# 1 Special relativity

## 1.1 Fundamental principles of special relativity (SR) theory

The way in which special relativity is taught at an elementary undergraduate level – the level at which the reader is assumed competent – is usually close in spirit to the way it was first understood by physicists. This is an algebraic approach, based on the Lorentz transformation (§ 1.7 below). At this basic level, we learn how to use the Lorentz transformation to convert between one observer's measurements and another's, to verify and understand such remarkable phenomena as time dilation and Lorentz contraction, and to make elementary calculations of the conversion of mass into energy.

This purely algebraic point of view began to change, to widen, less than four years after Einstein proposed the theory.[1] Minkowski pointed out that it is very helpful to regard $(t, x, y, z)$ as simply four coordinates in a four-dimensional space which we now call space-time. This was the beginning of the geometrical point of view, which led directly to general relativity in 1914–16. It is this geometrical point of view on special relativity which we must study before all else.

As we shall see, special relativity can be deduced from two fundamental postulates:

(1) *Principle of relativity* (Galileo): No experiment can measure the absolute velocity of an observer; the results of any experiment performed by an observer do not depend on his speed relative to other observers who are not involved in the experiment.
(2) *Universality of the speed of light* (Einstein): The speed of light relative to any unaccelerated observer is $c = 3 \times 10^8 \, \mathrm{m\,s^{-1}}$, regardless of the motion of the light's source relative to the observer. Let us be quite clear about this postulate's meaning: two different unaccelerated observers measuring the speed of the *same photon* will each find it to be moving at $3 \times 10^8 \, \mathrm{m\,s^{-1}}$ relative to themselves, regardless of their state of motion relative to each other.

As noted above, the principle of relativity is not at all a modern concept; it goes back all the way to Galileo's hypothesis that a body in a state of uniform motion remains in that state unless acted upon by some external agency. It is fully embodied in Newton's second

---

[1] Einstein's original paper was published in 1905, while Minkowski's discussion of the geometry of spacetime was given in 1908. Einstein's and Minkowski's papers are reprinted (in English translation) in *The Principle of Relativity* by A. Einstein, H. A. Lorentz, H. Minkowski, and H. Weyl (Dover).

law, which contains only accelerations, not velocities themselves. Newton's laws are, in fact, all invariant under the replacement

$$\boldsymbol{v}(t) \to \boldsymbol{v}'(t) = \boldsymbol{v}(t) - \boldsymbol{V},$$

where $\boldsymbol{V}$ is any *constant* velocity. This equation says that a velocity $\boldsymbol{v}(t)$ relative to one observer becomes $\boldsymbol{v}'(t)$ when measured by a second observer whose velocity relative to the first is $\boldsymbol{V}$. This is called the Galilean law of addition of velocities.

By saying that Newton's laws are *invariant* under the Galilean law of addition of velocities, we are making a statement of a sort we will often make in our study of relativity, so it is well to start by making it very precise. Newton's first law, that a body moves at a constant velocity in the absence of external forces, is unaffected by the replacement above, since if $\boldsymbol{v}(t)$ is really a constant, say $\boldsymbol{v}_0$, then the new velocity $\boldsymbol{v}_0 - \boldsymbol{V}$ is also a constant. Newton's second law

$$\boldsymbol{F} = m\boldsymbol{a} = m\,\mathrm{d}\boldsymbol{v}/\mathrm{d}\,t,$$

is also unaffected, since

$$\boldsymbol{a}' = \mathrm{d}\boldsymbol{v}'/\mathrm{d}\,t = \mathrm{d}(\boldsymbol{v} - \boldsymbol{V})/\mathrm{d}\,t = \mathrm{d}\boldsymbol{v}/\mathrm{d}\,t = \boldsymbol{a}.$$

Therefore, the second law will be valid according to the measurements of both observers, provided that we add to the Galilean transformation law the statement that $\boldsymbol{F}$ and $m$ are themselves invariant, i.e. the same regardless of which of the two observers measures them. Newton's third law, that the force exerted by one body on another is equal and opposite to that exerted by the second on the first, is clearly unaffected by the change of observers, again because we assume the forces to be invariant.

So there is no absolute velocity. Is there an absolute acceleration? Newton argued that there was. Suppose, for example, that I am in a train on a perfectly smooth track,[2] eating a bowl of soup in the dining car. Then, if the train moves at constant speed, the soup remains level, thereby offering me no information about what my speed is. But, if the train changes its speed, then the soup climbs up one side of the bowl, and I can tell by looking at it how large and in what direction the acceleration is.[3]

Therefore, it is reasonable and useful to single out a class of preferred observers: those who are unaccelerated. They are called *inertial observers*, and each one has a constant velocity with respect to any other one. These inertial observers are fundamental in special relativity, and when we use the term 'observer' from now on we will mean an inertial observer.

The postulate of the universality of the speed of light was Einstein's great and radical contribution to relativity. It smashes the Galilean law of addition of velocities because it says that if $\boldsymbol{v}$ has magnitude $c$, then so does $\boldsymbol{v}'$, regardless of $\boldsymbol{V}$. The earliest direct evidence for this postulate was the Michelson–Morely experiment, although it is not clear whether Einstein himself was influenced by it. The counter-intuitive predictions of special relativity all flow from this postulate, and they are amply confirmed by experiment. In fact it is probably fair to say that special relativity has a firmer experimental basis than any other of

---

[2] Physicists frequently have to make such idealizations, which often are far removed from common experience !
[3] For Newton's discussion of this point, see the excerpt from his *Principia* in Williams (1968).

our laws of physics, since it is tested every day in all the giant particle accelerators, which send particles nearly to the speed of light.

Although the concept of relativity is old, it is customary to refer to Einstein's theory simply as 'relativity'. The adjective 'special' is applied in order to distinguish it from Einstein's theory of gravitation, which acquired the name 'general relativity' because it permits us to describe physics from the point of view of both accelerated and inertial observers and is in that respect a more general form of relativity. But the real physical distinction between these two theories is that special relativity (SR) is capable of describing physics only in the absence of gravitational fields, while general relativity (GR) extends SR to describe gravitation itself.[4] We can only wish that an earlier generation of physicists had chosen more appropriate names for these theories !

## 1.2 Definition of an inertial observer in SR

It is important to realize that an 'observer' is in fact a huge information-gathering system, not simply one man with binoculars. In fact, we shall remove the human element entirely from our definition, and say that an inertial observer is simply a coordinate system for spacetime, which makes an observation simply by recording the location $(x, y, z)$ and time $(t)$ of any event. This coordinate system must satisfy the following three properties to be called *inertial*:

(1) The distance between point $P_1$ (coordinates $x_1, y_1, z_1$) and point $P_2$ (coordinates $x_2, y_2, z_2$) is independent of time.
(2) The clocks that sit at every point ticking off the time coordinate $t$ are synchronized and all run at the same rate.
(3) The geometry of space at any constant time $t$ is Euclidean.

Notice that this definition does not mention whether the observer accelerates or not. That will come later. It will turn out that only an unaccelerated observer can keep his clocks synchronized. But we prefer to start out with this geometrical definition of an inertial observer. It is a matter for experiment to decide whether such an observer can exist: it is not self-evident that any of these properties *must* be realizable, although we would probably expect a 'nice' universe to permit them! However, we will see later in the course that a gravitational field does generally make it impossible to construct such a coordinate system, and this is why GR is required. But let us not get ahead of the story. At the moment we are assuming that we *can* construct such a coordinate system (that, if you like, the gravitational fields around us are so weak that they do not really matter). We can envision this coordinate system, rather fancifully, as a lattice of rigid rods filling space, with a clock at every intersection of the rods. Some convenient system, such as a collection of GPS

---

[4] It is easy to see that gravitational fields cause problems for SR. If an astronaut in orbit about Earth holds a bowl of soup, does the soup climb up the side of the bowl in response to the gravitational 'force' that holds the spacecraft in orbit? Two astronauts in different orbits accelerate relative to one another, but neither *feels* an acceleration. Problems like this make gravity special, and we will have to wait until Ch. 5 to resolve them. Until then, the word 'force' will refer to a nongravitational force.

satellites and receivers, is used to ensure that all the clocks are synchronized. The clocks are supposed to be very densely spaced, so that there is a clock next to every event of interest, ready to record its time of occurrence without any delay. We shall now define how we use this coordinate system to make observations.

An *observation* made by the inertial observer is the act of assigning to any event the coordinates $x$, $y$, $z$ of the location of its occurrence, and the time read by the clock at $(x, y, z)$ when the event occurred. It is *not* the time $t$ on the wrist watch worn by a scientist located at $(0, 0, 0)$ when he first learns of the event. A *visual* observation is of this second type: the eye regards as simultaneous all events it *sees* at the same time; an inertial observer regards as simultaneous all events that *occur* at the same time as recorded by the clock nearest them when the events occurred. This distinction is important and must be borne in mind. Sometimes we will say 'an observer sees . . .' but this will only be shorthand for 'measures'. We will never mean a *visual* observation unless we say so explicitly.

An inertial observer is also called an *inertial reference frame*, which we will often abbreviate to 'reference frame' or simply 'frame'.

## 1.3 New units

Since the speed of light $c$ is so fundamental, we shall from now on adopt a new system of units for measurements in which $c$ simply has the value 1! It is perfectly okay for slow-moving creatures like engineers to be content with the SI units: m, s, kg. But it seems silly in SR to use units in which the fundamental constant $c$ has the ridiculous value $3 \times 10^8$. The SI units evolved historically. Meters and seconds are not fundamental; they are simply convenient for human use. What we shall now do is adopt a new unit for time, the meter. One meter of time is the time it takes light to travel one meter. (You are probably more familiar with an alternative approach: a year of distance – called a 'light year' – is the distance light travels in one year.) The speed of light in these units is:

$$
\begin{aligned}
c &= \frac{\text{distance light travels in any given time interval}}{\text{the given time interval}} \\
&= \frac{1 \text{ m}}{\text{the time it takes light to travel one meter}} \\
&= \frac{1 \text{ m}}{1 \text{ m}} = 1.
\end{aligned}
$$

So if we consistently measure time in meters, then $c$ is not merely 1, it is also dimensionless! In converting from SI units to these 'natural' units, we can use any of the following relations:

$$
\begin{aligned}
3 \times 10^8 \text{ m s}^{-1} &= 1, \\
1 \text{ s} &= 3 \times 10^8 \text{ m}, \\
1 \text{ m} &= \frac{1}{3 \times 10^8} \text{ s}.
\end{aligned}
$$

The SI units contain many 'derived' units, such as joules and newtons, which are defined in terms of the basic three: m, s, kg. By converting from s to m these units simplify considerably: energy and momentum are measured in kg, acceleration in $m^{-1}$, force in $kg\,m^{-1}$, etc. Do the exercises on this. With practice, these units will seem as natural to you as they do to most modern theoretical physicists.

## 1.4 Spacetime diagrams

A very important part of learning the geometrical approach to SR is mastering the spacetime diagram. In the rest of this chapter we will derive SR from its postulates by using spacetime diagrams, because they provide a very powerful guide for threading our way among the many pitfalls SR presents to the beginner. Fig. 1.1 below shows a two-dimensional slice of spacetime, the $t - x$ plane, in which are illustrated the basic concepts. A single point in this space[5] is a point of fixed $x$ *and* fixed $t$, and is called an *event*. A line in the space gives a relation $x = x(t)$, and so can represent the position of a particle at different times. This is called the particle's *world line*. Its slope is related to its velocity,

$$\text{slope} = \mathrm{d}t/\mathrm{d}x = 1/v.$$

Notice that a light ray (photon) always travels on a 45° line in this diagram.
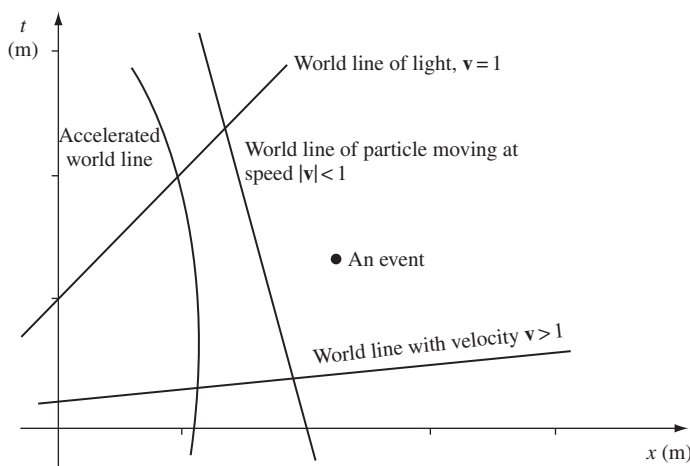


Figure 1.1    A spacetime diagram in natural units.

[5] We use the word 'space' in a more general way than you may be used to. We do not mean a Euclidean space in which Euclidean distances are necessarily physically meaningful. Rather, we mean just that it is a set of points that is continuous (rather than discrete, as a lattice is). This is the first example of what we will define in Ch. 5 to be a 'manifold'.

We shall adopt the following notational conventions:

(1) *Events* will be denoted by cursive capitals, e.g. $\mathcal{A}, \mathcal{B}, \mathcal{P}$. However, the letter $\mathcal{O}$ is reserved to denote observers.
(2) The *coordinates* will be called $(t, x, y, z)$. Any quadruple of numbers like $(5, -3, 2, 10^{16})$ denotes an event whose coordinates are $t = 5$, $x = -3$, $y = 2$, $z = 10^{16}$. Thus, we always put $t$ first. All coordinates are measured in meters.
(3) It is often convenient to refer to the coordinates $(t, x, y, z)$ as a whole, or to each indifferently. That is why we give them the *alternative names* $(x^0, x^1, x^2, x^3)$. These superscripts are *not* exponents, but just labels, called indices. Thus $(x^3)^2$ denotes the square of coordinate 3 (which is $z$), not the square of the cube of $x$. *Generically*, the coordinates $x^0$, $x^1$, $x^2$, and $x^3$ are referred to as $x^\alpha$. A *Greek index* (e.g. $\alpha, \beta, \mu, \nu$) will be assumed to take a value from the set (0, 1, 2, 3). If $\alpha$ is not given a value, then $x^\alpha$ is *any* of the four coordinates.
(4) There are occasions when we want to distinguish between $t$ on the one hand and $(x, y, z)$ on the other. We use *Latin indices* to refer to the spatial coordinates alone. Thus a Latin index (e.g. $a, b, i, j, k, l$) will be assumed to take a value from the set (1, 2, 3). If $i$ is not given a value, then $x^i$ is *any* of the three spatial coordinates. Our conventions on the use of Greek and Latin indices are by no means universally used by physicists. Some books reverse them, using Latin for {0, 1, 2, 3} and Greek for {1, 2, 3}; others use $a, b, c, \ldots$ for one set and $i, j, k$ for the other. Students should always check the conventions used in the work they are reading.

## 1.5 Construction of the coordinates used by another observer

Since any observer is simply a coordinate system for spacetime, and since all observers look at the same events (the same spacetime), it should be possible to draw the coordinate lines of one observer on the spacetime diagram drawn by another observer. To do this we have to make use of the postulates of SR.

Suppose an observer $\mathcal{O}$ uses the coordinates $t, x$ as above, and that another observer $\bar{\mathcal{O}}$, with coordinates $\bar{t}, \bar{x}$, is moving with velocity $v$ in the $x$ direction relative to $\mathcal{O}$. Where do the coordinate axes for $\bar{t}$ and $\bar{x}$ go in the spacetime diagram of $\mathcal{O}$?

$\bar{t}$ *axis*: This is the locus of events at constant $\bar{x} = 0$ (and $\bar{y} = \bar{z} = 0$, too, but we shall ignore them here), which is the locus of the origin of $\bar{\mathcal{O}}$'s spatial coordinates. This is $\bar{\mathcal{O}}$'s world line, and it looks like that shown in Fig. 1.2.

$\bar{x}$ *axis*: To locate this we make a construction designed to determine the locus of events at $\bar{t} = 0$, i.e. those that $\bar{\mathcal{O}}$ measures to be simultaneous with the event $\bar{t} = \bar{x} = 0$.

Consider the picture in $\bar{\mathcal{O}}$'s spacetime diagram, shown in Fig. 1.3. The events on the $\bar{x}$ axis all have the following property: A light ray emitted at event $\mathcal{E}$ from $\bar{x} = 0$ at, say, time $\bar{t} = -a$ will reach the $\bar{x}$ axis at $\bar{x} = a$ (we call this event $\mathcal{P}$); if reflected, it will return to the point $\bar{x} = 0$ at $\bar{t} = +a$, called event $\mathcal{R}$. The $\bar{x}$ axis can be *defined*, therefore, as the locus of
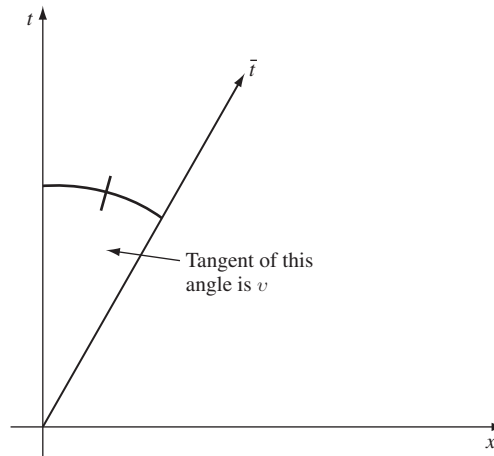
Figure 1.2

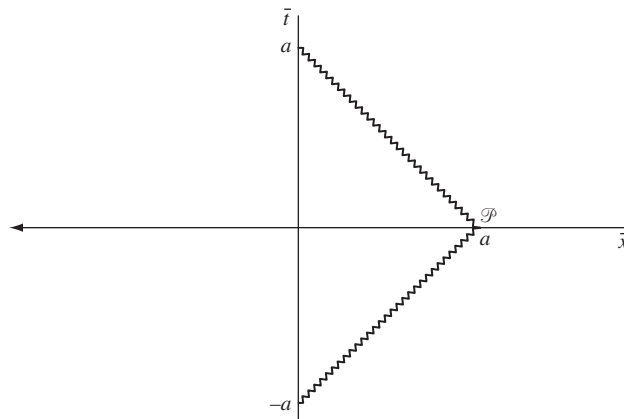The time-axis of a frame whose velocity is $v$.



Figure 1.3

Light reflected at $a$, as measured by $\bar{\mathcal{O}}$.

events that reflect light rays in such a manner that they return to the $\bar{t}$ axis at $+a$ if they left it at $-a$, for any $a$. Now look at this in the spacetime diagram of $\mathcal{O}$, Fig. 1.4.

We know where the $\bar{t}$ axis lies, since we constructed it in Fig. 1.2. The events of emission and reception, $\bar{t} = -a$ and $\bar{t} = +a$, are shown in Fig. 1.4. Since $a$ is arbitrary, it does not matter where along the negative $\bar{t}$ axis we place event $\mathcal{E}$, so no assumption need yet be made about the calibration of the $\bar{t}$ axis relative to the $t$ axis. All that matters for the moment is that the event $\mathcal{R}$ on the $\bar{t}$ axis must be as far from the origin as event $\mathcal{E}$. Having drawn them in Fig. 1.4, we next draw in the same light beam as before, emitted from $\mathcal{E}$, and traveling on a 45° line in this diagram. The reflected light beam must arrive at $\mathcal{R}$, so it is the 45° line with negative slope through $\mathcal{R}$. The intersection of these two light beams must be the event of reflection $\mathcal{P}$. This establishes the location of $\mathcal{P}$ in our diagram. The line joining it with the origin – the dashed line – must be the $\bar{x}$ axis: it does
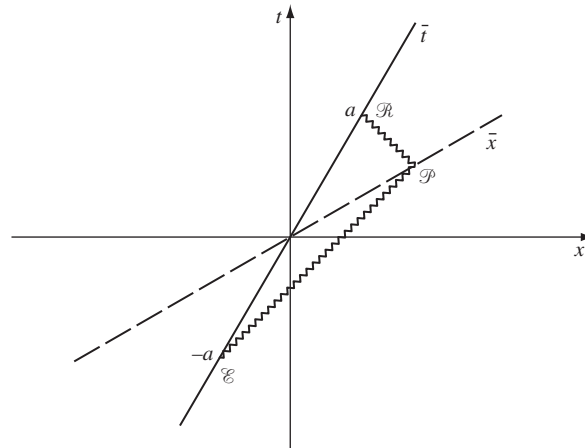
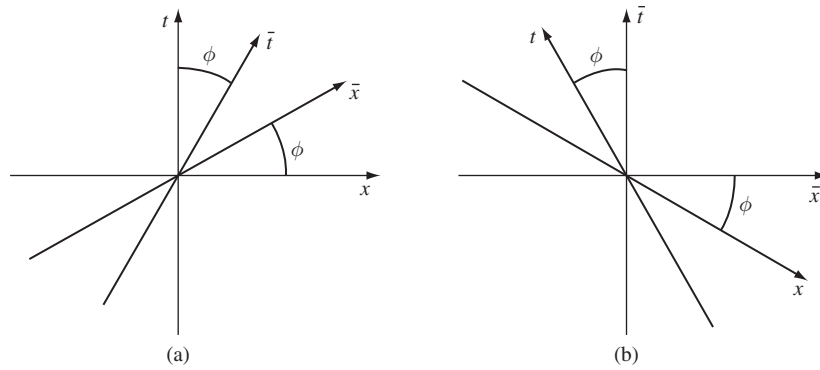Figure 1.4          The reflection in Fig. 1.3, as measured $\mathcal{O}$.



Figure 1.5          Spacetime diagrams of $\mathcal{O}$ (left) and $\bar{\mathcal{O}}$ (right).

*not* coincide with the $x$ axis. If you compare this diagram with the previous one, you will see why: in both diagrams light moves on a 45° line, while the $t$ and $\bar{t}$ axes change slope from one diagram to the other. This is the embodiment of the second fundamental postulate of SR: that the light beam in question has speed $c = 1$ (and hence slope = 1) with respect to *every* observer. When we apply this to these geometrical constructions we immediately find that the events simultaneous to $\bar{\mathcal{O}}$ (the line $\bar{t} = 0$, his $x$ axis) are not simultaneous to $\mathcal{O}$ (are not parallel to the line $t = 0$, the $x$ axis). This *failure of simultaneity* is inescapable.

The following diagrams (Fig. 1.5) represent the same physical situation. The one on the left is the spacetime diagram $\mathcal{O}$, in which $\bar{\mathcal{O}}$ moves to the right. The one on the right is drawn from the point of view of $\bar{\mathcal{O}}$, in which $\mathcal{O}$ moves to the left. The four angles are all equal to arc tan $|v|$, where $|v|$ is the relative speed of $\mathcal{O}$ and $\bar{\mathcal{O}}$.

## 1.6 Invariance of the interval

We have, of course, not quite finished the construction of $\bar{\mathcal{O}}$'s coordinates. We have the position of his axes but not the length scale along them. We shall find this scale by proving what is probably the most important theorem of SR, the invariance of the interval.

Consider two events on the world line of the same light beam, such as $\mathcal{E}$ and $\mathcal{P}$ in Fig. 1.4. The differences $(\Delta t, \Delta x, \Delta y, \Delta z)$ between the coordinates of $\mathcal{E}$ and $\mathcal{P}$ in some frame $\mathcal{O}$ satisfy the relation $(\Delta x)^2 + (\Delta y)^2 + (\Delta z)^2 - (\Delta t)^2 = 0$, since the speed of light is 1. But by the universality of the speed of light, the coordinate differences between the same two events in the coordinates of $\bar{\mathcal{O}}(\Delta \bar{t}, \Delta \bar{x}, \Delta \bar{y}, \Delta \bar{z})$ also satisfy $(\Delta \bar{x})^2 + (\Delta \bar{y})^2 + (\Delta \bar{z})^2 - (\Delta \bar{t})^2 = 0$. We shall define the *interval* between *any* two events (not necessarily on the same light beam's world line) that are separated by coordinate increments $(\Delta t, \Delta x, \Delta y, \Delta z)$ to be[6]

$$\Delta s^2 = -(\Delta t)^2 + (\Delta x)^2 + (\Delta y)^2 + (\Delta z)^2. \tag{1.1}$$

It follows that if $\Delta s^2 = 0$ for two events using their coordinates in $\mathcal{O}$, then $\Delta \bar{s}^2 = 0$ for the same two events using their coordinates in $\bar{\mathcal{O}}$. What does this imply about the relation between the coordinates of the two frames? To answer this question, we shall assume that the relation between the coordinates of $\mathcal{O}$ and $\bar{\mathcal{O}}$ is *linear* and that we choose their origins to coincide (i.e. that the events $\bar{t} = \bar{x} = \bar{y} = \bar{z} = 0$ and $t = x = y = z = 0$ are the same). Then in the expression for $\Delta \bar{s}^2$,

$$\Delta \bar{s}^2 = -(\Delta \bar{t})^2 + (\Delta \bar{x})^2 + (\Delta \bar{y})^2 + (\Delta \bar{z})^2,$$

the numbers $(\Delta \bar{t}, \Delta \bar{x}, \Delta \bar{y}, \Delta \bar{z})$ are linear combinations of their unbarred counterparts, which means that $\Delta \bar{s}^2$ is a *quadratic* function of the unbarred coordinate increments. We can therefore write

$$\Delta \bar{s}^2 = \sum_{\alpha=0}^{3} \sum_{\beta=0}^{3} M_{\alpha\beta} (\Delta x^\alpha)(\Delta x^\beta) \tag{1.2}$$

for some numbers $\{M_{\alpha\beta}; \alpha, \beta = 0, \ldots, 3\}$, which may be functions of $\boldsymbol{v}$, the relative velocity of the two frames. Note that we can suppose that $M_{\alpha\beta} = M_{\beta\alpha}$ for all $\alpha$ and $\beta$, since only the sum $M_{\alpha\beta} + M_{\beta\alpha}$ ever appears in Eq. (1.2) when $\alpha \neq \beta$. Now we again suppose that $\Delta s^2 = 0$, so that from Eq. (1.1) we have

$$\Delta t = \Delta r, \quad \Delta r = [(\Delta x)^2 + (\Delta y)^2 + (\Delta z)^2]^{1/2}.$$

---

[6] The student to whom this is new should probably regard the notation $\Delta s^2$ as a single symbol, not as the square of a quantity $\Delta s$. Since $\Delta s^2$ can be either positive or negative, it is not convenient to take its square root. Some authors do, however, call $\Delta s^2$ the 'squared interval', reserving the name 'interval' for $\Delta s = \sqrt{(\Delta s^2)}$. Note also that the notation $\Delta s^2$ *never* means $\Delta(s^2)$.

(We have supposed $\Delta t > 0$ for convenience.) Putting this into Eq. (1.2) gives

$$\Delta \bar{s}^2 = M_{00}(\Delta r)^2 + 2 \left( \sum_{i=1}^{3} M_{0i} \Delta x^i \right) \Delta r$$

$$+ \sum_{i=1}^{3} \sum_{j=1}^{3} M_{ij} \Delta x^i \Delta x^j. \tag{1.3}$$

But we have already observed that $\Delta \bar{s}^2$ must vanish if $\Delta s^2$ does, and this must be true for *arbitrary* $\{\Delta x^i, i = 1, 2, 3\}$. It is easy to show (see Exer. 8, § 1.14) that this implies

$$M_{0i} = 0 \quad i = 1, 2, 3 \tag{1.4a}$$

and

$$M_{ij} = -(M_{00}) \delta_{ij} \quad (i, j = 1, 2, 3), \tag{1.4b}$$

where $\delta_{ij}$ is the Kronecker delta, defined by

$$\delta_{ij} = \begin{cases} 1 & \text{if } i = j, \\ 0 & \text{if } i \neq j. \end{cases} \tag{1.4c}$$

From this and Eq. (1.2) we conclude that

$$\Delta \bar{s}^2 = M_{00}[(\Delta t)^2 - (\Delta x)^2 - (\Delta y)^2 - (\Delta z)^2].$$

If we define a function

$$\phi(\boldsymbol{v}) = -M_{00},$$

then we have proved the following theorem: *The universality of the speed of light implies that the intervals $\Delta s^2$ and $\Delta \bar{s}^2$ between any two events as computed by different observers satisfy the relation*

$$\Delta \bar{s}^2 = \phi(\boldsymbol{v}) \Delta s^2. \tag{1.5}$$

We shall now show that, in fact, $\phi(\boldsymbol{v}) = 1$, which is the statement that the interval is independent of the observer. The proof of this has two parts. The first part shows that $\phi(\boldsymbol{v})$ depends only on $|\boldsymbol{v}|$. Consider a rod which is oriented perpendicular to the velocity $\boldsymbol{v}$ of $\bar{\mathcal{O}}$ relative to $\mathcal{O}$. Suppose the rod is at rest in $\mathcal{O}$, lying on the $y$ axis. In the spacetime diagram of $\mathcal{O}$ (Fig. 1.6), the world lines of its ends are drawn and the region between shaded. It is easy to see that the square of its length is just the interval between the two events $\mathcal{A}$ and $\mathcal{B}$ that are simultaneous in $\mathcal{O}$ (at $t = 0$) and occur at the ends of the rod. This is because, for these events, $(\Delta x)_{\mathcal{AB}} = (\Delta z)_{\mathcal{AB}} = (\Delta t)_{\mathcal{AB}} = 0$. Now comes the key point of the first part of the proof: the events $\mathcal{A}$ and $\mathcal{B}$ are simultaneous as measured by $\bar{\mathcal{O}}$ as well. The reason is most easily seen by the construction shown in Fig. 1.7, which is the same spacetime diagram as Fig. 1.6, but in which the world line of a clock in $\bar{\mathcal{O}}$ is drawn. This line is perpendicular to the $y$ axis and parallel to the $t - x$ plane, i.e. parallel to the $\bar{t}$ axis shown in Fig. 1.5(a).

Suppose this clock emits light rays at event $\mathcal{P}$ which reach events $\mathcal{A}$ and $\mathcal{B}$. (Not every clock can do this, so we have chosen the one clock in $\bar{\mathcal{O}}$ which passes through the $y$ axis