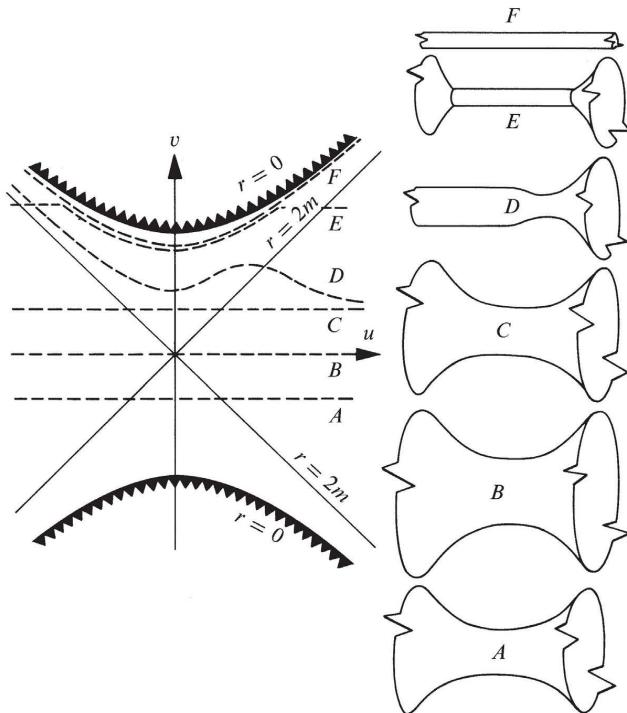
Another choice is of concern to the analyst, especially one doing his analysis on a digital computer. He is in the course of determining, via (21.114–21.115) written as difference equations, what happens on a lattice work of points, typified by $x = \dots, 73, 74, 75, 76, 77, \dots$, etc. He finds that the curvatures are developing most strongly in a localized region in the range around $x = 83$ to $x = 89$. He wants to increase the density of coverage of his tracer points in this region. He does so by causing points at lesser and greater x values to drift into this region moment by moment as t increases: $t = \dots, 122, 123, 124, \dots$ He makes the tracer points at lesser x -values start to move to the right (N_1 positive) and points at greater x -values move to the left (N_1 negative). In other words, the choice of the three shift functions $N_i(t, x, y, z)$ is just as much the responsibility of the analyst as is the choice of the lapse function N . The equations will never tell him what to pick. He has to tell the equations.

These options, far from complicating dynamic equations (21.114–21.115), make them flexible and responsive to the wishes of the analyst in following the course of whatever geometrodynamical process is in his hands for study.

The freedom that exists in general relativity in the choice of the four functions N, N_i , is illuminated from another side by comparing it with the freedom one has in electrodynamics to pick the one function $\phi(t, x, y, z)$, the scalar potential. In no way do the dynamic Maxwell equations (21.103) and (21.107), as generalized in exercise 21.12 determine ϕ . Instead they demand that it be determined (by the analyst) as the price for predicting the time-development of the vector potential A_i . An altered choice of $\phi(t, x, y, z)$ in its dependence on position and time means altered results from the dynamic equations for the development of the three A_i in time and space. However, the physically significant quantities, the electric and magnetic fields themselves on successive hypersurfaces, come out the same (exercise 21.14) regardless of this choice of ϕ . Similarly in geometrodynamics: an altered choice for the four

Lapse and shift chosen to push forward the integration in time as one finds most convenient

Same 4-geometry regardless of lapse and shift options

**Figure 21.4.**

Some of the many ways to make distinct spacelike slices through one and the same $(^4)\mathcal{G}$, the complete Schwarzschild 4-geometry.

functions N, N_i , means (a) an altered laying down of coordinates in spacetime, and therefore (b) altered results for the intrinsic metric $(^3)g_{ij}$ and extrinsic curvature K^{ij} of successive spacelike hypersurfaces, but yields the same 4-geometry $(^4)\mathcal{G}$ (Figure 21.4) regardless of this choice of coordinatization (exercise 21.15).

§21.9. THE INITIAL-VALUE PROBLEM IN THE THIN-SANDWICH FORMULATION

Initial-value data: what is freely disposable? and what is thereby fixed?

Given appropriate initial-value data, one can integrate the dynamic equations ahead in time and determine the evolution of the geometry; but what are “appropriate initial-value data”? They are six functions $(^3)g_{ij}(x, y, z)$ plus six more functions $\pi^{ij}(x, y, z)$ or $K^{ij}(x, y, z)$ that together satisfy the four initial-value equations (21.116) and (21.117). To be required to give coordinates and momenta accords with the familiar plan of Hamiltonian mechanics; but to have consistency conditions or “constraints” imposed on such data is less familiar. A particle moving in two-dimensional space is catalogued by coordinates x, y , and coordinates p_x, p_y ; but a particle forced to remain on the circle $x^2 + y^2 = a^2$ satisfies the constraint $xp_x + yp_y = 0$. Thus the existence of a “constraint” is a signal that the system possesses fewer degrees

of freedom than one would otherwise suppose. Fully to analyze the four “initial-value” or “constraint” conditions (21.116) and (21.117) is thus to determine (1) how many dynamic degrees of freedom the geometry possesses and (2) what these degrees of freedom are; that is to say, precisely what “handles” one can freely adjust to govern completely the geometry and its evolution with time. The counting one can do today, with the conclusion that the geometry possesses the same count of true degrees of freedom as the electromagnetic field. The identification of the “handles,” or freely adjustable features of the dynamics, is less advanced for geometry than it is for electromagnetism (Box 21.2), but most instructive so far as it goes.

By rights the identification of the degrees of freedom of the field, whether that of Einstein or that of Faraday and Maxwell, requires nothing more than knowing what must be fixed on initial and final spacelike hypersurfaces to make the appropriate variation principle well-defined. One then has the option whether (1) to give that quantity on both hypersurfaces or (2) to give that quantity and its dynamic conjugate on one hypersurface or (3) to give the quantity on both hypersurfaces, as in (1), but go to the limit of an infinitely thin sandwich, so that one ends up specifying the quantity and its time rate of change on one hypersurface. This third “thin sandwich” procedure is simplest for a quick analysis of the initial-value problem in both electrodynamics and geometrodynamics. Take electrodynamics first, as an illustration.

Give the divergence-free magnetic field and its time-rate of change: on an arbitrary smooth spacelike hypersurface in curved spacetime in the general case; on the hypersurface $t = 0$ in Minkowski spacetime in the present illustrative treatment,

$$\mathcal{B}^i(0, x, y, z) \text{ given,} \quad (21.118)$$

$$\dot{\mathcal{B}}^i(0, x, y, z) = \left(\frac{\partial \mathcal{B}^i}{\partial t} \right) \text{ also given.} \quad (21.119)$$

In electromagnetism, give magnetic field and its rate of change as initial data

These quantities together contain four and only four independent data per space point. How is one now to obtain the momenta $\pi^i \sim -\mathcal{E}^i$ so that one can start integrating the dynamic equations (21.103) and (21.107) forward in time? (1) Find a set of three functions $A_i(0, x, y, z)$ such that their curl gives the three specified \mathcal{B}^i . That this can be done at all is guaranteed by the vanishing of the divergence $\mathcal{B}^i_{,i}$. However, the choice of the A_i is not unique. The new set of potentials $A_{i,\text{new}} = A_i + \partial\lambda/\partial x^i$ with arbitrary smooth λ , provide just as good a solution as the original A_i . No matter. Pick one solution and stick to it. (2) Similarly, find a set of three $\dot{A}_i(0, x, y, z)$ such that their curl gives the specified $\dot{\mathcal{B}}^i(0, x, y, z)$, and resolve all arbitrariness of choice by *fiat*. (3) Recall that the electric field (negative of the field momentum) is given by

$$\mathcal{E}_i = -\dot{A}_i - \partial\phi/\partial x^i \quad (21.120)$$

(formula valid without amendment only in flat space). The initial-value or constraint equation $\mathcal{E}^i_{,i} = 0$ translates to the form

$$\nabla^2\phi = -\eta^{ij}\dot{A}_{i,j}. \quad (21.121)$$

Box 21.2 COUNTING THE DEGREES OF FREEDOM OF THE ELECTROMAGNETIC FIELD
A. First Approach: Number of "Field Coordinates" per Spacepoint

Superficial tally of the degrees of freedom of the source-free electromagnetic field gives three field coordinates $A_i(x, y, z)$ per spacepoint on the initial simultaneity Σ , plus three field momenta $\pi_{\text{true}}^i = \pi^i/4\pi$ [with $\pi^i = -\mathcal{E}^i(x, y, z)$] per spacepoint.

Closer inspection reveals that the number of coordinate degrees of freedom per spacepoint is not three but two. Thus the change in vector potential $A_i \rightarrow A_i + \partial\lambda/\partial x^i$ makes no change in the actual physics, the magnetic field components,

$$B^i = \frac{1}{2} [ijk] (\partial A_k / \partial x^j - \partial A_j / \partial x^k).$$

Moreover, though those components are three in number, they satisfy one condition per spacepoint, $\mathcal{B}^i,_i = 0$, thus reducing the effective net number of coordinate degrees of freedom per spacepoint to two.

The momentum degrees of freedom per spacepoint are likewise reduced from three to two by the one condition per spacepoint $\mathcal{E}^i,_i = 0$.

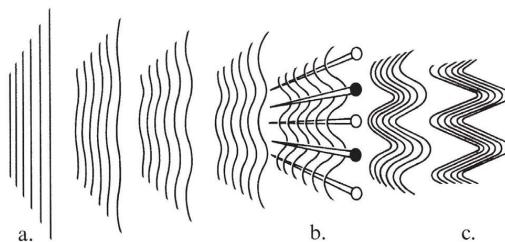
B. Alternative Approach: Count Fourier Coefficients

In textbooks on field theory [see, for example, Wentzel (1949)], attention focuses on flat spacetime. The electromagnetic field is decomposed by Fourier analysis into individual running waves. Instead of counting degrees of freedom per point in coordinate space, one does the equivalent: counts up degrees of freedom per point in wave-number space. Thus for each (k_x, k_y, k_z) , there are two independent states of polarization. Each state of polarization requires for its description an amplitude ("coordinate") and time-rate of change of amplitude ("momentum") at the initial time, t'_0 . Thus the number of degrees of freedom per point in wave-number space is two for coordinates and two for momenta, in accord with what one gets by carrying out the count in coordinate space.

In curved spacetime, Fourier analysis is a less convenient way of identifying the degrees of freedom of the electromagnetic field [for such a Fourier analysis, see Misner and Wheeler (1957), especially their Table X and following text] than direct analysis in space, as above.

C. Another Alternative: Analyze "Deformation of Structure"

Still a third way to get a handle on the degrees of freedom of a divergence-free field, whether \mathcal{E} or \mathcal{B} , rests on the idea of deformation of structure [diagram from Wheeler (1964)]. Represent the



magnetic field by Faraday's picture of lines of force (a) continuing through space without ever ending, automatic guarantee that $\mathcal{B}^i,_i$ is everywhere zero. Insert "knitting needles" (b) into the spaghetti-like structure of the lines of force and move these needles as one will. Sliding the "knitting needles" along a line of force causes no movement of the line of force. (c) With the help of two knitting needles perpendicular to each other and to the line of force, one can give any given line of force any small displacement one pleases perpendicular to its length: again two degrees of freedom per spacepoint. Granted any non-zero field to begin with, no matter how small, one can build it up by a sequence of such small deformations to agree with any arbitrary field pattern of zero divergence, no matter what its complexity and strength may be.

Solve for ϕ . Then (4) equation (21.120) gives the initial-value electric field, or electrodynamic field momentum $\pi^i \sim -\mathcal{E}^i$, required (along with the field coordinate A_i) for starting the integration of the dynamic equations (21.103) and (21.107). [Misner and Wheeler (1957) deal with the additional features that come in when the space is multiply connected. Each wormhole or handle of the geometry is able to trap electric lines of force. The flux trapped in any one wormhole defines the classical electric charge q_w associated with that wormhole. One has to specify all these charges once and for all in addition to the data (21.118) and (21.119) in order to determine fully the dynamic evolution of the electromagnetic field. There is no geometrodynamical analog to electric charge, according to Unruh (1971).] (5) In this integration, the scalar potential ϕ at each subsequent time step is not to be calculated; it is to be chosen. Only when one has made this free choice definite do the dynamic equations come out with definite results for the A_i and the π^i or \mathcal{E}^i at these successive steps.

In the thin-sandwich formulation of the initial-value problem of electrodynamics, to summarize, one gives \mathcal{B}^i and $\dot{\mathcal{B}}^i$ (equivalent to \mathcal{B} on two nearby hypersurfaces). One chooses the A_i and \dot{A}_i with much arbitrariness to represent these initial-value data. The arbitrariness having been seized on to give the initial A_i and \dot{A}_i , there is no arbitrariness left in the initial ϕ . However, at all subsequent times the situation is just the other way around. All the arbitrariness is sopped up in the choice of the ϕ , leaving no arbitrariness whatever in the three A_i (as given by the integration of the dynamic equation).

The situation is quite similar in geometrodynamics. One gives the beginnings of a 1-parameter family of spacelike hypersurfaces; namely,

Scalar potential: fixed at start; freely disposable later

In ADM treatment, give 3-geometry and its time-rate of change

$${}^{(3)}\mathcal{G}(0) \text{ given,} \quad (21.122)$$

$${}^{(3)}\dot{\mathcal{G}}(0) = \frac{\partial {}^{(3)}\mathcal{G}}{\partial t} \text{ given,} \quad (21.123)$$

Then (1) one picks a definite set of coordinates $x^i = (x, y, z)$ and in terms of those coordinates finds the unique metric coefficients $g_{ij}(x, y, z)$ that describe that 3-geometry. The existence of a solution is guaranteed by the circumstance that ${}^{(3)}\mathcal{G}$ is a Riemannian geometry. However, one could have started with different coordinates and ended up with different metric coefficients for the description of the same 3-geometry. No matter. Pick one set of coordinates, take the resulting metric coefficients, and stick to them as giving half the required initial-value data. (2) Similarly, to describe the 3-geometry ${}^{(3)}\mathcal{G} + {}^{(3)}\dot{\mathcal{G}} dt$ at the value of the parameter $t + dt$, make use of coordinates $x^i + \dot{x}^i dt$ and arrive at the metric coefficients $g_{ij} + \dot{g}_{ij} dt$. The arbitrariness in the x^i having thus been resolved by fiat, and the ${}^{(3)}\mathcal{G}$ being given as definite initial physical data, the g_{ij} are thereby completely fixed. (3) Recall that the components of the extrinsic curvature K_{ij} or the momenta π^{ij} are given in terms of the g_{ij} and \dot{g}_{ij} and the lapse and shift functions N and N_i by (21.67) or by (21.67) plus (21.91) or by (21.114). The four initial-value or “constraint” equations (21.116) and (21.117) thus become four conditions for finding the four

quantities N, N_i . One can shorten the writing of these conditions by introducing the abbreviations

$$\gamma_{ij} = \frac{1}{2} [N_{i|j} + N_{j|i} - \partial g_{ij}/\partial t] \quad (21.124)$$

and

$$\gamma_2 = \begin{pmatrix} \text{"shift"} \\ \text{"anomaly"} \end{pmatrix} = (\text{Tr } \gamma)^2 - \text{Tr } \gamma^2 \quad (21.125)$$

(both for functions of x, y, z on the initial simultaneity). Then one has

$${}^{(3)}R + \gamma_2/N^2 = 16\pi T_{nn} = 16\pi T^{nn} \quad (21.126)$$

for the one initial-value equation; and for the other three,

$$\left[\frac{\gamma_i^k - \delta_i^k \text{Tr } \gamma}{N} \right]_{|k} = -8\pi T_i^n. \quad (21.127)$$

Lapse and shift initially determinate; thereafter freely disposable

In summary, one chooses the g_{ij} and \dot{g}_{ij} with much arbitrariness (because of the arbitrariness in the coordinates, not by reason of any arbitrariness in the physics) to represent the given initial-value data, ${}^{(3)}g$ and ${}^{(3)}\dot{g}$. The arbitrariness at the initial time all having been soaked up in this way, one expects no arbitrariness to be left in the initial N and N_i as obtained by solving (21.126) and (21.127). However, on all later spacelike slices, the award of the arbitrariness is reversed. The lapse and shift functions are freely disposable, but, with them once chosen, there is no arbitrariness whatever in the six g_{ij} (and the six K^{ij} or π^{ij}) as given by the integration of the dynamic equations (21.114) and (21.115). The analogy with electrodynamics is clear. There the one "gauge-controlled" function ϕ was fixed at the start by the elliptic equation (21.121), but was thereafter free. Here the four lapse and shift functions are fixed at the start by the four equations (21.126) and (21.127), but are thereafter free.

Exercise 21.16 applies the initial-value equation (21.126) to analyze the whole evolution in time of any Friedmann universe in which one knows the equation $p = p(\rho)$ connecting pressure with density. Exercise 21.17 looks for a variation principle on the spacelike hypersurface Σ equivalent in content to the elliptic initial-value equation (21.121) for the scalar potential ϕ . Exercises 21.18 and 21.19 look for similar variation principles to determine the lapse and shift functions.

Counting initial-value data

How many degrees of freedom, or how many "handles," are there in the specification of the 4-geometry that one will obtain? The metric coefficients of the initial 3-geometry provided six numbers per space point. However, they were arbitrary to the extent of a coordinate transformation, specified by three functions of position,

$$\begin{aligned} x &= x(x', y', z'), \\ y &= y(x', y', z'), \\ z &= z(x', y', z'). \end{aligned}$$

The net number of quantities per space point with any physical information was therefore $6 - 3 = 3$. One can visualize these three functions as the three diagonal components of the metric in a coordinate system in which g_{ij} has been transformed to diagonal form. Ordinarily it is not useful to go further and actually spell out the analysis in any such narrowly circumscribed coordinate system.

Now think of the $(^3)\mathcal{G}$ in question as imbedded in the $(^4)\mathcal{G}$ that comes out of the integrations. Moreover, think of that $(^4)\mathcal{G}$ as endowed with the lumps, bumps, wiggles, and waves that distinguish it from other generic 4-geometries and that make Minkowski geometry and special cosmologies so unrepresentative. The $(^3)\mathcal{G}$ is a slice in that $(^4)\mathcal{G}$. It partakes of the lumps, bumps, wiggles, and waves present in all those regions of the $(^4)\mathcal{G}$ that it intersects. To the extent that the $(^4)\mathcal{G}$ is generic, it does not allow the $(^3)\mathcal{G}$ to be moved to another location without becoming a different $(^3)\mathcal{G}$. If one tries to push the $(^3)\mathcal{G}$ “forward in time” a little in a certain locality, leaving it unchanged in location elsewhere, one necessarily changes the $(^3)\mathcal{G}$. By this circumstance, one sees that the $(^3)\mathcal{G}$ “carries information about time” [Sharp (1960); Baierlein, Sharp, and Wheeler (1962)]. Moreover, this “forward motion in time” demands for its description one number per space point. It is possible to think of this number in concrete terms by imagining an arbitrary coordinate system $\bar{t}, \bar{x}, \bar{y}, \bar{z}$ laid down in the $(^4)\mathcal{G}$. Then the hypersurface can be conceived as defined by the value $\bar{t} = \bar{t}(\bar{x}, \bar{y}, \bar{z})$ at which it cuts the typical line $\bar{x}, \bar{y}, \bar{z}$. A forward movement carries it to $\bar{t}(\bar{x}, \bar{y}, \bar{z}) + \delta\bar{t}(\bar{x}, \bar{y}, \bar{z})$, and changes shape and metric coefficients on $(^3)\mathcal{G}$ accordingly. It is usually better not to tie one’s thinking down to such a concrete model, but rather to recognize as a general point of principle (1) that the location of the $(^3)\mathcal{G}$ in spacetime demands for its specification one datum per spacepoint, and (2) that this datum is already willy-nilly present in the three data per spacepoint that mark any $(^3)\mathcal{G}$.

In conclusion, there are only two data per spacepoint in a $(^3)\mathcal{G}$ that really tell anything about the $(^4)\mathcal{G}$ in which it is imbedded, or to be imbedded (as distinguished from where the $(^3)\mathcal{G}$ slices through that $(^4)\mathcal{G}$). Similarly for the other $(^3)\mathcal{G}$ that defines the other “face of the sandwich,” whether thick or thin. Thus one concludes that the specification of $(^3)\mathcal{G}$ and $(^3)\mathcal{G}$ actually gives four net pieces of dynamic information per spacepoint about the $(^4)\mathcal{G}$ (all the rest of the information being “many-fingered time,” telling where the 3-geometries are located in that $(^4)\mathcal{G}$). According to this line of reasoning, geometrodynamics has the same number of dynamic degrees of freedom as electrodynamics. One arrives at the same conclusion in quite another way through the weak-field analysis (§35.3) of gravitational waves on a flat spacetime background: the same ranges of possible wave numbers as for Maxwell waves; and for each wave number two states of polarization; and for each polarization one amplitude and one phase (the equivalent of one coordinate and one momentum).

In electrodynamics in a prescribed spacetime manifold, one has a clean separation between the one time-datum per spacepoint (when one deals with electromagnetism in the context of many-fingered time) and the two dynamic variables per spacepoint; but not so in the superspace formulation of geometrodynamics. There the two kinds of quantities are inextricably mixed together in the one concept of 3-geometry.

Four pieces of
geometrodynamic information
per space point on initial
simultaneity

Turn from initial- and final-value data to the action integral that is determined by (1) these data and (2) the principle that the action be an extremum,

$$I = I_{\text{extremum}} = S.$$

The action depends on the variables on the final hypersurface, according to the formula

$$S = S(\Sigma, \mathcal{B}) \quad (21.128)$$

in electrodynamics, but according to the formula

$$S = S^{(3)\mathcal{G}} \quad (21.129)$$

in geometrodynamics. In each case, there are three numbers per spacepoint in the argument of the functional (one in Σ ; two in a divergence-free magnetic field; three in $(3)\mathcal{G}$).

This mixing of the one many-fingered time and the two dynamic variables in a 3-geometry makes it harder in general relativity than in Maxwell theory to know when one has in hand appropriate initial value data. Give Σ and give \mathcal{B} and $\dot{\mathcal{B}}$ on Σ : that was enough for electrodynamics. For geometrodynamics, to give the six $g_{ij}(x, y, z)$ and the six $\dot{g}_{ij}(x, y, z)$ is not necessarily enough. For example, let the time parameter t be a fake, so that dt , instead of leading forward from a given hypersurface Σ to a new hypersurface $\Sigma + d\Sigma$, merely recoordinatizes the present hypersurface:

$$\begin{aligned} x^i &\longrightarrow x^i - \xi^i dt, \\ g_{ij} &\longrightarrow g_{ij} + (\xi_{i|j} + \xi_{j|i}) dt. \end{aligned} \quad (21.130)$$

A first inspection may make one think that one has adequate data in the six g_{ij} and the six

$$\dot{g}_{ij} = \xi_{i|j} + \xi_{j|i}, \quad (21.131)$$

but in the end one sees that one has not both faces of the thin sandwich, as required, but only one. Thus one must reject, as improperly posed data in the generic problem of dynamics, any set of six \dot{g}_{ij} that let themselves be expressed in the form (21.131) [Belasco and Ohanian (1969)].

Similar difficulties occur when the two faces of the thin sandwich, instead of coinciding everywhere, coincide in a limited region, be it three-dimensional, two-dimensional, or even one-dimensional (“crossover of one face from being earlier than the other to being later”). Thus it is enough to have (21.131) obtaining even on only a curved line in Σ to reject the six g_{ij} as inappropriate initial-value data.

That one can impose conditions on the g_{ij} and \dot{g}_{ij} which will guarantee existence and uniqueness of the solution $N(x, y, z)$, $N_i(x, y, z)$ of the initial-value equations (21.126) and (21.127) is known as the “thin-sandwich conjecture,” a topic on which there has been much work by many investigators, but so far no decisive theorem.

Problem in assuring completeness and consistency of initial data

The “thin sandwich conjecture”

To presuppose existence and uniqueness is to make the first step in giving mathematical content to Mach's principle that the distribution of mass-energy throughout space determines inertia (§21.12).

§21.10. THE TIME-SYMMETRIC AND TIME-ANTISYMMETRIC INITIAL-VALUE PROBLEMS

Turn from the general initial-value problem to two special initial-value problems that lend themselves to detailed treatment, one known as the time-symmetric initial-value problem, the other as the time-antisymmetric problem.

A 4-geometry is said to be time-symmetric when there exists a spacelike hypersurface Σ at all points of which the extrinsic curvature vanishes. In this case the three initial value equations (21.127) are automatically satisfied, and the fourth reduces to a simple requirement on the three-dimensional scalar curvature invariant,

$$R = 16\pi\rho. \quad (21.132)$$

Still further simplifications result when one limits attention to empty space. Simplest of all is the case of spherical symmetry in which (21.132) yields at once the full Schwarzschild geometry at the moment of time symmetry (two asymptotically flat spaces connected by a throat), as developed in exercise 21.20.

Consider a 3-geometry with metric

$$ds_1^2 = g_{(1)ik} dx^i dx^k. \quad (21.133)$$

Call it a "base metric." Consider another 3-geometry with metric

$$ds_2^2 = \psi^4(x^i) ds_1^2. \quad (21.134)$$

Angles are identical in the two geometries. On this account they are said to be conformally equivalent. The scalar curvature invariants of the two 3-geometries are related by the formula [Eisenhart (1926)]

$$R_2 = -8\psi^{-5} \nabla_1^2 \psi + \psi^{-4} R_1, \quad (21.135)$$

where

$$\nabla_1^2 \psi = \psi_{|i}^{||i} = g_1^{-1/2} (\partial/\partial x^i) [g_1^{1/2} g^{ik} (\partial\psi/\partial x^k)] \quad (21.136)$$

Demand that the scalar curvature invariant R_2 vanish, and arrive [Brill (1959)] at the "wave equation"

$$\nabla_1^2 \psi - (R_1/8)\psi = 0 \quad (21.137)$$

for the conformal correction factor ψ . Brill takes the base metric to have the form suggested by Bondi,

$$ds_1^2 = e^{2Aq_1(\rho, z)} (dz^2 + d\rho^2) + \rho^2 d\phi^2, \quad (21.138)$$

and takes the conformal correction factor ψ also to possess axial symmetry. In the application:

- $q_1(\rho, z)$ measures the “distribution of gravitational wave amplitude,” assumed for simplicity to vanish outside $r = (\rho^2 + z^2)^{1/2} = a$;
- A measures the “amplitude of the distribution of gravitational wave amplitude”;
- $\psi(\rho, z)$ is the conformal correction factor, which varies with position at large distances as $1 + (m/2r)$. The quantity $m(\text{cm})$ is uniquely determined by the condition that the geometry be asymptotically flat. It measures the mass-energy of the distribution of gravitational radiation.

Wave amplitude determines mass-energy: $m = m(A)$

“Time-antisymmetric” initial-value data

The mass m of the gravitational radiation is proportional to A^2 for small values of the amplitude A . It is inversely proportional to the reduced wavelength $\lambda = (\text{effective wavelength}/2\pi)$ that measures the scale of rapid variations in the gravitational wave amplitude $q_1(\rho, z)$ in the “active zone.” Thus the metric is dominated by wiggles, proportional in amplitude to A , in the active zone, and at larger distances dominated by something close to a Schwarzschild $(1 + 2m/r)$ factor in the metric. When the amplitude A is increased, a critical value is attained, $A = A_{\text{crit}}$, at which m goes to infinity and the geometry curves up into closure (“universe closed by its own content of gravitational-wave energy”). Further analysis and examples will be found in Wheeler (1964a), pp. 399–451, also in Wheeler (1964c).

Brill has carried out a similar analysis [Brill (1961)] for the vacuum case of what he calls time-antisymmetric initial-value conditions, sketched below as amended by York (1973). (1) The initial slice is maximal, $\text{Tr } \mathbf{K} = 0$. (2) This slice is conformally flat,

$$g_{ij} = \psi^4 \delta_{ij}. \quad (21.139)$$

(3) Work in the “base space” with metric δ_{ij} and afterwards transform to the geometry (21.139). Three of the initial-value equations become

$$K_{\text{base}, j}^{ij} = 0. \quad (21.140)$$

To solve these equations, (1) take any localized trace-free symmetric tensor B_{km} ; (2) solve the flat-space Laplace equation $\nabla^2 A = (3/2) \partial^2 B_{km} / \partial x^k \partial x^m$ for A ; (3) define the six potentials $A_{km} = B_{km} + \frac{1}{3} A \delta_{km}$; and (4) calculate

$$K_{\text{base}}^{ij} = [ik\ell][jmn] \partial^2 A_{km} / \partial x^\ell \partial x^n, \quad (21.141)$$

that automatically satisfy (21.140) and give $\text{Tr } \mathbf{K}_{\text{base}} = 0$. Then $K^{ij} = \psi^{-10} K_{\text{base}}^{ij}$ also automatically satisfies these conditions, but now in the *curved* geometry (21.139). The final initial-value equation becomes a quasilinear elliptic equation, in the flat base space, for the conformal factor ψ ,

$$8\nabla_{\text{base}}^2 \psi + \psi^{-7} \sum_{i,j} (K_{\text{base}, ij})^2 = 0. \quad (21.142)$$

The asymptotic form of ψ reveals that the mass of the wave is positive.

In addition to the time-symmetric and time-antisymmetric cases, there are at least two further cases where the initial-value problem possess special simplicity. One is the case of a geometry endowed with a symmetry, as, for example, for the Friedmann universe of Chapter 27 or the mixmaster universe of Chapter 30 or cylindrical gravitational waves in the treatment of Kuchař (1971a). One starts with a spacelike slice on which the g_{ij} and π^{ij} have a special symmetry, and makes all future spacelike slices in a way that preserves this symmetry. The geometry on any one of these simultaneities, though almost entirely governed by these symmetry considerations, still typically demands some countable number of parameters for its complete determination, such as the radius of the Friedmann universe, or the three principal radii of curvature of the mixmaster universe. These parameters and the momenta conjugate to them define a miniphase space. In this miniphase space, the dynamics runs its course as for any other problem of classical dynamics [see, for example, Box 30.1 and Misner (1969) for the mixmaster universe; Kuchař (1971a) and (1972) for waves endowed with cylindrical symmetry; and Gowdy (1973) for waves with spherical symmetry]. Even the evidence for the existence of many-fingered time, most characteristic feature of general relativity, is suppressed as the price for never having to give attention to any spacelike slice that departs from the prescribed symmetry.

Finite dimensional dynamics
for geometries endowed with
high symmetry

Exercise 21.16. POOR MAN'S WAY TO DO COSMOLOGY

Consider a spacetime with the metric

$$ds^2 = -dt^2 + a^2(t)[d\chi^2 + \sin^2\chi(d\theta^2 + \sin^2\theta d\phi^2)],$$

corresponding to a 3-geometry with the form of a sphere of radius $a(t)$ changing with time. Show that the tensor of extrinsic curvature as expressed in a local Euclidean frame of reference is

$$\mathbf{K} = -a^{-1}(da/dt)\mathbf{1},$$

where $\mathbf{1}$ is the unit tensor. Show that the initial value equation (21.77) reduces to

$$(6/a^2)(da/dt)^2 + (6/a^2) = 16\pi\rho(a)$$

[for the value of the second term on the left, see exercise 14.3 and Boxes 14.2 and 14.5], and explain why it is appropriate to write the term on the right as $6a_0/a^3$ for a “dust-filled model universe.” More generally, given any equation of state, $p = p(\rho)$, explain how one can find $\rho = \rho(a)$ from

$$d(\rho a^3) = -p d(a^3);$$

and how one can thus forecast the history of expansion and recontraction, $a = a(t)$.

**Exercise 21.17. THIN-SANDWICH VARIATIONAL PRINCIPLE FOR
THE SCALAR POTENTIAL IN ELECTRODYNAMICS**

(a) Choose the unknown U^m in the expression

$$\frac{1}{8\pi} g^{mn} \frac{\partial\phi}{\partial x^m} \frac{\partial\phi}{\partial x^n} + U^m \frac{\partial\phi}{\partial x^n}$$

EXERCISES

in such a way that this expression, multiplied by the volume element $g^{1/2} d^3x$, and integrated over the simultaneity Σ , is extremized by a ϕ , and only by a ϕ , that satisfies the initial-value equation (21.108) of electrodynamics.

(b) Show that the resulting variational principle, instead of having to be invented “out of the blue,” is none other than what follows directly from the action principle build on the Lagrangian density (21.100) of electrodynamics (independent variation of ϕ and the three A_i everywhere between the two faces of a sandwich to extremize I , subject only to the prior specification of the A_i on the two faces of the sandwich, in the limit where the thickness of the sandwich goes to zero).

Exercise 21.18. THIN-SANDWICH VARIATIONAL PRINCIPLE FOR THE LAPSE AND SHIFT FUNCTIONS IN GEOMETRODYNAMICS

- (a) Extremize the action integral

$$I_3 = \int \{ [R - (\text{Tr}K)^2 + \text{Tr}K^2 - 2T_{nn}^*]N - 2T_n^{*k}N_k \} g^{1/2} d^3x$$

with respect to the lapse and shift functions, and show that one arrives in this way at the four initial-value equations of geometrodynamics. It is understood that one has given the six g_{ij} and the six $\partial g_{ij}/\partial t$ on the simultaneity where the analysis is being done. The extrinsic curvature is considered to be expressed as in (21.67) in terms of these quantities and the lapse and shift. The energy density and energy flow are referred to a unit normal vector n and three arbitrary coordinate basis vectors e_i within the simultaneity, as earlier in this chapter, and the asterisk is an abbreviation for an omitted factor of 8π .

(b) Derive this variational principle from the ADM variational principle by going to the limit of an infinitesimally thin sandwich [see derivation in Wheeler (1964)].

Exercise 21.19. CONDENSED THIN-SANDWICH VARIATIONAL PRINCIPLE

- (a) Extremize the action I_3 of the preceding exercise with respect to the lapse function N .
(b) What is the relation between the result and the principle that “3-geometry is a carrier of information about time”?
(c) By elimination of N , arrive at a “condensed thin-sandwich variational principle” in which the only quantities to be varied are the three shift functions N_i .

Exercise 21.20. POOR MAN’S WAY TO SCHWARZSCHILD GEOMETRY

On curved empty space evolving deterministically in time, impose the conditions (1) that it possess a moment of time-symmetry, a spacelike hypersurface, the extrinsic curvature of which, with respect to the enveloping spacetime, is everywhere zero, and (2) that this spacelike hypersurface be endowed with spherical symmetry. Write the metric of the 3-geometry in the form

$$ds^2 = \psi^4(\bar{r})(d\bar{r}^2 + \bar{r}^2 d\theta^2 + \bar{r}^2 \sin^2\theta d\phi^2).$$

From the initial-value equation (21.127), show that the conformal factor ψ up to a multiplicative factor must have the form $\psi = (1 + m/2\bar{r})$. Show that the proper circumference $2\pi\bar{r}\psi^2(\bar{r})$ assumes a minimum value at a certain value of \bar{r} , thus defining the *throat* of the 3-geometry. Show that the 3-geometry is mirror-symmetric with respect to reflection in this throat in the sense that the metric is unchanged in form under the substitution $r' = m^2/4\bar{r}$. Find the transformation from the conformal coordinate \bar{r} to the Schwarzschild coordinate r .

§21.11. YORK'S "HANDLES" TO SPECIFY A 4-GEOMETRY

On a simultaneity—or on *the* simultaneity—of extremal proper volume, give the conformal part of the 3-geometry and give the two inequivalent components of the dynamically conjugate momentum in order (1) to have freely specifiable, but also complete, initial-value data and thus (2) to determine completely the whole generic four-dimensional spacetime manifold. This in brief is York's extension (1971, 1972b) to the generic case of what Brill did for special cases (see the preceding section). York and Brill acknowledge earlier considerations of Lichnerowicz (1944) and Bruhat (1962 and earlier papers cited there on conformal geometry and the initial-value problem). But why conformal geometry, and why pick such a special spacelike hypersurface on which to give the four dynamic data per spacepoint?

Few solutions of Maxwell's equations are simpler than an infinite plane monochromatic wave in Minkowski's flat spacetime, and few look more complex when examined on a spacelike slice cut through that spacetime in an arbitrary way, with local wiggles and waves, larger-scale lumps and bumps, and still larger-scale general curvatures. No one who wants to explore electrodynamics in its evolution with many-fingered time can avoid these complexities; and no one will accept these complexities of many-fingered time who wants to see the degrees of freedom of the electromagnetic field in and by themselves exhibited in their neatest form. He will pick the simplest kind of timelike slice he can find. On that simultaneity, there are two and only two field coordinates, and two and only two field momenta per spacepoint. Similarly in geometrodynamics.

When one wants to untangle the degrees of freedom of the geometry, as distinct from analyzing the dynamics of the geometry, one therefore retreats from the three items of information per spacepoint that are contained in a 3-geometry [or in any other way of analyzing the geometrodynamics, as especially seen in the "extrinsic time" formulation of Kuchař (1971b and 1972)] and following York (1) picks the simultaneity to have maximal proper volume and (2) on this simultaneity specifies the two "coordinate degrees of freedom per spacepoint" that are contained in the conformal part of the 3-geometry.

An element of proper volume $g^{1/2} d^3x$ on the spacelike hypersurface Σ undergoes, in the next unit interval of proper time as measured normal to the hypersurface, a fractional increase of proper volume [see Figure 21.3 and equations 21.59 and 21.66] given by

$$-\text{Tr } \boldsymbol{\kappa} = -\frac{1}{2} g^{-1/2} \text{Tr } \boldsymbol{\pi}. \quad (21.143)$$

For the volume to be extremal this quantity must vanish at every point of Σ . This condition is satisfied in a Friedmann universe (Chapter 27) and in a Taub universe (Chapter 30) at that value of the natural time-coordinate t at which the universe switches over from expansion to recontraction. It is remarkable that the same condition on the choice of simultaneity, Σ , lets itself be formulated in the same natural way,

$$\text{Tr } \boldsymbol{\kappa} = 0 \text{ or } \text{Tr } \boldsymbol{\pi} = 0, \quad (21.144)$$

The degrees of freedom of the geometry in brief

Pick hypersurface of extremal proper volume

Case of open 3-geometry

for a closed universe altogether deprived of any symmetry whatsoever. Alternatively, one can deal with a spacetime that is topologically the product of an open 3-space by the real line (time). Then it is natural to think of specifying the location in it of a bounding spacelike 2-geometry S with the topology of a 2-sphere. Then one has many ways to fill in the interior of S with a spacelike 3-geometry Σ ; but of all these Σ 's, only the one that is extremal, or only the ones that are extremal, satisfy (21.144).

Who is going to specify this 2-geometry with the topology of a 2-sphere? The choice of that 2-geometry is not a matter of indifference. In a given 4-geometry, distinct choices for the bounding 2-geometry will ordinarily give distinct results for the extremizing 3-geometry, and therefore different choices for the “initial-value simultaneity,” Σ . No consideration immediately thrusts itself forward that would give preference to one choice of 2-geometry over another. However, no such infinity of options presents itself when one limits attention to a closed 3-geometry. Therefore it will give concreteness to the following analysis to consider it applied to a closed universe, even though the analysis surely lets itself be made well-defined in an open region by appropriate specification of boundary values on the closed 2-geometry that bounds that open region. In brief, by limiting attention to a closed 3-geometry, one lets the obvious condition of closure take the place of boundary conditions that are not obvious.

York's analysis remains simple when his extrinsic time

$$\tau = \frac{2}{3} g^{-1/2} \text{Tr } \boldsymbol{\pi} = \frac{4}{3} \text{Tr } \boldsymbol{K}$$

has any constant value on the hypersurface, not only the value $\tau = 0$ appropriate for the hypersurface of extremal proper volume.

On the simultaneity Σ specified by the condition of constant extrinsic time, $\tau = \text{constant}$, begin by giving the conformal 3-geometry,

$$< = {}^{(3)}< = \left(\begin{array}{l} \text{the equivalence class of all those positive definite} \\ \text{Riemannian three-dimensional metrics that are} \\ \text{equivalent to each other under (1) diffeomorphism} \\ (\text{smooth sliding of the points over the manifold to}) \\ \text{(2) changes of scale that vary} \\ \text{smoothly from point to point, leaving fixed all} \\ \text{local angles (ratios of local distances), but} \\ \text{changing local distances themselves or (3) both.} \end{array} \right) \quad (21.145)$$

Meaning of conformal 3-geometry

The conformal 3-geometry is a geometric object that lends itself to definition and interpretation quite apart from the specific choice of coordinate system and even without need to use any coordinates at all. The conformal 3-geometry (*on the hypersurface Σ where $\tau = \text{constant}$*) may be regarded much as one regards the magnetic field in electromagnetism. The case of conformally flat 3-geometry,

$$ds^2 = \psi^4(x, y, z) ds_{\text{base}}^2 \quad (21.146)$$

(with $g_{ij\text{base}} = \delta_{ij}$), is analogous to those initial-value situations in electromagnetism where the magnetic field is everywhere zero (the time-antisymmetric initial-value problem of Brill); but now we consider the case of general ds^2_{base} .

The six metric coefficients g_{ij} of the conformal 3-geometry, subject to being changed by change of the three coordinates x^i , and undetermined at any one point up to a common position-dependent multiplicative factor, carry $6 - 3 - 1 = 2$ pieces of information per spacepoint. In this respect, they are like the components of the divergenceless magnetic field \mathcal{B} . The corresponding field momentum $\pi_{EM}^i \propto \mathcal{E}^i$ (Box 21.1, page 496) has its divergence specified by the charge density, and so also carries

two pieces of information (in addition to the prescribed information
about the density of charge) per spacepoint. (21.147)

The comparison is a little faulty between the components of \mathcal{B} and the metric coefficients. They are more like potentials than like components of the physically relevant field.

The appropriate measure of the "field" in geometrodynamics is the curvature tensor; but how can one possibly define a curvature tensor for a geometry that is as rudimentary as a conformal 3-geometry? York (1971) has raised and answered this question. The Weyl conformal-curvature tensor [equation (13.50) and exercise 13.13] is independent [in the proper (\S) representation], in spaces of higher dimensionality, of the position-dependent factor ψ^4 with which one multiplies the metric coefficients, but vanishes identically in three-dimensional space (exercise 21.21). One arrives at a non-zero conformally invariant measure of the curvature only when one goes to one higher derivative (exercise 21.22). In this way, one comes to *York's curvature tensor*

York's curvature tensor

$$Y^{ab} = Y^{ba} \text{ (symmetric);}$$

$$Y_a^a = 0 \text{ (traceless);}$$

$$Y^{ab}_{|b} = 0 \text{ (transverse);}$$

Y^{ab} invariant with respect to position-dependent
changes in the conformal scale factor;

$Y^{ab} = 0$ when and only when the 3-geometry is conformally flat. (21.148)

Y^{ab} provides what York calls the pure spin-two representation of the 3-geometry intrinsic to Σ . It is the analog of the field \mathcal{B} of electrodynamics on the spacelike initial-value simultaneity. It directly carries physical information about the conformal 3-geometry.

In addition to the conformal geometry ${}^{(3)}<\!\!$, specified by the "potentials" $g_{ij}/g^{1/3}$, and measured by the "field components" Y^{ij} , one must also specify on Σ the corresponding conjugate momenta:

The associated momenta

$$\begin{aligned}\tilde{\pi}^{ab} &= \tilde{\pi}^{ab} \text{ (symmetric);} & \tilde{\pi}_a^a &= 0 \text{ (traceless);} \\ \tilde{\pi}^{ab}_{\mid b} &= 0 \text{ (transverse) in case there is no flow of energy in} \\ &\quad \text{space; otherwise} \\ \tilde{\pi}^{ab}_{\mid b} &= 8\pi \text{ (density of flow of energy)}^a;\end{aligned}$$

two pieces of information (in addition to the prescribed information
about the flow of energy) per spacepoint. (21.149)

It might appear to be essential to specify with respect to which of the 3-geometries, distinguished from one another by different values of the conformal factor one calculates the covariant derivatives of tensor densities of weight 5/3 (see §21.2) in (21.148) and (21.149). However, York has shown that the conditions (21.149) do not in any way depend on the value of the conformal factor ψ^4 .

These equations (21.149) for what York calls the “momentum density of weight 5/3,”

$$\tilde{\pi}^{ab} = g^{1/3} \left(\pi^{ab} - \frac{1}{3} g^{ab} \operatorname{Tr} \boldsymbol{\pi} \right), \quad (21.150)$$

are linear, and therefore lend themselves to analysis by standard methods. It is a great help in this enterprise that York (1973a,b) has provided a “conformally invariant orthogonal decomposition of symmetric tensors on Riemannian manifolds” that allows one to generate solutions of these requirements (“transverse traceless,” “conformal Killing,” and “trace” parts, respectively, measure deformation of conformal part of geometry, mere recoordinatization, and change of scale). It is a further assistance, as York notes, that one has the same $\tilde{\pi}^{ab}$ for an entire conformal equivalence class of metrics; that is, for a given

$$\tilde{g}_{ab} = g^{-1/3} g_{ab}, \quad (21.151)$$

Unique solution for
conformal factor

no matter how different the g_{ab} and ψ themselves may be.

The conformal 3-geometry and the “momentum density of weight 5/3” once picked, the remaining initial-value equation (21.116) then becomes the “scale” equation,

$$8\nabla^2\psi - {}^{(3)}R\psi + M\psi^{-7} + Q\psi^{-3} - \frac{3}{8}\tau^2\psi^5 = 0, \quad (21.152)$$

for the determination of the conformal factor ψ . Here ∇^2 stands for the Laplacian

$$\nabla^2\psi \equiv g^{-1/2}(\partial/\partial x^a)g^{1/2}g^{ab}(\partial\psi/\partial x^b). \quad (21.153)$$

It, like ${}^{(3)}R$, M , and Q , refers to the base space. It is interesting that

$$\nabla^2 - \frac{1}{8} {}^{(3)}R$$

is a conformally invariant wave operator, whereas ∇^2 itself is not. The quantity M in York’s analysis is an abbreviation for

$$M \equiv g^{-5/3}g_{ac}g_{bd}\tilde{\pi}^{ab}\tilde{\pi}^{cd}, \quad (21.154a)$$

and

$$Q \equiv 16\pi\rho_{\text{base}} (= 16\pi\psi^8\rho = 16\pi\psi^8\rho_{\text{in final 3-geometry}}). \quad (21.154b)$$

One seeks a solution ψ that is continuous over the closed manifold and everywhere real and positive. When does such a solution ψ of the elliptic equation (21.152) exist? When is it unique? *Always* (when $M > 0$ and $\tau \neq 0$), is the result of O'Murchadha and York (1973); see also earlier investigations of Choquet-Bruhat (1972). Some of the physical considerations that come into this kind of problem have been discussed by Wheeler (1964a, pp. 370–381).

§21.12. MACH'S PRINCIPLE AND THE ORIGIN OF INERTIA

In my opinion the general theory of relativity can only solve this problem [of inertia] satisfactorily if it regards the world as spatially self-enclosed.

ALBERT EINSTEIN (1934), p. 52.

On June 25, 1913, two years before he had discovered the geometrodynamic law that bears his name, Einstein (1913b) wrote to Ernst Mach (Figure 21.5) to express his appreciation for the inspiration that he had derived for his endeavors from Mach's ideas. In his great book, *The Science of Mechanics*, Mach [(1912), Chapter 2, section 6] had reasoned that it could not make sense to speak of the acceleration of a mass relative to absolute space. Anyone trying to clear physics of mystical ideas would do better, he reasoned, to speak of acceleration relative to the distant stars. But how can a star at a distance of 10^9 light-years contribute to inertia in the here and the now? To make a long story short, one can say at once that Einstein's theory (1) identifies gravitation as the mechanism by which matter there influences inertia here; (2) says that this coupling takes place on a spacelike hypersurface [in what one, without a closer examination, might mistakenly think to be a violation of the principle of causality; see Fermi (1932) for a discussion and clarification of the similar apparent paradox in electrodynamics; see also Einstein (1934), p. 84: “Moreover I believed that I could show on general considerations a law of gravitation invariant in relation to any transformation of coordinates whatever was inconsistent with the principle of causation. These were errors of thought which cost me two years of excessively hard work, until I finally recognized them as such at the end of 1915”]; (3) supplies in the initial-value equations of geometrodynamics a mathematical tool to describe this coupling; (4) demands closure of the geometry in space [one conjectures; see Wheeler (1959, 1964c) and Hönl (1962)], as a boundary condition on the initial-value equations if they are to yield a well-determined [and, we know now, a unique] 4-geometry; and (5) identifies the collection of local Lorentz frames near any point in this resulting spacetime as what one means quantitatively by speaking of inertia at that point. This is how one ends up with inertia here determined by density and flow of mass-energy there.

There are many scores of papers in the literature on Mach's principle, including many—even one by Lenin (English translation, 1927)—one could call anti-Machian; and many of them make interesting points [see especially the delightful dialog by Weyl (1924a) on “inertia and the cosmos,” and the article (1957) and book (1961) of Sciama]. However, most of them were written before one had anything like the understanding of the initial-value problem that one possesses today. Therefore no

No violation of causality,
despite appearances

An enormous literature

(continued on page 546)

Figure 21.5.

Einstein's appreciation of Mach, written to Ernst Mach June 25, 1913, while Einstein was working hard at arriving at the final November 1915 formulation of standard general relativity. Regarding confirmation at a forthcoming eclipse: "If so, then your happy investigations on the foundations of mechanics, Planck's unjustified criticism notwithstanding, will receive brilliant confirmation. For it necessarily turns out that inertia originates in a kind of interaction between bodies, quite in the sense of your considerations on Newton's pail experiment. The first consequence is on p. 6 of my paper. The following additional points emerge: (1) If one accelerates a heavy shell of matter S , then a mass enclosed by that shell experiences an accelerative force. (2) If one rotates the shell relative to the fixed stars about an axis going through its center, a Coriolis force arises in the interior of the shell; that is, the plane of a Foucault pendulum is dragged around (with a practically unmeasurably small angular velocity)." Following the death of Mach, Einstein (1916a) wrote a tribute to the man and his work. Reprinted with the kind permission of the estate of Albert Einstein, Helen Dukas and Otto Nathan, executors.

Zinsel. 25. VI 13

Hoch geschätzter Herr Kollege!

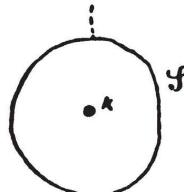
Diesen Tage haben Sie wohl meine neue Arbeit über Relativität und Gravitation erhalten, die nach unendlicher Mühe und quälendem Zweifel nun endlich fertig geworden ist. Nächstes Jahr bei der Sonnenfinsternis soll sich zeigen, ob die Lichtstrahlen aus der Sonne gekrümmt werden, ob in a. W. die zugegrundete gelegte fundamentale Annahme vor der Äquivalenz von Beschleunigung des Bezugssystems einerseits und Schwerkraft andererseits wirklich getroffen.

Wein ja, so erfahren Ihre genauen Untersuchungen über die Grundlagen der Mechanik-Planck's ungerichtfertigter Kritik zum Trotz - wie

glänzende Bestätigung. Denn es ergibt sich mit Notwendigkeit, dass die Trägheit an einer Art Wechselwirkung der Körper ihren Ursprung habe, ganz im Sinne Ihrer Überlegungen zum Newton'schen Einer - Versuch.

Eine erste Konsequenz in diesem Sinne finden Sie oben auf Seite 6 der Arbeit. Es hat sich ferner folgendes ergeben:

- 1) Beschleunigt man eine kreisförmige Kugelschale S, so erfährt nach der Theorie ein von der eingeschlossenen Körper aus beschleunigende Kraft



- 2) Rotiert die Schale S um eine durch ihren Mittelpunkt gehende Achse (relativ zum System der Fixsterne ("Restsystem"), so entsteht im Innern der Schale ein Coriolis - Tid., d. h. (des Torsionals-) Pendels wird (mit ferner allerdings unwesentlich kleinen "Geschwindigkeitsst.) aufgenommen

Es ist nur eine grosse Tendenz, Ihnen dies mittleren zu können, jemal eine kritische Plancks vor schon immer höchst ungerecht fertigt erscheinen will.

Mit grösster Hochachtung grüsst
Ihr ergebener d. K. weiter.

Ich dankte Ihnen zugleich für
die Übersendung Ihres Buches

attempt will be made to summarize or analyze the literature, which would demand a book in itself. Moreover, Mach's principle as presented here is more sharply formulated than Einstein ever put it in the literature [except for his considerations arguing that the universe must be closed; see Einstein's book (1950), pp. 107–108]; and Mach would surely have disowned it, for he could never bring himself to accept general (or even special) relativity. Nevertheless, it is a fact that Mach's principle—that matter there governs inertia here—and Riemann's idea—that the geometry of space responds to physics and participates in physics—were the two great currents of thought which Einstein, by means of his powerful equivalence principle, brought together into the present-day geometric description of gravitation and motion.

Mach's principle updated and spelled out

“Specify everywhere the distribution and flow of mass-energy and thereby determine the inertial properties of every test particle everywhere and at all times”. Spelled out, this prescription demands (1) a way of speaking about “everywhere”: a spacelike hypersurface Σ . Let one insist—in conformity with Einstein—(2) that it be a closed 3-geometry, and for convenience, not out of necessity, (3) that τ be independent of position on Σ . (4) Specify this 3-geometry to the extent of giving the conformal metric; without the specification of at least this much 3-geometry, there would be no evident way to say “where” the mass-energy is to be located. (5) Give density ρ_{base} as a function of position in this conformal 3-geometry. (6) Recognize that giving the mass-energy only of fields other than gravity is an inadequate way to specify the distribution of mass-energy throughout space. Formally, to be sure, the gravitational fields does not and cannot make any contribution to the source term that stands on the righthand side of Einstein's field equation. However, the analysis of gravitational waves (Chapters 18 and 35) shows that perturbations in the geometry of scale small compared to the scale of observation have to be regarded as carrying an effective content of mass-energy. Moreover, one has in a geon [Wheeler (1955); Brill and Hartle (1964); for more on gravitational-wave energy, see §35.14] an object built out of gravitational waves (or electromagnetic waves, or neutrinos, or any combination of the three) that holds itself together for a time that is long in comparison to the characteristic period of vibration of the waves. It looks from a distance like any other mass, even though nowhere in its interior can one put a finger and say “here is mass.” Therefore it, like any other mass, must have “its influence on inertia.” But to specify this mass, one must give enough information to characterize completely the gravitational waves on the simultaneity Σ . For this, it is not enough merely to have given the two “wave-coordinates” per spacepoint that one possesses in ${}^{(3)}<$. One must give in addition (7) the two “wave-momenta” per spacepoint that appear in York's “momentum density of weight 5/3,” $\tilde{\pi}^{ab}$; and at the same time, as an inextricable part of this operation, one must (8) specify the density of flow of field energy. (9) Solve for the conformal factor ψ . (10) Then one has complete initial-value data that satisfy the initial-value equations of general relativity. (11) These data now known, the remaining, dynamic, components of the field equation determine the 4-geometry into the past and the future. (12) In this way, the inertial properties of every test particle are determined everywhere and at all times, giving concrete realization to Mach's principle.

Much must still be done to spell out the physics behind these equations and to

see this physics in action. Some significant progress had already been made in this direction before the present stage in one's understanding of the initial-value equations. Especially interesting are results of Thirring (1918) and (1921) and of Thirring and Lense (1918), discussed by Einstein (1950) in the third edition of his book, *The Meaning of Relativity*.

Consider a bit of solid ground near the geographic pole, and a support erected there, and from it hanging a pendulum. Though the sky is cloudy, the observer watches the track of the Foucault pendulum as it slowly turns through 360° . Then the sky clears and, miracle of miracles, the pendulum is found to be swinging all the time on an arc fixed relative to the far-away stars. If "mass there governs inertia here," as envisaged by Mach, how can this be?

Enlarge the question. By the democratic principle that equal masses are created equal, the mass of the earth must come into the bookkeeping of the Foucault pendulum. Its plane of rotation must be dragged around with a slight angular velocity, ω_{drag} , relative to the so-called "fixed stars." How much is ω_{drag} ? And how much would ω_{drag} be if the pendulum were surrounded by a rapidly spinning spherical shell of mass M and radius R_{shell} , turning at angular velocity ω_{shell} ?

Einstein's theory says that inertia is a manifestation of the geometry of spacetime. It also says that geometry is affected by the presence of matter to an extent proportional to the factor $G/c^2 = 0.742 \times 10^{-28}$ cm/g. Simple dimensional considerations leave no room except to say that the rate of drag is proportional to a expression of the form

$$\omega_{\text{drag}} = k \frac{G}{c^2} \frac{m_{\text{shell, conv}}}{R_{\text{shell}}} \omega_{\text{shell}} = k \frac{m_{\text{shell}}}{R_{\text{shell}}} \omega_{\text{shell}}. \quad (21.155)$$

Here k is a numerical factor to be found only by detailed calculation. Lense and Thirring [(1918) and (1921)], starting with a flat background spacetime manifold, calculated in the weak-field approximation of Chapter 18 the effect of the moving current of mass on the metric. Expressed in polar coordinates, the metric acquires a non-zero coefficient $g_{\phi t}$. Inserted into the equation of geodesic motion, this off-diagonal metric coefficient gives rise to a precession. This precession (defined here about an axis parallel to the axis of rotation, not about the local vertical) is given by an expression of the form (21.155), where the precession factor k has the value

$$k = 4/3. \quad (21.156)$$

There is a close parallelism between the magnetic component of the Maxwell field and the precession component of the Einstein field. In neither field does a source at rest produce the new kind of effect when acting on a test particle that is also at rest. One designs a circular current of charge to produce a magnetic field; and a test charge, in order to respond to this magnetic field, must also be in motion. Similarly here: no pendulum vibration means no pendulum precession. Moreover, the direction of the precession depends on where the pendulum is, relative to the rotating shell of mass. The precession factor k has the following values:

The Foucault pendulum

The dragging of the inertial frame

- $k = 4/3$ for pendulum anywhere inside rotating shell of mass;
 $k = 4/3$ for pendulum at North or South pole; (21.157)
 $k = -2/3$ for pendulum just outside the rotating shell at its equator.

This position-dependence of the drag, ω_{drag} , makes still more apparent the analogy with magnetism, where the field of a rotating charged sphere points North at the center of the sphere, and North at both poles, but South at the equator.

Whether the Foucault pendulum is located in imagination at the center of the earth or in actuality at the North pole, the order of magnitude of the expected drag is

$$\begin{aligned}\omega_{\text{drag}} &\sim \frac{m_{\text{earth}}}{R_{\text{earth}}} \omega_{\text{earth}} \sim \frac{0.44 \text{ cm}}{6 \times 10^8 \text{ cm}} \frac{1 \text{ radian}}{13700 \text{ sec}} \\ &\sim 5 \times 10^{-14} \text{ rad/sec},\end{aligned}\quad (21.158)$$

too small to allow detection, let alone actual measurement, by any device so far built—but perhaps measurable by gyroscopes now under construction (§40.7). By contrast, near a rapidly spinning neutron star or near a black hole endowed with substantial angular momentum, the calculated drag effect is not merely detectable; it is even important (see Chapter 33 on the physics of a rotating black hole).

The distant stars must influence the natural plane of vibration of the Foucault pendulum as the nearby rotating shell of matter does, provided that the stars are not so far away ($r \sim$ radius of universe) that the curvature of space begins to introduce substantial corrections into the calculation of Thirring and Lense. In other words, no reason is apparent why all masses should not be treated on the same footing, so that (21.158) more appropriately, if also somewhat symbolically, reads

$$\omega_{\substack{\text{plane of} \\ \text{vibration} \\ \text{of Foucault} \\ \text{pendulum}}} \sim \frac{m_{\text{shell}}}{R_{\text{shell}}} \omega_{\text{shell}} + \sum_{\substack{\text{far-away} \\ \text{"stars"}}} \frac{m_{\text{"star"}}}{r_{\text{"star"}}} \omega_{\text{"star"}}. \quad (21.159)$$

Moreover, when there is no nearby shell of matter, or when it has negligible effects, the plane of vibration of the pendulum, if experience is any guide, cannot turn with respect to the frame defined by the far-away “stars.” In this event ω_{Foucault} must be identical with ω_{stars} ; or the “sum for inertia,”

$$\sum_{\substack{\text{far-away} \\ \text{"stars"}}} \frac{m_{\text{"star"}}}{r_{\text{"star"}}} \sim \frac{m_{\text{universe}}}{r_{\text{universe}}}, \quad (21.160)$$

must be of the order of unity. Just such a relation of approximate identity between the mass content of the universe and its radius at the phase of maximum expansion is a characteristic feature of the Friedman model and other simple models of a closed universe (Chapters 27 and 30). In this respect, Einstein’s theory of Mach’s principle exhibits a satisfying degree of self-consistency.

The “sum for inertia”

At phases of the dynamics of the universe other than the stage of maximum expansion, r_{universe} can become arbitrarily small compared to m_{universe} . Then the ratio (21.160) can depart by powers of ten from unity. Regardless of this circumstance, one has no option but to understand that the *effective* value of the “sum for inertia” is still unity after all corrections have been made for the dynamics of contraction or expansion, for retardation, etc. Only so can ω_{Foucault} retain its inescapable identity with $\omega_{\text{far-away stars}}$. Fortunately, one does not have to pursue the theology of the “sum for inertia” to the uttermost of these sophistications to have a proper account of inertia. Mach’s idea that mass there determines inertia here has its complete mathematical account in Einstein’s geometrodynamic law, as already spelled out. For the first strong-field analysis of the dragging of the inertial reference system in the context of relativistic cosmology, see Brill and Cohen (1966) and Cohen and Brill (1967); see also §33.4 for dragging by a rotating black hole.

Still another clarification is required of what Mach’s principle means and how it is used. The inertial properties of a test particle are perfectly well-determined when that particle is moving in ideal Minkowski space. “Point out, please,” the anti-Machian critic says, “the masses that are responsible for this inertia.” In answer, recall that Einstein’s theory includes not only the geometrodynamic law, but also, in Einstein’s view, the boundary condition that the universe be closed. Thus the section of spacetime that is flat is to be viewed, not as infinite, but as part of a closed universe. (For a two-dimensional analog, fill a rubber balloon with water and set it on a glass tabletop and look at it from underneath). The part of the universe that is curved acquires its curvature by reason of its actual content of mass-energy or—if animated only by gravitational waves—by reason of its effective content of mass-energy. This mass-energy, real or effective, is to be viewed as responsible for the inertial properties of the test particle that at first sight looked all alone in the universe.

It in no way changes the qualitative character of the result to turn attention to a model universe where the region of Minkowski flatness, and all the other linear dimensions of the universe, have been augmented tenfold (“ten times larger balloon; ten times larger face”). The curvature and density of the curved part of the model universe are down by a factor of 100, the volume is up by a factor of 1,000, the mass is up by a factor of 10; but the ratio of mass to radius, or the “sum for inertia” (the poor man’s substitute for a complete initial-value calculation) is unchanged.

Einstein acknowledged a debt of parentage for his theory to Mach’s principle (Figure 21.5). It is therefore only justice that Mach’s principle should in return today owe its elucidation to Einstein’s theory.

Minkowski geometry as limit
of a closed 3-geometry

Exercise 21.21. WHY THE WEYL CONFORMAL CURVATURE TENSOR VANISHES

EXERCISES

How many independent components does the Riemann curvature tensor have in three-dimensional space? How many does the Ricci curvature tensor have? Show that the two tensors are related by the formula

$$\begin{aligned} R^d_{abc} &= \delta_b^d R_{ac} - \delta_c^d R_{ab} + g_{ac} R^d_b - g_{ab} R^d_c \\ &+ \frac{1}{2} R (\delta_c^d g_{ab} - \delta_b^d g_{ac}) \end{aligned}$$

with no need of any Weyl conformal-curvature tensor to specify (as in higher dimensions) the further details of the Riemann tensor. Show that the Weyl tensor, from an n -dimensional modification of equation (13.50) as in exercise 13.13, vanishes for $n = 2$.

Exercise 21.22. YORK'S CURVATURE

[York (1971)]. (a) Define the tensor [Eisenhart (1926)]

$$R_{abc} = R_{ab|c} - R_{ac|b} + \frac{1}{4} (g_{ac} R_{|b} - g_{ab} R_{|c}).$$

(b) Show that a 3-geometry is conformally flat when and only when $R_{abc} = 0$.

(c) Show that the following identities hold and reduce to five the number of independent components of R_{abc} :

$$R^a_{ac} = g^{ab} R_{bac} = 0;$$

$$R_{abc} + R_{acb} = 0;$$

$$R_{abc} + R_{cab} + R_{bca} = 0.$$

(d) Show that York's curvature

$$\begin{aligned} Y^{ab} &= g^{1/3}[ae] \left(R_f^b - \frac{1}{4} \delta_f^b R \right)_{|e} \\ &= -\frac{1}{2} g^{1/3}[ae] g^{bm} R_{mef} \end{aligned}$$

is conformally invariant and has the properties listed in equations (21.148).

Exercise 21.23. PULLING THE POYNTING FLUX VECTOR "OUT OF THE AIR"

From the condition that the Hamilton-Jacobi functional $S(g_{ij}, A_m)$ (extremal of the action integral) for the combined Einstein and Maxwell fields, ostensibly dependent on the six metric coefficients $g_{ij}(x, y, z)$ and the three potentials $A_m(x, y, z)$, shall actually depend only on the 3-geometry of the spacelike hypersurface and the distribution of magnetic field strength on this hypersurface, show that the geometrodynamic field momentum $\pi^{ij} = \delta S / \delta g_{ij}$ satisfies a condition of the form

$$\pi^{ij}_{|j} = c[imn] \mathcal{E}_m \mathcal{B}_n,$$

and evaluate the coefficient c in this equation [Wheeler (1968b)]. Hint: Note that the transformation

$$x^i \rightarrow x^i - \xi^i, g_{ij} \rightarrow g_{ij} + \xi_{i|j} + \xi_{j|i}$$

in no way changes the 3-geometry itself, and therefore the corresponding induced change in S ,

$$\delta S = \int \left[\frac{\delta S}{\delta g_{ij}} \delta g_{ij} + \frac{\delta S}{\delta A_m} \delta A_m \right] d^3x$$

must vanish identically for arbitrary choice of the $\xi^i(x, y, z)$, which measure the equivalent of the sliding of a ruled transparent rubber sheet over an automobile fender.

Exercise 21.24. THE EXTREMAL ACTION ASSOCIATED WITH THE HILBERT ACTION PRINCIPLE DEPENDS ON CONFORMAL 3-GEOMETRY AND EXTRINSIC TIME [K. Kuchař (1972) and J. York (1972)]

Show that the data demanded by the Hilbert action principle $\delta \int^{(4)} R (-^{(4)}g)^{1/2} d^4x = 0$ on each of the two bounding spacelike hypersurfaces consist of (1) the conformal 3-geometry ${}^{(3)}<$ of the hypersurface plus (2) the extrinsic time variable defined by

$$\tau = \frac{2}{3} g^{-1/2} \text{Tr } \mathbf{n} = \frac{4}{3} \text{Tr } \mathbf{K},$$

conveniently represented by the pictogram , measured by one number per spacepoint, and independent of the conformal factor in the metric of the 3-geometry. This done, explain in a few words why in this formulation of geometrodynamics the Hamilton-Jacobi function (\hbar times the phase of the wave function in the semiclassical or JWKB approximation) is appropriately expressed in the form

$$S = S({}^{(3)}<, \text{pic}).$$

§21.13. JUNCTION CONDITIONS

The intrinsic and extrinsic curvatures of a hypersurface, which played such fundamental roles in the initial-value formalism, are also powerful tools in the analysis of “junction conditions.”

Recall the junction conditions of electrodynamics: across any surface (e.g., a capacitor plate), the tangential part of the electric field, $\mathbf{E}_{||}$, and the normal part of the magnetic field, \mathbf{B}_{\perp} , must be continuous; thus,

$$\begin{aligned} [\mathbf{E}_{||}] &\equiv (\text{discontinuity in } \mathbf{E}_{||}) \\ &\equiv (\mathbf{E}_{||} \text{ on } "+" \text{ side of surface}) - (\mathbf{E}_{||} \text{ on } "-" \text{ side of surface}) \\ &\equiv \mathbf{E}_{||}^+ - \mathbf{E}_{||}^- = 0, \end{aligned} \tag{21.161a}$$

$$[\mathbf{B}_{\perp}] \equiv \mathbf{B}_{\perp}^+ - \mathbf{B}_{\perp}^- = 0; \tag{21.161b}$$

Junction conditions for
electrodynamics

while the “jump” in the parts \mathbf{E}_{\perp} and $\mathbf{B}_{||}$ must be related to the charge density (charge per unit area) σ , the current density (current per unit area) \mathbf{j} , and the unit normal to the surface \mathbf{n} by the formulas

$$[\mathbf{E}_{\perp}] = \mathbf{E}_{\perp}^+ - \mathbf{E}_{\perp}^- = 4\pi\sigma\mathbf{n}, \tag{21.161c}$$

$$[\mathbf{B}_{||}] = \mathbf{B}_{||}^+ - \mathbf{B}_{||}^- = 4\pi\mathbf{j} \times \mathbf{n}. \tag{21.161d}$$

Recall also that one derives these junction conditions by integrating Maxwell’s equations over a “pill box” that is centered on the surface.

Similar junction conditions, derivable in a similar manner, apply to the gravitational field (spacetime curvature), and to the stress-energy that generates it.* Focus

*The original formulation of gravitational junction conditions stemmed from Lanczos (1922, 1924). The formulation given here, in terms of intrinsic and extrinsic curvature, was developed by Darmois (1927), Misner and Sharp (1964), and Israel (1966). For further references to the extensive literature, see Israel.

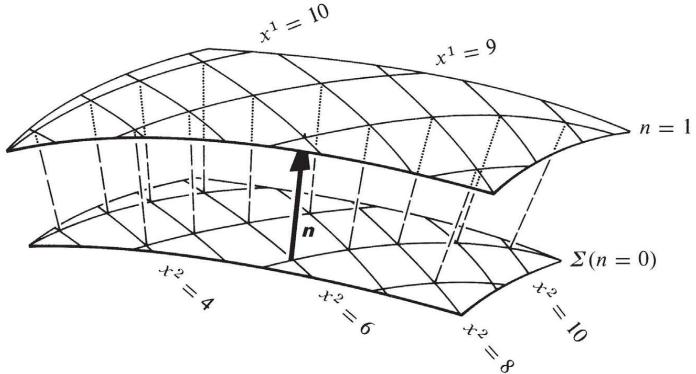


Figure 21.6.

Gaussian normal coordinates in the neighborhood of a 3-surface Σ . The metric in Gaussian normal coordinates has the form

$$ds^2 = (\mathbf{n} \cdot \mathbf{n})^{-1} dn^2 + g_{ij} dx^i dx^j$$

with $\mathbf{n} = \partial/\partial n$, $(\mathbf{n} \cdot \mathbf{n}) = -1$ if the surface is spacelike, and $(\mathbf{n} \cdot \mathbf{n}) = 1$ if it is timelike. (See exercise 27.2.) The extrinsic curvature of the surfaces $n = \text{constant}$ is $K_{ij} = -\frac{1}{2} \partial g_{ij}/\partial n$, and the Einstein field equations written in “3 + 1” form are (21.162).

Einstein equation in “3 + 1” form

attention on a specific three-dimensional slice through spacetime—the 3-surface Σ on Figure 21.6. Let the surface be either spacelike [unit normal \mathbf{n} timelike; $(\mathbf{n} \cdot \mathbf{n}) = -1$] or timelike [\mathbf{n} spacelike; $(\mathbf{n} \cdot \mathbf{n}) = +1$]. The null case will be discussed later. As an aid in deriving junction conditions, introduce Gaussian normal coordinates in the neighborhood of Σ [see the paragraph preceding equation (21.82)]. In terms of the intrinsic and extrinsic curvatures of Σ and of neighboring 3-surfaces $n = \text{constant}$, the Einstein tensor and Einstein field equation have components

$$G^n_n = -\frac{1}{2} {}^{(3)}R + \frac{1}{2}(\mathbf{n} \cdot \mathbf{n})^{-1}\{\text{Tr } \mathbf{K}\}^2 - \text{Tr } (\mathbf{K}^2) = 8\pi T^n_n, \quad (21.162a)$$

$$G^n_i = -(\mathbf{n} \cdot \mathbf{n})^{-1}\{K_i^m|_m - (\text{Tr } \mathbf{K})|_i\} = 8\pi T^n_i, \quad (21.162b)$$

$$G^i_j = {}^{(3)}G^i_j + (\mathbf{n} \cdot \mathbf{n})^{-1}\{(K^i_j - \delta^i_j \text{Tr } \mathbf{K})_{,n}$$

$$- (\text{Tr } \mathbf{K})K^i_j + \frac{1}{2} \delta^i_j (\text{Tr } \mathbf{K})^2 + \frac{1}{2} \delta^i_j \text{Tr } (\mathbf{K}^2)\} = 8\pi T^i_j. \quad (21.162c)$$

[See equations (21.77), (21.81), (21.76), and (21.82).]

Surface stress-energy tensor

Suppose that the stress-energy tensor T^α_β contains a “delta-function singularity” at Σ —i.e., suppose that Σ is the “world tube” of a two-dimensional surface with finite 4-momentum per unit area (analog of surface charge and surface current in electrodynamics). Then define the *surface stress-energy tensor on Σ* to be the integral of T^α_β with respect to proper distance (n), measured perpendicularly through Σ :

$$S^\alpha_\beta = \lim_{\epsilon \rightarrow 0} \left[\int_{-\epsilon}^{+\epsilon} T^\alpha_\beta dn \right]. \quad (21.163)$$

To discover the effect of this surface layer on the spacetime geometry, perform a “pill-box integration” of the Einstein field equation (21.162)

$$\lim_{\epsilon \rightarrow 0} \left[\int_{-\epsilon}^{+\epsilon} G^{\alpha}_{\beta} dn \right] = 8\pi S^{\alpha}_{\beta}. \quad (21.164)$$

Examine the integral of G^{α}_{β} . If the 3-metric g_{ij} were to contain a delta function or a discontinuity at Σ , then Σ would not have any well-defined 3-geometry—a physically inadmissible situation, even in the presence of surface layers. Absence of delta functions, $\delta(n)$, in g_{ij} means absence of delta functions in ${}^{(3)}R$; absence of discontinuities in g_{ij} means absence of delta functions in $K_{ij} = -\frac{1}{2}g_{ij,n}$. Thus, equations (21.162) when integrated say

$$\int G^n_n dn = 0 = 8\pi S^n_n, \quad (21.165a)$$

$$\int G^n_i dn = 0 = 8\pi S^n_i, \quad (21.165b)$$

$$\int G^i_j \cdot dn = (\mathbf{n} \cdot \mathbf{n})(\gamma^i_j - \delta^i_j \text{Tr } \boldsymbol{\gamma}) = 8\pi S^i_j, \quad (21.165c)$$

where γ^i_j is the “jump” in the components of the extrinsic curvature

$$\boldsymbol{\gamma} \equiv [\mathbf{K}] \equiv (\mathbf{K} \text{ on } "n = +\epsilon \text{ side}" \text{ of } \Sigma) - (\mathbf{K} \text{ on } "n = -\epsilon \text{ side}" \text{ of } \Sigma) \quad (21.166)$$

$$\equiv \mathbf{K}^+ - \mathbf{K}^-.$$

In the absence of a delta-function surface layer, the above junction conditions say, simply, that $\boldsymbol{\gamma} \equiv [\mathbf{K}] = 0$. In words: if one examines how Σ is embedded in the spacetime above its “upper” face, and how it is embedded in the spacetime below its “lower” face, one must discover identical embeddings—i.e., identical extrinsic curvatures \mathbf{K} . Of course, the intrinsic curvature of Σ must also be the same, whether viewed from above or below. More briefly:

$$(\text{absence of surface layers}) \iff (\text{"continuity" of } g_{ij} \text{ and } K_{ij}). \quad (21.167)$$

If a surface layer is present, then Σ must be the world tube of a two-dimensional layer of matter, and the normal to Σ must be spacelike, $(\mathbf{n} \cdot \mathbf{n}) = +1$. The junction conditions (21.165a,b) then have the simple physical meaning

$$\mathbf{s}(\mathbf{n}, \dots) = 0 \iff \begin{cases} \text{the momentum flow is entirely in } \Sigma; \\ \text{i.e., no momentum associated with the} \\ \text{surface layer flows out of } \Sigma; \text{ i.e., } \Sigma \\ \text{is the world tube of the surface layer} \end{cases}, \quad (21.168a)$$

which tells one nothing new. The junction condition (21.165c) says that the surface stress-energy generates a discontinuity in the extrinsic curvature (different embedding in spacetime “above” Σ than “below” Σ), given by

$$\gamma^i_j - \delta^i_j \text{Tr } \boldsymbol{\gamma} = 8\pi S^i_j. \quad (21.168b)$$

Of course, the intrinsic geometry of Σ must be the same as seen from above and below,

$$g_{ij} \text{ continuous across } \Sigma. \quad (21.169)$$

Derivation of junction conditions

Junction conditions in absence of surface layers

Junction conditions for a surface layer

In analyzing surface layers, one uses not only the junction conditions (21.168a) to (21.169), but also the four-dimensional Einstein field equation applied on each side of the surface Σ separately, and also an equation of motion for the surface stress-energy. The equation of motion is derived by examining the jump in the field equation $G^n_i = 8\pi T^n_i$ (equation 21.162b); thus $[G^n_i] = 8\pi[T^n_i]$ says

$$(\gamma_i^m - \delta_i^m \operatorname{Tr} \gamma)_{|m} = -8\pi[T^n_i];$$

and when reexpressed in terms of S_i^m by means of the junction condition (21.168b), it says

$$S^{im}_{|m} + [T^{in}] = 0. \quad (21.170)$$

Equation of motion for a surface layer

Gravitational-wave shock fronts

[For intuition into this equation of motion, see Exercises 21.25 and 21.26. For applications of the “surface-layer formalism” see exercise 21.27; also Israel (1966), Kuchař (1968), Papapetrou and Hamoui (1968).]

When one turns attention to junction conditions across a *null* surface Σ , one finds results rather different from those in the spacelike and timelike cases. A “pill-box” integration of the field equations reveals that even in vacuum the extrinsic curvature may be discontinuous. A discontinuity in K_{ij} across a null surface, without any stress-energy to produce it, is the geometric manifestation of a *gravitational-wave shock front* (analog of a shock-front in hydrodynamics). For quantitative details see, e.g., Pirani (1957), Papapetrou and Treder (1959, 1962), Treder (1962), and especially Choquet-Bruhat (1968b).

That a discontinuity in the curvature tensor can propagate with the speed of light is a reminder that all gravitational effects, like all electromagnetic effects, obey a causal law. The initial-value data on a spacelike initial-value hypersurface uniquely determine the resulting spacetime geometry [see the work of Cartan, Stellmacher, Lichnerowicz, and Bruhat (also under the names Fourès-Bruhat and Choquet-Bruhat) and others cited and summarized in the article of Bruhat (1962)] but determine it in a way consistent with causality. Thus a change in these data throughout a limited region of the initial value 3-geometry makes itself felt on a slightly later hypersurface solely in a region that is also limited, and only a little larger than the original region.

When one turns from classical dynamics to quantum dynamics, one sees new reason to focus attention on a spacelike initial-value hypersurface: the observables at different points on such a hypersurface commute with one another; i.e., are in principle simultaneously observable.

Not every four-dimensional manifold admits a global singularity-free spacelike hypersurface. Those manifolds that do admit such a hypersurface have more to do with physics, it is possible to believe, than those that do not.

Even in a manifold that does admit a spacelike hypersurface, attention has been given sometimes, in the context of classical theory, to initial-value data on a hypersurface that is not spacelike but “characteristic,” in the sense that it accommodates null geodesics [see, for example, Sachs (1964) and references cited there]. It is typical in such situations that one can predict the future but not the past, or predict the past but not the future.

Children of light and children of darkness is the vision of physics that emerges from this chapter, as from other branches of physics. The children of light are the differential equations that predict the future from the present. The children of darkness are the factors that fix these initial conditions.

Exercise 21.25. EQUATION OF MOTION FOR A SURFACE LAYER
EXERCISES

(a) Let \mathbf{u} be the “mean 4-velocity” of the matter in a surface layer—so defined that an observer moving with 4-velocity \mathbf{u} sees zero energy flux. Let σ be the total mass-energy per unit proper surface area, as measured by such a “comoving observer.” Show that the surface stress-energy tensor can be expressed in the form

$$\mathbf{S} = \sigma \mathbf{u} \otimes \mathbf{u} + \mathbf{t}, \text{ where } (\mathbf{t} \cdot \mathbf{u}) = 0, \quad (21.171)$$

and where \mathbf{t} is a symmetric stress tensor.

(b) Show that the component along \mathbf{u} of the equation of motion (21.170) is

$$d\sigma/d\tau = -\sigma u^j_{|j} + u_j t^{jk}_{|k} + u_j [T^{jn}], \quad (21.172)$$

where $d/d\tau = \mathbf{u}$. Give a physical interpretation for each term.

(c) Let a_j be that part of the 4-acceleration of the comoving observer which lies in the surface layer Σ . By projecting the equation of motion (21.170) perpendicular to \mathbf{u} , show that

$$\sigma a_j = -P_{ja} \{ t^{ab}_{|b} + [T^{an}] \}, \quad (21.173)$$

where P_{ja} is the projection operator

$$P_{ja} = g_{ja} + u_j u_a. \quad (21.174)$$

Give a physical interpretation for each term of equation (21.182).

Exercise 21.26. THIN SHELLS OF DUST

For a thin shell of dust surrounded by vacuum ($[T^{jn}] = 0$, $\mathbf{t} = 0$), derive the following equations

$$d\sigma/d\tau = -\sigma u^b_{|b}, \quad (21.175a)$$

$$\mathbf{a}^+ + \mathbf{a}^- = 0, \quad (21.175b)$$

$$\mathbf{a}^+ - \mathbf{a}^- = (4\pi\sigma)\mathbf{n} \quad (21.175c)$$

$$\gamma = 8\pi\sigma \left(\mathbf{u} \otimes \mathbf{u} + \frac{1}{2}\mathbf{g} \right). \quad (21.175d)$$

Here \mathbf{a}^+ and \mathbf{a}^- are the 4-accelerations as measured by accelerometers that are fastened onto the outer and inner sides of the shell, and \mathbf{g} is the 3-metric of the shell. Show that the first of these equations is the law of “conservation of rest mass.”

Exercise 21.27. SPHERICAL SHELL OF DUST

Apply the formalism of exercise 21.25 to a collapsing spherical shell of dust [Israel (1967b)]. For the metric inside and outside the shell, take the flat-spacetime and vacuum Schwarzschild expressions (Chapter 23),

$$ds^2 = -dt^2 + dr^2 + r^2(d\theta^2 + \sin^2\theta d\phi^2) \text{ inside,} \quad (21.176a)$$

$$ds^2 = -\left(1 - \frac{2M}{r}\right)dt^2 + \frac{dr^2}{1 - 2M/r} + r^2(d\theta^2 + \sin^2\theta d\phi^2) \text{ outside.} \quad (21.176b)$$

Let the “radius” of the shell, as a function of proper time measured on the shell, be

$$R \equiv \frac{1}{2\pi} \times (\text{proper circumference of shell}) = R(\tau). \quad (21.176c)$$

Show that the shell’s mass density varies with time as

$$\sigma(\tau) = \mu/4\pi R^2(\tau), \quad \mu = \text{constant} = \text{“total rest mass”}; \quad (21.176d)$$

and derive and solve the equation of motion

$$M = \mu \left\{ 1 + \left(\frac{dR}{d\tau} \right)^2 \right\}^{1/2} - \frac{\mu}{2R}. \quad (21.176e)$$

CHAPTER 22

THERMODYNAMICS, HYDRODYNAMICS, ELECTRODYNAMICS, GEOMETRIC OPTICS, AND KINETIC THEORY

§22.1. THE WHY OF THIS CHAPTER

Astrophysical applications of gravitation theory are the focus of the rest of this book, except for Chapters 41–44. Each application—stars, star clusters, cosmology, collapse, black holes, gravitational waves, solar-system experiments—can be pursued by itself at an elementary level, without reference to the material in this chapter. But deep understanding of the applications requires a prior grasp of thermodynamics, hydrodynamics, electrodynamics, geometric optics, and kinetic theory, all in the context of curved spacetime. Hence, most Track-2 readers will want to probe these subjects at this point.

§22.2. THERMODYNAMICS IN CURVED SPACETIME*

Consider, for concreteness and simplicity, the equilibrium thermodynamics of a perfect fluid with fixed chemical composition (“simple perfect fluid”—for example, the gaseous interior of a collapsing supermassive star. The thermodynamic state of a fluid element, as it passes through an event \mathcal{P}_0 , can be characterized by various thermodynamic potentials, such as n , ρ , p , T , s , μ . The numerical value of each potential at \mathcal{P}_0 is measured in the proper reference frame (§13.6) of an observer who moves with the fluid element—i.e., in the fluid element’s “rest frame.” Despite

This chapter is entirely Track 2. No earlier Track-2 material is needed as preparation for it, but Chapter 5 (stress-energy tensor) will be helpful.

§22.5 (geometric optics) is needed as preparation for Chapter 34 (singularities and global methods). The rest of the chapter is not needed as preparation for any later chapter; but it will be extremely helpful in most applications of gravitation theory (Chapters 23–40).

*For more detailed treatments of this subject see, e.g., Stueckelberg and Wanders (1953), Kluitenberg and de Groot (1954), Meixner and Reik (1959), and references cited therein; see also the references on hydrodynamics cited at the beginning of §22.3, and the references on kinetic theory cited at the beginning of §22.6.

Thermodynamic potentials are defined in rest frame of fluid

this use of rest frame to measure the potentials, the potentials are frame-independent functions (scalar fields). At the chosen event \mathcal{P}_0 , a given potential (e.g., n) has a unique value $n(\mathcal{P}_0)$; so n is a perfectly good frame-independent function.

The values of n , ρ , p , T , s , μ measure the following quantities in the rest frame of the fluid element:

Definitions of thermodynamic potentials

n , baryon number density; i.e., number of baryons per unit three-dimensional volume of rest frame, with antibaryons (if any) counted negatively.

ρ , density of total mass-energy; i.e., total mass-energy (including rest mass, thermal energy, compressional energy, etc.) contained in a unit three-dimensional volume of the rest frame.

p , isotropic pressure in rest frame.

T , temperature in rest frame.

s , entropy per baryon in rest frame. (The entropy per unit volume is ns .)

μ , chemical potential of baryons in rest frame [see equation (22.8) below].

Definition of “simple fluid”

The chemical composition of the fluid (number density of hydrogen molecules, number density of hydrogen atoms, number density of free protons and electrons, number density of photons, number density of ^{238}U nuclei, number density of Λ hyperons . . .) is assumed to be fixed uniquely by two thermodynamic variables—e.g., by the total number density of baryons n and the entropy per baryon s . In this sense the fluid is a “simple fluid.” Simple fluids occur whenever the chemical abundances are “frozen” (reaction rates too slow to be important on the time scales of interest; for example, in a supermassive star except during explosive burning and except at temperatures high enough for $e^- - e^+$ pair production). Simple fluids also occur in the opposite extreme of complete chemical equilibrium (reaction rates fast enough to maintain equilibrium despite changing density and entropy; for example, in neutron stars, where high pressures speed up all reactions). When one examines nuclear burning in a nonconvecting star, or explosive nuclear burning, or pair production and neutrino energy losses at high temperatures, one must usually treat the fluid as “multicomponent.” Then one introduces a number density n_j and a chemical potential μ_j for each chemical species with abundance not fixed by n and s . For further details see, e.g., Zel'dovich and Novikov (1971).

Law of baryon conservation

The most fundamental law of thermodynamics—even more fundamental than the “first” and “second” laws—is *baryon conservation*. Consider a fluid element whose moving walls are attached to the fluid so that no baryons flow in or out. As the fluid element moves through spacetime, deforming along the way, its volume V changes. But the number of baryons in it must remain fixed, so

$$\frac{d}{d\tau} (nV) = 0. \quad (22.1)$$

The changes in volume are produced by the flow of neighboring bits of fluid away from or toward each other—explicitly (exercise 22.1)

$$dV/d\tau = (\nabla \cdot \mathbf{u})V, \quad (22.2)$$

where $\mathbf{u} = d/d\tau$ is the 4-velocity of the fluid. Consequently, baryon conservation [equation (22.1)] can be reexpressed as

$$0 = \frac{dn}{d\tau} + \frac{n}{V} \frac{dV}{d\tau} = \nabla_{\mathbf{u}} n + n(\nabla \cdot \mathbf{u}) = \mathbf{u} \cdot \nabla n + n(\nabla \cdot \mathbf{u}) = \nabla \cdot (n\mathbf{u});$$

i.e.,

$$\nabla \cdot \mathbf{S} = 0, \quad (22.3)$$

$$\mathbf{S} = n\mathbf{u} = \text{baryon number-flux vector} \quad (22.4)$$

(see §5.4 and exercise 5.3.) Moreover, this abstract geometric version of the law must be just as valid in curved spacetime as in flat (equivalence principle).

Note the analogy with the law of charge conservation, $\nabla \cdot \mathbf{J} = 0$, in electrodynamics (exercise 3.16) and with the local law of energy-momentum conservation, $\nabla \cdot \mathbf{T} = 0$ (§§5.9 and 16.2). In a very deep sense, the forms of these three laws are dictated by the theorem of Gauss (§5.9, and Boxes 5.3, 5.4).

The second law of thermodynamics states that, in flat spacetime or in curved, entropy can be generated but not destroyed. Apply this law to a fluid element of volume V containing a fixed number of baryons N . The entropy it contains is

Second law of
thermodynamics

$$S = Ns = nsV.$$

Entropy may flow in and out across the faces of the fluid element (“heat flow” between neighboring fluid elements); but for simplicity assume it does not; or if it does, assume that it flows too slowly to have any significance for the problem at hand. Then the entropy in the fluid element can only increase:

$$d(nsV)/d\tau \geq 0 \quad \text{when negligible entropy is exchanged between neighboring fluid elements;}$$

i.e. [combine with equation (22.1)]

$$ds/d\tau \geq 0 \quad (\text{no entropy exchange}). \quad (22.5)$$

So long as the fluid element remains in thermodynamic equilibrium, its entropy will actually be conserved [= in equation (22.5)]; but at a shock wave, where equilibrium is momentarily broken, the entropy will increase (conversion of “relative kinetic energy” of neighboring fluid elements into heat). [For discussions of heat flow in special and general relativity, see Exercise 22.7. For discussion of shock waves, see Taub (1948), de Hoffman and Teller (1950), Israel (1960), May and White (1967), Zel'dovich and Rayzer (1967), Lichnerowicz (1967, 1971), and Thorne (1973a).]

Shock waves and heat flow

The first law of thermodynamics, in the proper reference frame of a fluid element, is identical to the first law in flat spacetime (“principle of equivalence”); and in flat spacetime the first law is merely the law of energy conservation:

First law of thermodynamics

$$d\left(\begin{array}{l} \text{energy in a volume element containing} \\ \text{a fixed number, } A, \text{ of baryons} \end{array}\right) = -p \, d(\text{volume}) + T \, d(\text{entropy});$$

i.e.,

$$d(\rho A/n) = -p d(A/n) + T d(As);$$

i.e.,

$$d\rho = \frac{\rho + p}{n} dn + nT ds.$$

Query: what kind of a “ d ” appears here? For a simple fluid, the values of two potentials, e.g., n and s , fix all the others uniquely; so *any* change in ρ must be determined uniquely by the changes in n and s . It matters not whether the changes are measured along the world line of a given fluid element, or in some other direction. Thus, the “ d ” in the first law can be interpreted as an exterior derivative

$$\mathbf{d}\rho = \frac{\rho + p}{n} \mathbf{d}n + nT \mathbf{d}s; \quad (22.6)$$

and the changes along a given direction in the fluid (along a given tangent vector \mathbf{v}) can be written

$$\begin{aligned} \nabla_{\mathbf{v}}\rho &\equiv \langle \mathbf{d}\rho, \mathbf{v} \rangle = \frac{\rho + p}{n} \langle \mathbf{d}n, \mathbf{v} \rangle + nT \langle \mathbf{d}s, \mathbf{v} \rangle \\ &= \frac{\rho + p}{n} \nabla_{\mathbf{v}}n + nT \nabla_{\mathbf{v}}s. \end{aligned}$$

Pressure and temperature calculated from $\rho(n, s)$

Equation (22.6) lends itself to interpretation in two opposite senses: as a way to deduce the density of mass-energy of the medium from information about pressure (as a function of n and s) and temperature (as a function of n and s); and conversely, as a way to deduce the two functions $p(n, s)$ and $T(n, s)$ from the one function $\rho(n, s)$. It is natural to look at the second approach first; who does not like a strategy that makes an intellectual profit? Regarding ρ as a known (or calculable) function of n and s , one deduces from (22.6)

$$\begin{aligned} \frac{\rho + p}{n} &= \left(\frac{\partial \rho}{\partial n} \right)_s, \\ nT &= \left(\frac{\partial \rho}{\partial s} \right)_n, \end{aligned}$$

and thence pressure and temperature individually,

$$p(n, s) = n \left(\frac{\partial \rho}{\partial n} \right)_s - \rho, \quad (22.7a)$$

$$T(n, s) = \frac{1}{n} \left(\frac{\partial \rho}{\partial s} \right)_n \quad (22.7b)$$

(“two equations of state from one”). The analysis simplifies still further when the fluid, already assumed to be everywhere of the same composition, is also everywhere

endowed with the same entropy per baryon, s , and is in a state of adiabatic flow (no shocks or heat conduction). Then the density $\rho = \rho(n, s)$ reduces to a function of one variable out of which one derives everything (ρ, p, μ) needed for the hydrodynamics and the gravitation physics of the system (next chapter). Other choices of the “primary thermodynamic potential” are appropriate under other circumstances (see Box 22.1).

If differentiation leads from $\rho(n, s)$ to $p(n, s)$ and $T(n, s)$, it does not follow that one can take any two functions $p(n, s)$ and $T(n, s)$ and proceed “backwards” (by integration) to the “primary function”, $\rho(n, s)$. To be compatible with the first law of thermodynamics (22.6), the two functions must satisfy the consistency requirement [“Maxwell relation”; equality of second partial derivatives of ρ]

Maxwell relation

$$(\partial p / \partial s)_n = n^2 (\partial T / \partial n)_s. \quad (22.7c)$$

Box 22.1 PRINCIPAL ALTERNATIVES FOR “PRIMARY THERMODYNAMIC POTENTIAL” TO DESCRIBE A FLUID

Primary thermodynamic potential and quantities on which it is most appropriately envisaged to depend

“Secondary” thermodynamic quantities obtained by differentiation of primary with or without use of

Conditions under which convenient, appropriate, and relevant

$$d\left(\frac{\rho}{n}\right) + pd\left(\frac{1}{n}\right) - T ds = 0$$

“Density”; total amount of mass-energy (rest + thermal + ...) per unit volume

$$\rho = \rho(n, s)$$

$$p(n, s) = n \left(\frac{\partial \rho}{\partial n} \right)_s - \rho$$

Conditions of adiabatic flow (no shocks or heat conduction), so that s stays constant along streamline

$$T(n, s) = \frac{1}{n} \left(\frac{\partial \rho}{\partial s} \right)_n$$

$$\mu(n, s) = \frac{p + \rho}{n} = \left(\frac{\partial \rho}{\partial n} \right)_s$$

“Physical free energy”

$$a(n, T) = \frac{\rho}{n} - Ts$$

$$p(n, T) = n^2 \left(\frac{\partial a}{\partial n} \right)_T$$

Know or can calculate a (or the “sum over states” of statistical mechanics) for conditions of specified volume per baryon and temperature

$$s(n, T) = - \left(\frac{\partial a}{\partial T} \right)_n$$

$$\rho(n, T) = -nT^2 \left[\frac{\partial(a/T)}{\partial T} \right]_n$$

“Chemical free energy”

$$f(p, T) = \frac{\rho + p}{n} - Ts$$

$$1/n(p, T) = (\partial f / \partial p)_T$$

$$s(p, T) = -(\partial f / \partial T)_p$$

Relevant for determining equilibrium when pressure and temperature are specified

$$\rho(p, T) = \frac{f - T(\partial f / \partial T)_p}{(\partial f / \partial p)_T} - p$$

“Chemical potential” (“energy to inject” expressed on a “per baryon” basis)

$$\mu(p, s) = \frac{p + \rho}{n}$$

$$1/n(p, s) = (\partial \mu / \partial p)_s$$

$$T(p, s) = (\partial \mu / \partial s)_p$$

$$\rho(p, s) = \frac{\mu}{(\partial \mu / \partial p)_s} - p$$

When injection energy [= Fermi energy for an ideal Fermi gas, relativistic or not; see exercise 22.3] is the center of attention

Chemical potential equals “injection energy” at fixed entropy per baryon and total volume

The chemical potential μ is also a unique function of n and s . It is defined as follows. (1) Take a sample of the simple fluid in a fixed thermodynamic state (fixed n and s). (2) Take, separately, a much smaller sample of the same fluid, containing δA baryons in the same thermodynamic state as the large sample (same n and s). (3) Inject the smaller sample into the larger one, holding the volume of the large sample fixed during the injection process. (4) The total mass-energy injected,

$$\delta M_{\text{injected}} = \rho \times (\text{volume of injected fluid}) = \rho(\delta A/n),$$

plus the work required to perform the injection

$$\begin{aligned}\delta W_{\text{injection}} &= \left(\begin{array}{l} \text{work done against pressure of large sample} \\ \text{to open up space in it for the injected fluid} \end{array} \right) \\ &= p(\text{volume of injected fluid}) = p(\delta A/n),\end{aligned}$$

is equal to $\mu \delta A$:

$$\mu \delta A = \delta M_{\text{injected}} + \delta W_{\text{injection}} = \frac{\rho + p}{n} \delta A.$$

Stated more briefly:

$$\begin{aligned}\mu &= \left(\begin{array}{l} \text{total mass-energy required, per baryon, to “create” and} \\ \text{inject a small additional amount of fluid into a given} \\ \text{sample, without changing } s \text{ or volume of the sample} \end{array} \right) \\ &= \frac{\rho + p}{n} = \left(\frac{\partial \rho}{\partial n} \right)_s. \\ &\quad \uparrow \\ &\quad [\text{by first law of thermodynamics (22.6)}]\end{aligned}\tag{22.8}$$

All the above laws and equations of thermodynamics are the same in curved spacetime as in flat spacetime; and the same in (relativistic) flat spacetime as in classical nonrelativistic thermodynamics—except for the inclusion of rest mass, together with all other forms of mass-energy, in ρ and μ . The reason is simple: the laws are all formulated as scalar equations linking thermodynamic variables that one measures in the rest frame of the fluid.

§22.3. HYDRODYNAMICS IN CURVED SPACETIME*

Laws of hydrodynamics for simple fluid without heat flow or viscosity:

A simple perfect fluid flows through spacetime. It might be the Earth’s atmosphere circulating in the Earth’s gravitational field. It might be the gaseous interior of the Sun at rest in its own gravitational field. It might be interstellar gas accreting onto a black hole. But whatever and wherever the fluid may be, its motion will be governed by the curved-spacetime laws of thermodynamics (§22.2) plus the local

*For more detailed treatments of this subject see, e.g., Ehlers (1961), Taub (1971), Ellis (1971), Lichnerowicz (1967), Cattaneo (1971), and references cited therein; see also the references on kinetic theory cited at the beginning of §22.6.

law of energy-momentum conservation, $\nabla \cdot \mathbf{T} = 0$. The chief objective of this section is to reduce the equation $\nabla \cdot \mathbf{T} = 0$ to usable form. The reduction will be performed in the text using abstract notation; the reader is encouraged to repeat the reduction using index notation.

The stress-energy tensor for a perfect fluid, in curved spacetime as in flat (equivalence principle!), is

$$\mathbf{T} = (\rho + p)\mathbf{u} \otimes \mathbf{u} + p\mathbf{g}. \quad (22.9)$$

(See §5.5.) Its divergence is readily calculated using the chain rule; using the “compatibility relation between \mathbf{g} and ∇ ,” $\nabla \mathbf{g} = 0$; using the identity $(\nabla p) \cdot \mathbf{g} = \nabla p$ (which one readily verifies in index notation); and using

$$\begin{aligned} 0 &= \nabla \cdot \mathbf{T} = [\nabla(\rho + p) \cdot \mathbf{u}] \mathbf{u} + [(\rho + p) \nabla \cdot \mathbf{u}] \mathbf{u} + [(\rho + p) \mathbf{u}] \cdot \nabla \mathbf{u} + (\nabla p) \cdot \mathbf{g} \\ &\quad \uparrow \text{[divergence on first slot]} \\ &= [\nabla_{\mathbf{u}}\rho + \nabla_{\mathbf{u}}p + (\rho + p)\nabla \cdot \mathbf{u}] \mathbf{u} + (\rho + p) \nabla_{\mathbf{u}}\mathbf{u} + \nabla p. \end{aligned} \quad (22.10)$$

The component of this equation along the 4-velocity is especially simple (recall that $\mathbf{u} \cdot \nabla_{\mathbf{u}}\mathbf{u} = \frac{1}{2}\nabla_{\mathbf{u}}\mathbf{u}^2 = 0$ because $\mathbf{u}^2 \equiv -1$):

$$\begin{aligned} 0 &= \mathbf{u} \cdot (\nabla \cdot \mathbf{T}) = -[\nabla_{\mathbf{u}}\rho + \nabla_{\mathbf{u}}p + (\rho + p)\nabla \cdot \mathbf{u}] + \nabla_{\mathbf{u}}p \\ &= -\nabla_{\mathbf{u}}\rho - (\rho + p)\nabla \cdot \mathbf{u}. \end{aligned}$$

Combine this with the equation of baryon conservation (22.3); the result is

$$\frac{d\rho}{d\tau} = \frac{(\rho + p)}{n} \frac{dn}{d\tau}. \quad (22.11a)$$

(2) Local energy conservation: adiabaticity of flow

Notice that this is identical to the first law of thermodynamics (22.6) applied along a flow line, plus the assumption that the entropy per baryon is conserved along a flow line

$$ds/d\tau = 0. \quad (22.11b)$$

There is no reason for surprise at this result. To insist on thermodynamic equilibrium and to demand that the entropy remain constant is to require zero exchange of heat between one element of the fluid and another. But the stress-energy tensor (22.9) recognizes that heat exchange is absent. Any heat exchange would show up as an energy flux term in \mathbf{T} (Ex. 22.7); but no such term is present. Consequently, when one studies local energy conservation by evaluating $\mathbf{u} \cdot (\nabla \cdot \mathbf{T}) = 0$, the stress-energy tensor reports that no heat flow is occurring—i.e. that $ds/d\tau = 0$.

Three components of $\nabla \cdot \mathbf{T} = 0$ remain: the components orthogonal to the fluid’s 4-velocity. One can pluck them out of $\nabla \cdot \mathbf{T} = 0$, leaving behind the component along \mathbf{u} , by use of the “projection tensor”

$$\mathbf{P} \equiv \mathbf{g} + \mathbf{u} \otimes \mathbf{u}. \quad (22.12)$$

Box 22.2 THERMODYNAMICS AND HYDRODYNAMICS FOR A SIMPLE PERFECT FLUID IN CURVED SPACETIME
A. Ten Quantities Characterize the Fluid

Thermodynamic potentials all measured in rest frame

n , baryon number density

ρ , density of total mass-energy

p , pressure

T , temperature

s , entropy per baryon

μ , chemical potential per baryon

Four components of the fluid 4-velocity

Equation for chemical potential

$$\mu = (\rho + p)/n, \quad (4)$$

which can be combined with $\rho(n, s)$ and $p(n, s)$ to give $\mu(n, s)$.

Law of baryon conservation

$$dn/d\tau \equiv \nabla_u n = -n \nabla \cdot u. \quad (5)$$

Conservation of energy along flow lines, which (assuming no energy exchange between adjacent fluid elements) means “adiabatic flow”

$ds/d\tau = 0$ except in shock waves, where

$$ds/d\tau > 0. \quad (6)$$

[Shock waves are not treated in this book; see Taub (1948), de Hoffman and Teller (1950), Israel (1960), May and White (1967), Zel'dovich and Rayzer (1967); Lichnerowicz (1967, 1971); and Thorne (1973a).]

Euler equations

$$(\rho + p) \nabla_u u = -(\mathbf{g} + \mathbf{u} \otimes \mathbf{u}) \cdot \nabla p, \quad (7), (8), (9)$$

which determine the flow lines to which \mathbf{u} is tangent.

Normalization of 4-velocity

$$\mathbf{u} \cdot \mathbf{u} = -1. \quad (10)$$

(See exercise 22.4.) Contracting \mathbf{P} with $\nabla \cdot \mathbf{T} = 0$ [equation (22.10)] gives

(3) Euler equation

$$(\rho + p) \nabla_u u = -\mathbf{P} \cdot (\nabla p) \equiv -[\nabla p + (\nabla_u p) \mathbf{u}]. \quad (22.13)$$

This is the “Euler equation” of relativistic hydrodynamics. It has precisely the same form as the corresponding flat-spacetime Euler equation:

$$\left(\begin{array}{l} \text{inertial mass} \\ \text{per unit volume} \\ \text{[exercise 5.4]} \end{array} \right) \times \left(\begin{array}{l} \text{4-acceleration} \\ \text{of fluid} \end{array} \right) = - \left(\begin{array}{l} \text{pressure gradient} \\ \text{in the 3-surface} \\ \text{orthogonal to 4-velocity} \end{array} \right). \quad (22.13')$$

The pressure gradient, not “gravity,” is responsible for all deviation of flow lines from geodesics.

Box 22.2 reorganizes and summarizes the above laws of thermodynamics and hydrodynamics.

Exercise 22.1. DIVERGENCE OF FLOW LINES PRODUCES VOLUME CHANGES EXERCISES

Derive the equation $dV/d\tau = (\nabla \cdot \mathbf{u})V$ [equation (22.2)] for the rate of change of volume of a fluid element. [Hint: Pick an event \mathcal{P}_0 , and calculate in a local Lorentz frame at \mathcal{P}_0 which momentarily moves with the fluid (“rest frame at \mathcal{P}_0 ”).] [Solution: At events near \mathcal{P}_0 the fluid has a very small ordinary velocity $v^j = dx^j/dt$. Consequently a cube of fluid at \mathcal{P}_0 with edges $\Delta x = \Delta y = \Delta z = L$ changes its edges, after time δt , by the amounts

$$\begin{aligned}\delta(\Delta x) &= [(dx/dt)\delta t]_{\text{at front face}} - [(dx/dt)\delta t]_{\text{at back face}} \\ &= (\partial v^x/\partial x)L\delta t, \\ \delta(\Delta y) &= (\partial v^y/\partial y)L\delta t, \\ \delta(\Delta z) &= (\partial v^z/\partial z)L\delta t.\end{aligned}$$

The corresponding change in volume is

$$\delta(\Delta x \Delta y \Delta z) = (\partial v^j/\partial x^j)L^3\delta t;$$

so the rate of change of volume is

$$\partial V/\partial t = V(\partial v^j/\partial x^j).$$

But in the local Lorentz rest frame at and near \mathcal{P}_0 (where $x^\alpha = 0$), the metric coefficients are $g_{\mu\nu} = \eta_{\mu\nu} + O(|x^\alpha|^2)$, and the ordinary velocity is $v^j = O(|x^\alpha|)$; so

$$\begin{aligned}u^0 &= \frac{dt}{d\tau} = \frac{dt}{(-g_{\mu\nu}dx^\mu dx^\nu)^{1/2}} = 1 + O(|x^\alpha|^2), \\ u^j &= \frac{dx^j}{d\tau} = v^j + O(|x^\alpha|^3).\end{aligned}$$

Thus, the derivatives $\partial V/\partial t$ and $V(\partial v^j/\partial x^j)$ at \mathcal{P}_0 are

$$\begin{aligned}\partial V/\partial t &= u^\alpha \partial V/\partial x^\alpha = u^\alpha V_{,\alpha} = dV/d\tau \\ &= V(\partial v^j/\partial x^j) = V(\partial u^\alpha/\partial x^\alpha) = Vu^\alpha_{;\alpha} = V(\nabla \cdot \mathbf{u}). \quad \text{Q.E.D.}\end{aligned}$$

[Note that by working in flat spacetime, one could have inferred more easily that $\partial V/\partial t = dV/d\tau$ and $\partial v^j/\partial x^j = \nabla \cdot \mathbf{u}$; one would then have concluded $dV/d\tau = (\nabla \cdot \mathbf{u})V$; and one could have invoked the equivalence principle to move this law into curved spacetime.]

Exercise 22.2. EQUATION OF CONTINUITY

Show that in the nonrelativistic limit in flat spacetime the equation of baryon conservation (22.3) becomes the “equation of continuity”

$$\frac{\partial n}{\partial t} + \frac{\partial}{\partial x^j}(nv^j) = 0.$$

Exercise 22.3. CHEMICAL POTENTIAL FOR IDEAL FERMI GAS

Show that the chemical potential of an ideal Fermi gas, nonrelativistic or relativistic, is (at zero temperature) equal to the Fermi energy (energy of highest occupied momentum state) of that gas.

Exercise 22.4. PROJECTION TENSORS

Show that contraction of a tangent vector \mathbf{B} with the “projection tensor” $\mathbf{P} \equiv \mathbf{g} + \mathbf{u} \otimes \mathbf{u}$ projects \mathbf{B} into the 3-surface orthogonal to the 4-velocity vector \mathbf{u} . [Hint: perform the

calculation in an orthonormal frame with $\mathbf{e}_\hat{\alpha} = \mathbf{u}$, and write $\mathbf{B} = B^\alpha \mathbf{e}_\hat{\alpha}$; then show that $\mathbf{P} \cdot \mathbf{B} = B^j \mathbf{e}_j$.] If \mathbf{n} is a unit spacelike vector, show that $\mathbf{P} \equiv \mathbf{g} - \mathbf{n} \otimes \mathbf{n}$ is the corresponding projection operator. Note: There is no *unique* concept of “the projection orthogonal to a null vector.” Why? [Hint: draw pictures in flat spacetime suppressing one spatial dimension.]

Exercise 22.5. PRESSURE GRADIENT IN STATIONARY GRAVITATIONAL FIELD

A perfect fluid is at rest (flow lines have $x^j = \text{constant}$) in a stationary gravitational field (metric coefficients are independent of x^0). Show that the pressure gradient required to “support the fluid against gravity” (i.e., to make its flow lines be $x^j = \text{constant}$ instead of geodesics) is

$$\frac{\partial p}{\partial x^0} = 0, \quad \frac{\partial p}{\partial x^j} = -(\rho + p) \frac{\partial \ln \sqrt{-g_{00}}}{\partial x^j}. \quad (22.14)$$

Evaluate this pressure gradient in the Newtonian limit, using the coordinate system and metric coefficients of equation (18.15c).

Exercise 22.6. EXPANSION, ROTATION, AND SHEAR

Let a field of fluid 4-velocities $\mathbf{u}(\mathcal{P})$ be given.

- (a) Show that $\nabla \mathbf{u}$ can be decomposed in the following manner:

$$u_{\alpha;\beta} = \omega_{\alpha\beta} + \sigma_{\alpha\beta} + \frac{1}{3} \theta P_{\alpha\beta} - a_\alpha u_\beta, \quad (22.15a)$$

where \mathbf{a} is the *4-acceleration* of the fluid

$$a_\alpha \equiv u_{\alpha;\beta} u^\beta, \quad (22.15b)$$

θ is the “*expansion*” of the fluid world lines

$$\theta \equiv \nabla \cdot \mathbf{u} = u^\alpha_{;\alpha}, \quad (22.15c)$$

$P_{\alpha\beta}$ is the *projection tensor*

$$P_{\alpha\beta} \equiv g_{\alpha\beta} + u_\alpha u_\beta, \quad (22.15d)$$

$\sigma_{\alpha\beta}$ is the *shear tensor* of the fluid

$$\sigma_{\alpha\beta} \equiv \frac{1}{2} (u_{\alpha;\mu} P^\mu_\beta + u_{\beta;\mu} P^\mu_\alpha) - \frac{1}{3} \theta P_{\alpha\beta}, \quad (22.15e)$$

and $\omega_{\alpha\beta}$ is the *rotation 2-form* of the fluid

$$\omega_{\alpha\beta} \equiv \frac{1}{2} (u_{\alpha;\mu} P^\mu_\beta - u_{\beta;\mu} P^\mu_\alpha). \quad (22.15f)$$

(b) Each of the component parts of this decomposition has a simple physical interpretation in the local rest frames of the fluid. The interpretation of the 4-acceleration \mathbf{a} in terms of accelerometer readings should be familiar. Exercise 22.1 showed that the expansion $\theta = \nabla \cdot \mathbf{u}$ describes the rate of increase of the volume of a fluid element,

$$\theta = (1/V)(dV/d\tau). \quad (22.15g)$$

Exercise 22.4 explored the meaning and use of the projection tensor \mathbf{P} . Verify that in a local Lorentz frame ($g_{\alpha\beta} = \eta_{\alpha\beta}$, $\Gamma^{\hat{\alpha}}_{\beta\gamma} = 0$) momentarily moving with the fluid ($u^{\hat{\alpha}} = \delta^{\alpha}_0$), $\sigma_{\hat{\alpha}\hat{\beta}}$ and $\omega_{\hat{\alpha}\hat{\beta}}$ reduce to the classical (nonrelativistic) shear and rotation of the fluid. [See, e.g., §§2.4 and 2.5 of Ellis (1971) for both classical and relativistic descriptions of shear and rotation.]

Exercise 22.7. HYDRODYNAMICS WITH VISCOSITY AND HEAT FLOW.*

(a) In §15 of Landau and Lifshitz (1959), one finds an analysis of viscous stresses for a classical (nonrelativistic) fluid. By carrying that analysis over directly to the local Lorentz rest frame of a relativistic fluid, and by then generalizing to frame-independent language, show that the contribution of viscosity to the stress-energy tensor is

$$\mathbf{T}^{(\text{visc})} = -2\eta\sigma - \xi\theta\mathbf{P}, \quad (22.16a)$$

where $\eta \geq 0$ is the “coefficient of dynamic viscosity”; $\xi \geq 0$ is the “coefficient of bulk viscosity”; and $\sigma, \theta, \mathbf{P}$ are the shear, expansion, and projection tensor of the fluid.

(b) An idealized description of heat flow in a fluid introduces the *heat-flux 4-vector* \mathbf{q} with components in the local rest-frame of the fluid,

$$q^{\hat{0}} = 0, \quad q^j = \begin{pmatrix} \text{energy per unit time crossing unit} \\ \text{surface perpendicular to } \mathbf{e}_j \end{pmatrix}. \quad (22.16b)$$

By generalizing from the fluid rest frame to frame-independent language, show that the contribution of heat flux to the stress-energy tensor is

$$\mathbf{T}^{(\text{heat})} = \mathbf{u} \otimes \mathbf{q} + \mathbf{q} \otimes \mathbf{u}. \quad (22.16c)$$

Thereby conclude that, in this idealized picture, the stress-energy tensor for a viscous fluid with heat conduction is

$$T^{\alpha\beta} = \rho u^\alpha u^\beta + (p - \xi\theta) P^{\alpha\beta} - 2\eta\sigma^{\alpha\beta} + q^\alpha u^\beta + u^\alpha q^\beta. \quad (22.16d)$$

(c) Define the entropy 4-vector \mathbf{s} by

$$\mathbf{s} \equiv ns\mathbf{u} + \mathbf{q}/T. \quad (22.16e)$$

By calculations in the local rest-frame of the fluid, show that

$$\begin{aligned} \nabla \cdot \mathbf{s} &= \left(\begin{array}{l} \text{rate of increase of entropy} \\ \text{in a unit volume} \end{array} \right) - \left(\begin{array}{l} \text{rate at which heat and fluid} \\ \text{carry entropy into a unit volume} \end{array} \right) \\ &= \left(\begin{array}{l} \text{rate at which entropy is being} \\ \text{generated in a unit volume} \end{array} \right). \end{aligned} \quad (22.16f)$$

Thereby arrive at the following form of the *second law of thermodynamics*:

$$\nabla \cdot \mathbf{s} \geq 0. \quad (22.16g)$$

(d) Calculate the law of local energy conservation, $\mathbf{u} \cdot \nabla \cdot \mathbf{T} = 0$, for a viscous fluid with heat flow. Combine with the first law of thermodynamics and with the law of baryon conservation to obtain

$$T \nabla \cdot \mathbf{s} = \xi\theta^2 + 2\eta\sigma_{\alpha\beta}\sigma^{\alpha\beta} - q^\alpha(T_{,\alpha}/T + a_\alpha). \quad (22.16h)$$

Interpret each term of this equation as a contribution to entropy generation (*example*: $2\eta\sigma_{\alpha\beta}\sigma^{\alpha\beta}$ describes entropy generation by viscous heating). [Note: The term $q^\alpha a_\alpha$ is relativistic in origin. It is associated with the inertia of the flowing heat.]

(e) When one takes account of the inertia of the flowing heat, one obtains the following generalization of the classical law of heat conduction:

$$q^\alpha = -\kappa P^{\alpha\beta}(T_{,\beta} + Ta_\beta) \quad (22.16i)$$

*Exercise supplied by John M. Stewart.

(Eckart 1940). Here κ is the *coefficient of thermal conductivity*. Use this equation to show that, for a fluid at rest in a stationary gravitational field (Exercise 22.5),

$$q_0 = 0, \quad q_j = -\frac{\kappa}{\sqrt{-g_{00}}} (T\sqrt{-g_{00}})_{,j}. \quad (22.16j)$$

[Thus, thermal equilibrium corresponds not to constant temperature, but to the redshifted temperature distribution $T\sqrt{-g_{00}} = \text{constant}$; Tolman (1934a), p. 313.] Also, use the idealized law of heat conduction (22.16i) to reexpress the rate of entropy generation as

$$T \nabla \cdot \mathbf{s} = \xi \theta^2 + 2\eta \sigma_{\alpha\beta} \sigma^{\alpha\beta} + (\kappa/T) P^{\alpha\beta} (T_{,\alpha} + Ta_\alpha) (T_{,\beta} + Ta_\beta) \geq 0. \quad (22.16k)$$

[For further details about heat flow and for discussions of the limitations of the above idealized description, see e.g., §4.18 of Ehlers (1971); also Marle (1969), Anderson (1970), Stewart (1971), and papers cited therein.]

§22.4. ELECTRODYNAMICS IN CURVED SPACETIME

Electric and magnetic fields

In a local Lorentz frame in the presence of gravity, an observer can measure the electric and magnetic fields \mathbf{E} and \mathbf{B} using the usual Lorentz force law for charged particles. As in special relativity, he can regard \mathbf{E} and \mathbf{B} as components of an electromagnetic field tensor,

$$F^{\hat{\alpha}\hat{\beta}} = -F^{\hat{\beta}\hat{\alpha}} = E^{\hat{\beta}}, \quad F^{\hat{\beta}\hat{k}} = \epsilon^{\hat{\beta}\hat{k}\hat{l}} B^{\hat{l}},$$

he can regard the charge and current densities as components of a 4-vector $J^{\hat{\alpha}}$, and he can write Maxwell's equations and the Lorentz force equation in the special relativistic form,

$$\begin{aligned} F^{\hat{\alpha}\hat{\beta}}_{,\hat{\beta}} &= 4\pi J^{\hat{\alpha}}, & F_{\hat{\alpha}\hat{\beta},\hat{\gamma}} + F_{\hat{\beta}\hat{\gamma},\hat{\alpha}} + F_{\hat{\gamma}\hat{\alpha},\hat{\beta}} &= 0, \\ ma^{\hat{\alpha}} &= F^{\hat{\alpha}\hat{\beta}} qu_{\hat{\beta}} & \left(\begin{array}{l} m = \text{mass of particle}, q = \text{charge}, \\ u^{\hat{\alpha}} = 4\text{-velocity}, a^{\hat{\alpha}} = 4\text{-acceleration} \end{array} \right). \end{aligned}$$

In any other frame these equations will have the same form, but with commas replaced by semicolons

$$F^{\alpha\beta}_{;\beta} = 4\pi J^\alpha, \quad (22.17a)$$

$$F_{\alpha\beta;\gamma} + F_{\beta\gamma;\alpha} + F_{\gamma\alpha;\beta} = 0, \quad (22.17b)$$

$$ma^\alpha = F^{\alpha\beta} qu_\beta. \quad (22.17c)$$

These are the basic equations of electrodynamics in the presence of gravity. From them follows everything else. For example, as in special relativity, so also here (exercise 22.9), they imply the equation of charge conservation

Charge conservation

$$J^\alpha_{;\alpha} = 0; \quad (22.18a)$$

and for an electromagnetic field interacting with charged matter (exercise 22.10) they imply vanishing divergence for the sum of the stress-energy tensors

$$(T^{(\text{EM})\alpha\beta} + T^{(\text{MATTER})\alpha\beta})_{;\beta} = 0. \quad (22.18b)$$

Local conservation of energy-momentum

As in special relativity, so also here, one can introduce a vector potential A^μ . Replacing commas by semicolons in the usual special-relativistic expression for $F^{\mu\nu}$ in terms of A^μ , one obtains

$$F_{\mu\nu} = A_{\nu;\mu} - A_{\mu;\nu}. \quad (22.19a)$$

Vector potential

If all is well, this equation should guarantee (as in special relativity) that the Maxwell equations (22.17b) are satisfied. Indeed, it does, as one sees in exercise 22.8. To derive the wave equation that governs the vector potential, insert expression (22.19a) into the remaining Maxwell equations (22.17a), obtaining

$$-A^{\alpha;\beta}_{\beta} + A^{\beta;\alpha}_{\beta} = 4\pi J^\alpha; \quad (22.19b)$$

then commute covariant derivatives in the first term using the identity (16.6c), to obtain

$$-A^{\alpha;\mu}_{\mu} + A^{\mu}_{;\mu} ;^{\alpha} + R^\alpha_\mu A^\mu = 4\pi J^\alpha. \quad (22.19b')$$

Finally, adopting the standard approach of special relativity, impose the Lorentz gauge condition

$$A^\mu_{;\mu} = 0, \quad (22.19c) \quad \text{Lorentz gauge condition}$$

thereby bringing the wave equation (22.19b') into the form

$$(\Delta_{dR} A)^\alpha \equiv -A^{\alpha;\beta}_{\beta} + R^\alpha_\beta A^\beta = 4\pi J^\alpha. \quad (22.19d)$$

Wave equation for vector potential

The “de Rham vector wave operator” Δ which appears here is, apart from sign, a generalized d’Alambertian for vectors in curved spacetime. Mathematically it is more powerful than $-A^{\alpha;\beta}_{\beta}$, and than any other operator that reduces to (minus) the d’Alambertian in special relativity. [For a discussion, see de Rham (1955).]

Although the electrodynamic equations (22.17a)–(22.19b) are all obtained from special relativity by the comma-goes-to-semicolon rule, the wave equation (22.19d) for the vector potential is not (“curvature coupling”; see Box 16.1). Nevertheless, when spacetime is flat (so $R^\alpha_\beta = 0$), (22.19d) does reduce to the usual wave equation of special relativity.

Exercise 22.8. THE VECTOR POTENTIAL FOR ELECTRODYNAMICS

Show that in any coordinate frame the connection coefficients cancel out of both equations (22.19a) and (22.17b), so they can be written

$$F_{\mu\nu} = A_{\nu,\mu} - A_{\mu,\nu}, \quad (22.20a)$$

$$F_{\alpha\beta,\gamma} + F_{\beta\gamma,\alpha} + F_{\gamma\alpha,\beta} = 0. \quad (22.20b)$$

EXERCISES

(In the language of differential forms these equations are $\mathbf{F} = d\mathbf{A}$, $d\mathbf{F} = 0$.) Then use this form of the equations to show that equation (22.19a) implies equation (22.17b), as asserted in the text.

Exercise 22.9. CHARGE CONSERVATION IN THE PRESENCE OF GRAVITY

Show that Maxwell's equations (22.17a,b) imply the equation of charge conservation (22.18a) when gravity is present, just as they do in special relativity theory. [Hints: Use the antisymmetry of $F^{\alpha\beta}$; and beware of the noncommutation of the covariant derivatives, which must be handled using equations (16.6). Alternatively, show that in coordinate frames, equation (22.17a) can be written as

$$\frac{1}{\sqrt{|g|}} \frac{\partial}{\partial x^\beta} (\sqrt{|g|} F^{\alpha\beta}) = 4\pi J^\alpha \quad (22.17a')$$

and (22.18a) as

$$J^\alpha_{;\alpha} \equiv \frac{1}{\sqrt{|g|}} \frac{\partial}{\partial x^\alpha} (\sqrt{|g|} J^\alpha) = 0, \quad (22.18a')$$

and carry out the demonstration in a coordinate frame.]

**Exercise 22.10. INTERACTING ELECTROMAGNETIC FIELD
AND CHARGED MATTER**

As in special relativity, so also in the presence of gravity ("equivalence principle"), the stress-energy tensor for an electromagnetic field is

$$T^{(\text{EM})\alpha\beta} = \frac{1}{4\pi} \left(F_{\alpha\mu} F_\beta^\mu - \frac{1}{4} F_{\mu\nu} F^{\mu\nu} g_{\alpha\beta} \right). \quad (22.21)$$

Use Maxwell's equations (22.17a,b) in the presence of gravity to show that

$$T^{(\text{EM})\alpha\beta}_{;\beta} = -F^{\alpha\beta} J_\beta. \quad (22.22)$$

But $F^{\alpha\beta} J_\beta$ is just the Lorentz 4-force per unit volume with which the electromagnetic field acts on the charged matter [see the Lorentz force equation (22.17c); also equation (5.43)]; i.e., it is $T^{(\text{MATTER})\alpha\beta}_{;\beta}$. Consequently, the above equation can be rewritten in the form (22.18b) cited in the text.

§22.5. GEOMETRIC OPTICS IN CURVED SPACETIME*

Radio waves from the quasar 3C279 pass near the sun and get deflected by its gravitational field. Light rays emitted by newborn galaxies long ago and far away propagate through the cosmologically curved spacetime of the universe, and get focused (and redshifted) producing curvature-enlarged (but dim) images of the galaxies on the Earth's sky.

*Based in part on notes prepared by William L. Burke at Caltech in 1968. For more detailed treatments of geometric optics in curved spacetime, see, e.g., Sachs (1961), Jordan, Ehlers, and Sachs (1961), and Robinson (1961); also references discussed and listed in §41.11.

These and most other instances of the propagation of light and radio waves are subject to the laws of geometric optics. This section derives those laws, in curved spacetime, from Maxwell's equations.

The fundamental laws of geometric optics are: (1) light rays are null geodesics; (2) the polarization vector is perpendicular to the rays and is parallel-propagated along the rays; and (3) the amplitude is governed by an adiabatic invariant which, in quantum language, states that the number of photons is conserved.

The conditions under which these laws hold are defined by conditions on three lengths: (1) the typical reduced wavelength of the waves,

$$\lambda \equiv \frac{\lambda}{2\pi} = \left(\text{“classical distance of closest approach for a photon with one unit of angular momentum”} \right), \quad (22.23a)$$

as measured in a typical local Lorentz frame (e.g., a frame at rest relative to nearby galaxies); (2) the typical length \mathcal{L} over which the amplitude, polarization, and wavelength of the waves vary, e.g., the radius of curvature of a wave front, or the length of a wave packet produced by a sudden outburst in a quasar; (3) the typical radius of curvature \mathcal{R} of the spacetime through which the waves propagate,

$$\mathcal{R} \equiv \left| \begin{array}{l} \text{typical component of } \mathbf{Riemann} \text{ as measured} \\ \text{in typical local Lorentz frame} \end{array} \right|^{-1/2}. \quad (22.23b)$$

Geometric optics is valid whenever the reduced wavelength is very short compared to each of the other scales present,

$$\lambda \ll \mathcal{L} \quad \text{and} \quad \lambda \ll \mathcal{R}, \quad (22.23c)$$

so that the waves can be regarded *locally* as plane waves propagating through spacetime of negligible curvature.

Mathematically one exploits the geometric-optics assumption, $\lambda \ll \mathcal{L}$ and $\lambda \ll \mathcal{R}$, as follows. Focus attention on waves that are highly monochromatic over regions $\lesssim \mathcal{L}$. (More complex spectra can be analyzed by superposition, i.e., by Fourier analysis.) Split the vector potential of electromagnetic theory into a rapidly changing, real phase,

$$\theta \sim (\text{distance propagated})/\lambda,$$

and a slowly changing, complex amplitude (i.e. one with real and imaginary parts),

$$\mathbf{A} = \text{Real part of} \{ \text{amplitude} \times e^{i\theta} \} \equiv \Re \{ \text{amplitude} \times e^{i\theta} \}.$$

Imagine holding fixed the scale of the amplitude variation, \mathcal{L} , and the scale of the spacetime curvature, \mathcal{R} , while making the reduced wavelength, λ , shorter and shorter. The phase will get larger and larger ($\theta \propto 1/\lambda$) at any fixed event in spacetime, but the amplitude as a function of location in spacetime can remain virtually unchanged,

$$\text{Amplitude} = \left[\begin{array}{l} \text{dominant part,} \\ \text{independent of } \lambda \end{array} \right] + \left[\begin{array}{l} \text{small corrections (deviations from} \\ \text{geometric optics) due to finite wavelength} \end{array} \right].$$

Overview of geometric optics

Conditions for validity of geometric optics

The “two-length-scale” expansion underlying geometric optics

This circumstance allows one to expand the amplitude in powers of λ :*

$$\text{Amplitude} = \mathbf{a} + \mathbf{b} + \mathbf{c} + \dots$$

\uparrow \uparrow \uparrow
 independent $\propto \lambda$ $\propto \lambda^2$
 of λ

[Actually, the expansion proceeds in powers of the dimensionless number

$$\lambda / (\text{minimum of } \mathcal{L} \text{ and } \mathcal{R}) \equiv \lambda / L. \quad (22.24)$$

Applied mathematicians call this a “two-length-scale expansion”; see, e.g., Cole (1968). The basic short-wavelength approximation here has a long history; see, e.g., Liouville (1837), Rayleigh (1912). Following a suggestion of Debye, it was applied to Maxwell’s equations by Sommerfeld and Runge (1911). It is familiar as the WKB approximation in quantum mechanics, and has many other applications as indicated by the bibliography in Keller, Lewis, and Seckler (1956). The contribution of higher order terms is considered by Kline (1954) and Lewis (1958). See especially the book of Fröman and Fröman (1965).]

It is useful to introduce a parameter ϵ that keeps track of how rapidly various terms approach zero (or infinity) as λ/L approaches zero:

The vector potential in
geometric optics

$$A_\mu = \Re\{(a_\mu + \epsilon b_\mu + \epsilon^2 c_\mu + \dots) e^{i\theta/\epsilon}\}. \quad (22.25)$$

Any term with a factor ϵ^n in front of it varies as $(\lambda/L)^n$ in the limit of very small wavelengths [$\theta \propto (\lambda/L)^{-1}$; $c_\mu \propto (\lambda/L)^2$; etc.]. By convention, ϵ is a dummy expansion parameter with eventual value unity; so it can be dropped from the calculations when it ceases to be useful. And by convention, all “post-geometric-optics corrections” are put into the amplitude terms $\mathbf{b}, \mathbf{c}, \dots$; none are put into θ .

Note that, while the phase θ is a real function of position in spacetime, the amplitude and hence the vectors $\mathbf{a}, \mathbf{b}, \mathbf{c}, \dots$ are complex. For example, to describe monochromatic waves with righthand circular polarization, propagating in the z direction, one could set $\theta = \omega(z - t)$ and $\mathbf{a} = 1/\sqrt{2}a(\mathbf{e}_x + i\mathbf{e}_y)$ with a real; so

$$\mathbf{A} = \Re\left\{\frac{1}{\sqrt{2}}a(\mathbf{e}_x + i\mathbf{e}_y)e^{i\omega(z-t)}\right\} = \frac{1}{\sqrt{2}}a[\cos[\omega(z-t)]\mathbf{e}_x - \sin[\omega(z-t)]\mathbf{e}_y]$$

Basic concepts of geometric optics:

The assumed form (22.25) for the vector potential is the mathematical foundation of geometric optics. All the key equations of geometric optics result from inserting this vector potential into the source-free wave equation $\Delta \mathbf{A} = 0$ [equation (22.19d)] and into the Lorentz gauge condition $\nabla \cdot \mathbf{A} = 0$ [equation (22.19c)]. The resulting equations (derived below) take their simplest form only when expressed in terms of the following:

*The equations for \mathbf{A} are linear. Therefore the analysis would proceed equally well assuming, instead of an amplitude independent of λ , a dominant term $\mathbf{a} \propto \lambda^n$, with $\mathbf{b} \propto \lambda^{n+1}$, $\mathbf{c} \propto \lambda^{n+2}$, etc. The results are independent of n . Choosing $n = 1$ would give field strengths $F_{\mu\nu}$ and energy densities $T_{\mu\nu} \propto F^2 \propto A^2/\lambda^2 \propto \text{constant as } \lambda \rightarrow 0$.

“wave vector,” $\mathbf{k} \equiv \nabla\theta$; (22.26a) (1) wave vector

“scalar amplitude,” $a \equiv (\mathbf{a} \cdot \bar{\mathbf{a}})^{1/2} = (a^\mu \bar{a}_\mu)^{1/2}$; (22.26b) (2) scalar amplitude

“polarization vector,” $\mathbf{f} \equiv \mathbf{a}/a = “unit complex vector along \mathbf{a}”$. (22.26c) (3) polarization vector

(Here $\bar{\mathbf{a}}$ is the complex conjugate of \mathbf{a} .) *Light rays* are defined to be the curves $\mathcal{P}(\lambda)$ normal to surfaces of constant phase θ . Since $\mathbf{k} \equiv \nabla\theta$ is the normal to these surfaces, the differential equation for a light ray is

$$\frac{dx^\mu}{d\lambda} = k^\mu(x) = g^{\mu\nu}(x)\theta_{,\nu}(x). \quad (22.26d)$$

Box 22.3, appropriate for study at this point, shows the polarization vector, wave vector, surfaces of constant phase, and light rays for a propagating wave; the scalar amplitude, not shown there, merely tells the length of the vector amplitude \mathbf{a} . Insight into the complex polarization vector, if not familiar from electrodynamics, can be developed later in Exercise 22.12.

So much for the foundations. Now for the calculations. First insert the geometric-optics vector potential (22.25) into the Lorentz gauge condition:

Derivation of laws of
geometric optics

$$0 = A^\mu_{;\mu} = \Re \left\{ \left[\frac{i}{\epsilon} k_\mu (a^\mu + \epsilon b^\mu + \dots) + (a^\mu + \epsilon b^\mu + \dots)_{;\mu} \right] e^{i\theta/\epsilon} \right\}. \quad (22.27)$$

The leading term (order $1/\epsilon$) says

$$\mathbf{k} \cdot \mathbf{a} = 0 \quad (\text{amplitude is perpendicular to wave vector}); \quad (22.28)$$

or, equivalently

$$\mathbf{k} \cdot \mathbf{f} = 0 \quad (\text{polarization is perpendicular to wave vector}). \quad (22.28')$$

The post-geometric-optics breakdown in this orthogonality condition is governed by the higher-order terms $[0(1), 0(\epsilon), 0(\epsilon^2), \dots]$ in the gauge condition (22.27); for example, the $0(1)$ terms say

$$\mathbf{k} \cdot \mathbf{b} = i \nabla \cdot \mathbf{a}.$$

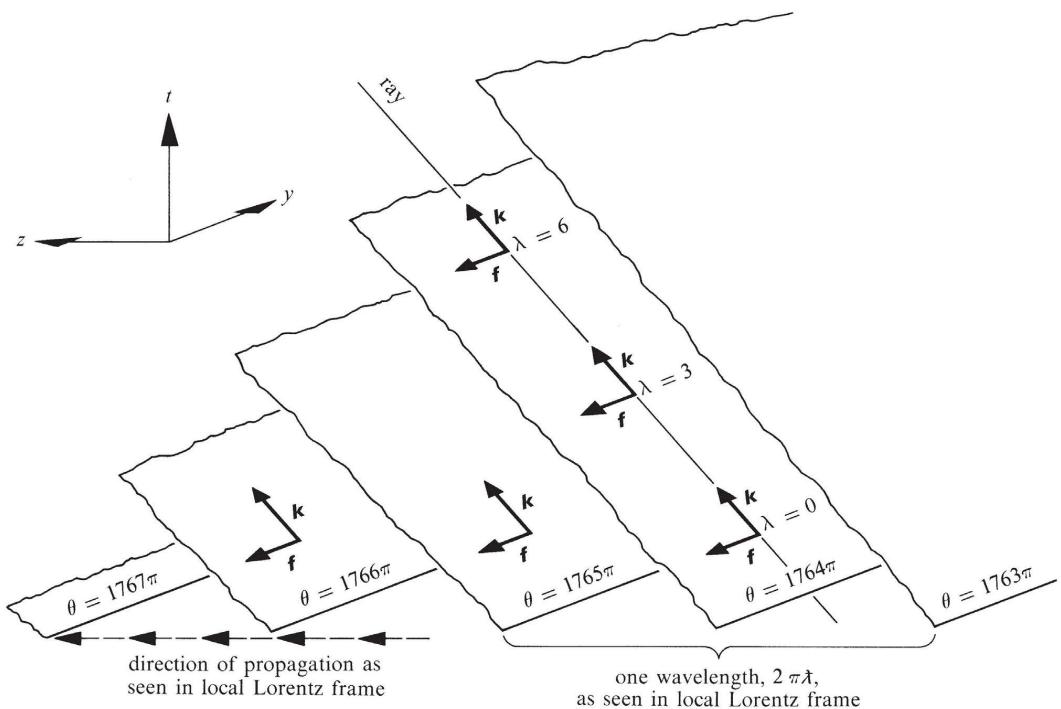
Next insert the vector potential (22.25) into the source-free wave equation (22.19d):

$$\begin{aligned} 0 &= (\mathcal{A}_{dR}\mathbf{A})^\alpha = -A^\alpha_{;\beta} + R^\alpha_\beta A^\beta \\ &= \Re \left\{ \left[\frac{1}{\epsilon^2} k^\beta k_\beta (a^\alpha + \epsilon b^\alpha + \epsilon^2 c^\alpha + \dots) - 2 \frac{i}{\epsilon} k^\beta (a^\alpha + \epsilon b^\alpha + \dots)_{;\beta} \right. \right. \\ &\quad \left. \left. - \frac{i}{\epsilon} k^\beta_{;\beta} (a^\alpha + \epsilon b^\alpha + \dots) - (a^\alpha + \dots)_{;\beta}^\beta + R^\alpha_\beta (a^\beta + \dots) \right] e^{i\theta/\epsilon} \right\}. \quad (22.29) \end{aligned}$$

Collect terms of order $1/\epsilon^2$ and $1/\epsilon$ (terms of order higher than $1/\epsilon$ govern post-geometric-optics corrections):

(continued on page 576)

Box 22.3 GEOMETRY OF AN ELECTROMAGNETIC WAVE TRAIN



The drawing shows surfaces of constant phase, $\theta = \text{constant}$, emerging through the “surface of simultaneity”, $t = 0$, of a local Lorentz frame. The surfaces shown are alternately “crests” ($\theta = 1764\pi, \theta = 1766\pi, \dots$) and “troughs” ($\theta = 1765\pi, \theta = 1767\pi, \dots$) of the wave train. These surfaces make up a 1-form, $\tilde{\mathbf{k}} = d\theta$. The “corresponding vector” $\mathbf{k} = \nabla\theta$ is the “wave vector.” The wave vector is null, $\mathbf{k} \cdot \mathbf{k} = 0$, according to Maxwell’s equations [equation (22.30)]. Therefore it lies in a surface of constant phase:

$$\left(\begin{array}{l} \text{number of surfaces} \\ \text{pierced by } \mathbf{k} \end{array} \right) = \langle d\theta, \mathbf{k} \rangle = \langle \tilde{\mathbf{k}}, \mathbf{k} \rangle = \mathbf{k} \cdot \mathbf{k} = 0.$$

But not only does it lie in a surface of constant phase; it is also perpendicular to that surface! Any vector \mathbf{v} in that surface must satisfy $\mathbf{k} \cdot \mathbf{v} = \langle \tilde{\mathbf{k}}, \mathbf{v} \rangle = \langle d\theta, \mathbf{v} \rangle = 0$ because it pierces no surfaces.

Geometric optics assumes that the reduced wavelength λ , as measured in a typical local Lorentz frame, is small compared to the scale \mathcal{L} of inhomogeneities in the wave train and small compared to the radius of curvature of spacetime, \mathcal{R} . Thus, over regions much larger than λ but smaller than \mathcal{L} or \mathcal{R} , the waves are plane-fronted

and monochromatic, and there exist Lorentz reference frames (Riemann normal coordinates). In one of these “extended” local Lorentz frames, the phase must be

$$\theta = k_\alpha x^\alpha + \text{constant};$$

no other expression will yield $\nabla\theta = \mathbf{k}$. The corresponding vector potential [equation (22.25)] will be

$$A^\mu = \Re\{a^\mu \exp[i(\mathbf{k} \cdot \mathbf{x} - k^0 t)]\} + (\text{“post-geometric-optics corrections”});$$

hence,

$$k^0 = 2\pi/(\text{period of wave}) = 2\pi\nu = \omega \equiv (\text{angular frequency}),$$

$$|\mathbf{k}| = 2\pi/(\text{wavelength of wave}) = 1/\lambda = \omega,$$

\mathbf{k} points along direction of propagation of wave.

At each event in spacetime there is a wave vector; and these wave vectors, tacked end-on-end, form a family of curves—the “light rays” or simply “rays”—whose tangent vector is \mathbf{k} . The rays, like their tangent vector, lie both in and perpendicular to the surfaces of constant phase.

The affine parameter λ of a ray (not to be confused with wavelength $= 2\pi\lambda$) satisfies $\mathbf{k} = d/d\lambda$; therefore it is given by

$$\lambda = t/k^0 + \text{constant} = t/\omega + \text{constant},$$

where t is proper time along the ray as measured, not by the ray itself (its proper time is zero!), but by the local Lorentz observer who sees angular frequency ω . Thus, while ω is a frame-dependent quantity and t is also a frame-dependent quantity, their quotient t/ω when measured along the ray (*not off the ray*) is the frame-independent affine parameter. For a particle it is possible and natural to identify the affine parameter λ with proper time τ . For a light ray this identification is unnatural and impossible. The lapse of proper time along the ray is identically zero. The springing up of λ to take the place of the vanished τ gives one a tool to do what one might not have suspected to be possible. Given a light ray shot out at event \mathcal{A} and passing through event \mathcal{B} , one can give a third event \mathcal{C} along the same null world line that is twice as “far” from \mathcal{A} as \mathcal{B} is “far,” in a new sense of “far” that has nothing whatever directly to do with proper time (zero!), but is defined by equal increments of the affine parameter ($\lambda_{\mathcal{C}} - \lambda_{\mathcal{B}} = \lambda_{\mathcal{B}} - \lambda_{\mathcal{A}}$). The “affine parameter” has a meaning for any null geodesic analyzed even in isolation. In this respect, it is to be distinguished from the so-called “luminosity distance” which is sometimes introduced in dealing with the propagation of radiation through curved spacetime, and which is defined by the spreading apart of two or more light rays coming from a common source.

Maxwell’s equations as explored in the text [equation (22.28’)] guarantee that the complex polarization vector \mathbf{f} is perpendicular to the wave vector \mathbf{k} and that, therefore, it lies in a surface of constant phase (see drawing). Intuition into the polarization vector is developed in exercise 22.12.

$$0\left(\frac{1}{\varepsilon^2}\right): \quad k^\beta k_\beta a^\alpha = 0 \\ \implies \mathbf{k} \cdot \mathbf{k} = 0 \text{ (wave vector is null);} \quad (22.30)$$

$$0\left(\frac{1}{\varepsilon}\right): \quad \underbrace{k^\beta k_\beta b^\alpha - 2i\left(k^\beta a^\alpha_{;\beta} + \frac{1}{2}k^\beta_{;\beta} a^\alpha\right)}_{[=0]} = 0 \\ \implies \nabla_{\mathbf{k}} \mathbf{a} = -\frac{1}{2}(\nabla \cdot \mathbf{k}) \mathbf{a} \text{ (propagation equation for vector amplitude).} \quad (22.31)$$

These equations (22.30, 22.31) together with equation (22.28) are the basis from which all subsequent results will follow. As a first consequence, one can obtain the geodesic law from equation (22.30). Form the gradient of $\mathbf{k} \cdot \mathbf{k} = 0$,

$$0 = (k^\beta k_\beta)_{;\alpha} = 2k^\beta k_{\beta;\alpha},$$

and use the fact that $k_\beta \equiv \theta_{,\beta}$ is the gradient of a scalar to interchange indices, $\theta_{;\beta\alpha} = \theta_{;\alpha\beta}$ or

$$0 = k^\beta k_{\beta;\alpha} = k^\beta k_{\alpha;\beta}.$$

The main laws of geometric optics:

(1) Light rays are null geodesics

The result is

$$\nabla_{\mathbf{k}} \mathbf{k} = 0 \text{ (propagation equation for wave vector).} \quad (22.32)$$

Notice that this is the geodesic equation! Combined with equation (22.30), it is the statement, derived from Maxwell's equations in curved spacetime, that *light rays are null geodesics*, the first main result of geometric optics.

Turn now from the propagation vector $\mathbf{k} = \nabla\theta$ to the wave amplitude $\mathbf{a} = a\mathbf{f}$, and obtain separate equations for the magnitude a and polarization \mathbf{f} . Use equation (22.31) to compute

$$2a \partial_{\mathbf{k}} a = 2a \nabla_{\mathbf{k}} a = \nabla_{\mathbf{k}} a^2 = \nabla_{\mathbf{k}} (\mathbf{a} \cdot \bar{\mathbf{a}}) = \bar{\mathbf{a}} \cdot \nabla_{\mathbf{k}} \mathbf{a} + \mathbf{a} \cdot \nabla_{\mathbf{k}} \bar{\mathbf{a}} \\ = -\frac{1}{2}(\nabla \cdot \mathbf{k})(\bar{\mathbf{a}} \cdot \mathbf{a} + \mathbf{a} \cdot \bar{\mathbf{a}}) = -a^2 \nabla \cdot \mathbf{k};$$

so

$$\partial_{\mathbf{k}} a = -\frac{1}{2}(\nabla \cdot \mathbf{k})a \text{ (propagation equation for scalar amplitude).} \quad (22.33)$$

Next write $\mathbf{a} = a\mathbf{f}$ in equation (22.31) to obtain

$$0 = \nabla_{\mathbf{k}}(a\mathbf{f}) + \frac{1}{2}(\nabla \cdot \mathbf{k})a\mathbf{f} = a\nabla_{\mathbf{k}}\mathbf{f} + \mathbf{f}\left[\nabla_{\mathbf{k}}a + \frac{1}{2}(\nabla \cdot \mathbf{k})a\right] = a\nabla_{\mathbf{k}}\mathbf{f}$$

or

$$\nabla_{\mathbf{k}} \mathbf{f} = 0 \text{ (propagation equation for polarization vector).} \quad (22.34)$$

This together with equation (22.28'), constitutes the second main result of geometric optics, that *the polarization vector is perpendicular to the rays and is parallel-propagated along the rays*. It is now possible to see that these results, derived from equations (22.30) and (22.31) are consistent with the gauge condition (22.28). The vectors \mathbf{k} and \mathbf{f} , specified at one point, are fixed along the entire ray by their propagation equations. But because both propagation equations are parallel-transport laws, the conditions $\mathbf{k} \cdot \mathbf{k} = 0$, $\mathbf{f} \cdot \bar{\mathbf{f}} = 1$, and $\mathbf{k} \cdot \mathbf{f} = 0$, once imposed on the vectors at one point, will be satisfied along the entire ray.

The equation (22.33) for the scalar amplitude can be reformulated as a conservation law. Since $\partial_{\mathbf{k}} \equiv (\mathbf{k} \cdot \nabla)$, one rewrites the equation as $(\mathbf{k} \cdot \nabla)a^2 + a^2 \nabla \cdot \mathbf{k} = 0$, or

$$\nabla \cdot (a^2 \mathbf{k}) = 0. \quad (22.35)$$

Consequently the vector $a^2 \mathbf{k}$ is a “conserved current,” and the integral $\int a^2 k^\mu d^3 \Sigma_\mu$ has a fixed, unchanging value for each 3-volume cutting a given tube formed of light rays. (The tube must be so formed of rays that an integral of $a^2 \mathbf{k}$ over the walls of the tube will give zero.) What is conserved? To remain purely classical, one could say it is the “number of light rays” and call $a^2 k^0$ the “density of light rays” on an $x^0 = \text{constant}$ hypersurface. But the proper correspondence and more concrete physical interpretation make one prefer to call equation (22.35) *the law of conservation of photon number*. It is the third main result of geometric optics. Photon number, of course, is not always conserved; it is an adiabatic invariant, a quantity that is not changed by influences (e.g., spacetime curvature, $\sim 1/\mathcal{R}^2$) which change slowly ($\mathcal{R} \gg \lambda$) compared to the photon frequency.

Box 22.4 summarizes the above equations of geometric optics, along with others derived in the exercises.

(2) polarization vector is perpendicular to ray and is parallel propagated along ray

(3) conservation of “photon number”

Exercise 22.11. ELECTROMAGNETIC FIELD AND STRESS ENERGY

Derive the equations given in part D of Box 22.4 for \mathbf{F} , \mathbf{E} , \mathbf{B} , and \mathbf{T} .

EXERCISES

Exercise 22.12. POLARIZATION

At an event \mathcal{P}_0 through which geometric-optics waves are passing, introduce a local Lorentz frame with z -axis along the direction of propagation. Then $\mathbf{k} = \omega(\mathbf{e}_0 + \mathbf{e}_z)$. Since the polarization vector is orthogonal to \mathbf{k} , it is $\mathbf{f} = f^0(\mathbf{e}_0 + \mathbf{e}_z) + f^1 \mathbf{e}_x + f^2 \mathbf{e}_y$; and since $\mathbf{f} \cdot \bar{\mathbf{f}} = 1$, it has $|f^1|^2 + |f^2|^2 = 1$.

(a) Show that the component f^0 of the polarization vector has no influence on the electric and magnetic fields measured in the given frame; i.e., show that one can add a multiple of \mathbf{k} to \mathbf{f} without affecting any physical measurements.

(continued on page 581)

Box 22.4 GEOMETRIC OPTICS IN CURVED SPACETIME
 (Summary of Results Derived in Text and Exercises)

A. Geometric Optics Assumption

Electromagnetic waves propagating in a source-free region of spacetime are locally plane-fronted and monochromatic (reduced wavelength $\lambda \ll$ scale \mathcal{L} over which amplitude, wavelength, or polarization vary; and $\lambda \ll \mathcal{R}$ = mean radius of curvature of spacetime).

B. Rays, Phase, and Wave Vector (see Box 22.3)

Everything (amplitude, polarization, energy, etc.) is transported along *rays*; and the quantities on one ray do not influence the quantities on any other ray. The rays are null geodesics of curved spacetime, with tangent vectors (“wave vectors”) \mathbf{k} :

$$\nabla_{\mathbf{k}} \mathbf{k} = 0.$$

The rays both lie in and are perpendicular to surfaces of constant phase, $\theta = \text{const.}$; and their tangent vectors are the gradient of θ :

$$\mathbf{k} = \nabla\theta.$$

In a local Lorentz frame, k^0 is the “angular frequency” and $k^0/2\pi$ is the ordinary frequency of the waves, and

$$\mathbf{n} = \mathbf{k}/k^0$$

is a unit 3-vector pointing along their direction of propagation.

C. Amplitude and Polarization Vector

The waves are characterized by a real amplitude a and a complex polarization vector \mathbf{f} of unit length, $\mathbf{f} \cdot \bar{\mathbf{f}} = 1$. (Of the fundamental quantities θ , \mathbf{k} , a , \mathbf{f} , all are real except \mathbf{f} . See exercise 22.12 for deeper understanding of \mathbf{f} .)

The polarization vector is everywhere orthogonal to the rays, $\mathbf{k} \cdot \mathbf{f} = 0$; and is parallel-transported along them, $\nabla_{\mathbf{k}} \mathbf{f} = 0$.

The propagation law for the amplitude is

$$\partial_{\mathbf{k}} a = -\frac{1}{2}(\nabla \cdot \mathbf{k})a.$$

This propagation law is equivalent to a *law of conservation of photons* (classically: of rays); $a^2\mathbf{k}$ is the “conserved current” satisfying $\nabla \cdot (a^2\mathbf{k}) = 0$; and $(8\pi\hbar)^{-1} \int a^2 k^0 \sqrt{|g|} d^3x$ is the number of photons (rays) in the 3-volume of integration on any $x^0 = \text{constant}$ hypersurface, and is constant as this volume is carried along the rays.

The propagation law holds separately on each hypersurface of constant phase. There it can be interpreted as conservation of a $a^2\mathcal{A}$, where \mathcal{A} is a two-dimensional cross-sectional area of a pulse of photons or rays. See exercise 22.13.

D. Vector Potential, Electromagnetic Field, and Stress-Energy-Momentum

At any event the vector potential in Lorentz gauge is

$$\mathbf{A} = \Re\{ae^{i\theta}\mathbf{f}\},$$

where \Re denotes the real part.

The electromagnetic field tensor is orthogonal to the rays, $\mathbf{F} \cdot \mathbf{k} = 0$, and is given by

$$\mathbf{F} = \Re\{iae^{i\theta}\mathbf{k} \wedge \mathbf{f}\}.$$

The corresponding electric and magnetic fields in any local Lorentz frame are

$$\mathbf{E} = \Re\{iak^0e^{i\theta}(\text{projection of } \mathbf{f} \text{ perpendicular to } \mathbf{k})\},$$

$$\mathbf{B} = \mathbf{n} \times \mathbf{E}, \text{ where } \mathbf{n} \equiv \mathbf{k}/k^0.$$

The stress-energy tensor, averaged over a wavelength, is

$$\mathbf{T} = (1/8\pi)a^2\mathbf{k} \otimes \mathbf{k},$$

corresponding to an energy density in a local Lorentz frame of

$$T^{00} = (1/8\pi)(ak^0)^2$$

and an energy flux of

$$T^{0j} = T^{00}n^j,$$

so that energy flows along the rays (in $\mathbf{n} = \mathbf{k}/k^0$ direction) with the speed of light. This is identical with the stress-energy tensor that would be produced by a beam of photons with 4-momenta $\mathbf{p} = \hbar\mathbf{k}$.

Conservation of energy-momentum $\nabla \cdot \mathbf{T} = 0$ follows from the ray conservation law $\nabla \cdot (a^2\mathbf{k}) = 0$ and the geodesic law $\nabla_{\mathbf{k}}\mathbf{k} \equiv (\mathbf{k} \cdot \nabla)\mathbf{k} = 0$:

$$8\pi \nabla \cdot \mathbf{T} = \nabla \cdot (a^2\mathbf{k} \otimes \mathbf{k}) = [\nabla \cdot (a^2\mathbf{k})]\mathbf{k} + a^2(\mathbf{k} \cdot \nabla)\mathbf{k} = 0.$$

Box 22.4 (continued)

The adiabatic (geometric optics) invariant “ray number” $a^2 k^0$ or “photon number” $(8\pi\hbar)^{-1} a^2 k^0$ in a unit volume is proportional to the energy, $(8\pi)^{-1} a^2 (k^0)^2$, divided by the frequency, k^0 —corresponding exactly to the harmonic oscillator adiabatic invariant E/ω [Einstein (1912), Ehrenfest (1916), Landau and Lifshitz (1960)].

E. Photon Reinterpretation of Geometric Optics

The laws of geometric optics can be reinterpreted as follows. This reinterpretation becomes a foundation of the standard quantum theory of the electromagnetic field (see, e.g., Chapters 1 and 13 of Baym (1969)); and the classical limit of that quantum theory is standard Maxwell electrodynamics.

Photons are particles of zero rest mass that move along null geodesics of spacetime (the null rays).

The 4-momentum of a photon is related to the tangent vector of the null ray (wave vector) by $\mathbf{p} = \hbar \mathbf{k}$. A renormalization of the affine parameter,

$$(\text{new parameter}) = (1/\hbar) \times (\text{old parameter}),$$

makes \mathbf{p} the tangent vector to the ray.

Each photon possesses a polarization vector, \mathbf{f} , which is orthogonal to its 4-momentum ($\mathbf{p} \cdot \mathbf{f} = 0$), and which it parallel-transports along its geodesic world line ($\nabla_{\mathbf{p}} \mathbf{f} = 0$).

A swarm of photons, all with nearly the same 4-momentum \mathbf{p} and polarization vector \mathbf{f} (as compared by parallel transport), make up a classical electromagnetic wave. The scalar amplitude a of the wave is determined by equating the stress-energy tensor of the wave

$$\mathbf{T} = \frac{1}{8\pi} a^2 \mathbf{k} \otimes \mathbf{k} = \frac{1}{8\pi} \left(\frac{a}{\hbar} \right)^2 \mathbf{p} \otimes \mathbf{p}$$

to the stress-energy tensor of a swarm of photons with number-flux vector \mathbf{S} ,

$$\mathbf{T} = \mathbf{p} \otimes \mathbf{S}$$

[see equation (5.18)]. The result:

$$\mathbf{S} = \frac{1}{8\pi} \left(\frac{a}{\hbar} \right)^2 \mathbf{p} = \frac{1}{8\pi\hbar} a^2 \mathbf{k}$$

or, in any local Lorentz frame,

$$a = (8\pi\hbar^2 S^0/p^0)^{1/2} = (8\pi)^{1/2} \hbar \left(\frac{\text{number density of photons}}{\text{energy of one photon}} \right)^{1/2}.$$

(b) Show that the following polarization vectors correspond to the types of polarization listed:

$$\mathbf{f} = \mathbf{e}_x, \text{ linear polarization in } x \text{ direction;}$$

$$\mathbf{f} = \mathbf{e}_y, \text{ linear polarization in } y \text{ direction;}$$

$$\mathbf{f} = \frac{1}{\sqrt{2}}(\mathbf{e}_x + i\mathbf{e}_y), \text{ righthand circular polarization;}$$

$$\mathbf{f} = \frac{1}{\sqrt{2}}(\mathbf{e}_x - i\mathbf{e}_y), \text{ lefthand circular polarization;}$$

$$\mathbf{f} = \alpha\mathbf{e}_x + i(1 - \alpha^2)^{1/2}\mathbf{e}_y, \text{ righthand elliptical polarization.}$$

(c) Show that the type of polarization (linear; circular; elliptical with given eccentricity of ellipse) is the same as viewed in any local Lorentz frame at any event along a given ray. [Hint: Use pictures and abstract calculations rather than Lorentz transformations and component calculations.]

Exercise 22.13. THE AREA OF A BUNDLE OF RAYS

Write equation (22.31) in a coordinate system in which one of the coordinates is chosen to be $x^0 = \theta$, the phase (a retarded time coordinate).

(a) Show that $g^{00} = 0$ and that no derivatives $\partial/\partial\theta$ appear in equation (22.33); so propagation of a can be described within a single $\theta = \text{constant}$ hypersurface.

(b) Perform the following construction (see Figure 22.1). Pick a ray \mathcal{C}_0 along which a is to be propagated. Pick a bundle of rays, with two-dimensional cross section, that (i) all lie in the same constant-phase surface as \mathcal{C}_0 , and (ii) surround \mathcal{C}_0 . (The surface is three-di-

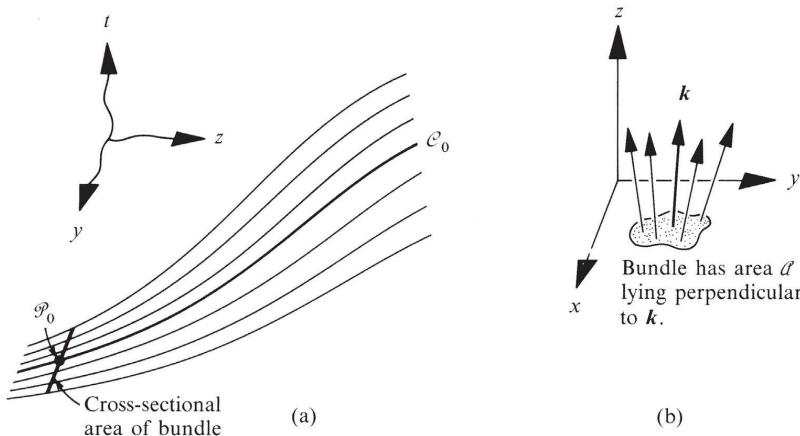


Figure 22.1.

Geometric optics for a bundle of rays with two-dimensional cross section, all lying in a surface of constant phase, $\theta = \text{const}$. Sketch (a) shows the bundle, surrounding a central ray \mathcal{C}_0 , in a spacetime diagram with one spatial dimension suppressed. Sketch (b) shows the bundle as viewed on a slice of simultaneity in a local Lorentz frame at the event \mathcal{P}_0 . Slicing the bundle turns each ray into a “photon”; so the bundle becomes a two-dimensional surface filled with photons. The area \mathcal{A} of this photon-filled surface obeys the following laws (see exercises 22.13 and 22.14): (1) \mathcal{A} is independent of the choice of Lorentz frame; it depends only on location \mathcal{P}_0 along the ray \mathcal{C}_0 . (2) The amplitude a of the waves satisfies

$$\mathcal{A}a^2 = \text{constant all along the ray } \mathcal{C}_0$$

(“conservation of photon flux”). (3) \mathcal{A} obeys the “propagation equation” (22.36).