

Yu.VLADIMIROV,  
N.MITSKIÉVICH, J.HORSKÝ

# Space Time Gravitation

Mir Publishers





# Space Time Gravitation

**Ю. С. ВЛАДИМИРОВ,  
Н. В. МИЦКЕВИЧ, Я. ХОРСКИ**

**ПРОСТРАНСТВО  
ВРЕМЯ  
ГРАВИТАЦИЯ**

**ИЗДАТЕЛЬСТВО «НАУКА» МОСКВА**



Yu. VLADIMIROV, N. MITSKIÉVICH, J. HORSKÝ

# Space Time Gravitation

Translated from Russian by  
A. G. Zilberman

Edited by F. I. Fedorov,  
Mem. Belorussian Acad. Sc.



MIR PUBLISHERS MOSCOW

FIRST PUBLISHED 1987  
REVISED FROM THE 1984 RUSSIAN EDITION

*На английском языке*

*Printed in the Union of Soviet Socialist Republics*

© Издательство «Наука», 1984

© English translation, Mir Publishers, 1987

There is a wide variety of literature for the general reader on the subjects of space, time and gravitation, or on general relativity, as Einstein put it (though Einstein's terminology has recently been criticised). However, this book by Dr. Yu. S. Vladimirov of Moscow State University, Prof. N. V. Mitskiévich of the Partrice Lumumba Peoples' Friendship University, and Prof. J. Horský of J. E. Purkyně University, Brno, Czechoslovakia, will still find its own place.

The usual attitude today is to try to explain the basics of modern gravitation theory in very simple terms, so as to bring it home to a layman. Usually a book like this is addressed to a reader with a minimum amount of training in physics and mathematics, and the text would contain few or no formulae. Prof. Horský and his colleagues have adopted a different attitude for this book. They have included a considerable number of formulae but, none the less, it can be said to be for the general reader. Indeed, the notions of derivative and integral have long been taught at school, and the mathematics in this book does not go much beyond this level. Hence, this book should in fact be comprehensible to a reader with just a high school education. What other criterion could be a better yardstick for its suitability for a general reader? It should be noted, however, that this book, especially Chapters Two and Three, does demand more of the reader than is usual.

Although the book is comparatively small, the authors have been quite scrupulous in bringing to the reader's attention the major points of the modern theory of gravitation and many of these points have been ingeniously presented. Chapter One, "Yesterday", is a review of the ideas leading to the theory of general relativity. In Chapter Two,

"Today", the basic established findings in modern gravitation theory are explained. Although this chapter is more traditional in style it still contains some topics which the reader would not easily find in monographs. These include subjects such as gravitational lenses, or the dragging phenomenon in the gravitational fields of Kerr rotating black holes, or monadic method in reference frame theory.

Chapter Three is called "Tomorrow" and takes up almost half the book. It looks at how the science of space, time and gravitation may develop. It covers the theory of gravitational waves and the modern techniques for detecting them; elements of relativistic astrophysics (including black holes); and Penrose diagrams and cosmological singularities. Various attempts at generalization of Einstein's classical gravitation theory are discussed as well as the problems encountered when quantizing gravitational fields, and even the principal points of space-time dimensionality, including the five- and six-dimension theories.

The long lecturing and research experience of the authors contributed to the quality of this book, which, I believe, will be interesting and useful both to students specializing in physics, mathematics or technology and to research physicists including those working on gravitation theory. It should also be instructive to a lecturer giving a course on this subject.

*F. Fedorov*

# CONTENTS

5	Foreword
8	Introduction
	Chapter ONE. YESTERDAY. The Evolution of the Space and Time as Concepts. The Main Steps Towards the General Theory of Relativity
11	1.1. Formation of the Notions of Relativity and the Universe
12	1.2. The Law of Universal Gravitation
17	1.3. From Euclid to Lobachevski
20	1.4. From Riemann to Einstein
26	1.5. The Special Relativity
32	1.6. The Creation of General Relativity (the Geometrization of Gravitational Interaction)
39	1.7. Some Typical Features and Properties of General Relativity
44	Chapter TWO. TODAY. A Review of the Basic Results of Modern Relativity Theory
51	2.1. The Principle of Equivalence and Gravitational Redshift
51	2.2. Schwarzschild's Space-Time
56	2.3. Perihelion Advance and the Solution of Mercury's "Abnormal" Precession
64	2.4. Deflection of Light by Gravitational Fields
67	2.5. Gravitational Lenses
72	2.6. Space-Time Around Rotating Bodies
79	2.7. Dragging in the Kerr Field
87	2.8. Reference Frames in the General Relativity
93	2.9. Study of the Universe as a Whole (Cosmology)
102	Chapter THREE. TOMORROW. Modern Problems in the Theory of Gravitation. The Prospects for the Study of Space and Time
111	3.1. Gravitational Waves
114	3.2. Black Holes and Relativistic Astrophysics
127	3.3. Generalizations of Einstein's Gravitation Theory
161	3.4. Gravitation and Quantum Physics
174	3.5. Dimensionality of Physical Space-Time
186	Conclusion
207	References
211	

## INTRODUCTION

We have tried in this book to present both the evolution of the ideas leading to, and our current understanding of space and time. This branch of science has advanced considerably over the last century. Many perennial philosophical problems, such as whether the Universe is finite or infinite and the relationship between space and time, have been mathematically formulated and thus made accessible for physical investigation.

In our age of the scientific and technological revolution interest in the profound problems concerning our physical understanding of the Universe that we shall consider has grown. Some theoretical findings have already found practical applications and today we are coming to grips with new discoveries which may be even more important for our outlook. These will certainly induce radical changes in science and technology.

When we prepared this book we made use of our experience in lecturing on gravitation theory at universities in our two countries and of our work on books and films about science for the general public. Interest in the space-time problem has recently risen and this gives us reason to believe that this book is really quite topical.

The book was meant to be simple enough to be suitable for a large audience, and so as to include students and lecturers at university level, engineers, and school leavers interested in the subject. Naturally, it is difficult to meet the requirements of readers with such diverse backgrounds in mathematics and physics. Therefore, we have incorporated information at various levels of complexity. Certain places, which are very difficult, such as the paragraphs dealing with the mathematics or illustrated by formulae, or sections

in the final part of the book on the generalizations of gravitation theory, should be skipped at the first reading. The reader can always return to omitted passages for the second and more careful reading. Some of the mathematical points have been deliberately included so that this book may serve as a primer for study of the general theory of relativity. Such readers can use this book alongside texts on Riemannian geometry and general relativity.

A group of readers we wish especially to mention are those, who have had no special training, and yet who try to solve the fundamental problems of the Universe. They send us and our colleagues vast numbers of letters containing "theories" and "solutions" to an impressive variety of the fundamental problems of theoretical physics. In many cases these solutions answer every problem at once, from the nature of gravity to nuclear forces, and the structure of the Universe. Clearly, these ideas have no scientific value. When reading such a "work" you can't help noting that the author uses, say, an 18th century (Lomonosov's time) approach, while another may have gone as far as Faraday. It's only seldom that we come across an amateur whose reasoning is at an early 20th century level. Many of these authors sincerely believe in their ideas and their ability to solve the problems of the Universe (perhaps, that they have solved them), but they only waste their energies and their spare time because there is no way they can contribute to science without having acquired the knowledge already accumulated by mankind and without modern research techniques. This has been demonstrated by the history of science. These readers should study, systematically and seriously, those branches of physics and mathematics they are fond of.

Finally, we would like to emphasize that in modern gravitation theory, as in any other science, there are, on every important question, a variety of viewpoints and differing ideas, between which there are constant clashes. Science cannot progress without this debate, nor could two scientists or research teams be found who approach a broad range of problems in the same way. We, the authors, are no exception and on many issues touched upon in this book we do not have the same opinion. We have however attempted to strike a compromise that reflects each of our positions when presenting debatable aspects. It would however be beyond the scope of this book to go into all the details and the slight

differences of opinion. An interested reader can look these up in our monographs [49, 70, 74, 109, 110] on gravitation theory.

We have divided the text into three chapters. In Chapter One, "Yesterday", we discuss the evolution of the science of space, time and gravitation and the way it has led to the general theory of relativity. Chapter Two, "Today", presents the established facts and concepts in the physics of space-time, that is those that completely agree with the theory and have been proved experimentally. Of course, the chapter includes some unproven theoretical findings but whose truth is unquestionable in the framework of general relativity. Some findings that could have been included in this chapter, though they remain debatable having not been rigorously verified, but which nevertheless are important for physics, have been placed in Chapter Three, "Tomorrow". It is here that the prospects and generalizations of the theory of space, time, and gravitation are discussed.



## YESTERDAY

---

### The Evolution of the Space and Time as Concepts. The Main Steps Towards the General Theory of Relativity

There are many articles and books devoted to the history of scientific picture of the Universe and our modern understanding of the space-time and gravitation. Is it worth digging into the past, again, to expose the faults and misjudgements of our predecessors? Should we not leave this to historians and instead go straight into the problems of modern physics? The evolution of science is a complex process in which every step is determined by what has already occurred. Therefore, it is difficult to come to a proper understanding of today's problems without a knowledge of the past.

It is not enough, however, to restate that the future originates in the past. The process of understanding the external world has objectively and dialectically led to a framework of living and interacting ideas, images and concepts. It would be rash to claim that some of these old ideas are wrong and obsolete. In reality, it is for there to be several competing ideas at any one time. At a given moment (and in a given approximation) one prevails over the others, which are discarded. With time, however, a seemingly contradictory idea may become dominant. This can be illustrated by quite a number of cases in the history of physics, such as the heated debate about the short-range or long-range nature of physical interactions,\* about the continuity or discontinuity of the Universe, or about the corpuscular or wave nature of matter.

That is why some ideas that will be needed for the development of physics will undoubtedly come back from the

---

\* Strictly, the debate concerned 'action at a distance', that is whether a force requires a medium through which to act. Our terminology should not be confused with the range over which different forces act, a topic we shall not be discussing in this book.

past. There is even a saying that the new is the thoroughly forgotten past.

The history of physics in fact shows that its evolution involves many different ideas which quite unexpectedly interweave to emerge in new combinations and as new theories. In this chapter we will try to show which chains of physical thought were involved in the development of the notions of space-time and gravitation and how they finally led to modern gravitation theory. Physics today, as at any other time in its history, can be compared to a tangle of ideas. Which elements of the chain should be connected so as to advance to the next phase of our knowledge? This is our task today and tomorrow, but it can only be completed if we start from what was achieved in the past. This book is not limited to description of the main current of the developments in physics, we also include several of the side-streams of thinking which are, to a certain extent, present in modern research and may well become dominant tomorrow.

It is important to note that learning is a process of reproducing knowledge, and everyone who learns relives the history of humanity, but in a very compressed sort of way. We see our task as helping the younger generation walk along this tortuous but highly rewarding path with as little pain and as quickly as possible.

### **1.1. FORMATION OF THE NOTIONS OF RELATIVITY AND THE UNIVERSE**

Sometimes Aristotle (384-322 B.C.) is referred to as the father of physics. In modern terms, physics is the science of the basic rules ("the primary causation", as Aristotle put it), the principles ("the origin") of Nature and its "elements". The method of cognition employed by Aristotle is very different from that used nowadays. His *The Physics* was more philosophical than a guide to the advance of natural science.

Aristotle hesitated between materialism and idealism. Recognizing the objective existence of matter, he rejected Plato's Universe of *eidos* (Greek for idea, image) and propounded a theory of four types of cause: (1) a material cause or matter, (2) a formal cause or form, (3) an efficient cause, and (4) a final cause or purpose. Matter, according to

Aristotle, is the "first substrate" of each object [3]. And each object is shaped matter, i.e. the form is the essence of existence. The potential inherent in matter is realized due to the effect of form. Aristotle assumed that each phenomenon in nature contained an intrinsic *entelécheia* (Greek for objective, completion). Aristotle believed that what animated all natural phenomena was the "final" purpose. The evolution of nature was, he felt, a process in which matter was shaped, in which its potential was turned into reality. This is a surmise about the unity of form and matter, and yet Aristotle distinguished between these notions.

Progress consists in the destruction and creation of things, of their growth or change. Let us discuss the concept of a "local progress", that is translation. In all conceivable circumstances local progress is understood to be movement from place to place. According to Aristotle, there are two concepts of "place". The first implies that a place is the internal surface of the environment of a body. This concept is practically useless. Suppose, for example, that a boat anchored in a river is washed by changing currents of water, it would then be changing its place, whereas a boat carried along by the river would be stationary. According to the second concept, a place is determined with respect to a body at rest, but calls for definitions of absoluteness and relativity, and of what the immobility ascribed to the reference body is.

In *The Physics*, Aristotle brings up the concept of natural motion as a basic concept of dynamics. Bodies that occupy places that are not natural for them are impelled to move, whereas those in their natural places are not. No natural motion can exist, according to Aristotle, if every point of space is equivalent. He also denied the existence of a vacuum. In a vacuum, i.e. in a uniform space, he pointed out that nobody can say why a body set in motion would stop, for why should it stop in one place and not another? This Aristotelian conception of motion contradicts the viewpoint of the Greek atomists Democritus and Epicurus who claimed that there is an infinite outer space which contains things and is the arena for their movement.

In the same manner that he denied the existence of the vacuum, Aristotle refused to accept the notion that time is independent of events. His argument was that time cannot exist without change, that is if the present were not differ-

ent in each situation but remained the same, there would be no time.

At the foundation of Aristotelian dynamics and cosmology there lies the concept of an active causation that sustains every kind of motion. Thus, radial motions (the Earth was believed to be the centre of the Universe) were considered absolute. As a result of these motions the states of bodies and their role in the Universe harmony were changed. Relative motions, on the other hand (circular paths around the Earth, the centre of the Universe), Aristotle held, did not affect the static harmony of the centre and the spheres. The cosmological systems of Aristotle and Ptolemy (2nd cen. A.D.) were different but not in essence. Aristotle's geocentric Universe, however, was distinctly different from Aristarchus's (4th-3d cen. B.C.) heliocentric system of the world.\*

As mentioned above, Aristotle defined form as the force that makes matter be what it is. Moreover, he believed the *nús* (Greek for intellect, universal spirit) to be the "form of all forms". Thus his idealism had spilled over into theology. The universal *nús* is the "primary power" that governs the expedient motion of the whole Universe. It was this aspect of Aristotle's philosophy that interested scholastics in the Middle Ages such as St. Thomas Aquinas.

In the Middle Ages Aristotle's teachings were canonized by the Church, eventually becoming an obstacle to the development of physics and astronomy. But we should not jump to the other extreme and reject all the reasonable aspects of Aristotle's philosophy. Lenin once said that Aristotle was empirically minded but his empiricism was intellectual rather than limited. He had worked out questions of dialectics and logic. Aristotle is a founder of formal logic, the science of the laws and forms of thought.

In 1543, Nicolaus Copernicus (1473-1543) published *De revolutionibus orbium coelestium* (On the Revolutions of the Heavenly Orbs). This introduced a new theory concerning the Universe, and this led inevitably to fresh concepts of relative motion and physical relativity. Copernicus adduced the following analogy: It seems to the sailors aboard a ship sailing in calm weather that everything around the ship

---

\* We mention also the world model of Pythagoras the centre of which was the central fire, while the Sun revolved, like the other celestial bodies and the Earth, about this fire. As to the role of Ptolemy, see *The Crime of Claudius Ptolemy* by Robert R. Newton.

is moving, as if reflecting its movement, whilst they believe that they themselves as well as everything close to them are at rest. Obviously, Copernicus continued the argument, the same occurs for a moving Earth and it seems to us that the rest of the Universe rotates around us [18].

So, in terms of kinematics, it is equally possible for either the observer or the observed to be in motion, either one could be at rest. However, considerations of an astronomical and philosophical nature led Copernicus to the view that immobility of the Earth was illusory and that it moved around the Sun. Copernicus's ideas were later taken up by Galileo.

**Giordano Bruno** (1548-1600) developed Copernicus's teaching philosophically. He criticised the doctrines of Aristotle and Ptolemy, becoming in fact a proponent of the philosophy of Democritus and Epicurus. He rejected Aristotle's teaching of finiteness of the Universe, of the opposition of the earth and the heavens, and that there was an absolutely fixed centre in the Universe. In his interpretation of the relativity of motion and rest, Bruno shared the ideas of **Nicolaus Cusanus** (1401-1464) whose writings were known to Copernicus.

Whilst advocating Copernicus's teachings, **Galilei Galileo** (1564-1642) proposed his own principle of relativity, and to illustrate it, he described the motions to be observed in a closed cabin on board a ship at rest and those on board a ship at sea. This description is contained in his *Diálogo sopra i due massimi sistemi del mundo: tolemaico e copernico* (Dialogue About Two Basic Systems of the World: Ptolemy's and Copernicus's). Having described what happened when the ship was at rest he pointed out that if the ship was then to move at some speed and provided its motion was uniform (no rolling or pitching) then none of the phenomena would be changed in the slightest nor could an observer distinguish whether the ship was moving or standing still [41].

This formulation contains a very important physical principle, the Galilean principle of relativity. No mechanical test will reveal whether a system is at rest or is moving uniformly in a straight line. Any movement within these two reference frames is identical.

**Albert Einstein** paid a great deal of attention to Galileo's scientific heritage. He noted, in particular, the striking similarity between the contributions of Faraday and Maxwell

in their epoch and those of Galileo and Newton in the 17th century. In both cases the first member of the pair qualitatively described a fundamental law while the second produced the exact mathematical formulation and applied it quantitatively.

Isaac Newton's (1643-1727) concepts of absolute space were the culmination of a long historical process. Aristotle's tennet that "Nature abhors a vacuum" dominated thinking for many centuries. After a prolonged debate the concepts of the atomists of antiquity regained their place in science, moreover, after the discovery of the vacuum in the 17th century they started to acquire more and more proponents.

Newton considered space to be a void arena of things and phenomena. It was three-dimensional, continuous, static, infinite, uniform, and isotropic. He believed that absolute space, in its own nature and with regard to anything external, always remains similar and unmovable.

Newtonian time was also absolute and independent. He regarded it to be "receptacle of events" and that the course of events did not affect the flow of time. Time was thus unidimensional, continuous, homogeneous and infinite.

Newton's view of motion was similar. In a reference frame stationary with respect to absolute space, Newton's three laws must hold: (1) the law of inertia, (2) the law of motion, and (3) the law of action and reaction.

The force  $F$  in the second law is due to the interaction between bodies. A gravitational force is an example of such a force.

An absolute frame of reference, fixed with respect to absolute space is an inertial frame; and the transition from one such frame to another is accomplished by a Galilean transformation, i.e.

$$t' = t, \quad \mathbf{x}' = \mathbf{x} + \mathbf{v}t. \quad (1.1)$$

Suppose a particle with a mass  $m$  is moving in an absolute reference system  $S$  according to the law  $m d^2\mathbf{x}/dt^2 = \mathbf{F}$ . If now we consider another inertial reference frame  $S'$ , and because according to Newtonian mechanics,  $\mathbf{F}$  and  $m$  are absolute quantities, i.e. they are the same in both reference frames ( $\mathbf{F}' = \mathbf{F}$ ,  $m' = m$ ), then using the Galilean transformations, it can be shown that in the new reference system,  $m' d^2\mathbf{x}'/dt'^2 = \mathbf{F}'$ . This means that Newton's second law is

invariant with respect to these transformations. Thus, all inertial reference frames are equivalent and there is no way of detecting absolute space.

## 1.2. THE LAW OF UNIVERSAL GRAVITATION

The discovery of the law of universal gravitation became possible as a result of evolution of ideas. In Aristotle's philosophy, for example, the problem of gravitation simply did not exist. The celestial bodies were "naturally placed", and their motion was taken for granted without the need for any force. The other bodies were always "attracted" to the centre of the Universe, that is, the centre of the Earth. It was also believed that the velocity of a falling body was proportional to its mass.\*

Copernicus's view was a significant step forward towards an understanding of gravitation, for it stated that gravity not only existed on the Earth, but affected the other celestial bodies. Next, it was necessary to eliminate the delusion that the velocity of a falling body depended on its mass. Galileo is said to have started experimenting with different weights released simultaneously from the top of the leaning tower of Pisa in about 1589.

Johannes Kepler's (1571-1630) contribution was also important. His first scientific study, *The Cosmographic Enigma*, which was published in 1596, was essentially a search for a numerical relationship between the various characteristics of the planetary orbits in the Solar system. In 1602, Kepler discovered the second law of planetary motion, viz. the radius vector from the Sun to any planet sweeps equal areas in equal intervals of time. In 1602, Kepler discovered the law later called the "first", viz. the orbits of the planets are ellipses with the Sun at a focus.

It is thought that Newton discovered the universal law of gravitation

$$F = G \frac{m_1 m_2}{r^2} \quad (1.2)$$

between 1667 and 1670 [41, 60], but he did not publish his discovery for a long time. Independently and at about the

---

\* In a vacuum the free fall velocity, according to Aristotle, should be infinite.—N.M.

same time **Robert Hooke** (1635-1703), **Giovanni Borelli** (1608-1679) and **Christian Huygens** (1629-1695) all came close to discovering the law, too. Hooke published an essay on the Earth's motion in 1674 in which he formulated the idea of universal gravitation qualitatively. He assumed, however, that the force was inversely proportional to the first power of the distance. In 1680, in a letter to Newton, Hooke gave the correct form for the law, i.e. that the force was inversely proportional to the square of the distance. When in 1686 Newton presented his *Principia*, containing the law of universal gravitation, to the Royal Society, Hooke claimed priority for the discovery. Newton rejoined that he had known the law for 20 years and referred to a letter he had sent Huygens via the Secretary of the Royal Society [41, 86].

The discovery of the law of universal gravitation enabled other scientific concepts to be clearly formulated. We shall consider two of these in view of their significance for the development of physics. The first one is the concept of mass. On the one hand, mass can be determined by measuring the force with which a body is attracted to a standard body. The value obtained will characterize the gravitational properties of the test body, its attractability to the standard. Thus, the "gravitating mass"  $m_{gr}$  can be determined by applying the universal law of gravitation. On the other hand, the mass of a body can be determined using Newton's second law by measuring the test body's acceleration when acted upon by a standard force. This value will characterize the body's inertial properties and its ability to keep its velocity, and is called inertial mass  $m_{in}$ . Galileo found that the acceleration of a freely falling body is constant and independent of its mass. We can get a similar result by assuming that  $m_{in} = m_{gr}$ , given that  $m_{in}a = m_{gr}g$ . (In general, the masses need only be proportional for us to be able to get the same conclusion.) Newton experimented with pendulums and proved that the periods of their oscillations were independent of their mass. This experimental evidence enabled Newton to equate the masses in his second law and in the law of gravitation. That the masses should be equated only follows from experiment. Moreover, the physical sense of the two concepts is profoundly different: a gravitating mass is essentially the gravitational charge of a body, whereas its inertial mass is a measure of its "resistance" to the action of a force.



It should be noted that neither Galileo's nor Newton's experiments were very accurate. Later, similar experiments were done and repeated many times with ever greater precision. That the two masses are equal has been established with a very high degree of accuracy (see Sec. 2.1). Jumping ahead we may note that the equality of the gravitating and the inertial masses resulted in the principle of equivalence, which in turn happened to be extremely important for the development of the general theory of relativity.

The second concept was related to the nature of interaction between two gravitating bodies. It can be said that Newton's works initiated a centuries-long discussion as to whether bodies can affect each other at a distance (long-range interaction), or whether a medium is needed in order to mediate the action (short-range interaction). Newton himself was not a consistent advocate of either point of view. However, some of the opponents of the long-range interaction theory (e.g. Faraday, Maxwell, Thompson) would cite Newton, who in a letter to Bently wrote that he could not imagine a situation in which an inanimate substance could directly, without a material intermediary, influence another substance other than by touching it. Were an intermediary unnecessary, as Epicurus thought, then gravity would be a property of matter. Then Newton wrote that he wished that Bently had not credited him with the doctrine that gravity originated from matter, for to do so was to allow that a body can act on another over a distance, across a void, without any intermediary to transmit the action and force. This view, Newton felt, was so absurd that no one who could think philosophically would make such an error. Gravity, he felt, must be caused by a concrete "factor" which acts according to specific laws [60]. The opponents of action at a distance would usually cut short their quotations at this point, but Newton went on to write that he left it to his readers to discover whether the "factor" he had spoken of was material or not. We can thus conclude that Newton was inclined to believe in a certain divine intervention [86]. The mathematician Cotes, an associate of Newton's in Cambridge, is believed by some to have first introduced the theory of long-range interaction. He did so in his preface to the second edition of *Principia* in 1713. It is also believed that Newton never read this preface and therefore could not have either consented or dissented. All in all, the long-range theory took

root in physics and prevailed for a long time, apparently until the 19th century. The first theories concerning electricity were also influenced by it. Then, after the success of Maxwell's theory of the electromagnetic field it was replaced with the short-range theory, which has come down to us now in the form of modern field theory. The general theory of relativity is of a short-range nature. This does not mean, however, that the long range interaction concept turned out to be wrong. We shall show later that it can be used to formulate both theories of electromagnetic and gravitational interactions, and in a way that is not inferior to the generally accepted theories developed by Maxwell and Einstein. These theories are now known as the theories of direct interparticle interaction (see Sec. 3.3.4). Only the future will show which of the two concepts will finally win.

### 1.3. FROM EUCLID TO LOBACHEVSKI

Difficult was the path mankind traversed to its realization that space, or to be more exact space-time, is curved. To be able to appreciate the grandeur of the achievements of science and try to anticipate the future let us consider the main events in this chain of human thought, i.e. the development of ideas about the geometry of the world.

We shall start with Euclid's fifth postulate. A monumental 13-volume work known as Euclid's *Elements* was compiled by Euclid early in the 3d cen. B.C. He summarized all the geometry of the ancients and presented as demanded by Aristotle's logic. In the history of mathematics this book, with its impeccable logic and depth of insight, has never been challenged by any other single work.

*Elements* remained the foundation of geometrical instruction practically to the beginning of the 20th century. It explained geometry on an axiomatic basis in the form of propositions or theorems derived from a limited number of basic axioms (that is postulated without proof and regarded as self-evident). Euclid's axioms were [42, 64].

- (1) Given two points there is an interval that joins them.
- (2) An interval can be prolonged indefinitely.
- (3) A circle can be constructed when its centre and a point on it are given.
- (4) All right angles are equal.
- (5) If a straight line falling on two straight lines makes the

interior angles on the same side less than two right angles, the two straight lines, if produced indefinitely, meet on that side on which the angles are less than two right angles.

At first glance it becomes obvious that the fifth postulate differs in clarity from the first four. For over 2000 years many mathematicians believed this postulate could be logically derived from the other Euclidean axioms. Some believe that Euclid himself hesitated when including the fifth postulate in his list of axioms. That *Elements* consists of two parts, i.e. theorems that are proved without the fifth postulate (absolute geometry) and the theorems that are based on the fifth postulate (Euclidean geometry proper), supports this idea. Euclid may have chosen to separate his material this way after unsuccessful attempts to prove the fifth postulate.

In any case, for 2000 years many attempts were made to prove the fifth postulate. During the history of mathematics many proofs have been suggested, for example, by Posidonius (1st cen. B.C.), Ptolemy (2nd cen. A.D.), Proclus (410-485), Nasir Eddin (1201-1274), J. Wallis (1616-1703), G. Saccheri (1667-1733), Lambert (1728-1777), A. M. Legendre (1752-1833), and F. Bolyai (1775-1856) [20, 42]. When closely inspected these purported proofs either contain logical faults or are based on assumptions that were taken to be self-evident but which in reality were equivalent to the fifth postulate. For example, the fifth axiom is equivalent to the following formulations [52]:

“Through a point  $C$ , lying on a given straight line  $AB$ , only one straight line parallel to  $AB$  can be drawn” (that is, a line lying in the same plane with the given straight line and not intersecting it);

“two parallel lines are equidistant”;

“three non-collinear points always lie on a circle”;

“the sum of the angles in any plane triangle is equal to two right angles”.

The history of human race does not abound in problems on which so much effort has been spent as the proof of Euclid's fifth axiom. Perhaps, only the search for the philosophers' stone or the innumerable attempts to create a perpetuum mobile in the Middle Ages are comparable. Sometimes those efforts took on a dramatic tone. For example, there is the letter Farkas Bolyai wrote to his son János Bolyai, who took up the work on this problem [64]:

"You must not try to conquer the theory of parallel lines; I know this path, I have followed it to the end, I have lived through an endless night, and buried all the light and all the joy in my life. I beg of you, leave the parallel-line axiom alone. Avoid it as you should avoid lust, for it will rob you of your health, your time, your patience; it will destroy the happiness in your life. This bottomless abyss of darkness can engulf a thousand such giants as Newton. There will never be light in the world, for mankind will never find the absolute truth, nor will it reach it in geometry. It was a terrible and endless wound in my soul; God protect you from this obsession, which has so powerfully seized you. It will rob you of happiness not only in geometry but in your earthly life, too. I was ready to become a martyr for truth, to present humanity with a geometry without this blemish. I have done a gigantic, hard job, I have achieved much more than has been obtained by others, and yet I am not content.

"Learn from my example: I remain unknown because of my desire to overcome the parallel-line postulate. It has sucked dry my blood, it has devoured all my time. It was the root of all my subsequent mistakes. If I could have resolved the enigma of the parallel-line postulate, even if no one else were to have known I had done so, I would have thought myself an angel.

"It is inconceivable that this impenetrable darkness, this eternal blackness, this cloud, this black spot on a virgin, pure truth should still exist in geometry.... To eradicate the blemish is to commit yourself to a Herculean task, hold on, or you will perish!"

The Herculean task that had faced humanity for over twenty centuries was only solved in the first half of the 19th century. This important achievement in the history of thought is associated with the names of **Nikolai Lobachevski** (1792-1856), **János Bolyai** (1802-1860), and **Carl Gauss** (1777-1855). The history of their discovery has been recounted more than once [20, 42, 52, 60, 64], and it teaches us many lessons. We would like to emphasize a few of them.

First, it is a classical illustration of a discovery made when the time was ripe for it. Typically, a mature idea occurs almost simultaneously and independently to several people. The problem was solved by three mathematicians strikingly close in time. Lobachevski made his famous presentation

*On the Fundamentals of Geometry* to a session of the Scientific Council of the Physics and Mathematics Department, Kazan University, on February 23, 1826, and published it in 1829. J. Bolyai after five years of study published his *Appendix* (to a voluminous work by his father) in 1832. We know from his notebooks that Gauss had also developed aspects of the new geometry in the 1820s. But the list is not ended. Professor **Schweickart** (1780-1859), a law lecturer at Kharkov University between 1812 and 1816 who moved to Germany in 1817, discussed his ideas with Gauss and is also known to have developed a non-Euclidean geometry. In 1824, similar ideas were communicated in a letter to Gauss by **F. Taurinus** (1794-1874), Schweickart's nephew and also a lawyer. Other names could be included such as **Wachter** and **DeTilli** [52].

Despite all the differences between the people who made this great discovery—in their temper and nationality, in their attitude to their results—they all faced one thing in common: an almost complete misunderstanding and even hostility on the part of their colleagues and the general public. Lobachevski's interest in non-Euclidean geometry caused him to be viewed in Russia as a crank, at best. Worse, he was attacked in a humiliating and ignorant article in *The Son of Fatherland* periodical, and there were mocking and rude remarks by distinguished contemporaries. All of Lobachevski's students turned their backs on him. At his funeral, when it is common to praise a deceased's deeds, nothing was said about the subject that was the main thing in his life—non-Euclidean geometry.

J. Bolyai also had a bitter life. He died in 1860, and his burial ceremony resembled a ritual of oblivion. Only three people were present to see his remains placed in a nameless mass grave, and the entry in the church register read: "His life was passed uselessly" [64].

Carl Gauss, the greatest European mathematician of the time, was an example of common sense. He clearly realized the scale of the perturbation in geometry (and not only in geometry) that would be occasioned by the discovery of non-Euclidean geometry, but he realized what the reaction of his colleagues and contemporaries would be to the discovery itself and to those who would dare to support it openly. He preferred to retain his status in society, he chose a quiet life and did not publish the results of his work. To give him

credit, however, it should be admitted that he did not betray his ideas in science. In a letter to Bessel in 1829, Gauss wrote that he would probably not adapt his lengthy investigations on this subject for publication. He thought he might even decide never to do so because he was afraid of the clamour of the Boeotians that would have been raised if he were to publish all his ideas. In a letter to Herling he wrote that he was very glad that Herling had had the courage to express himself in such a way as to admit the possibility that the parallel-line theory and, therefore, that the whole of geometry (Euclidean) might be incorrect. But he warned Herling to look out for the hornets whose nest he was destroying [64].

But let us now turn to the essence of the discovery. Despite the slight differences in method, in depth and scale of the research, the basic finding was the same. Basically the mathematicians investigated what would happen if they were to disregard the fifth postulate and assume the opposite, that is, that through a point  $C$  not on a given straight line  $AB$ , not one but two (and, consequently, an infinite number of) lines parallel to  $AB$  can be drawn. The task was to construct a geometry based on this new axiom. The idea was that if the fifth postulate was really a theorem, then, sooner or later, the new geometry would contain logical contradictions, which would mean that the initial assumption was wrong, and the fifth postulate would thus be proven. But after constructing this new geometry the mathematicians could find no contradictions. Moreover, they discovered that they had a new and elegant geometry with a number of interesting and unique characteristics. In the new geometry, the sum of the angles of a triangle was less than  $180^\circ$  and would indeed depend on the triangle's linear dimensions. A certain parameter with a dimension of length (and called the space-constant) arises in the theory; and system's geometrical properties depend on the relation of its size to this parameter. No similar figures exist in this theory. In a very small space (in comparison with this new parameter), the new geometry was practically Euclidean, but in a large one the two were essentially different. Lobachevski called his geometry "imaginary" (or "pangeometry"); Schweickart called it an "astral" geometry. The problem, however, was not in the name but in the difference between this new geometry and the Euclidean one.

It is important to note that Lobachevski and Gauss did not confine themselves to the mathematical aspect of the discovery, they pondered about how the new geometry was related to the physical world. They wanted to know which of the two geometries describes real space. To answer this question, Gauss measured the sum of the angles in a triangle formed by three mountain peaks, while Lobachevski preferred a much larger triangle—that formed by two positions of the Earth in its orbit and a distant star, he measured the parallaxes of the stars. But neither Gauss's measurements nor Lobachevski's astronomical observations answered the question (and now we know that at that time they simply could not, for accuracy of astronomical measurements was far too small).\*

Even though Lobachevski, J. Bolyai and the others were convinced of the correctness of their geometry they did not finally prove its logical consistency. For it is one thing to have no contradictions in a theoretical construction, even if it is very advanced, but it is quite another to claim that there will never be any contradictions in the new theory. A final proof of the consistency of Lobachevski's geometry only came in the 1870s after work by **Eugenio Beltrami** (1835-1900), an Italian geometer, and **Felix Klein** (1849-1925), a German mathematician. The main idea of the proof was to generalize the first non-Euclidean geometry, which was originally constructed for a plane, to a geometry on a three-dimensional hypersurface with a constant negative curvature (three-dimensional hyperboloid) in the framework of four-dimensional Euclidean geometry, whose consistency was already known. It was only necessary to replace the notions of straight lines (the shortest lines in Euclidean world) by those of geodesics (extremal curves) on a hypersurface. Then, all the statements regarding the straight lines in Lobachevski's geometry would be converted into

---

\* It is interesting to note that even if Lobachevski had possessed instruments to measure the star's parallax with enough accuracy he would have been disappointed for the sum of the angles in his triangle would have been in excess of  $180^\circ$ , instead of less than  $180^\circ$ . Neither his imaginary geometry nor the Euclidean one would have been correct. As Arifov and Kadyiev showed in their presentation to the 5th International Gravitational Conference in Tbilisi in 1968, the negative parallaxes obtained by astronomers (angle sums greater than  $180^\circ$ ) were due to mechanisms described by the theory of general relativity. See Sec. 2.4.

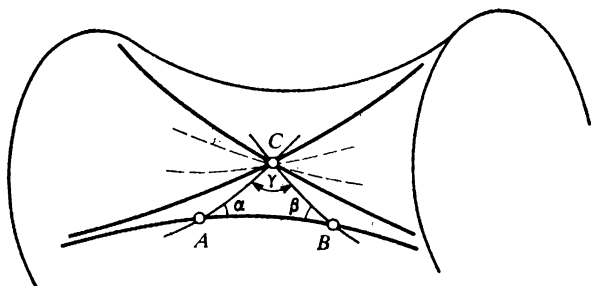


Fig. 1

corresponding statements for geodesics on a hyperboloid. Since it is impossible to visualize a hyperbolically transformed three-dimensional world, we shall illustrate the idea using lines, called hyperbolas, on a two-dimensional hyperboloid. Figure 1, for example, is an explanation of a generalized version of Euclid's fifth postulate. At a point  $C$  not on a given hyperbola  $AB$ , two hyperbolas meet which do not intersect  $AB$ . Therefore, all the other hyperbolas (broken lines) will not intersect  $AB$ . A triangle formed by three intersecting hyperbolas is also shown in the figure. It is easy to see that the sum of the angles is  $\alpha + \beta + \gamma < 180^\circ$ . For this reason the first non-Euclidean geometry (Lobachevski's geometry) is often referred to as the hyperbolic geometry. In this geometry, the space-constant acquires the sense of the radius of curvature of a three-dimensional hyperboloid. Now it is easy to understand that the properties of geometrical figures depend on their size.

#### 1.4. FROM RIEMANN TO EINSTEIN

The next substantial step in the development of the geometry of space was made by **Bernhard Riemann** (1826-1866), a German mathematician, in 1854. Incidentally, this contribution by a young scientist of the generation that followed Lobachevski and Bolyai was also connected, to some extent, with Gauss. In order to secure his position as assistant professor at Göttingen University, Riemann had to deliver a test lecture. Following the established procedure, he presented three alternate topics to the Department. The first two were related to problems then current among mathema-



ticians, whereas the third, and least prepared by Riemann, was devoted to the foundations of geometry. Riemann did not expect the latter to be chosen. But, **Wilhelm Weber** (1804-1891), a German mathematician, later wrote that Gauss had intentionally picked this topic. He had admitted that he was eager to hear how such a young man would manage to cope with this rather difficult game [64]. Riemann did indeed deliver this lecture, and it was later published as *About the Hypotheses Lying at the Foundation of Geometry* [85]. Apparently, Riemann prepared the lecture exclusively for Gauss's consumption. Riemann was successful. When the lecture was over, Gauss rose in silence and sauntered out. According to Weber the lecture went beyond all of Gauss's expectations. It reduced him to a state of utter amazement, and as he left the session, he extolled Riemann to the skies, which was very rare for him [64].

Gauss was amazed by the approach Riemann used to a non-Euclidean geometry, since it was quite different from those used by his predecessors. Apparently Riemann knew nothing about Lobachevski or Bolyai and had just a vague idea of Gauss's interest in the subject. He succeeded, however, in incorporating into his study two extremely fruitful ideas: Gauss's mathematical apparatus for describing the geometry of two-dimensional curved surfaces and his own new concept of a multidimensional manifold (multiply extended geometrical objects). A surface is thus a two-dimensional manifold, a space is a three-dimensional manifold, and that is all the difference there is. All the ideas and methods used to describe two-dimensional surfaces can be directly applied to three-dimensional curved spaces. Among notions used the most important one is the metric, i.e. the quadratic form for the differences between coordinates, which describes the length of the interval between two neighbouring points in a curved manifold. A historian aware of the interest shown by Gauss himself can only wonder why the idea did not strike Gauss before Riemann.

This successful integration of ideas enabled Riemann to advance when constructing both particular cases of non-Euclidean spaces and a theory of arbitrarily curved spaces. To paraphrase **Albert Einstein** (1879-1955), Riemann's credit for the development of our ideas about the relationship between geometry and physics is twofold. Firstly, Riemann discovered a spherical (elliptic) geometry which was the

opposite to the hyperbolic geometry of Lobachevski. Thus, he was the first to indicate the possibility of a finite geometrical space. The idea immediately took root and brought about the question as to whether our physical space is finite. Secondly, he had the courage to build much more general geometries than Euclid's or the narrowly non-Euclidean ones [30].

The first point deserves an explanation. The professional literature often refers to "Riemannian geometry" (in the narrow sense of the word) as the second non-Euclidean geometry

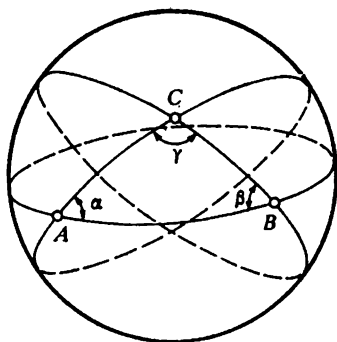


Fig. 2

of spaces of constant positive curvature and it corresponds to geometry on a three-dimensional hypersphere. The essential property of this three-dimensional space is that its volume is finite, so that if a point moves in the same direction it may eventually return to the starting point. Instead of straight lines in Euclidean space, we have in Riemannian spherical geometry geodesics, or the arcs of great circles. From a two-dimensional illustration of geometry on the sphere (Fig. 2), it is clear that the notion of parallel lines as given by Euclidean fifth postulate in this case has no sense at all, because any arc of a large circle that passes through a point  $C$ , not lying on  $AB$ , will necessarily intersect  $AB$ , and even at two points. Figure 2 also indicates that the sum of the angles of a triangle formed by three intersecting arcs of three great circles is always more than  $180^\circ$ .

Einstein emphasized this contribution of Riemann pointing out that he had arrived at the challenging idea that the geometrical relationships between bodies might depend on a physical cause, that is, on forces. Thus, by way of a purely mathematical analysis he discovered that geometry and physics might be combined. This was actually realized seventy years later in the general theory of relativity, which combined geometry and the gravitation theory [30]. Riemann could not, however, have perceived which forces would be related to the non-Euclidean properties of geometry. It is inter-

esting that he did speculate about the nature of gravitation [85], but leaving aside his geometrical ideas.

In his memoir *About the Hypotheses Lying at the Foundation of Geometry*, Riemann made some interesting speculations about the nature of space which have not lost their value even now, and which have become sources of new schools of thoughts. Among them we find the idea that the metric relations in infinitely small spaces may not correspond to the adopted geometrical assumptions, the idea of discreteness, the physical nature of the metric relations, higher dimensional manifolds, etc.

We should also note the substantial contribution to the subject by William Clifford (1845-1879), a British mathematician, who developed the idea that the physical properties of matter and the properties of curved space were related. In *The Common Sense of Exact Sciences* which was published after his death in 1885 and edited by Karl Pearson (1857-1936) [16, 17], he astoundingly anticipated the basic elements of general relativity, which came much later. Clifford wondered whether we could be wrongly interpreting as physical things that in reality arose due to the curvature of space. Might it be that all or some of the forces that we call physical originate from the geometry of our space. He considered there to be three types of curvature in space that we might accept:

I. Our space may have a curvature that changes from point to point, but which we cannot detect either because we are aware only of a small section of space, or because we confuse insignificant changes in it with changes in our physical existence, and we fail to connect the latter to changes in our position.

II. Our space may be everywhere identical (having the same curvature), but the value of this curvature may change over time. If this is the case, our geometry, which is based on the fact that space is identical, will be valid everywhere. The changes in the curvature, however, may produce a number of consecutive visible changes in space.

III. We may conceive of our space as having approximately the same curvature, but slight changes in the curvature may exist from point to point, the curvature itself being a function of time. These changes in curvature over time may produce effects which we, not too unnaturally, ascribe to physical forces independent of the geometry of our space.

In a note to his work, Clifford suggested, for instance, that heat, light and electromagnetic fields might be such physical phenomena. Let us note in passing that the true physical manifestation of curvature is gravitation which became clear after Einstein but which had escaped Clifford. He did, however, hypothesize a connection between the electromagnetic field and the geometry of the space. Clifford should obviously be credited for the first idea of the geometrization of the electromagnetic field.

Clifford was much more definite about possible physical manifestation of the curvature of space than Riemann had been. Later developments showed (and we shall demonstrate this in Chapter Two) that all three types of curvature Clifford had indicated were explained in the general theory of relativity. Jumping ahead, we should note that Clifford's first type of curvature concerns the curvature of space (and time) around gravitating bodies, such as the Sun and Earth. It is this kind of curvature that explains Newton's law of universal gravitation. Clifford's second category is the variation in time of the curvature of space, the curvature being the same at all points; this curvature is incorporated in modern cosmological theories (see Sec. 2.10). The third type of curvature may occur due to a gravitational wave, a search for which is under way.

Let us note also that Clifford translated Riemann's memoir into English. From the reminiscences of Einstein's contemporaries and from his biographers we know that he encountered Clifford's works when he was in Bern (1902 to 1909) [59].

Our scientific heritage from the past not only contains the generally recognized main-line achievements, but also the secondary branches of their thought, which normally reach far ahead of their time and start to become important only much later. This is very true about Clifford. He should be regarded as the originator of the concept of space as a substance, a concept now being advocated by theoretical physicists. Clifford wrote, for example, that what actually occurs when matter—either matter with mass or ether—moves is that the curvature of space changes. In the physical world, he continued, there is nothing except this change, which is (probably) subject to the law of continuity [17]. This statement is nothing but an expression of the complete geometrization of matter. J. Wheeler, a distinguished American

theorist, has said things of this sort. His school's programme is to make "mass without mass", "charge without charge", etc., that is to obtain all the properties of matter from the properties of "empty" space (and time).

Ernst Mach, an Austrian physicist (1838-1916), played an important role in the preparation of the conditions for the creation of the general theory of relativity. Einstein himself wrote that Mach clearly understood the weaker aspects of classical mechanics and was close to the general theory of relativity. This was some fifty years before Einstein's publication of the theory. It is possible that Mach might have developed general relativity if in his time physicists had been concerned with the interpretation of the velocity of light [28].

Even in 1903, i.e. on the very threshold of the creation of general relativity, Mach gave a detailed analysis of the mathematical and physical aspects of the geometry of space in a paper called *Space and Geometry from the Viewpoint of the Natural Sciences* [65, 66]. In the same article he cited and explained the contributions of the appropriate scientists such as Lobachevski, Bolyai, Riemann, and Gauss. His idea was that geometry is brought to life through mathematics applied to our experience gained with relation to space. In the same paper he predicted that the developments, which led to the revolution in our understanding of geometry, would be accepted as part of a healthy and inevitable evolution. Prepared by centuries of work, and recently substantially enhanced, this advance should not in any way be considered over. On the contrary, we must expect it to yield more rich fruit, particularly in the theory of cognition, and not only for mathematics and geometry but for other sciences as well. It was true, he pointed out, that the development was due to certain outstanding individuals, but in general it was called into being by the common need! This is clear from the simple fact that people of many different occupations took part in the movement. Philosophers or mathematicians, they all contributed to these studies and the approaches used by different researchers turned out to be very similar [66].

It is interesting that Einstein's creative genius was influenced by Mach's ideas. Whilst he was working on the general theory of relativity, he was convinced that he was realizing Mach's ideas. Einstein rarely quoted anybody, but

he referred to Mach in the majority of his works at that time.

In retrospect we can say that early in the 20th century all the groundwork necessary for the formulation of the general theory of relativity had been laid. The scientific community of the time was ready to assimilate the physical manifestations of the curvature of space. A number of outstanding geometers, e.g. Sophus Lie (1842-1899), Elwin Christoffel (1829-1900), Gregorio Ricci-Curbastro (1853-1925), and Tullio Levi-Civita (1873-1941), had developed the necessary mathematics of curved multidimensional manifolds (Riemannian geometry) but two elements were still missing [52]. First, there was the unification of space and time in the framework of a four-dimensional manifold, which was completed in the special theory of relativity (see next section).

The publication of the special theory was almost concurrent with the Mach work. Secondly, the theory of gravitation was to catch the eye of physicists. Finally, both elements came into being. Henri Poincaré (1854-1912) in his *About the Dynamics of the Electron*, published in 1906 [83], took the first step since Newton towards a real gravitation theory; he attempted to build it into the framework of the space-time of the special theory of relativity. In 1907, A. Einstein began to work on the theory of gravitation [27]. G. Nordström and M. Abraham (1912) followed suit, but they considered the subject in the framework of flat space-time. Finally, by 1913 the seeds had ripened, and the formulation of the general theory of relativity proper began.

## 1.5. THE SPECIAL RELATIVITY

The unification step of space and time in a four-dimensional manifold occurred as a result of another chain of ideas, hypotheses, experiments and disillusionments. The discovery is another example of conditions maturing before the final step is taken concurrently by several scientists.

As we mentioned above, Newton's equations of motion are invariant under the Galilean transformations. Thus, we can say that these transformations formulate the (mechanical) principle of relativity. No other physical equations capable of describing the motion of matter were known for a long time. In the second half of the 19th century, James Clerk Maxwell (1831-1879) formulated a system of equations that

described the electromagnetic field on the basis of results obtained by Michael Faraday (1791-1867). An electromagnetic field at that time was conceived of as being a strained state of ether. It was considered to be self-evident that Maxwell's field equations were only valid in a preferred inertial system, fixed relative to the absolute rest frame of ether. However, the equations were not invariant under the Galilean transformations. It is worth mentioning that the conceptual evolution of ether is another interesting chain of ideas which has perhaps not yet ended.

In the late 19th century, a series of experiments was launched to detect the motion of the Earth relative to the rest state of ether. There is no need to go into all the designs and the results of these famous experiments which Armand Fizeau (1819-1896) or Albert Michelson (1852-1931) conducted. The main finding was that the absolute motion of the Earth could not be detected. This result undermined the validity of the ether assumptions. A series of additional assumptions had to be introduced which brought the theory closer to our modern understanding step by step. We should make a special note of the hypothesis that a moving body contracts, which was proposed in 1892 by George Fitzgerald (1851-1901). A great deal of work was done in this field by Hendrik Antoon Lorentz (1853-1928). In 1904, Lorentz, and one year later Poincaré, developed the transformation under which Maxwell's vacuum equations were invariant. It is worth mentioning that this transformation, now known as the Lorentz transformation, had already been developed by W. Voigt (1850-1919) in 1887, but his discovery was premature.

The invariance of all natural laws under the Lorentz transformation was first suggested by Poincaré. Finally, in 1905 Einstein showed that this transformation corresponded to the properties of a unified four-dimensional space-time manifold (this fact was best expressed later by H. Minkowski). This development marks the start of the era of the special theory of relativity. The fundamentals and main results of this theory are now known to all high school graduates, but this does not mean, unfortunately, that they (and college students too) have fully understood the principles of the theory. It is not enough to learn the definitions and formulae, we must change the very nature of our physical thought, disentangle ourselves from conventional doctrines and rise

to a higher level of perception. From our own experience we know how difficult this is. We have to pass through a stage of setting up and resolving a number of "paradoxes" of the theory. Only by repeatedly resolving these paradoxes can we be convinced of the validity of this theory\*. This sequence of events occurred in physics throughout the first half of the 20th century: there were many hot discussions, doubts and even the harassment of the proponents of the theory of relativity. In the history of science there have been few barriers separating the comprehension levels of physical principles as difficult as that separating relativity theory from traditional physics. Perhaps, only the comprehension of quantum mechanics was as hard. Even the transition from the special theory of relativity to the general theory requires less effort.

Hermann Minkowski (1864-1909) made an important "geometrical" contribution to the refinement of the principles of special relativity and to its evolution into general relativity. Minkowski's space-time (or flat space-time) is a four-dimensional manifold in which the square of the interval between two neighbouring events in Cartesian coordinates  $s$  is ( $c = 3 \times 10^{10} \text{ cms}^{-1}$  being the velocity of light)

$$(\Delta s)^2 = (c\Delta t)^2 - (\Delta l)^2 = g_{\alpha\beta}^0 \Delta x^\alpha \Delta x^\beta, \quad (1.3)$$

the Greek indices here and from now on in our text can take one of four values, i.e. 0, 1, 2, and 3;  $c\Delta t = \Delta x^0$  and  $\Delta x^\alpha$  are the time interval and the difference between the Cartesian coordinates of the points at which the two events take place; and  $g_{\alpha\beta}^0$  is the metric tensor of Minkowski's space-time, which is defined thus

$$g_{\alpha\beta}^0 = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & -1 & 0 & 0 \\ 0 & 0 & -1 & 0 \\ 0 & 0 & 0 & -1 \end{bmatrix}. \quad (1.4)$$

Here and henceforth, we use the Einstein's summation convention which makes the mathematical expressions much simpler: over any two repeated (in the same summand) indices a summation is automatically performed within the

---

\* The principles of the special theory of relativity and a number of its "paradoxes" are discussed in [7, 8, 9, 23, 40, 43, 54, 63, 90, 102].



range of values that can be assumed by the indices. Minkowski's space-time is also called a quasi-Euclidean four-dimensional manifold (it would have been Euclidean had all the entries in the definition of the metric (1.4) had the same sign). In all inertial frames, the interval  $\Delta s$  remains invariant, but the  $\Delta t$  (time interval) or  $\Delta l$  (distance) may each differ. Clearly, the interval  $\Delta s$  can be written in terms of other sorts of coordinates, e.g. curvilinear ones. The Minkowski metric would then have to be taken for infinitesimal intervals, and formula (1.3) would only be applicable to the differentials of the coordinates. The structure of Minkowski's space-time would not be changed by changing the system of coordinates.

Now we shall go over some points of the special theory of relativity. An example of a Lorentz transformation (for Cartesian coordinates) is

$$x'^{\alpha} = L^{\alpha}_{\beta} x^{\beta}, \quad (1.5)$$

which describes the transition from one inertial frame to another which is moving relative to the first frame with a velocity  $v$  along the  $x^1$  axis, the other axes remaining parallel in both frames. The  $L^{\alpha}_{\beta}$  constants here are

$$\begin{aligned} L^0_0 &= L^1_1 = (1 - v^2/c^2)^{-1/2}; \\ L^0_1 &= L^1_0 = -(v/c) (1 - v^2/c^2)^{-1/2}; \\ L^2_2 &= L^3_3 = 1; \text{ the remaining } L^{\alpha}_{\beta} \text{ are zeros} \end{aligned} \quad (1.6)$$

(a hyperbolic rotation in the  $x^0, x^1$  plane).

Suppose that a process (the motion of a system or the propagation of a field) can be described in an inertial frame  $S$  by the equation

$$F \left( A, B, \dots, \frac{\partial A}{\partial x^{\alpha}}, \frac{\partial B}{\partial x^{\alpha}}, \dots \right) = 0, \quad (1.7)$$

where  $A, B, \dots$  are sets of quantities which characterize the system (the field). The special principle of relativity demands that the same process be described in another inertial frame  $S'$  by an equation with the same form but with new (primed) variables. This demand can be met if the equation is written in the form of the equality of two tensors in Minkowski's space-time.

This is the first time we have mentioned the word tensor. From our experience as teachers we know that at first it often frightens students of general relativity. Tensor analysis does present some difficulties for beginners but they only arise because they are unfamiliar objects. In fact, tensor analysis is not difficult. Tensors are the sets of  $n^k$  components which change, as defined below, when the coordinates are transformed. Here  $n$  is the dimension of the manifold in which the components are defined (we are interested in  $n = 4$ ), and  $k$  is the rank of the tensor. A tensor is usually notated by a symbol, such as  $B_{\alpha\beta}\dots$ , in which the number of indices,  $\alpha, \beta, \dots$  equals its rank. Tensors of the first rank are called vectors, and those of zero rank, scalars.

A reader will not need to know the rules for tensor transformations when reading through this book once; it is enough to know that it exists. For the benefit of more advanced readers, however, we shall explain what it is. Under an arbitrary transformation of coordinates

$$x'^\alpha = x'^\alpha(x^\beta) \equiv f^\alpha(x^\beta) \quad (1.8)$$

the tensors are transformed thus

$$B_{\sigma\lambda}^{\alpha\beta}\dots'(x') = \frac{\partial x'^\alpha}{\partial x^\mu} \frac{\partial x'^\beta}{\partial x^\nu} \dots \frac{\partial x^\kappa}{\partial x'^\sigma} \frac{\partial x^\rho}{\partial x'^\lambda} \dots B_{\kappa\rho}^{\mu\nu}\dots(x). \quad (1.9)$$

Tensors with superscripts are called contravariant, those with subscripts, covariant. A tensor may have both superscripts and subscripts simultaneously. The simplest illustration of a contravariant tensor of the first rank (vector) is the set of differentials of coordinates  $dx^\mu$ , for which we have  $dx'^\mu = (\partial x'^\mu/\partial x^\alpha) dx^\alpha$ , an example of a covariant vector is the set of components of a gradient  $\partial\varphi/\partial x^\alpha$ , which transform as  $\partial\varphi/\partial x'^\alpha = (\partial x^\mu/\partial x'^\alpha) \partial\varphi/\partial x^\mu$ .

For more detailed information we refer the reader to any textbook on tensor analysis or relativity theory [5, 25, 37, 62, 69, 84]. Note however that tensor analysis is an indispensable tool for the formulation of the laws of relativity theory, the general one in particular. This is the language of the theory and without it no progress can be made in the modern theory of gravitation.

It might be well to point out that in mechanics, vector analysis is used extensively and considered quite common. However, if you look through books published in the 1940s

you will see that at that time vector analysis was only just being introduced into mechanics and had both proponents and opponents. Before the 1940s mechanics could be studied without using vectors at all. Nowadays we have started to use tensor analysis to the practical exclusion of vector analysis. Indeed many people believe that even tensor analysis belongs to the past, and it is gradually being replaced in the general theory of relativity by coordinate-free formalisms like Cartan's exterior form analysis.

The squared interval (1.3) taken for two infinitesimally close points of space-time (we then write it as  $ds^2$ ), is an important feature of the relationship between these points. If  $ds^2 > 0$ , then the interval is *time-like*, and events at these points may be causally connected. If  $ds^2 = 0$ , the interval is called *light-like* (or *null*), and its endpoints are connected by the space-time path (*world line*) of a possible light signal. This means that there may be still causal connection between the two events. But if  $ds^2 < 0$ , the interval becomes *space-like*, and only superluminal "signals" may connect its endpoints, thus making them causally independent. If we take some fixed point and consider all the other points connected with it by null world lines (possible light signals), these lines being integral curves of the equation  $ds^2 = 0$ , we arrive at the notion of a *light cone*. The cone has two cavities, or interior regions, one in the future, and one in the past. For the both,  $ds^2 > 0$ , and all the points within these regions may be connected with the apex of the light cone by causally propagating signals; the same is for the hypersurface of the light cone. (The apex of the cone is sometimes referred to as the "*here and now*".) The exterior region (referred to as the "*elsewhere*") of the light cone represents all the events which cannot in principle be connected with the apex via causal interactions, that is such a point can neither cause nor be an effect of an event at the apex. If we observe by light the world around us, the more distant an object the older it is, since light takes more time to reach us from it. The whole set of these objects (let them fill space continuously) form the hypersurface of the light cone with its apex at the observation world point (our here and now). This intuitive explanation is nevertheless rigorous, and the light cone is a central notion in our picture of relativistic space-time.

We shall now demonstrate how tensor and vector expressions of the important equations of special relativity are the

same. Thus the mechanical motion equation can be presented in the following tensor form as

$$dp^\alpha/ds = F^\alpha, \quad (1.10)$$

where  $p^\alpha = mu^\alpha$  is the momentum of a particle,  $u^\alpha = dx^\alpha/ds$  are the components of the 4-velocity,  $F^\alpha$  are the components of the force (when transforming from one frame to another the force is changed by a vector transform). For purely mechanical forces we have  $F^\alpha u_\alpha = 0$ . The four-dimensional equation (1.10) can be split into a spatial three-dimensional part

$$\frac{1}{c^2 \sqrt{1-v^2/c^2}} \frac{dp^i}{dt} = F^i, \quad (1.11)$$

where

$$p^i = m_0 v^i / \sqrt{1-v^2/c^2}; \quad v^i = \frac{dx^i}{dt},$$

and a time part

$$\frac{1}{c^2 \sqrt{1-v^2/c^2}} \frac{dE}{dt} = F, \quad (1.12)$$

where

$$E = m_0 / \sqrt{1-v^2/c^2}; \quad F = F_0$$

(here the superscript  $i$  assumes three dimensions: 1, 2 or 3). The spatial part (1.11) can be presented in vector form, i.e.

$$\frac{1}{c^2 \sqrt{1-v^2/c^2}} \frac{d\mathbf{p}}{dt} = \mathbf{F}. \quad (1.13)$$

It is a relativistic generalization of Newton's equations of motion. Equation (1.12) essentially describes the kinematic energy variation of a particle.

It will be recalled that Maxwell's field equations can be given in vector form

$$\text{div } \mathbf{H} = 0, \quad \text{curl } \mathbf{E} = -\frac{1}{c} \frac{\partial \mathbf{H}}{\partial t}, \quad (1.14)$$

$$\text{div } \mathbf{E} = 4\pi\rho, \quad \text{curl } \mathbf{H} = \frac{1}{c} \frac{\partial \mathbf{E}}{\partial t} + \frac{4\pi}{c} \mathbf{j}, \quad (1.15)$$

where  $\mathbf{E}$  and  $\mathbf{H}$  are the strengths of the electric and magnetic fields, respectively,  $\rho$  and  $\mathbf{j}$  are the charge density and the

vector of the current density of the field source. By introducing a tensor for the electromagnetic field

$$F_{\mu\nu} = \begin{bmatrix} 0 & E_x & E_y & E_z \\ -E_x & 0 & -H_z & H_y \\ -E_y & H_z & 0 & -H_x \\ -E_z & -H_y & H_x & 0 \end{bmatrix}, \quad (1.16)$$

we can rewrite the first (1.14) and second (1.15) pairs of Maxwell's equations in tensor form

$$\frac{\partial F_{\mu\nu}}{\partial x^\alpha} + \frac{\partial F_{\nu\alpha}}{\partial x^\mu} + \frac{\partial F_{\alpha\mu}}{\partial x^\nu} = 0, \quad (1.17)$$

$$\frac{\partial F^{\mu\nu}}{\partial x^\nu} = -\frac{4\pi}{c} j^\mu. \quad (1.18)$$

Finally, by using the definition of the tensor of electromagnetic field, the equation of a moving charge  $q$  with a mass  $m$  in Minkowski's space-time (in Cartesian coordinates) could be given in tensor form as follows:

$$\frac{d^2 x^\mu}{ds^2} = \frac{q}{mc^2} F^\mu{}_\nu \frac{dx^\nu}{ds}. \quad (1.19)$$

## 1.6. THE CREATION OF GENERAL RELATIVITY (THE GEOMETRIZATION OF GRAVITATIONAL INTERACTION)

We have now arrived at the point of birth of the general theory of relativity. The emergence of the general theory was a painful process and took a whole decade, from 1907 to 1916. This decisive movement can be said to have started in section five of Einstein's article *Über das Relativitätsprinzip und die aus demselben gezogenen Folgerungen* (About the Principle of Relativity and Its Consequences)\* [27] in which he made a clearcut relationship between gravitation and the principle of equivalence and accelerated frames. We cannot say that other scientists had ignored this problem for, as we have mentioned, a number of prominent physicists, particularly from the group connected with the creation of the special theory of relativity, had started working on gravitation using relativity theory. A strange reversal of

---

\* For the life and work of Einstein see [46, 48, 59, 94].

roles then occurred in the history of physics between the theories of gravitation and electricity. Initially Newton's law of gravitation together with Poisson's equation for gravitational potential, which was first derived by **Pierre Laplace** (1749-1827) in 1782 for a region outside a source, had been the model for the development of electrostatics. However, the situation was opposite in the 1900s in that the theory of electromagnetism was used as the basis for a consistent relativistic theory for gravitation. We shall not discuss here this interesting stage in the scientific process because the reader can find enough information about it elsewhere (see [80, 106]).

We quoted Einstein's article here for two reasons. Firstly, it was Einstein who succeeded in bringing this line of research to a logical end. Secondly and on top of that, he united gravitation theory and the principles of relativity and equivalence. This was the approach of a physicist who had a sharp critical mind who was working on the most topical questions of his time. However, it took another five years after that article to find an adequate mathematical framework. At the end of this five-year period Einstein collaborated with **Marcel Grossmann** (1878-1936), someone he had known since they were students together. Their joint article, published in 1913, was the first to introduce the idea of the geometrization of gravitation. It was not, however, a final result. In *Entwurf einer verallgemeinerten Relativitätstheorie und einer Theorie der Gravitation* (The Draft of the Generalized Theory of Relativity and of the Theory of Gravitation) [33] they did not speak directly of Riemannian geometry or the curvature of space-time, but only mentioned the metric tensor. But they did discuss the main thing, that is the relationship between the new theory and the earlier geometries of Lobachevski, Riemann and Clifford.

Let us look at what was stated in this article from the viewpoint of Einstein and his contemporaries, and then, from the viewpoint of today's understanding of his ideas. The special relativity principle states that physical phenomena occur in the same way in all inertial frames. In other words, if in an arbitrary inertial frame we perform a physical experiment, its outcome will be the same (not only qualitatively, but also quantitatively) as if the same experiment had been performed in any other inertial frame. In the first place Einstein was worried by the limitation of

the principle to inertial frames even though there were reasonable indications that in Minkowski's world the noninertial frames possessed spatial and temporal inhomogeneities that show up as inertial forces, that depend on the specific characteristics of the reference frame. Obviously, the inertial forces have to have a noticeable effect on the physical processes in these reference frames. The problem was in fact whether these inhomogeneities could be accounted or compensated for so that all the frames could in a certain sense be made equivalent. This compensation should depend on the here and now in space-time because the properties of the reference frame were not uniform. This would be equivalent to the introduction of a field, but what sort of field? It would have to be indifferent to the nature of the objects it "affected", so that any mass in the field would be subjected to the same acceleration as every other mass in the field. This is the same as saying that in Newton's second law, the force on the right-hand side is proportional to the mass of the body it acts upon, and is cancelled together with the mass in the linear momentum term on the left-hand side. But if we identify this new field with gravitational, this is just a formulation of the good old principle of equivalence! Einstein always considered the principle of equivalence to be an obvious fact, and there was a reference to L. Eötvös's (1848-1919) famous experiments that verified it, in his joint article with Grossmann. We shall discuss this experiment in more detail at the beginning of Chapter Two.

From the viewpoint of today's understanding, on the other hand, the situation looks different. It has been known since Galileo's time that all bodies in the gravitational field of the Earth have the same acceleration, no matter what their individual properties (e.g. mass, substance, shape) are. Consequently, their acceleration depends only on the point in space where they happen to be. Can we, therefore, attribute the gravitational characteristics (acceleration) to the points in space, where the bodies are, rather than to the bodies themselves? However, Minkowski's flat space-time doesn't have the properties needed to implement this idea: it is homogeneous, that is, everywhere uniform and isotropic (the same in all directions). This means that the components<sup>0</sup> of the metric (the metric tensor)  $g_{\alpha\beta}$  in (1.4) are constant (their individual moduli are either zero or unity). Conse-

quently, we need a space-time whose metric tensor has components  $g_{\alpha\beta}(x)$  that change from point to point, i.e. the space-time should be curved. This enables us to consider geometrical properties of space-time that change at different points. The next problem is to determine the specific nature of the relationship between the values of the components  $g_{\alpha\beta}(x)$  and the properties of gravitational interactions. This was the task that Einstein and Grossmann formulated and began work on in their article. In the section "Physics", which Einstein wrote, he stated: "Thus we come to the conclusion that in a general case the gravitational field is characterized by ten space-time functions" [33]. These ten functions replace Newton's single gravitational potential  $\Phi_N$ .

In such a theory as special relativity the most important role is played by the (squared) interval. It is easy to see why this should not be taken as the interval between two arbitrarily selected events, as in (1.3), but as that between two infinitesimally close ones. Hence we have

$$ds^2 = g_{\alpha\beta} dx^\alpha dx^\beta, \quad (1.20)$$

where  $dx^\alpha$  are differentials of the coordinates of the neighbouring points. It is also clear that the metric tensor is symmetric,  $g_{\alpha\beta} = g_{\beta\alpha}$ , that is, in a general case, only ten of its components are independent. These components are the main "bricks" for building the general relativity mathematics.

One may ask why the interval is given in terms of a square. This is mainly due to the symmetry properties of the interval with respect to the direction between two adjacent points ( $AB$  is equivalent to  $BA$ ). In fact any even power with respect to  $dx$  would be suitable, but the quadratic form is the simplest. Perhaps, in the future, this simple form will also no longer be adequate. The possibilities of such generalizations were examined by Riemann himself. This topic is briefly discussed in Sec. 3.3.3.

How do bodies move in a curved space-time? It was quickly realised that test bodies (those with small masses) move along geodesics in curved space-time. We discussed geodesics when we talked about the Lobachevski and Riemann spaces (Secs. 1.3 and 1.4). These geodesics were discovered by mathematicians long time ago. In order to obtain a geodesic the extremal path between two points must be found by setting the variation of the path between the two points equal to



zero (if the ends of the path are fixed):

$$\delta \int_A^B ds = 0. \quad (1.21)$$

By substituting  $ds$  from (1.20), doing some straightforward manipulations, we obtain the equation of a geodesic, i.e.

$$\frac{d^2 x^\mu}{ds^2} = -\Gamma_{\alpha\beta}^\mu \frac{dx^\alpha}{ds} \frac{dx^\beta}{ds}, \quad (1.22)$$

where  $\Gamma_{\alpha\beta}^\mu$  is called the Christoffel symbols and denotes a combination of the first derivatives of the metric tensor  $g_{\alpha\beta}$

$$\Gamma_{\alpha\beta}^\mu = \frac{1}{2} g^{\mu\nu} \left( \frac{\partial g_{\alpha\nu}}{\partial x^\beta} + \frac{\partial g_{\beta\nu}}{\partial x^\alpha} - \frac{\partial g_{\alpha\beta}}{\partial x^\nu} \right). \quad (1.23)$$

Equation (1.22) is essentially the system of equations that describe the motion of test bodies in a curved space-time, or as we can say now, the equations of motion of a body in a gravitational field.

Let us look more closely at this equation. Firstly, it resembles the relativistic equation of motion of a charged particle in an electromagnetic field (1.19). On the left-hand side we have the "acceleration"  $d^2 x^\mu/ds^2$ . The first difference between the two equations is that in (1.22) two 4-velocities appear on the right-hand side. The Christoffel symbol has the sense of the intensity of the gravitational field since it is made up of gravitational "potentials"  $g_{\alpha\beta}$  in much the same way as the electromagnetic field  $F_{\mu\nu}$  is made up of the  $A$ -potentials  $A_\mu$ , i.e.  $F_{\mu\nu} = \partial A_\nu/\partial x^\mu - \partial A_\mu/\partial x^\nu$ . Secondly, in (1.22) the mass of the particle is not present, which means that the gravitational and inertial masses automatically cancel out. This is actually a realization of the principle of equivalence. Thirdly, the equation of a geodesic (1.22) yields as a basic approximation the radial equation for a particle moving in a central Newtonian gravitational field, i.e.  $\ddot{r} = -GM/r^2$ , if we use spherical coordinates and assume that

$$g_{00} = 1 - \frac{2GM}{c^2 r} \equiv 1 + \frac{2\Phi_N}{c^2} \quad (1.24)$$

( $\Phi_N$  is the Newtonian gravitational potential), the remaining components of the metric being the same as for Minkow-

ski's space-time. This is the first demonstration of how gravitation is related to curved space-time. If a particle is affected by nongravitational forces, they are accounted for by adding the corresponding vector expressions to the right-hand side of (1.22). For instance, this may be an electromagnetic force, as occurs on the right-hand side of equation (1.19). If there are no such forces, the particle is said to be moving freely, i.e. along a geodesic.

Several scientists, such as the outstanding physicists and mathematicians **Vladimir A. Fock** (1898-1974) and **John L. Synge** (b. 1897), have rejected the idea that gravitation and acceleration are equivalent [37, 101]. They are right in a sense. In the general relativity theory as formulated in its four-dimensional form and without any separation between physical space and time in reference frames, a test particle which has no structure moves along a geodesic. That is the relevant equation does not contain a mass factor, and so there is no need to discuss the equivalence of the two masses. At the same time, experiments conducted in the USSR and USA in the 1960s have been widely accepted and acclaimed as verifications of the principle of equivalence (see Sec. 2.1). These experiments indicate that the limits in which general relativity holds are very wide.

## 1.7. SOME TYPICAL FEATURES AND PROPERTIES OF GENERAL RELATIVITY

So far, we had not touched upon the equations of a gravitational field. It took almost three years more to complete the work Einstein and Grossmann started in 1913 and to derive these equations. Two papers were then published concurrently, after which general relativity became de facto. These were *Die Grundlagen der Physik* by **David Hilbert** (1862-1943) [47], and *Die Grundlagen der allgemeinen Relativitätstheorie* by Einstein [29].\* Hilbert's paper did not attract much attention of physicists because of its advanced mathematics, or perhaps because it concentrated on the theory of G. Mie, which quickly lost favour. **Gustav Mie** (1868-1957) started the series of attempts to create a unified field theory. It is also clear that it did not appeal to mathe-

---

\* In fact, general relativity had been finally though briefly formulated by Einstein in the autumn of 1915 in *Die Feldgleichungen der Gravitation*. This paper was preceded by a few days by Hilbert's.

maticians because the work was too physics-oriented (oddly enough, the giants of mathematics, such as Lobachevski, Gauss, Riemann and Clifford, had as much eagerness for physics as Hilbert). Hilbert believed that general physical principles could only be logically formulated in mathematical terms, whereas even the transition from the principle of relativity to that of variance of action seemed formal and unclear to the majority of physicists of that time. The physics of that period was mostly illustrative, and theoretical physicists had not assimilated the theories of groups, matrices, tensors and spinors. Maybe Hilbert's paper is clear and understandable to us now because we begin using its "language" (though with some technical simplification) from the second or third year in university. Einstein, in his article, not only substantiated and explained the new theory, he also taught his readers the ABC of geometrical language (Hilbert just used it). But even after being put on the right track, the majority of his contemporaries found it too difficult to follow, and for many years physicists questioned general relativity because they misunderstood its basic elements.

What are the equations for gravitational fields? We are not going to derive them here for that would be beyond the scope of this book. In addition, the basic equations of physics, such as the Maxwell, Einstein, Klein-Fock-Gordon, and Dirac equations, were not derived, they were rather discovered. Whenever we speak about the derivation of an equation we really mean either the introduction of a postulate or assumption that either is equivalent to the equation or leads to it, or we mean the formulation of propositions that explain the equation or make it more comprehensible. While Hilbert might have used the first approach, we shall follow the second.

To begin with, let us have a look at the 'building blocks' we have available. There are the ten components of the metric tensor  $g_{\alpha\beta}$  (and the contravariant components  $g^{\alpha\beta}$  that follow unambiguously from them) and the forty components of the first derivatives  $\partial g_{\alpha\beta}/\partial x^\mu$ , the latter could be combined more conveniently into the forty components of the Christoffel symbols,  $\Gamma_{\alpha\beta}^\mu = \Gamma_{\beta\alpha}^\mu$ , to which they are equivalent. But this is not all. The laws of nature have the feature that practically all basic equations (say, in physics) contain second derivatives. This is the case with the equations

describing mechanical motion: on the left-hand side we have acceleration ( $d^2x/dt^2$  or  $d^2x^\mu/ds^2$ ). Similarly, Maxwell's equations (1.18) are written using the first derivatives of  $F_{\mu\nu}$ , which is itself the first derivative of the vector potential  $A_\mu$ . Consequently, it is only natural to have in the equation of gravitational field the second derivatives of  $g_{\alpha\beta}$ . But how should we arrange them? Obviously, these equations should be tensor equations, and, therefore, it is important to know which tensors can be constructed from  $g_{\alpha\beta}$  and their derivatives.

The simplest tensor that can be constructed from the components of the metric tensor and its first and second derivatives, is the tensor of rank four,  $R^\lambda_{\alpha\nu\mu}$ , that is the curvature tensor or Riemann-Christoffel tensor, which plays an important role in differential geometry. For a detailed information about the properties of the curvature tensor, we refer the reader to textbooks on Riemannian geometry or on general relativity (see, for example, [25, 69, 84]). This tensor is

$$R^\lambda_{\alpha\nu\mu} = \frac{\partial \Gamma^\lambda_{\alpha\mu}}{\partial x^\nu} - \frac{\partial \Gamma^\lambda_{\alpha\nu}}{\partial x^\mu} + \Gamma^\sigma_{\alpha\mu} \Gamma^\lambda_{\sigma\nu} - \Gamma^\sigma_{\alpha\nu} \Gamma^\lambda_{\sigma\mu}. \quad (1.25)$$

It defines the change in a tensor quantity in a parallel transfer along a closed contour, that is a new property of curved manifold. From  $R^\lambda_{\alpha\mu\nu}$  one can construct tensors of lower rank,

$$R_{\alpha\mu} = R^\lambda_{\alpha\lambda\mu}; \quad R = R_{\alpha\mu} g^{\alpha\mu}$$

which are the Ricci tensor and scalar curvature, respectively.

Now, we can write and discuss Einstein's equations. The aim of these equations is to establish a relationship between the geometrical characteristics of space-time and the physical properties of matter which cause the curvature. It was not easy to select specific quantities which were to be connected in the equations in question. Firstly, there were several trial variants of the theory. Finally, Hilbert and Einstein found that the gravitational field equations must have the form:

$$R_{\mu\nu} - \frac{1}{2} g_{\mu\nu} R = \kappa T_{\mu\nu}, \quad (1.26)$$

where  $T_{\mu\nu}$  is the tensor for the energy-momentum of matter (for dust,  $T_{\mu\nu} = \rho c^2 u_\mu u_\nu$ , where  $\rho$  is the density of the dust). The  $T_{\mu\nu}$  term is multiplied by a dimensional coefficient  $\kappa = 8\pi G/c^4$ , which is now called Einstein's gravitational constant. All the textbooks on gravitational theory show how the Poisson equation,

$$\frac{\partial^2 \Phi_N}{\partial x^2} + \frac{\partial^2 \Phi_N}{\partial y^2} + \frac{\partial^2 \Phi_N}{\partial z^2} = 4\pi G\rho, \quad (1.27)$$

where  $\Phi_N$  is related to  $g_{00}$  as in (1.24), follows from Einstein's equations in the Newtonian approximation.

Thus, the curvature of the space-time has a physical meaning. We have long got used to the fact that many things in this world are different aspects of the same entity. This was the case with mass and energy, with inertial mass and gravitational mass, and now we see that the same is true of gravitational fields and the geometrical properties of space-time. Curvature is a geometrical entity whose components vary from point to point in any given system of coordinates. Thus curvature is a field, in the mathematical sense of the word, and at the same time a manifestation of a physical gravitational field. The credit for this discovery must go to Einstein.

Now, we shall discuss some of the assertions of general relativity. In the macroscopic world, gravitation is undoubtedly the most universal kind of interaction. It is universal because all macroscopic objects in physics generate a gravitational field (i.e. they form  $T_{\mu\nu}$ ) and interact this way (generally speaking they can also interact in other ways if they have charges). Gravitational charge (we mean here mass-energy and not simply rest mass) is common to all objects, at least macroscopic ones.

The most general form of existence of matter is thus a four-dimensional space-time. The word "form" is used here in the sense that if we do not consider the properties of space-time, the configuration and evolution of a specific type of matter cannot be described. But gravitation is characterized by the curvature of space-time and we could consider it natural that the universality of gravitational interaction and the universality of space-time as a form of the existence of matter are simply two sides of the same coin\*.

---

\* A discussion of the philosophical aspects of general relativity can be found in [1, 19, 38, 43, 55].

How does general relativity add to our understanding of a gravitational field? For Newton, gravitation was the instantaneous effect of two, or more, bodies interacting at a distance. He described it using a one-component potential whose gradient was proportional to a force. In electrostatics, this description fully reflects the behaviour of the electric (scalar) potential. However, we know that, on the one hand, the potential is only produced by a moving charged body at a distance after a delay ("retarded potentials"). This must also be reflected by the theory of gravitation, otherwise the principle of causality is broken, and information could be transmitted via gravitational interaction faster than the speed of light. On the other hand, if we measure the fields while moving in a purely electric field, we will find that a magnetic field is also present. In other words, it is only a combination of electric and magnetic fields, i.e. an electromagnetic field, that can occur. We may also say that the magnetic field is the relativistic effect of the electric field, i.e. it is caused by the observer's motion with respect to it. A similar effect should be present when an observer moves in a gravistatic field, and a quasi-magnetic gravitational field must exist as an independent reality in situations similar to those in electrodynamics when there are fields of mostly magnetic nature. These two features of the gravitational field in general relativity are implied by Einstein's equations (1.26). However, there is a new feature in them: the tensor gravitational field and this has no electromagnetic analogue. This new problem can be strictly solved by applying the techniques for describing reference frames (see Secs. 2.8 and 2.9). In the general case, the motion of a reference frame is described by three tensorial quantities, viz. the acceleration vector  $G_\mu$ , the antisymmetric tensor of angular velocity  $A_{\mu\nu}$ , and the rate-of-strain tensor  $D'_{\mu\nu}$ . It can be shown that in a number of situations, the acceleration  $G_\mu$  is an analogue to the vector describing the electric field, whereas the angular velocity tensor corresponds to the magnetic field. The rate-of-strain tensor for a reference frame (it is symmetric, i.e.  $D'_{\mu\nu} = D'_{\nu\mu}$ ) has no analogue in the theory of electromagnetic fields. In other cases, the  $D_{\mu\nu}$  tensor acts like electric field. Notice that some components of the curvature tensor are analogous to the electric and magnetic fields (see Sec. 3.1). However, the curvature tensor has more components than the electromagnetic field tensor.

Hence we see that although gravitation can be compared to electromagnetism, it has its own features. Each physical field has its own specific features, however, all of them are governed by a number of similar principles.

The nonlinearity of the gravitational field is understandable and is formally expressed in Einstein's equations as nonlinear combinations of potentials and their first derivatives. This nonlinearity arises in Einstein's equations because of the universality of the gravitational interaction. The gravitational field may possess energy and momentum (this problem is currently under discussion and still has no unambiguous and generally accepted solution). If so, energy and momentum must act as the sources of some secondary gravitational field, which, in turn, will produce energy and momentum, which again create a gravitational field, and so on and so forth! The nonlinearity of Einstein's equations is, in fact, a result of this sort of "self-generation": the simple superposition of two or more solutions does not produce a new solution. On the contrary, Maxwell's equations of a free electromagnetic field are linear and their solutions are in agreement with the superposition principle. The nonlinearity of the gravitational field equations, even though this property is natural and understandable, makes them very difficult to solve and difficult to construct a simple theory for gravitational field energy.

Gravitational interactions have been geometrized in general relativity. But how can we say where the boundary between physics and geometry lies? We can answer in several ways. Firstly, instead of seeking a direct answer we can put a counter-question, viz. must such a boundary exist, and hence is there any point in talking about any difference between physics and geometry? The only reason we regard geometry as an independent of physical reality is because we were educated in the spirit of Euclidean axioms. Think how absurd this viewpoint is: indeed, our Universe is whole and unique, and its properties can all be learned only from experience. The geometry of the real world should be regarded as one of the experimental physical sciences and its assertions and predictions must always be verified through experience or practice. The axioms of geometry are, in fact, no more than a refined form of human experience!

Secondly, in modern theoretical work we operate with two sorts of quantities. Some describe gravitational interac-

tions and have been geometrically interpreted in the framework of general relativity, while others, such as the electromagnetic field tensor or scalar field functions, are considered external to geometry. Our question can thus be put this way: Is it possible to construct a theory in which all the quantities have a geometrical interpretation? After this question was first asked by Clifford, various answers were offered by Weyl, Eddington and others. We shall return to these ideas later, in Chapter Three, but to anticipate matters we can say that the most elegant solution lies in the framework of higher-dimension geometry (see Sec. 3.5).

Finally, it is worth stressing that general relativity was not, even at the moment of its conception by Einstein and Hilbert, solely a theory of gravitation. It is valid for all fields and matter in curved space-time whose equations can be written in a covariant form.

Some of the other features of general relativity are discussed from various viewpoints elsewhere [6, 10, 44, 73, 78, 92, 103, 105].



## TODAY

---

### A Review of the Basic Results of Modern Relativity Theory

In the previous chapter we mainly discussed the development of our understanding of space and time. To highlight its present state, we shall now review the main results of general relativity (physical effects including). Although this book is classified as "popular science", only Chapter One was really written for an absolute layman. Chapter Three is devoted to more delicate scientific problems, while Chapter Two is intermediate between a popularizing style and a scientific presentation. We have included in this chapter both a qualitative description of the theory and experiments and a mathematical description (fairly strict while not very complicated).

#### 2.1. THE PRINCIPLE OF EQUIVALENCE AND GRAVITATIONAL REDSHIFT

We now return to the principle of equivalence which we discussed in Sec. 1.6 with respect to Einstein's general relativity. Firstly, remember that we are talking about the identity of the behaviours of test bodies (irrespective of their nature) when they are in similar conditions. Normally, test bodies are only assigned mass and no account is made of their angular momenta or their more sophisticated analogues (higher moments). According to the principle of equivalence, the actual mass is not important since it cancels out from both sides of the equations of motion (therefore, in the equation for a geodesic (1.22) the mass is not present). This approach is appropriate because the higher moments of test bodies under normal conditions make very insignificant contributions to their motion. In addition, the idea behind

the principle of equivalence is that we can analyze local effects, while the higher moments, from angular momentum upwards, are affected by the inhomogeneities in the fields (curvature and its derivatives) in the same way as the electric and magnetic moments are affected by inhomogeneities of the electromagnetic field in conventional electrodynamics. Hence, when we talk about the equivalence principle, we shall confine ourselves to comparisons of the inertial and gravitating (heavy) masses.

Since the inertial or gravitating masses of a body have the same absolute value or are related by a universal proportionality factor, it is only logical to include the latter in the gravitational constant in the equations of motion. This amounts to the assumption that the numerical values of the inertial and gravitating masses are the same. The principle of equivalence could thus be reduced to the equality of the inertial and gravitating masses (this is sometimes referred to as the weak equivalence principle) [10, 15, 21, 69]. We could not have anticipated this equality between the two masses. Our psychological readiness to accept it is due to the similarity of the words used rather than the experience: physically the inertial and gravitating masses of a body have as much in common as, say, mass and energy.

The first observations of equivalence (it was Einstein who actually raised it to a principle) were carried out by Galileo, who noted, as the story goes, the periods of the natural oscillations of the chandeliers in a cathedral, and the similarity between the free fall of balls made from different materials. The first quantitative verification of the principle was made by Newton who compared the oscillations of two different pendulums. The pendulum bobs were made from different materials but were equal in weight. They were set in the centres of gravity of equally sized wooden boxes to compensate for the difference in their aerodynamic drags. Then Newton studied the phases of oscillations with the same amplitude. The experimental technique was perfected by R. von Eötvös (1886), R. H. Dicke (1964) and finally, V. B. Braginsky (1971). R. von Eötvös attained a remarkable precision when he verified the equality of the inertial and gravitating masses: the difference was one part in  $10^8$ . Dicke experimented with gold and aluminium and reached three parts in  $10^{11}$  and Braginsky for platinum and aluminium reached 9 in  $10^{12}$ .

In the later of these experiments two balls made from different materials were suspended on a torsion balance. At equilibrium the following forces were in balance: the torque of the suspension filament, the moments due to the gravitational forces of the Earth and the Sun (the latter's role was the main), and the moment of the inertial forces that act on the weights as the apparatus rotates together with the Earth. Both the inertial and gravitating masses contribute to the moments of the gravitational and inertial forces. At different hours of the day the moments combine differently, but they must balance out for the two weights (provided that the torque on the suspension filament is constant) if the inertial and gravitating masses are equivalent. Otherwise, the torsion balance would oscillate with a 24-hour period and this can be detected electronically. No such oscillations have been recorded for any experimental technique such that a significant signal can be separated from the noise [11, 15, 21].

How can the results of these experiments, which rank among the most delicate in modern physics, be interpreted? They demonstrated the universality of all types of energy: the energy of any origin is equivalent to both inertial and gravitating masses. Indeed, materials (such as platinum and gold, aluminium or paraffin) are different not only in terms of their macroscopic characteristics—specific and molecular weight, etc.—but also in terms of their atomic and nuclear properties. All macroscopic properties are averaged manifestations of the fundamental interactions of physics, i.e. the electromagnetic, weak, strong, and (to a very small extent) gravitational interactions. The atoms of different chemical elements differ in the number of protons and neutrons in their nuclei and the number of electrons in the electron shells that surround the nuclei. To these differences we add the electromagnetic and nuclear fields that stabilize nuclei, atoms and molecules. These fields realise the interactions between the components of the nuclei and atoms, and as a whole lead to a mass defect for the composite systems, and in this way their stability (in the sense that this given state is energetically more favourable) is ensured. Hence, the complete mass (energy) of such a system consists of the rest mass of its components, the kinetic energy of the components, and the energy of the fields that bind them together. The experimental proof of the principle of equiv-

klence attests that the rest masses of any particles, their kinetic energies, and the energies of all the physical fields which ensure their interaction are equivalent and behave alike.

Since the various contributions to the total mass of a body differ by many orders of magnitude, the accuracies with which the principle of equivalence can be verified using the experiments of Eötvös, Dicke, and Braginsky are different. The experiments were the most inaccurate when estimating the equivalence of the gravitational field energy (the contribution of which to the mass defect is negligibly small under laboratory conditions). However, the gravitational mass defect rises (absolutely and relatively) as the total mass of the objects rises. For example, the relative shares of the strong, electromagnetic, weak, and gravitational interactions in the energy of an average atom are  $1 : 10^{-2} : 10^{-12} : 10^{-40}$ , for a laboratory macroscopic weight they are  $1 : 10^{-2} : 10^{-12} : 10^{-29}$ , and for an average star (the Sun, for example)  $1 : 10^{-2} : 10^{-12} : 10^{-6}$  [15]. It is clear that the gravitational share rises with mass. Indeed, gravitation is not important on the nuclear or atomic scale, but it becomes dominant on macroscopic and astronomical scales. This is because the other interactions have either a short range of action (nuclear interactions, for example) or rely on the presence of opposite charges which cancel in large bodies. Gravitational "charges", or masses, on the other hand, have the same sign for all particles and are simply added; while the gravitational field itself has a large range of action. This means that the principle of equivalence for gravitational energy should be verified on astronomical objects (though they should be small in comparison with the characteristic length of gravitational inhomogeneities). Appropriate objects are the massive planets, whose motion can be predicted by the laws of celestial mechanics (that is, Newton's theory, which is the nonrelativistic limit of general relativity). These predictions are equivalent to the experiments of Eötvös, Dicke, and Braginsky (see [15]).

Let us discuss now gravitational redshift, which is directly related to the principle of equivalence. Photons with different frequencies possess different energies (mass due to motion as opposed to rest mass, which is zero for photons), but their world lines are the same, all things being equal. This is, however, not the only manifestation of the prin-

ciple, another is that to conserve their total energy in potential gravitational fields their kinetic energies fall (since photons have no rest energy!) when they move away from a gravitating body (i.e. leave a gravitational potential well). The energy (and frequency) of the photons thus decreases and a redshift is observed. In general this effect is included in and inseparable from the Doppler effect.

Here, we shall show a simplified derivation of the redshift effect in time-independent weak fields. We begin with assuming the correspondence of Einstein's and Newton's theories at the level of the equations of motion. As the gravitational field is weak in this case, we can use a universal, almost Cartesian system of coordinates in which the axes are always straight and mutually orthogonal, i.e. the metric tensor has the form of (1.4). If there is, however, even weak gravitation, the Cartesian frame is only an approximation, and the true metric tensor becomes

$$g_{\mu\nu} = g_{\mu\nu}^0 + h_{\mu\nu}, \quad (2.1)$$

where  $h_{\mu\nu}$  is small,  $|h_{\mu\nu}| \ll 1$ . For a slowly moving test mass  $|dx^i/ds| \ll 1$  (note that the Latin indices run over the purely spatial values 1, 2, 3), we may suppose that  $dx^0/ds \simeq 1$ , where  $x^0 = ct$ . Then, from the spatial part of the geodesic equation (1.22), which can be more conveniently rewritten as

$$\frac{d}{ds} \left( g_{\mu\nu} \frac{dx^\nu}{ds} \right) = -\frac{1}{2} \frac{\partial g_{\kappa\lambda}}{\partial x^\mu} \frac{dx^\kappa}{ds} \frac{dx^\lambda}{ds}, \quad (2.2)$$

follows that

$$\frac{d^2 x^i}{dt^2} \simeq -\frac{c^2}{2} \frac{\partial g_{00}}{\partial x^i}. \quad (2.3)$$

In the standard vector form, we have

$$\frac{d^2 \mathbf{r}}{dt^2} \simeq -\frac{c^2}{2} \text{grad } g_{00}. \quad (2.4)$$

By comparing this equation with Newton's law of motion for a test mass in a gravitational field

$$d^2 \mathbf{r}/dt^2 = -\text{grad } \Phi \quad (2.5)$$

( $\Phi_N$  is the Newtonian gravitational potential; the inertial and gravitating masses have cancelled according to the prin-

ciple of equivalence), we have  $g_{00} \simeq 1 + 2\Phi_N/c^2$ , which is (1.24) again. Here the integration constant is taken to be unity so that the metric asymptotically (as  $r \rightarrow \infty$ ) approaches its Cartesian form (particularly, as  $g_{00} \rightarrow 1$ ). For a free photon  $E = h\nu$  ( $E$  is energy,  $\nu$  frequency, and  $h$  Planck's constant). If, however, the photon is in a gravitational field, its negative potential energy (the mass defect) should be added):  $E = h\nu + h\nu\Phi_N/c^2$ . This total energy is conserved when the photon moves from position 1 to position 2:

$$h\nu_1 + h\nu_1\Phi_N(1)/c^2 = h\nu_2 + h\nu_2\Phi_N(2)/c^2. \quad (2.6)$$

Hence

$$\frac{\nu_1}{\nu_2} = \frac{1 + \Phi_N(2)/c^2}{1 + \Phi_N(1)/c^2} \simeq \sqrt{\frac{g_{00}(2)}{g_{00}(1)}} \quad (2.7)$$

or

$$\frac{\nu_1 - \nu_2}{\nu_2} = \frac{\Phi_N(2) - \Phi_N(1)}{c^2 + \Phi_N(1)} \simeq \frac{\Phi_N(2)}{c^2} - \frac{\Phi_N(1)}{c^2}. \quad (2.8)$$

Thus, in conformity with the principle of equivalence, the relative gravitational shift of frequency is proportional to the difference between the gravitational potentials at the emission and absorption points of a photon. The redshift was then very accurately measured under terrestrial laboratory conditions (the astronomical effect is heavily masked by noise)\*. The experiments gave  $0.9990 \pm 0.0076$  of the figure predicted by the equivalence principle.

## 2.2. SCHWARZSCHILD'S SPACE-TIME

A static, spherically symmetric gravitational field in vacuum turned out to be one of the simplest from the viewpoint

---

\* The Mössbauer effect was used to obtain a beam of photons (gamma quanta) with very narrow range of frequencies. This could then be used to measure very accurately the variations in the frequency due to the photons rising about 20 m (or descending 20 m, in which case there would be a blueshift). The Mössbauer effect involves the distribution of recoil momentum among all the nuclei of a crystal lattice when a photon is emitted or absorbed. As a result, the recoil velocity of nucleus in the lattice due to an emitted or absorbed gamma quantum becomes negligibly small and the Doppler effect is eliminated as it would otherwise distort the results.

of calculation, but it was not easy to come to, as the history of science proves. This field is analogous, at least far away from its source, to Coulomb's electrostatic field, which may be considered as the simplest solution of Maxwell's equations. This gravitational field was named after **Karl Schwarzschild** (1873-1916), who found this solution to Einstein's equations with its zero right-hand side in 1915, almost immediately after the formulation of the general theory [91]. Since then Schwarzschild's solution has remained the most important example of a relativistic gravitational field, and not because we know no better (many other solutions have since been discovered), but because Schwarzschild space-time has so many implications for physics.

Here we shall discuss the most general characteristics of the Schwarzschild solution, but later (in Sec. 3.2) we shall return to it to consider some of its astronomical aspects. Typically, the solution is obtained by solving Einstein's equations using two assumptions: spherical symmetry and the absence of distributed sources. The equations then automatically imply that the gravitational field is static, so the system reduces to an ordinary differential equation. This can be easily solved, and the constants of integration are determined at infinity. This approach is somewhat complicated because the curvature components have to be calculated in a general form first and then substituted into the left-hand side of Einstein's equations. We cannot follow this path here, but we can give a semblance of the derivation given by **Arnold Sommerfeld** (1868-1951) [98], who took it from Lenz. A reader to whom our calculations seem difficult can pass directly to formula (2.15) and still be able to understand the remaining pages.

Assume that only one compact spherically symmetric mass exists in the Universe and that space-time is asymptotically characterized (at the spatial infinity) by the Minkowski metric, which in spherical coordinates can be given as

$$ds_{\infty}^2 = c^2 dt^2 - dr^2 - r^2 (d\theta^2 + \sin^2 \theta d\varphi^2). \quad (2.9)$$

We introduce four mutually orthogonal unit vectors (an orthonormal basis) parallel to the coordinate axes and located far from the mass, which is at rest at the origin. By designating the vectors by  $g^{(\alpha)}$ , where  $\alpha$  is the vector number ( $\alpha = 0$ ,

1, 2, 3), we have the following components:

$$\begin{aligned} g_{00}^{(0)} &= 1, & g_{11}^{(0)} &= g_{22}^{(0)} = g_{33}^{(0)} = 0; \\ g_{00}^{(1)} &= 0, & g_{11}^{(1)} &= 1, & g_{22}^{(1)} &= g_{33}^{(1)} = 0; \\ g_{00}^{(2)} &= g_{11}^{(2)} = 0, & g_{22}^{(2)} &= r, & g_{33}^{(2)} &= 0; \\ g_{00}^{(3)} &= g_{11}^{(3)} = g_{22}^{(3)} = 0, & g_{33}^{(3)} &= r \sin \theta. \end{aligned} \quad (2.10)$$

It is easier to operate on linear differential forms (covectors or 1-forms) rather than on separate components of these covariant vectors. These 1-forms are defined as

$$\omega^{(\alpha)} = g_{\mu}^{(\alpha)} dx^{\mu} \quad (2.11)$$

( $\mu = 0, 1, 2, 3$ , the summation is done over  $\mu$ ). In our case

$$\omega_{\infty}^{(0)} = c dt, \quad \omega_{\infty}^{(1)} = dr, \quad \omega_{\infty}^{(2)} = r d\theta, \quad \omega_{\infty}^{(3)} = r \sin \theta d\varphi, \quad (2.12)$$

so that for a distant observer, the coordinates  $t, r, \theta$  and  $\varphi$  have the sense of standard spherical coordinates (and can for instance be measured by standard methods). Suppose that this observer, his orthonormal basis, and all the instruments he needs, are inside a spacecraft that is freely falling toward the gravitating centre, the spherically symmetric mass. Sooner or later, as he continues to fall the gravitational field will be strong enough to matter and he will be travelling at a velocity  $v$  towards the origin. Of course, he will consider himself at rest, since according to the equivalence principle he will not feel any acceleration so long as the gravity field in his cabin remains relatively homogeneous. However, by measuring the length of a measuring rod and the speed of a clock which are both at rest relative to the central mass, he will notice that these outside standards will be gradually changing. According to the equivalence principle, he should ascribe the changes to the relative velocity of the measuring rod and the clock. We know from special relativity that for him the external clock (at rest relative to the central mass) should be slower and the length of the measuring rod will be shorter in the radial direction. The length of the rod perpendicular to his velocity will remain unchanged, of course. Since the time and space scales are essentially the basis for the frame relative to which the measurements are done the freely falling basis (2.12) carried by the observer from infinity is related to the basis of ref-



erence frame (the one the observer passes at a given instant) as follows:

$$\begin{aligned}\omega^{(0)} &= \sqrt{1 - v^2/c^2} \omega_{\infty}^{(0)}, & \omega^{(1)} &= \omega_{\infty}^{(1)} / \sqrt{1 - v^2/c^2}, \\ \omega^{(2)} &= \omega_{\infty}^{(2)}, & \omega^{(3)} &= \omega_{\infty}^{(3)}.\end{aligned}\quad (2.13)$$

The interval squared is invariant and can now be found in the static basis from (2.13) in the form

$$\begin{aligned}ds^2 &= \omega^{(0)}\omega^{(0)} - \omega^{(1)}\omega^{(1)} - \omega^{(2)}\omega^{(2)} - \omega^{(3)}\omega^{(3)} \\ &= (1 - v^2/c^2) \omega_{\infty}^{(0)}\omega_{\infty}^{(0)} - \frac{\omega_{\infty}^{(1)}\omega_{\infty}^{(1)}}{1 - v^2/c^2} - \omega_{\infty}^{(2)}\omega_{\infty}^{(2)} - \omega_{\infty}^{(3)}\omega_{\infty}^{(3)},\end{aligned}\quad (2.14)$$

assuming that each of the two bases is orthonormal in its own applicability region. The basis  $\omega^{(\alpha)}$  can be applied everywhere (at any  $r$ ) and turns into the basis  $\omega_{\infty}^{(\alpha)}$  at large distances (asymptotically) where the speed of the observer is zero. To formulate the metric finally, we must express  $v$  as a function of  $r$ . We shall do this by assuming that  $g_{00} = 1 + 2\Phi_N/c^2$  according to (1.24). Here again we shall apply the equivalence principle and assume the correspondence to Newtonian theory\*. Remembering that  $\omega_{\infty}^{(0)} = c dt$  and expressing the interval squared in terms of the coordinates used by a distant observer (which, generally speaking, change their sense at finite distances from the central mass), we have  $g_{00} = 1 - v^2/c^2$ , hence  $v^2 = -2\Phi_N$ . Assuming that the central mass is point-like, that is,  $\Phi_N = -Gm/r$ , we have

$$ds^2 = \left(1 - \frac{2Gm}{c^2 r}\right) c^2 dt^2 - \frac{dr^2}{1 - \frac{2Gm}{c^2 r}} - r^2 (d\theta^2 + \sin^2 \theta d\varphi^2). \quad (2.15)$$

---

\* How strong are the limitations that arise from assuming that the gravitational field is weak when the correspondence principle is to be used? We believe that there is no limitation in this case because our acceptance of  $g_{00} = 1 + 2\Phi_N/c^2$  before the metric has yet been determined means that we are using our freedom of choice to select the radial coordinate, and this is not a length itself but is associated with it through the metric. It does mean that the result may not represent an exact solution because we have only operated intuitively and without solving the field equations. However, a direct substitution of this solution into Einstein's equations confirms its correctness.

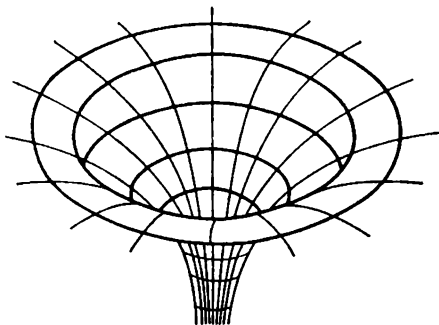


Fig. 3

This is in fact Schwarzschild's metric, the famous exact solution to Einstein's equations for a vacuum when the gravitational field vanishes at infinity and has spherical symmetry (this can be proved by substituting (2.15) into Einstein's equations). We should emphasize again that in this derivation we have only used the local validity of special relativity and the equivalence principle, and assumed that the theories of Einstein and Newton correspond.

Let us consider some of the basic properties of the space-time described by the metric in (2.15). Firstly, this metric turns asymptotically into the Minkowski metric (in spherical coordinates), that is, the world becomes flat at large distances from the origin. At finite distances, however, the space is not flat. It follows from (2.15), for example, that the circumferences of circles drawn around the origin increase with radius more slowly than they do in flat Euclidean space (Fig. 3.) Secondly, the components of the metric do not depend on the  $t$  coordinate. Can we unequivocally state that the metric is thus static? Yes, if  $t$  is physical time, that is, if the term with  $dt^2$  in (2.15) is positive, while the other terms remain negative. This is the case if  $r > 2Gm/c^2$ , and in this whole domain that extends to spatial infinity, the metric (2.15) actually describes a static spherically symmetric gravitational field. For  $r < 2Gm/c^2$ , however, the component  $g_{00}$  becomes negative, and  $g_{rr}$ , positive, so that in this domain, the role of time-like coordinate is played by  $r$ , whereas that of space-like coordinate, by  $t$ . Thus in this domain, the gravitational field depends significantly

on time ( $r$ ) and does not depend on the coordinate  $t$ . What happens on the boundary between the two domains? The surface  $r = 2Gm/c^2$  is special in the sense that the component  $g_{00}$  of the metric tensor vanishes on it, whereas  $g_{rr}$  becomes infinite. Another critical value is  $r = 0$ , which is not surprising in itself, because if we have a point source, there must always be a singularity at the origin as occurs, for example, in a Coulomb field in electrostatics. For gravitation, however, the singularity is peculiar in the sense that its location behaves differently from that of the Coulombian singularity which is simply at rest. It is not difficult to see that the Schwarzschild singularity at  $r = 0$  cannot stay at one place in principle and, moreover, it always moves faster than the light. This follows from the sign of  $ds^2$  for any coordinate differential taken along the world line of the singularity (it is more convenient to operate with points a little way off the world line, so as to work with finite elements). Any world line  $r = \text{const} < 2Gm/c^2$ ,  $d\theta = d\varphi = 0$  lies outside the light cone and, therefore, moves at a velocity greater than that of light. This statement is true at  $r = 0$ , too.

Let us get back to the "sphere"  $r = 2Gm/c^2$ . Inside the sphere the line  $r = \text{const}$  represents a superluminal velocity (the line is space-like,  $ds^2 < 0$ ); outside the sphere, this line is time-like (it represents normal subluminal velocity,  $ds^2 > 0$ ). This means that at  $r = 2Gm/c^2$ , the line is null ( $ds^2 = 0$ ): a point on a Schwarzschild sphere when  $d\theta = d\varphi = 0$  moves at the velocity of light. Thus, on this surface, we have at least one of the generatrices of the light cone. What happens to the other radial generatrix? The light cone is defined by  $ds^2 = 0$ , hence for the radial generatrices we have  $dr/dt = \pm(1 - 2Gm/c^2r)$ . If we have a sequence of light cones whose apexes lie progressively closer to the Schwarzschild sphere, we will see that they close up and each cone's generatrices merge as  $r \rightarrow 2Gm/c^2$ , so that the difference between all subluminal velocities and light velocities is lost. Does this have any sense? If so, one could only leave a Schwarzschild sphere in a space-like direction, that is, at a superluminal velocity, and arrive at it in the same way. All the other velocities (from zero to the light velocity) would be indistinguishable on this sphere. Why? Is this a law of nature, or a defect in the metric? It turns out to be a defect, not in the metric itself

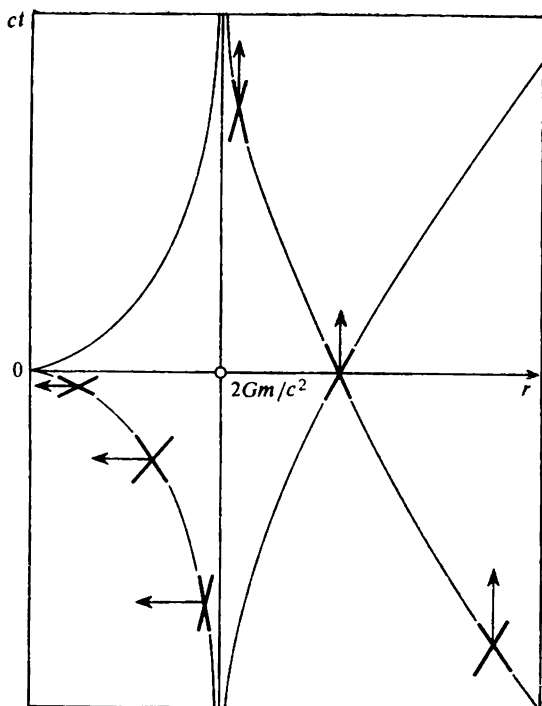


Fig. 4

but in the system of coordinates that we used. Remember what we meant by  $t$  and  $r$  at the very beginning. They meant conventional time and a conventional radial coordinate for an infinitely distant observer. The defect is due to the inapplicability of these coordinates on the Schwarzschild sphere. We cannot arrive at the sphere within a finite  $t$  value other than at a velocity above that of light, and a falling observer cannot attain this velocity. This rest mass is in principle nonzero, and so even to reach light velocity he would consume an infinite amount of energy. But any real energy source is limited. Even if the observer were falling from infinity with the maximum velocity for a real object (the velocity of light represents this limit), it would still not reach the Schwarzschild sphere within a finite time  $t$ . Indeed, the whole of a world line for a ray of light would

have to be on the sphere if even a single point of the line was on the sphere, and "arriving" light rays would have world lines that only asymptotically approach the Schwarzschild sphere (this discussion only concerns pure radial motion). Figure 4 shows local light cones whose apexes are located at various distances from  $r = 2Gm/c^2$  (the Schwarzschild sphere). The sphere itself is the vertical straight line representing  $r = 2Gm/c^2$  at all values of  $t$  ( $t$  is plotted along the ordinate axis). By accounting for the gradual increase of the generatrix slope as the light cone approaches the Schwarzschild sphere, we can easily draw lines of light escaping from the sphere to infinity, and falling on it from infinity. Thus our observer cannot reach the Schwarzschild sphere at any  $t$ , even less get inside. This means that the metric in (2.15) can only describe the world outside the sphere.

Let us now track our observer as he falls. He is falling freely, that is, he is moving along a geodesic, the equation of which has the form (2.2). This consists of four equations, two of which (corresponding to the angular components) are trivial for radial motion. Since the components of the metric do not depend on  $t$ , the time component of this geodesic equation ( $\mu = 0$ ) is

$$\frac{d}{ds} \left[ \left( 1 - \frac{2Gm}{c^2 r} \right) \frac{dt}{ds} \right] = 0. \quad (2.16)$$

This can be easily integrated and yields the following equality

$$c \left( 1 - \frac{2Gm}{c^2 r} \right) \frac{dt}{ds} = 1 \quad (2.17)$$

(this is the energy integral, and the integration constant is taken to be unity because the observer was at rest at  $r = \infty$ ). By expressing  $dt$  in terms of  $ds$ , we have after simple manipulations

$$ds = -\sqrt{c^2 r / 2Gm} dr, \quad (2.18)$$

where we choose the minus sign for the square root in order to describe the observer's fall (not escape). This equation can also be easily integrated, but an initial coordinate  $r_0$  (see Fig. 4) has to be used because it is impossible to get

anywhere from infinity at a finite velocity within a finite time lapse. As a result we have

$$\begin{aligned}\Delta s &= - \int_{r_0}^{2Gm/c^2} (2/3) \sqrt{c^2/2Gm} \, d(r^{3/2}) \\ &= \frac{r_0^{3/2}}{3} \sqrt{2c^2/Gm} - \frac{4Gm}{3c^2}.\end{aligned}\tag{2.19}$$

What does this mean? The proper time  $\Delta s/c$  that passed according to the falling observer's clock while he was freely falling from  $r_0$  to the Schwarzschild sphere turns out to be finite. At the same time, it follows from Fig. 4 that the time  $t$  when the observer arrives at the sphere having passed  $r_0$  at  $t = 0$  is infinitely long. The observer only reaches the sphere asymptotically. Hence no matter how long we prolong our study, we will never get to this point (along  $t$ ). This means that we can arrive at the Schwarzschild sphere, but to describe this trip we will need a coordinate other than  $t$ .

"Good" coordinates for the Schwarzschild field, which are valid both inside, outside, and on the sphere, exist and have been studied. In the following sections, however, we shall confine our discussion to phenomena that occur outside the Schwarzschild sphere where the metric (2.15) is valid.

### 2.3. PERIHELION ADVANCE AND THE SOLUTION OF MERCURY'S "ABNORMAL" PRECESSION

The laws which govern the motion of the planets were clearly explained by Newton's gravitational theory. This was a triumph for Newtonian mechanics, striking contemporaries by its universal applicability to phenomena both on Earth and in the heavens, which had been regarded as completely different to the "imperfect" bodies on Earth. Newton had brought astronomy "down to Earth", and every new fact concerning the behaviour of the planets and stars that his theory explained was a new stone in its foundation (Newton however undertook this "blasphemous" action with the utmost reluctance [86]). So great was the faith in the power of Newton's laws that astronomers interpreted visible deviations in the motions of the planets as indications of existence of unknown hitherto (unobserved) planets. For example, after studying deviations in the motion of Uranus **John Adams**

(1819-1892) in England and **Urbain Leverrier** (1811-1877) in France predicted in the 1840s the existence of Neptune and even computed its position. Adams began his calculations in 1843 and presented his results in October 1845 to **George Airy** (1801-1892), a British astronomer, who accepted them with some reluctance (no man is a prophet in his own country, and what is more, Adams was a student at the time!). Leverrier started his calculations, independently of Adams, in 1845 and was encouraged by **Dominique Arago** (1786-1853), a leading French astronomer. He finished by the summer of 1846, and his results came to the attention of British astronomers who then remembered their compatriot Adams, and began looking for the planet. But the first person to discover the new planet was **Halle**, an astronomer from Berlin, to whom Leverrier had sent up his calculations. Halle found the planet just  $1^\circ$  off its theoretically predicted position!

This story is not directly concerned with the general relativity theory, except to point out that the triumph of Newton's theory recruited such ardent proponents among physicists and astronomers that they began a search for new planets close to the Sun. In the mid-1800s some data appeared that indicated irregularities in Mercury's motion. Regular observations of Mercury's motion around the Sun allowed astronomers to determine its orbit very accurately, but over longer periods departures from the trajectory predicted by Newton's theory were found to exist. The same Leverrier set off to calculate the orbit and mass of the hypothetical planet which could cause these deviations.

Kepler's laws, which follow from Newtonian mechanics, can only be observed in their pure form if the Sun has a single planet in its system. It has, however, many planets, some of which have considerable masses (Jupiter, for example, has only one thousandth of the mass of the Sun; and the Earth is 333 000 times lighter than the Sun). The attraction of the planets to each other alters their trajectories around the Sun from being ideally elliptical (as they should be if there were a single Newtonian gravitational potential proportional to  $1/r$ ). Instead, the orbits are nearly elliptical and open. This means that they appear to be ellipses which are precessing or slowly rotating in their plane. Thus Mercury's perihelion, for example, advances by about  $575''$  (arc seconds) per century due to the action of all the other planets. We

should add to this figure a nondynamic (kinematic) advance of the perihelion which is due to the rotation (or, being more exact, to the precession) of our frame and amounts to  $5026''$  per century. These corrections are the consequences of Newton's theory, which was no less skilfully applied in the past than it is nowadays. At that time, the astronomers already found that Mercury's perihelion advances an extra  $40''$  on top of the explainable  $575'' + 5026'' = 5601''$ . That was the enigma that attracted Leverrier. However no evidence could be found for the hypothetical planet, neither transits across the disc of the Sun nor a planet's presence in the vicinity of the Sun during total solar eclipses. Attempts were made to generalize Newton's law of universal gravitation by adding summands to the potential or by changing the power of  $r$ , but this was like the addition of epicycles within the ancient astronomic models of Hipparchus and Ptolemy. Only the general theory of relativity finally resolved this problem.

As he was formulating his general relativity theory, Einstein showed that the planetary orbits cannot be closed and that they must resemble ellipses continuously rotating in their planes, their semimajor axes rotating at a constant velocity. We can quite easily derive the laws of energy and momentum conservation from the equation of a geodesic in the field of a point source (the Schwarzschild field) and show that motion occurs in the same plane (we shall refer to it as the equatorial plane). Since we want to know the trajectory of a planet's orbit, but not the four-dimensional world line, then we can combine the first integrals with the expression for the 4-interval, and arrive at first-order equation for the radial coordinate  $r$  as a function of the angle  $\varphi$ . This equation can be integrated by quadrature, but the result cannot be expressed through elementary functions (we get an elliptic integral). Hence, an approximate solution is used. This is not important for our basic analysis, and we shall confine ourselves to the most general characteristic of this solution. The first post-Newtonian approximation gives the equation for an orbit

$$r = \frac{a(1-\varepsilon^2)}{1+\varepsilon \cos \lambda\varphi}, \quad (2.20)$$

where  $a$  is the semimajor axis of an orbit,  $\varepsilon$  is its eccentricity,



and

$$\lambda = 1 - \frac{3Gm}{c^2 a (1 - \varepsilon^2)}.$$

The period of the function  $r(\varphi)$  is the same as the period of the cosine, that is, it equals  $2\pi/\lambda$ . Therefore,  $\varphi$  changes by more than  $2\pi$  during a single orbit. For example, the perihelion advances over one revolution by

$$\delta\varphi = \frac{2\pi}{\lambda} - 2\pi \simeq \frac{6\pi Gm}{c^2 a (1 - \varepsilon^2)}. \quad (2.21)$$

If we now substitute in the mass of the Sun ( $m_{\odot} = 1.99 \times 10^{33}$  g), the length of the semimajor axis of Mercury's orbit ( $a = 57.91 \times 10^{11}$  cm), the orbit's eccentricity ( $\varepsilon = 0.21$ ), and Mercury's orbital period around the Sun (87.97 days), we can easily calculate the angle by which Mercury's perihelion advances per terrestrial century, viz.  $43.03''$  (see Fig. 5). Observations from the surface of the Earth give us (after subtraction of kinematic and dynamic perturbations due to the other planets) a value of  $42.9 \pm 0.2''$ , a remarkable agreement with the prediction from general relativity [16] (discussions of this can be found in [6, 10, 92]).

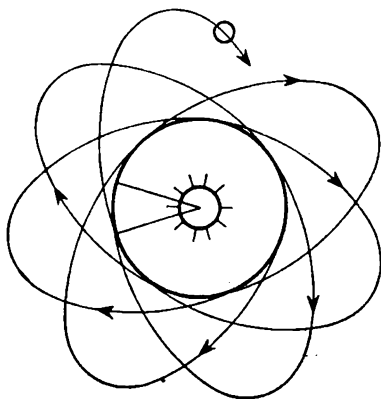


Fig. 5

Of course, similar perihelion advances must exist for other planets, too, but they are smaller because they are much farther away from the Sun than Mercury, and they move much more slowly.

#### 2.4. DEFLECTION OF LIGHT BY GRAVITATIONAL FIELDS

Gravitational interaction is universal and all bodies act on one another through their gravitational fields. Light also generates a gravitational field and hence logically "feels" the gravitational fields of other objects (let us call this an

analogue of Newton's third law). The gravitational field of light is negligibly small, and we cannot think of detecting its effect on the behaviour of macroscopic bodies. However, the world line of a light ray changes noticeably close to massive bodies such as the Sun. This deflection can be computed by solving the geodesic equation for light and the geodesic will be null, i.e.  $ds^2 = 0$ , thus differing from a planetary geodesic.

Let us consider the light propagating in the spherically symmetric gravitational field of a compact mass, that is, in a Schwarzschild field (2.15). We should, however, modify (2.15) by changing the radial coordinate in order to simplify the calculation. We shall not go through the whole procedure here, but we will simply use the final results and derive some corollaries. The equation for the trajectory of light (the geodesic equation solution in the first approximation) is

$$\frac{1}{r} = \frac{1 + \varepsilon \cos \varphi}{r_0 (1 + \varepsilon)}, \quad (2.22)$$

since the new radial coordinate here is a function of the old one (see (2.15)), we have designated it  $r$ , too, and in fact they coincide at great distances. Here  $r_0$  is the radius of the closest approach to the centre of attraction ( $r = r_0$  at  $\varphi = 0$ ; the trajectory is symmetric with respect to this position, i.e. it does not change under a change in the sign of the azimuthal angle  $\varphi$ ), and  $\varepsilon$  is the eccentricity of the "orbit". Also

$$\varepsilon = 1 + \frac{c^2 r_0}{2Gm} \gg 1. \quad (2.23)$$

Thus the ray's trajectory is a wide hyperbola.

We can see from Fig. 6 that if an observer looks from the left along the light ray, then the light source will seem to be located somewhere along the left asymptote, i.e. in the upper right quarter. In reality, the source is somewhere along the right asymptote (provided that the source is far enough away). Therefore, the change in the direction in which we see the source is determined by the angle  $\delta$  between the two asymptotes, and this can be easily calculated from Fig. 6 and equation (2.22). The values of the polar angle  $\varphi$  at which the radial coordinate  $r$  becomes infinitely large on the trajectory determine the directions of both asymptotes. The absolute value of both angles is  $(\delta + \pi)/2$  (remem-

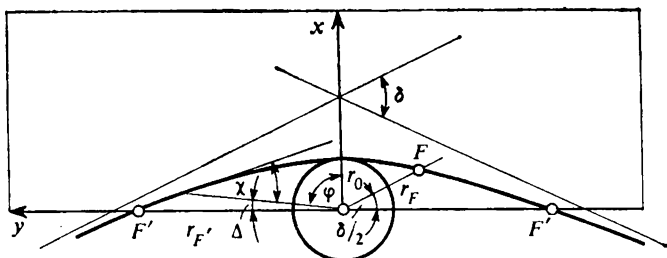


Fig. 6

ber the symmetry of the curve). The angle  $\delta$  is small enough for the sine to be approximated by the angle (in radians), and hence

$$\delta \simeq 2 \sin \frac{\delta}{2} = -2 \cos \left( \frac{\pi + \delta}{2} \right) = \frac{2}{\varepsilon}. \quad (2.24)$$

We used equation (2.22). Hence, by substituting the eccentricity from (2.23), in which the first summand (unity) should be ignored, because the second one is much larger, we obtain the main numerical results of general relativity for the deflection of light by a gravitational field of mass  $m$ , i.e.

$$\delta = Gm/c^2 r_0. \quad (2.25)$$

Taking the mass of the Sun to be  $m_{\odot} = 1.99 \times 10^{33}$  g and its radius to be  $r_0 = 7 \times 10^{10}$  cm (we take the radius of the closest approach), then for a ray tangent to the Sun's limb we have a deflection of  $1.75''$  (to verify this, don't forget to convert radians into arc seconds!).

Newton's theory also predicts that light will be deflected by a gravitational field (if light is considered to be a flux of particles), but the result is half the size. The calculation for Newton's theory was first made by Soldner, a German astronomer, in 1800 and Einstein did a similar calculation in 1911 using the preliminary variant of his theory. It was, therefore, important to verify by direct observation which theory should be chosen. The history of theories in the natural sciences, and particularly in physics, shows that in the end those theories finally become established (as most accurately and fully describing reality), which stand out as having a logical elegance and beauty. Beauty is subjective and has no rigorous value, but the first quality is objective to the extent that it can be used in practice. We should add

another characteristic, simplicity, but this, too, is not an objective quality. In this context, Einstein's gravitational theory is indisputably better than the others, and experience has indeed definitely shown it to be preferable.

During the first few decades after general relativity had appeared the only method of verifying the deflection of light was to observe the apparent positions of stars in the vicinity of the Sun's locus in the sky. They were observed first in the vicinity of the Sun itself, and then when the Sun was as far as possible. In the second case, it is very convenient to take photographs of the sky in standard conditions. Similar pictures (using the same camera and type of plates) should be taken in the first case but only during a full solar eclipse when the Sun is in the same part of the sky. This was difficult since total solar eclipses are quite rare and only take place in a narrow band of the Earth's surface. Costly field expeditions have thus to be sent to these areas, and several should be sent each time to offset the possibilities of foul weather. In addition, solar eclipses occur on the day side where the atmosphere is unstable because the air is in motion due to heat from the Sun. This motion seriously distorts the positions of the stars in the pictures, they "walk" around the frame because of local fluctuations in the refractive index of the atmosphere. Therefore, these observations, even though they support the predictions of general relativity have big quantitative errors. Information about these experiments can be found in [24, 39, 95] (about the first successful expedition by Sir Arthur Eddington in 1919) and in [15].

The same publications also contain the results of more accurate measurements that did not require observations during total solar eclipses. It turned out that close to the line of the ecliptic (the apparent annual locus of the Sun among the stars on the celestial sphere) there are several quasars. These extremely distant objects are powerful sources of radiation (in a number of cases radio sources), even "brighter" than the Sun although it is much nearer us. Since radio waves are not strongly attenuated by the Earth's atmosphere, radiointerferometry with a long baseline (the distance between two antennae in a large interferometer) can be used to pinpoint radio sources (quasars) on the celestial sphere, and detect deflection of radiation near the Sun. In this way, the precision of the experiments that confirm the theoret-

ically predicted deflection was raised to better than 1 percent. In addition it has been possible to eliminate partially the effect of radio wave refraction in plasma that surrounds the Sun (of which the corona is a part) using data on the variation in the refractive index of plasma with frequency. The quasar observations were conducted simultaneously in different wavebands. Of course it should also be possible to photograph the stars around the Sun from a spacecraft, masking the disc with a small screen.

The time-delays encountered when radar signals propagate close to the Sun are due to both the light being deflected and the gravitational redshift. If a narrow lobe-shaped radio signal is sent through the strong gravitational field of the Sun to Venus, Mercury or an unmanned space probe when they are about to go behind the Sun, and having rebounded from the body (triggering a transponder on the space ship) it returns through the strong gravitational field around the Sun to be detected by a receiver on the Earth, then the round-trip of the signal takes longer than the time theoretically calculated given a constant velocity of light. This does not however contradict the postulate that the velocity of light is constant, a postulate that is indispensable for general relativity. In fact we have measured the time by an Earth-bound clock, whereas in the medium where the signal was travelling, the time passes more slowly (due to the gravitational redshift). However, the path length is calculated according to a local scale closer to the Sun, so, the ratio of this length to the time measured by the Earth-bound clock cannot be equal to the standard velocity of light.

Another effect which is due entirely to the deflection of star light by the gravitational field of the Sun explains the well documented phenomenon of "negative parallax". It was found that the distances to some stars, if calculated by the parallax between observations taken from opposite points on the Earth's orbit, turned out formally to be greater than infinity. The parallax values for most stars are very small, so that even if the light from them is tangent to the Earth's orbit, and never comes near to the field source, i.e. the Sun, it may be deflected more than the change due to parallax (see the footnote at the end of Sec. 1.3).

The science of measuring cosmological distances, either using visible light or radio waves, the latter with long baseline interferometers, is important to modern astrono-

my. Over the last decade the precision of the measurements within the Solar system (i.e. measuring distances of the order of one astronomical unit,  $1.5 \times 10^8$  km or about 500 light seconds) has been raised so much that the errors are now some one to three orders of magnitude less than the Sun's gravitational radius (about 3 km). This means that even general relativistic effects such as perihelion advance or light deflection can be found even in routine astronomical observations (a detectable perihelion advance can accumulate within a few seconds!). This is why a few years ago the International Astronomical Union (IAU) ruled that all ephemerides had to be calculated using general relativity. The rule means that the equations of general relativity are now to be used not only to calculate the orbital motion of the planets themselves but also for the calculation of the light paths from the planets to an observer on the Earth or on board an orbital station, i.e. within the gravitational field filling the Solar system. The problems of relativistic astronomical measurement were discussed at the IAU symposium No. 114 that took place in Leningrad, USSR, in May 1985. Nowadays the experimental evidence cannot be explained using Newton's theory and to meet requirements of precise observation it has been substituted by the theory of general relativity.

## **2.5. GRAVITATIONAL LENSES**

Consider two photographs of the same part of the sky, one with the Sun in the centre and the other effectively without the Sun. If the pictures are superimposed, we will observe that the Sun has apparently pushed the stars away: by attracting their light it has "repelled" their images. In other words, the gravitational field of the Sun has acted like a lens, magnifying the image of the star field (though the stars in the pictures remained point images). The Sun's field becomes weaker with distance, and the "magnification" decreases the further the stars are away from the Sun in the angular sense. In this way, a concentrated mass forms a gravitational lens around itself, though by contrast with a standard optical lens it cannot "project" a sharp image onto a screen. This deficiency is because the "focal length" of the gravitational lens is directly proportional to the square of the radius of the light ray's closest approach to the gravi-

tating centre. Let us get back to Fig. 6 and use the formulae from Sec. 2.4 to determine the basic properties of a gravitational lens.

Consider two parallel rays coming from infinity and passing on either side of a mass  $m$  in the same plane with the same minimal radius  $r_0$ . They will be deflected by the gravitational field of the mass and intersect at some point. This point is called the focal point of the lens for the given mass  $m$  and given closest approach radius  $r_0$ . Let us designate the corresponding focal radius  $r_F$ . In addition, the rays emitted from a point farther from  $m$  than the focal one and passing the mass on either side at the same closest approach radius  $r_0$  will always have a conjugate point where the rays meet again. Thus, we have two conjugate focuses, and as one tends to infinity the other tends to  $r_F$ . Let us consider a special case, in which the two conjugate focuses are equidistant from the central mass, this distance being  $r_{F'}$ . All these points are shown in Fig. 6, where only one of the rays is drawn (the reader can add a second ray, which should be drawn corresponding to the choice of configuration). Obviously, the construction of the second ray will be different for focus  $F$  and for the conjugate focuses  $F'$ . Using Fig. 6 and (2.22), we have for  $F$

$$r_F = r_0^2 c^2 / 4Gm, \quad (2.26)$$

since at  $F$  we have  $\varphi = (\pi - \delta)/2$ , and for  $F'$

$$r_{F'} = r_0^2 c^2 / 2Gm, \quad (2.27)$$

since at the point  $F'$  we have  $\varphi = \pi/2$ . So, two symmetric conjugate focuses (both are identified as  $F'$ ) are twice as far from the central mass as the focus  $F$ , which has its conjugate point at infinity.

If central mass belongs to a distributed object of a spherical form (as our Sun is to a good approximation), we can evaluate the effect of its gravitational field on its own observed image. From Fig. 6 follows that due to gravitational deflection of light, this massive luminous object must appear larger than it is "in reality". Moreover, an observer stationed at some distance from the object does not see all of the half he faces but will see more or less than that depending on his position. Suppose the observer is at a distance  $r$  from the centre of mass (in Fig. 6 the position of the observer is connected with the origin by a straight line at the angle  $\Delta$

with respect to the  $y$  axis). From this point the observer will view a radius of the central body at an angle  $\chi$ , this radius being at an angle  $\varphi$  with the direction to the observer (see Fig. 6). Clearly, the angle  $\varphi$  can be computed from another form of equation (2.22)

$$\cos \varphi = \frac{1}{\varepsilon} [(1 + \varepsilon) r_0/r - 1], \quad (2.28)$$

and if it is only slightly different from a right angle, that is,  $\varphi = \pi/2 - \Delta$ , and  $\Delta \ll 1$ , then

$$\Delta = r_0/r - 2Gm/c^2 r_0 \quad (2.29)$$

(we also used (2.23) when getting this result). The first term describes the effect of flat geometry (in the absence of the gravitational deflection of light), and is simply related to the finite distance from the observer to the central body. The second term, on the other hand, is essentially due to gravitation, and is even independent of the position of the observer. Notice the different signs of the terms; since the light rays are deflected by gravitation, they tend to "show" the reverse side of the central body. As to the angle  $\chi$ , it can be presented, according to Fig. 6, in Cartesian coordinates  $x, y$  (as shown, the picture need only be turned to make the  $y$  axis vertical). Remember that a tangent to a curve has a slope whose tan is equal to the derivative  $dy/dx$ , where the function  $y(x)$  describes the curve. It is easy to see that

$$dy/dx = \tan(\chi + \varphi - \pi). \quad (2.30)$$

In Cartesian coordinates (2.22) can be rewritten as

$$\sqrt{x^2 + y^2} + \varepsilon x = r_0(1 + \varepsilon), \quad (2.31)$$

Table

Parameter	Gravitational lenses		
	neutron	Earth	Jupiter
$m, g$	$1.7 \times 10^{-24}$	$6 \times 10^{27}$	$1.9 \times 10^{30}$
$r_0, \text{ cm}$	$10^{-13}$	$6.4 \times 10^8$	$7 \times 10^9$
$r_g, \text{ cm}$	$2.5 \times 10^{-52}$	$0.89 \times 10^0$	$2.8 \times 10^2$
$r_F, \text{ cm}$	$2 \times 10^{25}$	$2.3 \times 10^{17}$	$0.9 \times 10^{17}$
$\Delta_F, \text{ rad}$	$2.5 \times 10^{-39}$	$1.4 \times 10^{-9}$	$4 \times 10^{-8}$
$\chi_F, \text{ rad}$	$5 \times 10^{-39}$	$2.8 \times 10^{-9}$	$8 \times 10^{-8}$



since  $r = \sqrt{x^2 + y^2}$  and  $x = r \cos \varphi$ . Differentiating (2.31) with respect to  $x$ , we have

$$\frac{x + y \, dy/dx}{r} + \varepsilon = 0 \quad (2.32)$$

or, rewriting in polar coordinates,

$$dy/dx = -(\varepsilon + \cos \varphi)/\sin \varphi. \quad (2.33)$$

A further calculation, in which the elementary properties of the trigonometric functions and their expansions in power series are used, brings us to the following expression (the angle  $\chi$  is small):

$$\chi \simeq \frac{r_0}{r} \left[ 1 + \frac{2Gm}{c^2} \left( \frac{1}{r_0} - \frac{1}{r} \right) + \frac{1}{2} \left( \frac{r_0}{r} \right)^2 \right]. \quad (2.34)$$

Anyone unwilling to do the trigonometry will have to take this relationship for granted. Has he not accepted the even more important relationship (2.22)?

The table shows the characteristics of various gravitational lenses. Consider the masses which generate the gravitational fields of these lenses: an elementary particle (neutron), our Earth, Sun, a neutron star (a pulsar, whose rotation we will not take into consideration here), and a galaxy (for simplicity, we will take the atypical case of a spherically symmetric mass distribution, so that the Schwarzschild solution can be used outside it). The first three lines of the table show masses, gravitational and geometrical radii of the objects (for a neutron and a galaxy these are typical sizes). The remaining lines contain the parameters of rays passing these objects along hyperbolae with the radii of the closest approach equal to the objects' radii. Since for symmetrically conjugate focuses

Gravitational lenses		
Sun	neutron star	galaxy
$2 \times 10^{33}$	$3 \times 10^{33}$	$10^{46}$
$7 \times 10^{10}$	$10^6$	$3 \times 10^{22}$
$3 \times 10^5$	$4.5 \times 10^5$	$1.5 \times 10^{18}$
$0.8 \times 10^{16}$	$1.12 \times 10^6$	$3 \times 10^{26}$
$4.3 \times 10^{-6}$	0.45	$5 \times 10^{-5}$
$8.3 \times 10^{-6}$	1.3	$10^{-4}$

$r_{F'} = 2r_F$ , the table only shows  $r_F$ . Note also that  $\Delta_\infty = -\Delta_F$  (the table contains the values for  $\Delta_F$ ; see Fig. 6). Clearly, a neutron is hopelessly inadequate as a gravitational lens. The planets of the Solar system have focal points farther than one hundred radii (semimajor axes) of the orbit of the farthest planet, Pluto\*. The distance to the closest focus of the Sun (a little more than ten such radii) is ten times shorter. For a neutron star, if it is 1.5 times more massive than the Sun, the focus is almost located on its surface. Since light rays which escape to infinity bend about the star, we might glimpse some of its opposite side,  $0.45 \text{ rad} = 26^\circ$ , that is, about 10% of its surface (all known pulsars, however, are so distant that we see them as points even in the most powerful telescopes). Notice that despite the increase of the area from which the radiation from the pulsar surface reaches us, its brightness remains the same; this fact can be easily explained if we think about the paths of all rays that leave it to infinity.

The focuses corresponding to parallel light rays that pass on either side of a central mass move out from  $r_F$  to infinity as the distance between the closest approach of the light rays gets larger. This is the basic difference between a gravitational lens and a normal optical one. If we look at a source of light that is far enough behind a large compact mass we will see the light (if we also are sufficiently far away from the mass) as a luminous circle whose centre is the mass creating the gravitational lens. The farther we go away from the mass the greater the radius of the circle. However, this is only valid if the light source, the mass creating the gravitational lens, and the observer are on the same straight line. Otherwise, the observer will not see a luminous circle around the mass, but will see separate images of the source lined up along a straight line, because in a spherically symmetric gravitational field (we shall stick to this approximation in this analysis) the trajectories of the rays are planar (the mass also lies in the same plane).

How many such images will the observer see? If the central massive body is opaque and has a radius significantly greater

---

\* 1 radian  $\simeq 57.5^\circ$ ; 1 astronomical unit (1 AU)  $\simeq 1.5 \times 10^{13} \text{ cm}$ ; 1 light-year (1 lt-yr)  $\simeq 10^{18} \text{ cm}$ . Pluto's semimajor axis  $\simeq 40 \text{ AU} \simeq 6 \times 10^{14} \text{ cm}$ . The distance to the nearest galaxy (M31, Andromeda Galaxy) is about  $2.2 \times 10^6 \text{ lt-yr}$ , and to a typical quasar, about  $10^{10} \text{ lt-yr}$ .

than that of a Schwarzschild sphere, it is clear that no more than two images can be seen: one apparently farther from the central mass than the other. The relative brightness of these images can be calculated and the one that is produced by least deflected light is the brightest (farther off the centre of the lens). If the central mass is somewhat transparent, three or more images will appear, the extra ones having been formed by the light that has passed through the gravitating medium. This situation may occur if the lens is created by a galaxy.

It is also possible for a lens to be generated by a ball, whose radius is comparable to that of Schwarzschild's. Rays which pass close to the surface of the central body are not simply bent by the gravitational field, they may even turn once or more times around the centre, and only then depart from it. Because of this, more but weaker images may appear at practically the same place as the central mass, if viewed by a faraway observer. A very small proportion of the light emitted by the source actually participates in creating these images. The source may, however, be away from the site (invisible mass constituting the lens) where these images appear. Of course, those rays that come "straight" to the observer must be much brighter.

Some idea about what happens when the mass creating a gravitational lens moves relative to a light source behind it (or vice versa, when the light source moves behind the lens) can be got by looking at the text on a sheet of printed paper through a convex lens. Some popularized descriptions of these phenomena can be found elsewhere [54], but rigorous analyses of gravitational lenses are scattered over many journals. However, to the best of our knowledge, there is no textbook that contains a clear general review of this problem. Those people who have dealt with it, however, differ in their interpretation of the effects, influenced by interference phenomena that significantly distort light patterns close to the "axis" of the lens.

Einstein predicted the gravitational-lens effect back in 1936 [32] and a number of researchers studied the details of the phenomenon. One may ask to what extent this effect has been proved by astronomical observations today. For many years it was believed that detecting it was a hopeless task. One star was put forward as a gravitational source, while other more remote stars were put forward as light

sources. Both kinds of stars had to be strong emitters, but the one we would call the "central mass", being relatively closer to us, should be much brighter than one posing as the light source, thus masking the effect. During the last few decades, however, we have witnessed the discovery of massive neutron stars, which have as much mass as the Sun, but only emit weakly in the optical waveband, and very remote and bright quasars, which emit specific sort of radiation. Thus the problem of masking has been resolved.

Quite recently, astronomers have discovered five potential gravitational lenses with a quasar as a light source and a galaxy as the mass creating the lens. In one case a central galaxy was even observed by telescope, and turned out to be a large elliptical galaxy with the mass of about  $10^{13}$  of that of the Sun. Quasars are fantastically far away from us. Moreover, they are moving rapidly away in a radial direction due to the expansion of the Universe (see Sec. 2.10). This shows up as a very large cosmological redshift. Since this parameter is related to the distance to the object, quasars with the same redshift must be equidistant from us. Any coincidence of redshift values, of course, is an extreme rarity. The spectra of quasars are even more diverse and each one emits a very individual radiation, although quasars are united into a very specific group of emitters. Therefore, astronomers were amazed at the discovery of two quasars in Ursa Major with practically the same spectra and redshifts. The observation showed that they were close to each other and "mimicked" each other's colour. This could not be a pure coincidence: either they were a unique quasar twin of the same origin, or simply two images of the same object viewed through a gravitational lens. Even more astonishing was the discovery of another similar object in another part of the sky; this time it was a "triple" quasar.

At present, we cannot categorically state that a gravitational lens has been actually observed, but the probability is really high. A final proof of the possibility would be, for example, if two quasar images duplicated each other's behaviour in time. As a matter of fact, quasars often change their brightness significantly (the period over which these changes occur indicates the approximate size of the quasar). If it is the same object, all the images must "wink" similarly, with only a time delay because the light, travelling by different paths, reaches us over substantially different times

(there was even a communication that a time delay has been already observed, and that it was about one and a half year).

There is another way of verifying the existence of gravitational lens. It is also associated with the individuality of the radiation but when it is measured over short time intervals, irrespective of the nature of its source. We have in mind here the coherence of rays that have travelled different paths and met again at some point. These paths should not, however, differ too much in length (or travel time). Generally speaking, therefore, if the light from different images of the same quasar is directly compared, it cannot be coherent, and no interference pattern would occur when the rays are superposed. There are ways, however, of storing electromagnetic signals which allow the signals to be superimposed with an artificially long-time delay in one of them. If the natural time delay is correct offset and the images are produced by the same source, the coherence will distinctly show up. This technique is rather simple for radio radiation, in which quasars are active, and the gravitational-lens effect is noticeable. This is a job for computers (before commissioning this project, the quasar "images" in radio band must be carefully recorded on tape for long periods of time). Of course, this task cannot be fulfilled manually because billions of comparisons of wave oscillations have to be completed.

## **2.6. SPACE-TIME AROUND ROTATING BODIES**

In Sec. 2.2 we learned what the Schwarzschild's solution is, i.e. a space-time in which there is only one spherically symmetric mass that does not rotate (rotation would have led to a preferred direction, that is, the axis of rotation, thus violating the symmetry of the system). However, the Schwarzschild case cannot be considered general enough for real objects. Therefore, ever since the emergence of general relativity many researchers have sought to describe a space-time containing at least one (for simplicity) rotating mass. In 1918, J. Lense and H. Thirring found an approximate solution of Einstein's field equations for this case and applied it to physical effects. They used linear approximation of Einstein's field equations assuming that the metric is close to a flat one (a gravitational field is weak), which is really

the case at distances much larger than the Schwarzschild radius. It was only in 1963 that R. Kerr [57] found the exact metric for such a space-time. Naturally, it transforms into the Lense-Thirring metric at great distances, and for a central body that is not rotating too fast. Here, we shall give a descriptive (but not rigorous) derivation of the Kerr metric by generalizing the consideration of the Schwarzschild metric we gave in Sec. 2.2, and then we shall write the approximate Lense-Thirring metric.

To begin with it is worth discussing some general properties of rotation. In order to describe the rotation of a rigid body an angular velocity  $\Omega$  is introduced in addition to the conventional (linear) velocity  $v$ , because the angular velocity is constant for a rigid rotation, whereas the linear velocity of any point of the body is proportional to the distance between the point and the axis of rotation (the distance is measured perpendicular to the axis, that is, it is a radial coordinate in a cylindrical system of coordinates). Thus, the relationship between angular and linear velocities is

$$v = \Omega \rho \quad (2.35)$$

in cylindrical coordinates, and

$$v = \Omega r \sin \theta \quad (2.36)$$

in spherical coordinates (because  $\rho = r \sin \theta$ ). However, a body can rotate not as a rigid one (for example, Jupiter's atmosphere rotates with different angular velocities at different latitudes as a result having different periods of rotation). The rotation period is related to angular velocity thus:  $T = 2\pi/\Omega$ . Hence the angular velocity may depend on position.

A reference frame (this notion will be discussed in detail in Sec. 2.8) may be rotating, too; though a rigid body rotation is even less natural for such a system than a rotation with different angular velocities at different points. Also, if a reference frame extended to infinity could rotate as a rigid body, that is, with a constant angular velocity  $\Omega$ , then a linear velocity at a finite distance from its axis (on a cylinder  $\rho = c/\Omega$ ) would reach the velocity of light  $c$ , and outside of this "light cylinder" would surpass it. Obviously, this kind of reference frame is impossible to simulate for any material bodies, therefore, the angular velocity of

any realistic reference frame must change with distance from the axis. The slowdown must not be less than inversely proportional to that distance. There should be a domain, well within the light cylinder, where the reference frame would rotate as a rigid body.

We live on a rotating planet and have to cope with the kinematic manifestations of rotation. We best understand two of these, that is the centrifugal force and the Coriolis force. The former is very noticeable and everyone has experienced it in a fast vehicle moving along a bend. The Coriolis force is less apparent, and to experience it it is not enough simply to sit or stand motionlessly in a car or bus while it turns, you have to move relative to the rotating frame of reference. While your bus is turning round a corner it is also rotating (usually, about an imaginary axis outside of the bus). The Coriolis force is perpendicular both to this axis and to the velocity of the object on which it acts:

$$\mathbf{F}_{\text{Cor}} = -2M\boldsymbol{\Omega}_{\text{r.f.}} \times \mathbf{v}_{\text{obs}}, \quad (2.37)$$

where  $\mathbf{v}_{\text{obs}}$  is the velocity of the observer,  $M$  is his/her mass,  $\boldsymbol{\Omega}_{\text{r.f.}}$  is the angular rotation velocity of the reference frame (its vector is directed along the axis of rotation according to the right-hand screw rule and has an absolute value equal to the angular velocity).

A well-known example of this force is the one which acts on the flow of rivers to produce a finding of **K. E. von Baer** (1792-1876). Each element of water flowing in a river is subjected to the Coriolis force, which is perpendicular to the direction of flow and to the Earth's rotation axis. The force would thus point in the direction a right-hand screw would move in if the screw were oriented at right angles both to the direction of the river and to the Earth's axis, and then rotated from the flow direction towards the North Pole (Fig. 7). Generally speaking, the Coriolis force does not act in a plane tangent to the Earth's surface (it would, if the river flowed along a meridian), hence the projection of the force on the plane is usually less than the complete Coriolis force. As a rule, we notice only this projection. In the northern hemisphere, the Coriolis force is directed towards the right bank of a river, whereas in the southern hemisphere it acts on the left one. Although it is not a large force, its action over a long period of time makes the river undermine one of its banks, which then becomes steeper than the

other (this is Baer's law in its original formulation). The Coriolis force may act along the vertical, and this is the case for rivers flowing along the equator. Those flowing eastwards will be subject to a Coriolis force directed upward, and those flowing westwards will feel it downward. Hence,

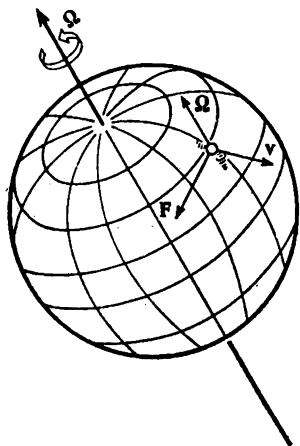


Fig. 7

all things being equal, on the equator the westward river will become deeper than the eastward one.

Both centrifugal and Coriolis forces are the varieties of inertial forces and are fictitious in that they are only the manifestation of the motion of the reference frame. They show how a simple description of the motion of bodies and particles in an inertial frame changes in a rotating frame. If the rotating frame is not a mathematical fiction but a physical object (for example, our Earth), these "fictitious" forces must be capable of converting one form of energy

into another. Baer's law, for example, is an illustration of this. Obviously, both a body that moves relative to the rotating frame, and the kinetic energy of the matter constituting the rotating body can act as the source of the energy.

It follows from the equivalence principle that both the centrifugal and Coriolis forces can be locally connected to gravitation. The gravitational field, however, will only exist if it cannot be compensated for by the choice of an appropriate reference frame (this compensation ought to be done in a finite domain or globally). Later, we shall see how these forces show up in dynamics, and not just in kinematics, in the form of a true gravitational field.

A rotating physical body possesses an angular momentum as a conserved characteristic, which in certain respects is related to energy and momentum, which are also subject to the conservation laws. Strictly speaking, energy and momentum are only conserved for isolated systems. Otherwise, they are controlled by the strict laws that determine the



exchange of energy, linear and angular momenta between the system and its surrounding (or the conversion of mechanical energy into other forms, such as thermal energy). In Newton's theory, energy is the source of a gravitational field (i.e., energy divided by the velocity of light squared, this energy in the form of mass), while linear and angular momenta have no such a role. In the general theory of relativity, however, a gravitational field is generated by a combination of distributions of energy, and linear and angular momenta (and the stress, too). Let us examine the angular momentum of an infinitely thin ring (which has, however, a finite mass) rotating around its axis. This angular momentum is a vector which is directed along the axis of rotation (using the right-hand screw rule) and has an absolute value of

$$L = mVR = m\Omega R^2,$$

where  $m$  is the mass of the ring,  $V$  its linear velocity,  $\Omega$  its angular velocity, and  $R$  is the radius of the ring.

Now, let us try to obtain a metric which describes the gravitational field around the rotating ring using a technique like the one we used in Sec. 2.2 for the Schwarzschild metric. We should warn the reader, however, that a rigorous derivation of this metric (the Kerr metric) is rather cumbersome. Even in lectures for students, it is normally presented without derivation in its final form. Here, we shall take the liberty of presenting a semiquantitative (accessible to the inquisitive mind of a qualified reader) "derivation" of the Kerr metric. A reader to whom this derivation may seem complicated can skip this passage on the first reading and continue from formula (2.47).

To account for rotation effects using the equivalence principle, we start from a rotating reference frame (it rotates not as a rigid body, but so that at large distances the effect of the rotation weakens; the nature of the frame rotation will be examined at the final stage of this analysis). In this rotating frame, we let a box with an observer fall towards the gravitating centre and in the box we take into account the slowing-down of the clock and the contraction of the scales in the direction of the fall, just as we did when deriving the Schwarzschild field. Assume that the box falls radially in the rotating frame. Then, we shall get back to the initial, nonrotating reference frame and consider the result.

We begin with the Minkowski's space-time in which we introduce spherical coordinates in a nonrotating frame; we assume the basis is (2.12), thus relative to it the flat space-time metric (2.9) will be

$$ds_{\infty}^2 = \omega_{\infty}^{(0)}\omega_{\infty}^{(0)} - \omega_{\infty}^{(1)}\omega_{\infty}^{(1)} - \omega_{\infty}^{(2)}\omega_{\infty}^{(2)} - \omega_{\infty}^{(3)}\omega_{\infty}^{(3)}. \quad (2.38)$$

A transition to a nonuniformly rotating reference frame is done by locally applying Lorentz transformations so that every point has its own speed of motion directed towards an increasing angle  $\varphi$ . The absolute value of this velocity is a function  $V$  which depends, generally speaking, on the coordinates  $r$  and  $\theta$ . Such a local Lorentz transformation is not equivalent to the transformation of the coordinates in the domain studied (in practice this domain is the whole of space) but is limited only to the transformation of the basis at each point. Thus, we have

$$\begin{aligned} \tilde{\omega}^{(0)} &= \left( \omega_{\infty}^{(0)} - \frac{V}{c} \omega_{\infty}^{(3)} \right) / \sqrt{1 - V^2/c^2}, & \tilde{\omega}^{(1)} &= \omega_{\infty}^{(1)}, \\ \tilde{\omega}^{(2)} &= \omega_{\infty}^{(2)}, & \tilde{\omega}^{(3)} &= \left( \omega_{\infty}^{(3)} - \frac{V}{c} \omega_{\infty}^{(0)} \right) / \sqrt{1 - V^2/c^2}. \end{aligned} \quad (2.39)$$

Now let the box with the observer be released from infinity. In this case we can write a new basis in which time has slowed down, and the lengths in radial direction have shortened. This is equivalent to the substitution of the  $\omega_{\infty}^{(\alpha)}$  in (2.13) by the basis covectors from (2.39)

$$\begin{aligned} \omega'^{(0)} &= \tilde{\omega}^{(0)} \sqrt{1 - v^2/c^2}, & \omega'^{(1)} &= \tilde{\omega}^{(1)} / \sqrt{1 - v^2/c^2}, \\ \omega'^{(2)} &= \tilde{\omega}^{(2)}, & \omega'^{(3)} &= \tilde{\omega}^{(3)}. \end{aligned} \quad (2.40)$$

Thus, we have assumed that the observer makes his measurements in the rotating frame and notices the relativistic changes in his observations. Now let us do the reverse transformation to the nonrotating reference frame by applying Lorentz transformations (inverse to (2.39)) to the basis (2.40):

$$\begin{aligned} \omega^{(0)} &= \left( \omega'^{(0)} + \frac{V}{c} \omega'^{(3)} \right) / \sqrt{1 - V^2/c^2}, & \omega^{(1)} &= \omega'^{(1)}, \\ \omega^{(2)} &= \omega'^{(2)}, & \omega^{(3)} &= \left( \omega'^{(3)} + \frac{V}{c} \omega'^{(0)} \right) / \sqrt{1 - V^2/c^2}. \end{aligned} \quad (2.41)$$

We now insert into (2.41) the  $\omega^{(\alpha)}$  basis, which is expressed in terms of the  $\tilde{\omega}^{(\alpha)}$  from (2.40), and then write this expression in terms of the  $\omega_{\infty}^{(\alpha)}$  from (2.39). We postulate, as we did for (2.14), that the resulting basis (2.41) remains orthonormalized. A few manipulations yield

$$ds^2 = \left(1 - \frac{v^2}{c^2 - V^2}\right) \omega_{\infty}^{(0)} \omega_{\infty}^{(0)} - (1 - v^2/c^2)^{-1} \omega_{\infty}^{(1)} \omega_{\infty}^{(1)} - \omega_{\infty}^{(2)} \omega_{\infty}^{(2)} \\ - \left(1 - \frac{v^2 V^2}{c^4 - c^2 V^2}\right) \omega_{\infty}^{(3)} \omega_{\infty}^{(3)} + \frac{2Vv^2}{c^3 - cV^2} \omega_{\infty}^{(0)} \omega_{\infty}^{(3)}. \quad (2.42)$$

The principle of correspondence with Newton's theory presupposes that (1.24) holds, as was the case for the derivation of the Schwarzschild field. Here the potential  $\Phi_N$  represents a solution of the Laplace equation, though under the new symmetry, that is rotational and not spherical. Therefore it is now worth considering oblique spheroidal coordinates in flat space. These coordinates,  $\rho$ ,  $\theta$ , and  $\varphi$  are defined as

$$x + iy = (\rho + ia) e^{i\varphi} \sin \theta, \quad z = \rho \cos \theta, \\ \frac{x^2 + y^2}{\varphi^2 + a^2} + \frac{z^2}{\rho^2} = 1, \quad r = \sqrt{x^2 + y^2 + z^2}. \quad (2.43)$$

We know that  $\Delta(1/r) = 0$  when  $r \neq 0$ , and this equality holds under any translation of coordinates. Let this translation be purely imaginary and directed along the  $z$  axis, i.e.  $x \rightarrow x$ ,  $y \rightarrow y$ , and  $z \rightarrow z - ia/c$ . Then we easily find that  $cr \rightarrow (c^2 r^2 - a^2 - 2iacz)^{1/2} = c\rho - ia \cos \theta$ . From here the expression for Newton's potential follows,

$$\Phi_N = -\frac{Gm}{c^2 r} \rightarrow \Phi_N = -\frac{Gm}{c} \operatorname{Re} \frac{1}{\rho c - ia \cos \theta} = \frac{Gm\rho}{c^2 \rho^2 + a^2 \cos^2 \theta}$$

since the Laplace equation is satisfied simultaneously by both the real and imaginary parts of the translated potential. Hence we can get with the help of (1.24)

$$\frac{v^2}{c^2 - V^2} = \frac{2Gm\rho}{c^2 \rho^2 + a^2 \cos^2 \theta}. \quad (2.44)$$

We determine the velocity  $V$  using a model of rotating ring of some radius  $\rho_0$  for the source of the Kerr field, this ring being stationary relative to the rotating reference frame (2.40). On the one hand,  $V = \Omega (x^2 + y^2)^{1/2} = (\rho^2 + a^2/c^2)^{1/2} \Omega \sin \theta$  corresponds to the relation (2.36). On the other hand, it is clear that the reference frame cannot

rotate as a rigid body, otherwise the frame wouldn't be extensible beyond the light cylinder as we dropped our box from infinity. Therefore the angular velocity  $\Omega$  has also to be a function of position. The ring lies naturally in the equatorial plane, so that its angular momentum is

$$L = mV \sqrt{\rho_0^2 + a^2/c^2} \left( \rho = \rho_0, \quad \theta = \frac{\pi}{2} \right).$$

We now introduce an important hypothesis which establishes a connection between the angular momentum and the Kerr parameter  $a$ , which is also a characteristic for spheroidal coordinates (2.43), namely we put  $a = L/m$ . These last three statements yield  $\Omega (\rho = \rho_0, \theta = \pi/2) = c^2 a / (c^2 \rho_0^2 + a^2)$ . If we now add a second hypothesis, that the field is independent of the choice of the ring radius (depending only on its angular momentum), then we get

$$\Omega = c^2 a / (c^2 \rho^2 + a^2)$$

and finally

$$V = ca (c^2 \rho^2 + a^2)^{-1/2} \sin \theta. \quad (2.45)$$

It only remains for us to choose the expression for a basis  $\omega_\infty^{(\alpha)}$  which would correspond to the assumed rotational symmetry (i.e., to the oblique spheroidal coordinates). The reader may substitute the coordinates  $x$ ,  $y$ , and  $z$  from (2.43) him/herself into the Minkowski squared interval,  $ds^2 = c^2 dt^2 - dx^2 - dy^2 - dz^2$ , hence getting a quadratic form with a nondiagonal term. This term, which contains  $d\rho d\varphi$ , can be excluded by a simple change of the azimuth angle:  $d\varphi \rightarrow d\varphi + ca (c^2 \rho^2 + a^2)^{-1} d\rho$  thus leading to a diagonal quadratic form. If now the square roots of the separate summands are taken, we get the final form of the initial basis  $\omega_\infty^{(\alpha)}$ :

$$\begin{aligned} \omega_\infty^{(0)} &= dt, & \omega_\infty^{(1)} &= \sqrt{(c^2 \rho^2 + a^2 \cos^2 \theta) / (c^2 \rho^2 + a^2)} d\rho, \\ \omega_\infty^{(2)} &= \sqrt{\rho^2 + a^2 \cos^2 \theta / c^2} d\theta, \\ \omega_\infty^{(3)} &= \sqrt{\rho^2 + a^2 / c^2} \sin \theta d\varphi. \end{aligned} \quad (2.46)$$

A mere substitution of these expressions into (2.42) yields the standard form of the Kerr metric in terms of the

Boyer-Lindquist coordinates,

$$\begin{aligned}
 ds^2 = & \left(1 - \frac{2Gm\rho}{c^2\rho^2 + a^2\cos^2\theta}\right) c^2 dt^2 - \frac{c^2\rho^2 + a^2\cos^2\theta}{c^2\rho^2 - 2Gm\rho + a^2} d\rho^2 \\
 & - \left(\rho^2 + \frac{a^2}{c^2}\cos^2\theta\right) d\theta^2 - \left(\rho^2 + a^2/c^2 + \frac{2Gm\rho a^2\sin^2\theta}{c^4\rho^2 + a^2c^2\cos^2\theta}\right) \\
 & \times \sin^2\theta d\varphi^2 + 2 \frac{2Gm\rho a\sin^2\theta}{c^2\rho^2 + a^2\cos^2\theta} dt d\varphi. \quad (2.47)
 \end{aligned}$$

Our approach does not, of course, guarantee that the resulting metric is automatically a solution of Einstein's gravitational field equations. That this is in fact the case we know from other arguments (from calculations that have been performed by many other workers, see e.g. [57]), but the method does give some hint as to how to understand the Kerr metric and its sources, and it lets us look at the structure of the latter.

We propose the reader perform the following exercise (simpler than that considered above). He/she is invited to assume in the calculations that  $(V/c)^2 \ll 1$ , thus dropping the corresponding terms in (2.39) and (2.41). This is the assumption of slow rotation (more exactly, of the smallness of  $L$ , the angular momentum of the source) and it leads to  $V = a \sin \theta/r$  instead of (2.45). Thus instead of the Kerr metric (2.47) he/she will get the approximate metric

$$\begin{aligned}
 ds^2 = & \left(1 - \frac{2Gm}{c^2r}\right) c^2 dt^2 - \left(1 - \frac{2Gm}{c^2r}\right)^{-1} dr^2 \\
 & - r^2 (d\theta^2 + \sin^2\theta d\varphi^2) + 2 \frac{2Gma}{c^2r} \sin^2\theta dt d\varphi, \quad (2.48)
 \end{aligned}$$

which is known as the Lense-Thirring metric. We have written in it  $r$  instead of  $\rho$  and taken into account the approximate sense of the expressions. This approximation will be used in the next section.

## 2.7. DRAGGING IN THE KERR FIELD

A reader for whom the calculations in this section seem too awkward can pass directly to the paragraphs discussing the results after formula (2.55).

The feature that immediately attracts our attention in the Kerr metric (and the Lense-Thirring metric) is the presence

of a nondiagonal term with the product  $dt d\varphi$ . This metric cannot be diagonalized by any coordinate transformation unless an explicit  $t$ -dependence of the metric tensor components is admitted. Space-times that possess this property are referred to as *stationary* as opposed to *static* ones (in the latter case, the metric can be obtained in a form that is both diagonal and  $t$ -independent, which means that the field source does not rotate). Thus the source of a Kerr metric (or Lense-Thirring metric) must rotate, as was assumed in the construction of the model. The nondiagonality of stationary metrics is the cause for a phenomenon which is known as dragging (sometimes, called the "dragging of local inertial frames"). Let us look into this effect and see if we can find analogues in other branches of physics which are more common to the reader.

Imagine that we are freely moving in a Lense-Thirring field at a constant (at least instantaneously) radial coordinate, that is, at this moment  $dr/dt = 0$ , but the second derivative of  $r$  with respect to time may be nonzero; the derivative is not determined by us, but by Nature itself (namely, by the geodesic equation). The value of the first derivative, however, defines the initial state of motion, which we can control (provided its velocity is less than the velocity of light). We shall base our analysis on one of the forms of the geodesic equation, namely (2.2). Its radial component is

$$\frac{d}{ds} \left( g_{11} \frac{dr}{ds} \right) = \frac{1}{2} \frac{\partial g_{\mu\nu}}{\partial r} \frac{dx^\mu}{ds} \frac{dx^\nu}{ds}$$

( $x^1 = r$ , the metric is diagonal with respect to the coordinate  $x^1$ ). Since the  $r$  value could be considered as instantaneously constant, we can derive the following expression for its second derivative with respect to  $s^*$

$$\begin{aligned} \frac{d^2 r}{ds^2} &= \frac{1}{2g_{11}} \frac{\partial g_{\mu\nu}}{\partial r} \frac{dx^\mu}{ds} \frac{dx^\nu}{ds} = \frac{1}{2g_{11}} \left[ \frac{\partial g_{tt}}{\partial r} \left( c \frac{dt}{ds} \right)^2 \right. \\ &\quad \left. + 2c \frac{\partial g_{t\varphi}}{\partial r} \frac{dt}{ds} \frac{d\varphi}{ds} + \frac{\partial g_{\varphi\varphi}}{\partial r} \left( \frac{d\varphi}{ds} \right)^2 + \frac{\partial g_{\theta\theta}}{\partial r} \left( \frac{d\theta}{ds} \right)^2 \right]. \quad (2.49) \end{aligned}$$

---

\* In order to understand the changes in  $r$  bear in mind that it can be differentiated with respect to either time variable because both the coordinate time  $t$  and the proper time  $s/c$  run in the same direction.

Note that the factors  $(d\varphi/ds)^2$  and  $(d\theta/ds)^2$  in the Lense-Thirring field (2.47) are the same as those that would have arisen in a Schwarzschild field or in a flat space-time metric written in terms of spherical coordinates. Hence these terms are not specific for the field of a rotating source, and our attention should only be focused on the contribution to  $d^2r/ds^2$  of the second term in the square brackets in (2.49). The coefficient for the  $(cdt/ds)^2$  term is the same as that for the Schwarzschild field, thus it is not of interest. Thus, the specific part of the radial acceleration (the second derivative of  $r$  with respect to  $s$ ) which is due to the source's rotation is

$$\left(\frac{d^2r}{ds^2}\right)_{\text{rot}} = \frac{c}{g_{rr}} \frac{\partial g_{t\varphi}}{\partial r} \frac{dt}{ds} \frac{d\varphi}{ds}.$$

By expanding the right-hand side and retaining only those terms with the same order of magnitude as  $Gm/c^2r$ , we get

$$\left(\frac{d^2r}{ds^2}\right)_{\text{rot}} = \frac{2Gma \sin^2 \theta}{c^2 r^2} \frac{dt}{ds} \frac{d\varphi}{ds}.$$

When the point which represents our observer moves slowly,  $ds$  differs sufficiently little from  $cdt$ , so we may assume that  $ds \simeq cdt$ . We designate  $\omega$  to be the instantaneous angular velocity of the observer about the  $z$  axis (i.e.  $d\varphi/dt$ ). Whence

$$\left(\frac{d^2r}{dt^2}\right)_{\text{rot}} = \frac{2Gma}{c^2 r^2} \sin^2 \theta \omega. \quad (2.50)$$

The additional force due to the rotation of the gravitational field source (dragging force), which causes this contribution to the radial acceleration, may thus be directed both towards or away from the centre, depending on the sign of the angular velocity  $\omega$ . If the observer moves in the same direction as the field source rotates, the dragging force is directed outward. Otherwise, if he moves counter to the field source's rotation the dragging force is inward, just like the Coriolis force in a river flowing along the equator (see page 82). The reader may ask how we know the direction of a source's rotation. The answer follows from the transformation (2.39) which was introduced to "accompany" the source (see the

comments made before the introduction of this transformation). The source therefore rotates in the direction of increasing  $\varphi$ . Look again at the

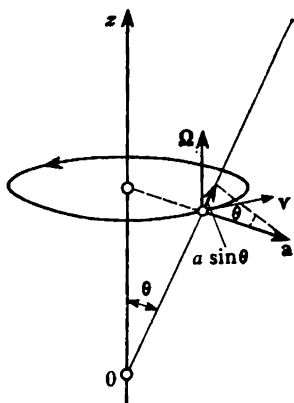


Fig. 8

Coriolis force (2.37) and its radial component, which is only influenced by the velocity component  $v_{\text{obs}}$  directed along angle  $\varphi$ , that is,  $\omega r \sin \theta$ . The force must be divided by the test particle's (observer's) mass to assess its contribution to acceleration (in our case, to the second derivative of  $r$ ). However, an acceleration due to the vector  $\Omega_{r.f.}$  directed along the  $z$  axis and to the velocity  $v_{\text{obs}}$  in the direction will not be oriented along the radius, but along a perpendicular to the  $z$  axis, its radial

projection resulting from usual multiplication by  $\sin \theta$  (see Fig. 8). In a vector product, on the other hand, the absolute values of  $\Omega_{r.f.}$  and  $v_{\text{obs}}$  will be simply multiplied (since the two vectors are mutually orthogonal). Consequently, we have

$$\left(\frac{d^2 r}{dt^2}\right)_{\text{Cor}} = 2\Omega_{r.f.} v_{\text{obs}} \sin \theta = 2\Omega_{r.f.} \omega r \sin^2 \theta. \quad (2.51)$$

This expression is in agreement with a corollary of the geodesic equation (2.50), if

$$\Omega_{r.f.} = \frac{Gma}{c^2 r^3} = \frac{Gm}{c^2 r} \Omega \quad (2.52)$$

is assumed as a residual correction for the rotation of the reference frame. In our determination of the metric in (2.47) we used the transformations (2.39) and (2.41), which would have cancelled each other out if there had not been the intermediary substitution (2.40). The factor  $\Omega$  in (2.52) is the same as that in (2.45) in the Lense-Thirring approximation. Let us verify this interpretation by a direct computation:

$$\begin{aligned} \omega^{(3)} &\simeq \omega'^{(3)} + \frac{V}{c} \omega'^{(0)} = \tilde{\omega}^{(3)} + \frac{V}{c} \sqrt{1-v^2/c^2} \tilde{\omega}^{(0)} \\ &\simeq \omega_{\infty}^{(3)} - \frac{V}{c} (1 - \sqrt{1-v^2/c^2}) \omega_{\infty}^{(0)}, \end{aligned} \quad (2.53)$$



from which it follows that the residual correction for the linear velocity is

$$v_{r.f.} = V (1 - \sqrt{1 - v^2/c^2}) \simeq \frac{1}{2} V v^2/c^2 = \frac{Gma}{c^2 r^2} \sin \theta, \quad (2.54)$$

in which (2.45) and the expression for  $v^2$  were approximately used. If we recollect now that  $\Omega = v/(r \sin \theta)$  and interpret the new angular velocity as  $\Omega_{r.f.}$ , we will have a relationship which exactly corresponds to (2.52)! This residual angular velocity shows up as a dragging effect, its form being like Coriolis acceleration. We could say that a massive rotating body entrains local inertial frames which results in effects similar to inertial forces. Note, however, that we have not found centrifugal-like forces in this study: they could not appear in this approximation due to their very structure. In classical mechanics, a centrifugal force is a triple vector product

$$\mathbf{F}_{\text{centr}} = -M\Omega_{r.f.} \times (\Omega_{r.f.} \times \mathbf{r}). \quad (2.55)$$

This means that it is proportional to  $(Gm/c^2 r)^2 (a/cr)^2$ , that is, to the small quantity of the second order that we ignored.

At the equator ( $\theta = \pi/2$ ), a force of the Coriolis type becomes purely radial in a Kerr field (if the motion occurs at a constant  $r$ ), and acts as a supplement to the centripetal force of gravitational attraction, which was introduced by Newton. This supplement may be both positive (directed outward, and weakening the force of gravity as a result) if the motion coincides with the direction of rotation of the central body, and negative (enhancing the force of gravity), if the motion is counter to the rotation. For a test particle moving along an equatorial orbit in a Kerr field, the rotation period is longer in the first case, and shorter in the second (remember that the period of a pendulum's oscillations decreases as the force acting on it increases). It is interesting that this change in the period of particle's orbit around a Kerr centre is a universal and structurally very simple quantity (either added to or subtracted from the standard Newtonian period, which dominates):

$$\Delta T = 2\pi \frac{L}{mc^3}. \quad (2.56)$$

Here, as before,  $L$  is the angular momentum of the centre,  $m$  is its mass, and  $c$  is the velocity of light. It is amazing that

$\Delta T$ , which appears from the exact solution of the geodesic equation in a Kerr field, is independent of both the gravitational constant and the orbital radius! This effect is really weak, and only becomes essential for the very rapidly rotating compact neutron stars (pulsars).

The Coriolis force (2.37), and its gravitational counterpart (2.50), which we examined for a Kerr field, are strikingly similar in structure to the Lorentz force, which determines the action of a magnetic field on a moving charge in electrodynamics:

$$\mathbf{F} = ev \times \mathbf{H}. \quad (2.57)$$

The only thing missing is a constant (equal to 2), and the apparent difference in sign is simply due to the skew-symmetry of the vector product relative to the order of the multiplicands; the mass  $M$  acts as the gravitational charge of the particle. Thus, the vector of the angular velocity (2.52), which describes the dragging effect, is analogous to the magnetic field, and once you know how to determine the direction in which a magnetic field acts on a charge, you will easily find out where the dragging effect pulls a test mass in a stationary gravitational field. On the other hand, the electric (Coulomb) field is analogous to usual Newtonian field of gravitation, but all the gravitational charges (masses) have the same sign, and are attracted to each other (in electrodynamics, charges may be either positive or negative, and like charges repel).

This analogy is not simply qualitative, it extends to many quantitative aspects of fields and the laws of motion. It therefore becomes useful in assessing specific physical effects and in planning research. Gravitational fields, however, have a far richer structure than electromagnetic ones. The metric tensor, for example, generally has ten independent components for a gravitational potential, whereas an electromagnetic potential only has four and so gravitational fields have some qualitatively new effects. The frozen-in field is another illustration of the analogy between electrodynamics and gravitation. In electrodynamics this phenomenon is observed in media with a high electrical conductivity for which the magnetic force lines seem to be anchored to the medium ("frozen in"). As the medium is compressed, the strength of the magnetic field may rise significantly even though the total magnetic flux through a given cross-

sectional area remains the same. This effect is used in physics to produce superstrong magnetic fields. A conducting medium with an induced magnetic field is placed in a container surrounded by an explosive. When compressed by an explosion the magnetic field reaches several million Gauss, but only for a short duration.

The same effect shows up on a cosmic scale. For example, when a supernova outbursts, its magnetic field, which for a normal star is rather weak, is compressed together with the star's core, which turns into a neutron star or pulsar. These extraordinary astrophysical objects can develop a magnetic field with the strength in excess of  $10^{12}$  Gauss! (The energy of such a magnetic field is so great that  $1 \text{ cm}^3$  of field contains even more energy than 100 g of purely electromagnetic mass.) In the case of gravitation, the situation is almost trivial, a simple demonstration of the law of conservation of angular momentum. If you sit in a rotating chair with your arms extended and holding two heavy books and then bring them towards you, you will rotate faster. Here, the angular velocity of rotation is analogous to the magnetic field and the magnetic flux (or magnetic moment) corresponds to angular momentum (they differ by a constant factor). Incidentally, this is the very reason that pulsars also have very high angular velocities, which are amazing for such massive bodies (they have as much mass as the Sun).

## **2.8. REFERENCE FRAMES IN THE GENERAL RELATIVITY**

We have so far used the notion of reference frame several times without defining it rigorously. In the last two sections, we even applied it to numerical calculations. Obviously, a reference frame is directly associated to the measurement and observation of physical effects. Any observation requires instruments (including our senses) whose readings will depend on how they move (how the observer moves, in particular). When general relativity first emerged, during the first stages of its development most efforts were spent on uniting space and time. As the importance of the notion of reference frame grew, however, the need to associate the four-dimensional theory to experiment demanded

that researchers find valid methods for splitting the four-dimensional manifold of special and general relativity into physical time and to the physical three-dimensional space of a given reference frame. This means that geometrical objects that would reflect the nature of the motion of the observer and his instruments had to be introduced into the theory alongside the four-dimensional geometry of space-time.

Modern observations are being carried out in three very different domains, and these are sketched out in Fig. 9. The

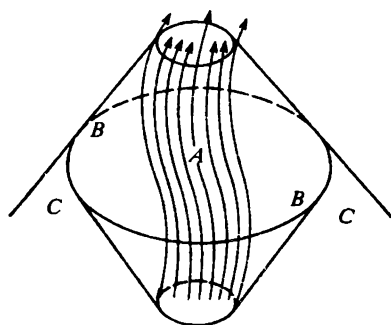


Fig. 9

world tube *A* in this figure represents the motion of a reference body, i.e. the laboratory of the observer (in a more general sense, we will consider a deformable "reference body"). In this laboratory, measurements can be carried out directly (e.g. by applying a ruler to the objects to be measured). The domain *A* can be extended, for example, to the Earth, or, at least, to that part of its surface acces-

sible to the observer, and where a triangulation network is available. The domain *B* comprises that part of space we can reach by active research, such as radar, i.e. we send signals along the light cone of the future and receive them back with the information we need along the light cone of the past. Thus, we can determine distances to a number of the bodies in the Solar system, and probe their properties. The third domain, *C*, lies beyond the reach of our active methods and we have to pump out information about it from the light cone of the past, and from inside the cone (e.g. observations by telescope, radiotelescope, charged particle counters and nuclear emulsions for cosmic rays). Since all these devices are in some state of motion (basically with the Earth), the information we collect is automatically related to a moving reference frame. The methods which we have mostly developed in the *A* domain are involuntarily transferred by us to other domains (*B* and *C*), even though to do so we have to apply special mathematical procedures.

Of course, in  $A$  we can use both direct measurement and all the other means of investigation which are applicable to the other domains. The results, however, are expressed as if they had been obtained by direct measurement. Therefore, with regard to domain  $A$  one can say that it is covered by an observer which permeates the entire domain. The observer can put point-like clocks and rulers everywhere in  $A$  (it is enough to have three mutually orthogonal rulers at each point). We thus have a field of a local basis (e.g. the basis  $\omega^{(\alpha)}$  in Sec. 2.2), which is sometimes referred to as a tetrad or vierbein (four orthonormal vectors) field. We can, however, do with clocks only, since the measurement of lengths can be reduced to the transmission and reception of light (radar) signals. Then, by postulating the constancy of the velocity of light (see the end of Sec. 2.4), we can determine the absolute value of any distance in the three-dimensional system, that is, we can estimate an observed length with relation to the given reference frame. Formally it is enough for us to know the metric of space-time and the world lines of the reference body. These world lines are the lines of the point-like instruments and observers, which can be completely described by the field of its unit tangent vector. Since the world lines are time-like, a tangent vector simply becomes a 4-velocity. The world lines of the reference body are said to form a congruence, which means that one and only one world line passes through every point of the domain in question (domain  $A$ ). If more than one world line passed through a point the description becomes ambiguous, at the very least. If a 4-velocity field is used to describe a reference frame, it is designated  $\tau^\mu$  and called a monad field. The intermediate fields between monad and tetrad ones are diad and triad fields (the fields of sets of two and three orthonormal vectors, one of which is time-like).

In the  $B$  and  $C$  domains we use accessible observational techniques and try to describe this world in the language of the reference frames following the procedure we used in the  $A$  domain. We must realize that in doing this our observations contain contributions due to the effect of the laws of physics (light propagation, etc.) and contributions due to the motion of the observer. The observations are not influenced by the fashion in which we extend the reference frame we used in  $A$  to the  $B$  and  $C$  domains, which are inaccessible for direct measurement. Nevertheless, the use

of monad and other reference frames in  $B$  and  $C$  is fruitful, and the techniques for extending the monad and tetrad fields from  $A$  to these domains have been developed.

When speaking about reference frames, we implicitly assumed that the reference bodies on which they are built are test bodies, i.e. they consist of test particles which do not perturb the geometry of space-time, and which do not influence either each other or the physical processes, i.e. they play an exclusively kinematic role. At the same time they are material points because their world lines must be time-like, otherwise they would not be able to simulate instruments and observers or behave according to the principle of causality. Thus the same phenomena can be viewed from different reference frames simultaneously (if the reference frames had not been built of test bodies, the phenomena observed in them would have been different because of the perturbations involved).

A reference frame is therefore an idealized model necessary to describe the measurement process. What role should be ascribed to a system of coordinates, which is often wrongly confused with a reference frame? It goes without saying that to obtain a quantitative prediction of any physical law we must use an appropriate system of coordinates, and the selection of a good system that accounts for the symmetry of a problem, say, may radically simplify the calculations. This does not mean, however, that the system of coordinates and the reference frame are identical. Academician Vladimir A. Fock (1898-1974), a Soviet theoretical physicist, wrote, "The notion of a physical reference frame (laboratory) is not identical, in general, to the notion of a system of coordinates, even if we disregarded all the properties of a laboratory with the exception of its motion as a whole" [38]. At the same time, the confusion is so widespread that even Einstein wrote that if such a coordinate system  $K$  is selected so that physical laws in it hold in their simplest form, the same laws hold in any other coordinate system  $K'$  that moves uniformly along a straight line relative to  $K$ . This postulate is called 'the special principle of relativity' [36]. In regard to the general relativity, Einstein wrote that the general laws of Nature must be expressed in terms of equations valid in all coordinate systems. These covariant equations must be relative to all kinds of substitutions (generally covariant). Obviously, the physics that meets

this postulate must meet the general relativity postulate, too, since the whole set of substitutions contains, in any case, those which correspond to all relative motions of (three-dimensional) coordinate systems [29]. Without underestimating these ideas, we would like to note that a covariant formulation can be applied practically to any equation in physics, so that it is not a principle of physical theory but a technique of formulation. By the same token, the form of the equivalence principle which locally equates gravity to acceleration depends on our ability at formulating the mathematical assertions of the theory (this formulation reduces to a transition from arbitrary coordinates to locally geodesic ones). We would here agree with Fock who asserted: "The true logical foundation of Einstein's gravitation theory is not the idea of general relativity, nor even the equivalence principle, but other two ideas: firstly, the unification of space and time into a single four-dimensional chronogeometrical manifold with indefinite metric (this was realized by Einstein in his theory in 1905, the special theory of relativity); and, secondly, the rejection of a "rigid" metric, which allowed him to relate it to the effect of gravitation, and thus to heavy matter (Einstein's field equations). The concept of the general covariance of equations (called the general relativity principle) and the kinematic interpretation of gravitation (called the equivalence principle), on the other hand, have only played a heuristic role" [37]. The same viewpoint is advocated by another authority on gravitational theory, John L. Synge [101].

The notion of reference frame in general relativity has undoubtedly been fruitful and reflects the relationship between four-dimensional theory and experience. In cases of strong gravitational fields, the application of this relationship is an absolute must. General relativity differs from special relativity in that special relativity admits the presence of privileged reference frames, namely the inertial frames, whereas general relativity allows, in general, no privileged frames whatsoever. Strictly speaking, however, these frames do appear in the presence of symmetries with respect to time (static and stationary space-times), but they are much poorer than those in special relativity. In this sense, the general theory of relativity is, Fock argues, more absolute than the special theory.

Getting back to the topic of the relation between reference

frames and systems of coordinates, we should say that a relationship can often be established, but only if special conventions are introduced.

Before we start the monad description of reference frames, we should warn the reader that the subject may prove somewhat difficult. If he does not feel like putting in the effort to wade through the argument, he can skip the section without damaging his perception of the subsequent ones.

The body of mathematics used to describe reference frames can be conveniently presented in five parts: (1) algebra; (2) the definition of the tensors describing the reference frames; (3) analysis (definition of the differential operators); (4) formulation of all the identities, equations and covariant relationships in the framework of this method (i.e., using the above quantities and operators); and (5) basic gauges of the method (i.e., the application of the conventions about the relationship between reference frames and coordinate systems). Below, we shall briefly describe the first two parts with reference to the monad technique.

**Algebra.** We mentioned above that the monad technique is based on setting a unit vector field  $\tau^\mu$  ( $\tau_\mu \tau^\mu = 1$ ) tangent to the congruence of the time-like world lines of points of a reference body. Hence, the physical sense of a monad is the 4-velocity vector field of the points of a reference body  $\tau^\mu = dx^\mu/ds$ . This vector field can be directly applied to one of the basic problems in the theory of reference frames, i.e. the determination of time intervals and the time components of tensors, as if the observer could only use clocks. Mathematically, this is done by contracting the corresponding tensors with the vector  $\tau^\mu$  (by projecting them onto  $\tau^\mu$ ). For example, the time interval relative to a given reference frame  $d\tau$  is given as  $d\tau = \tau_\mu dx^\mu$  (generally speaking, this is not a total differential!). The energy density is constructed by using the energy-momentum tensor  $T_{\mu\nu}$  as  $\varepsilon = T_{\mu\nu} \tau^\mu \tau^\nu$  in a flat space-time in Cartesian coordinates whose time coordinate lines coincide with the world lines of the reference body's points; this value is determined as  $T_0^0$ . Thus, the observed energy density is a true scalar, that is, it does not depend on the choice of a system of coordinates. In the theory of reference frames, all observed values have to be described by scalars (though not with respect to regauging the reference frame).

In reality, our observer is not such a simpleton that he



only uses clocks for his measurements. We can easily see that in his neighbourhood an observer may equally use rulers or light signals to measure distances (quantitatively the results are the same) without exceeding the limits of the monad technique, since it does not require any new mathematical objects other than monad field and metric tensor. From the  $g_{\mu\nu}$  tensor we can always extract the  $\tau_\mu\tau_\nu$  tensor. The remainder is designated

$$b_{\mu\nu} = \tau_\mu\tau_\nu - g_{\mu\nu}. \quad (2.58)$$

By definition, the projection of the tensor  $b_{\mu\nu}$  onto  $\tau^\mu$  is zero:  $\tau^\mu b_{\mu\nu} \equiv 0$ . This means that tensor  $b_{\mu\nu}$  is orthogonal to the time-like direction  $\tau^\mu$ , that is, it lies in a space-like hypersurface which corresponds to the 3-space of the given reference frame (see Fig. 10). These planes can be local ones, that is, the hypersurface may be nonholonom. This means that local area elements may not necessarily have a global 3-space envelope. At the same time it is easy to see that  $b_{\mu\nu}b^\nu_\lambda = -b_{\mu\lambda}$  and the tensor  $b_{\mu\nu}$  is a typical projector onto a local subspace which is an idempotent orthogonal to the physical time of the reference frame (i. e. to the congruence of the world lines of the points belonging to the reference body). This means that the projection of any tensor quantity onto the 3-space of the reference frame is done by contracting with  $b_{\mu\nu}$  (over one index). For example,  $T_{\mu\nu}b^\mu_\alpha b^\nu_\beta$  is a three-dimensional stress tensor. We can build mixed projections so that for one group of indices the projection is onto the temporal direction, while for another it is onto the 3-space orthogonal to the former;  $T_{\mu\nu}\tau^\mu b^\nu_\alpha = S_\alpha$  is the Poynting 3-vector of the energy flux density. These values are not yet observable, because they depend on which system of coordinates is chosen (as all usual components do), but their absolute values, expressed as scalar four-dimensional squares, are observable for a given monad reference frame. For example, the squared interval  $ds^2$  can be split into two parts: a squared observable time interval  $d\tau$  (we pointed out that it may not be a total differential) and a three-dimensional length (the measure of spatial translation relative to given

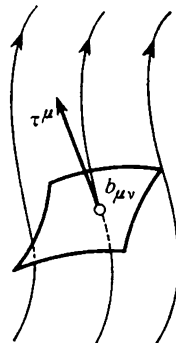


Fig. 10

reference frame)  $dl$ :

$$ds^2 = g_{\mu\nu} dx^\mu dx^\nu = (\tau_\mu dx^\mu)^2 - b_{\mu\nu} dx^\mu dx^\nu = d\tau^2 - dl^2. \quad (2.59)$$

This technique is sometimes called  $(1 + 3)$ -splitting of space-time (with respect to a reference frame). Also, the stress tensor is simultaneously the density of momentum flux; while the energy flux density, which we defined earlier, is the density of momentum for a given distribution of matter or fields. Following the established tradition, we will use the term "matter distribution" instead of "substance distribution", even though it does not survive a more general philosophical definition of matter. The reader can easily determine which quantities in flat space-time with Cartesian coordinates correspond to the flux densities defined above.

**Setting of Three Tensor (both Physical and Geometrical) Characteristics of a Reference Frame.** A reference frame is described, we have seen, by the congruence of the world lines of the points belonging to the reference body. In special relativity, we generally use inertial reference frames, the reference body of which behaves as a rigid body without acceleration or rotation (and, being rigid, it cannot of course be deformed). Therefore, the world lines of its points are straight lines. However, even in the framework of special relativity, we can introduce noninertial reference frames. What then is the difference between the motion of a reference body in an inertial and a noninertial frame? Firstly, the reference body's points which move noninertially possess acceleration. Moreover, a reference body may be rotating (see Secs. 2.6 and 2.7) and deforming (that is, it may not behave like a rigid body). The deformations, in turn, may be classified into dilatation (also called expansion; there may alternatively be compression if its sign is reversed) and shear. Under expansion, the shape of the initial configuration of reference body's points is not changed; under shear, on the other hand, if we had initially marked out (for example, painted) a set of points that form the surface of a sphere, then with time the sphere would turn into an ellipsoid with a certain orientation of its axes, but with the same size as the initial sphere. In general relativity too, any noninertial motion of a reference body is reduced just to these four (or three, if we do not subdivide the deformations into expansion or shear) characteristics.

Let us discuss rotation in more detail. In the case of a rigid body, the linear and angular velocities of its rotation are related as  $\mathbf{v} = \boldsymbol{\omega} \times \mathbf{r}$  (compare this to (2.36)). Taking the curl of this expression for a constant vector  $\boldsymbol{\omega}$  (the body is rigid), we have  $\boldsymbol{\omega} = (1/2) \text{curl } \mathbf{v}$ . This equation defines the local angular velocity of the rotation also in the case of nonrigid body (that is, when deformation is present). The transition to relativity presupposes the substitution of a conventional velocity by a four-dimensional one, and the generalization of the three-dimensional curl to the four-dimensional operator, projecting it onto the local 3-space of the reference frame. This situation is analogous to the four-dimensional form of electromagnetic field strength tensor expressed through a four-potential in Maxwell's electrodynamics. Thus, if we take the gradient  $\tau_{\mu\nu}$  (covariant differentiation!) of the four-velocity vector field of a reference body's points, that is, of the monad, then after skewsymmetrization and projecting it onto 3-space using the projector  $b_{\alpha}^{\mu}$ , we have the tensor of angular velocity of the reference frame rotation,  $A_{\alpha\beta}$ . Incidentally, it is not necessary to use a covariant derivative, a partial derivative will do in this case as well.

The acceleration vector of a reference frame is obtained if the gradient of the monad is skewsymmetrized and projected with respect to one index onto the physical time of the reference frame,  $\tau^{\mu}$  (for another index, this quantity will be automatically projected onto the 3-space because of skewsymmetry). It is also possible simply to project onto the monad the covariant gradient of the same monad with respect to the differentiation index. The deformation is characterized by the rate-of-strain tensor  $D_{\mu\nu}$  (the instantaneous state of deformation is not our concern because the reference body has no elasticity); this tensor by definition is the projection of the symmetrized covariant gradient of the monad onto the 3-space of the reference frame.

The concept of monad and the corresponding tensors which have been introduced in this section allow us to interpret physically in terms of the reference frame all the tensor equations and relations which are dealt with in general relativity.

## **2.9. STUDY OF THE UNIVERSE AS A WHOLE [COSMOLOGY]**

Cosmology, like astronomy, dates back to before recorded history. Man's inquisitive mind has been attempting to penetrate the mystery of the universe, the sense and purpose of his life, and his place in this world since time immemorial. This somewhat painful journey has led to truths that are more grandiose and, perhaps, even more fantastic than the old myths and legends which encapsulate kernels of our ancestors' view of the world. Cosmology has always been the battle ground for conflicting ideologies because it attempts to explain the world as a whole for all space and for all time. At first, humanity sought for the sense and purpose of its existence in this world, not even comprehending that these are entirely human concepts. Nor even was it realized or believed that man was not a central figure in the universe. But the development of cognition could not have unfolded otherwise, and it was inevitable that at the start humanity should have a childish egocentricity. Are we, in our majority, ready to accept that our planet is our most valuable treasure and that to be in harmony with it is our greatest virtue, which we should conscientiously and sincerely pursue freely, and without being constrained or attracted by any possible future benefit we or the next generation might accrue from a healthier habitat?

Millennia elapsed to be followed by centuries that matched them in import, which were followed in turn by decades of intense progress. Aristotle, who considered the Earth to be the centre of the universe and the focal point of all forces and goals was followed by Copernicus who placed the Sun at the centre of the universe while the Earth was reduced to the third planet orbiting around it. The winds of change that blew during the Renaissance spread the seeds of equality, both in the human and in cosmic sense. Giordano Bruno was hundreds of years ahead of his time when he put forward the idea that there were an infinite number of inhabited worlds, in a universe which had no holy centre. For this, he was burned at the stake by the Inquisition in 1600; human egocentricity had defended its citadel.

In the 19th century it was discovered that stars are far-away suns as large and hot as our own, and sometimes even larger and hotter. Even then it was obvious that some sort

of turbulent energy release must take place in stars (thermodynamics was just entering a boom period), but still our thought was conservative and we continued to adhere to the postulate of invariability and similarity of the world at all places and all times. This was the "absolute cosmological principle" and it held even Einstein at first (and Fred Hoyle and other prominent researchers in later decades). When in 1921, **Alexander A. Friedmann** (1888-1925) in Petrograd (now Leningrad) constructed the first nonstationary cosmological model, Einstein published a short comment in which he pointed out a mathematical error. His next response however was to acknowledge that he himself had been in error. In the late 1920s, **Edwin Hubble** (1889-1953) proved by observational data that the universe is expanding, as Friedmann had predicted. This theory was further developed by **Robertson and Walker**, and it is now referred to as the Friedmann-Robertson-Walker (FRW) cosmological model.

Since then, the cosmological principle has been interpreted in a narrower sense. Now we say that the universe, on a large scale, is homogeneous and isotropic, but that its properties are varying in time. However, if homogeneity and isotropy exist in one reference frame (on one section of space-time) but on another hypersurface corresponding to another reference frame, inhomogeneity and anisotropy appear, then the question arises as to in which frame is the universe the simplest and most symmetric? The answer is rather simple: homogeneity and isotropy are found in the frame that is comoving\* with all the matter and the fields in the universe. We know, however, that in our neighbourhood and on scales psychologically accessible to us, the world is far from homogeneous (but while inhomogeneity is sufficient for anisotropy, it is not necessary for it). The cosmological principle can only be true on a cosmological (and not on simply a macroscopic or even astronomical) scale.

How can these scales be measured? The determination of the distances to the observable astronomical objects—stars, galaxies, quasars—is one of the central problems of modern astronomy and cosmology. In cosmology, astronomers often use the photometric distance  $D$  to a source of radi-

---

\* 'Comoving' here means that the frame is fixed with respect to the average matter and fields of the universe.

ation, which is defined as  $D = \sqrt{L/4\pi l}$ . Here  $L$  is the absolute luminosity of the source (the luminosity which is observed from a standard distance, 1 parsec (pc), for example)\*, and  $l$  is its apparent luminosity (think about the sense of factor  $4\pi$ , how it is related to the geometry and what properties it reflects). The main difficulty in determining  $D$  in practice is the measurement of the  $L$  of the source. Hubble overcame this by finding variable stars in the galaxies closest to us (Cepheid variables). The brightness periods of these stars, the Cepheid variables, are unambiguously related to their absolute luminosities. He also found, using the Doppler effect to determine the radial (along our line of sight) velocities of these galaxies, that these velocities always show a component directed away from us which increases linearly with distance.

The observable universe (the distribution of the galaxies) starts to be isotropic in volumes of 30 cubic megaparsec ( $1 \text{ Mpc} = 10^6 \text{ pc}$ ). In volumes a thousand times larger, matter is distributed isotropically to a few percent. This isotropy of distribution holds for both optical and radio sources, though the situation is not yet clear with respect to quasars. The isotropy of background radiation (see Chapter Three) has now been measured to better than 0.1%.

Before solving Einstein's equations, the right-hand sides, i.e. the energy-momentum tensor of the gravitational field sources, must be specified. If the average matter distribution (of galaxies, or, to be more exact, their clusters) in the universe is likened to an isotropic fluid, the corresponding tensor  $T_{\mu\nu}$  will be

$$T_{\mu\nu} = (\rho + p/c^2) u_\mu u_\nu - g_{\mu\nu} p/c^2, \quad (2.60)$$

where the mass (energy) density  $\rho$  and pressure  $p$  depend only on cosmic time  $t$ , and in a frame referred to the matter  $u^0 = 1$ ,  $u^i = 0$ . The function  $\rho$  contains both the contribution of the mass density of the galaxies and their kinetic energy and those due to electromagnetic fields (including radiation), to neutrinos, and to gravitational binding energy, indeed to all forms of matter in general.

We must emphasize once more that the isotropy and homogeneity discussed in this section are natural in a spatial section of a privileged reference frame moving with all the

---

\*  $1 \text{ pc} = 3.262 \text{ lt-yr} = 3.086 \times 10^{18} \text{ cm}$ .

matter in the universe. If we include this symmetry of spatial sections in Einstein's field equations and write them in the (1 + 3)-splitting form (see previous sections about reference frames), we can show that the Ricci tensor  ${}^3R_{ik}$  of three-dimensional spatial sections is connected to the components of their metric tensor  $b_{ik}$  in a very remarkable way:

$${}^3R_{ik} = B b_{ik}, \quad (2.61)$$

where  $B(t)$  is a function which does not depend on spatial coordinates. There are three different types of spatial sections (3-spaces) depending on the sign of  $B$ :

(a) If  $B > 0$ , then the space has constant (in the sense of its dependence on spatial coordinates) and positive curvature.

(b) If  $B < 0$ , then the space has constant (in the same sense) and negative curvature.

(c) If  $B = 0$ , the space has zero curvature.

By solving equation (2.61), we get the particular forms of the 3-metric  $b_{ik}$  for each of these three cases. The forms can be combined as the following expression for a squared three-dimensional distance

$$dl^2 = R^2(t) \left[ \frac{dr^2}{1 - Kr^2} + r^2 (d\theta^2 + \sin^2 \theta d\varphi^2) \right], \quad (2.62)$$

where  $R(t)$  (the scale factor) is a function of cosmic time  $t$ ,  $K$  is an auxiliary parameter which assumes the values  $+1$ ,  $-1$ , and  $0$  for the positive, negative, and zero 3-space curvatures respectively.

Let us briefly discuss these three cases separately.

(a) The spatial metric of the space with constant positive curvature can be rewritten as:

$$dl^2 = R^2(t) [dx_1^2 + \sin^2 x_1 (d\theta^2 + \sin^2 \theta d\varphi^2)],$$

where the new coordinate  $x_1$  is related to the  $r$  in (2.62) as  $r = \sin x_1$ . The parameter  $x_1$  takes the values  $0 \leq x_1 \leq \pi$ . The length of a circumference and the area of a sphere whose centre is at the origin in terms of  $x_1$  and  $\theta = \pi/2$  are  $l = 2\pi R \sin x_1$  and  $S = 4\pi R^2 \sin^2 x_1$ . As  $x_1$  increases both quantities first rise to a maximum ( $l_{\max} = 2\pi R$ ,  $S_{\max} = 4\pi R^2$ ), and then fall to zero. The radius  $\tilde{r}(x_1)$  of the circum-

ference or the sphere corresponding to  $x_1$  is  $Rx_1$ . Its maximum value is  $\tilde{r}_{\max} = \pi R$ . The ratio of the lengths of the circumference to the radius is  $l/\tilde{r} = 2\pi (\sin x_1)/x_1 < 2\pi$ .

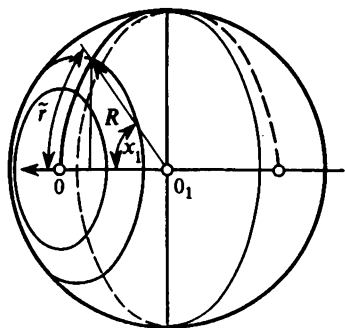


Fig. 11

This corresponds to the geometry on a three-dimensional sphere with a radius  $R$  in a four-dimensional Euclidean space (Fig. 11). The volume of such a three-dimensional space (3-sphere) is finite, viz.  $V = 2\pi^2 R^3$ , and hence these spaces are called closed or finite.

We recounted in Chapter One of this book how the geometry of these spaces was first developed by Riemann. It became the second example of a non-Euclidean geometry (after that of Lobachevski and Bolyai) and is often referred to as Riemannian geometry (in the narrow sense). Remember that Clifford and then Mach looked into the application of this spherical geometry to real spaces.

(b) The spatial metric of a constant negative curvature can be presented as

$$dl^2 = R^2(t) [dx_1^2 + \sinh^2 x_1 (d\theta^2 + \sin^2 \theta d\varphi^2)],$$

where the coordinate  $x_1$  is related to the  $r$  in (2.62) as  $r = \sinh x_1$ . Obviously,  $x_1$  may take any value from zero to infinity. The circumference and area of a sphere whose centre is at the origin in terms of  $x_1$  and  $\theta = \pi/2$  are  $l = 2\pi R \times \sinh x_1$  and  $S = 4\pi R^2 \sinh^2 x_1$ . As  $x_1$  rises they grow from zero to infinity. As before,  $\tilde{r} = Rx_1$  which is unbounded now, and  $l/\tilde{r} = 2\pi (\sinh x_1)/x_1 > 2\pi$ . The volume of the space is infinite. Therefore, spaces with a constant negative curvature are called open spaces. They are described by the Lobachevski geometry, which we discussed in Chapter One as the first illustration of a non-Euclidean geometry. Remember that the establishment of the relationship between the Lobachevski plane geometry and the geometry on hyperboloid by Felix Klein (1849-1925) and Eugenio Beltrami



(1835-1900) was the final proof of the consistency of this geometry.

(c) Spaces with zero curvature are described by Euclidean geometry. The form of this metric has the form of (2.62) with  $K = 0$ . Remember that in a Euclidean space  $l/r = 2\pi$ .

Thus each of the three types of geometry which have played an important historical role in our understanding of space and time (Euclid's, Lobachevski's, and Riemann's) are the three possibilities for spatial sections of our real four-dimensional universe, so far as the general theory of relativity is concerned.

If we include a temporal part, then the four-dimensional interval for a homogeneous cosmological model is

$$ds^2 = c^2 dt^2 - R^2(t) \left[ \frac{dr^2}{1 - Kr^2} + r^2 (d\theta^2 + \sin^2 \theta d\varphi^2) \right]. \quad (2.63)$$

This metric, in its most general form, is called in the literature the Friedmann-Robertson-Walker metric. Equation (2.63) must be substituted into Einstein's field equations with the source tensor (2.60) on the right-hand side using a pressure-density equation of state. We notice that at various stages of the expansion of the universe, the equation of state may be different as the interplay between matter and radiation changes when the universe evolves. The present stage is best described by the simplest equation of state for dust, that is, for a dust (incoherent fluid),  $p = 0$ . In the neighbourhood of  $R = 0$ , the predominant equation of state is  $p = \rho/3$  which is characteristic of radiation.

The solutions of Einstein's equations yield different expressions for  $R(t)$  for all three types of spatial sections. It can be shown that for closed worlds ( $\Lambda = 0$ ,  $p = 0$ ), the solution, which was first obtained by Friedmann, "oscillates" in time (Fig. 12, curve *a*). In parametric form this solution is

$$t = (R_0/c) (x^0 - \sin x^0), \quad R = R_0 (1 - \cos x^0), \quad (2.64)$$

where  $R_0$  is a constant.

Friedmann's model for an open world evolves in a qualitatively different way: it either monotonically expands or contracts (see curve *b* in Fig. 12 for the expansion). This behaviour is described by the following parametric formulae:

$$t = (R_0/c) (\sinh x^0 - x^0), \quad R = R_0 (\cosh x^0 - 1). \quad (2.65)$$

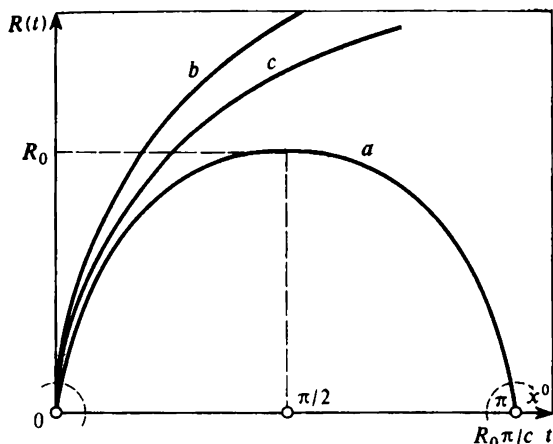


Fig. 12

We did not show the monotonically contracting model in Fig. 12 because it has been established experimentally that our universe is expanding.

The behaviour of an open model with zero-curvature spatial sections is shown by curve *c*. Note that when *R* vanishes, the density  $\rho$  becomes infinite. In the neighbourhood of these points (in the figure they are circled by broken lines), general relativity would appear to lose its applicability.

Einstein's equations contain the cosmological term  $\Lambda$ . It is the only additional term, which does not contradict their structure. Einstein had to introduce this term to obtain static solutions (the absolute cosmological principle!). Following the development of the FRW models, and after Hubble had discovered the expansion of the universe, this term was neglected for many years. Lately, however, it was found that its presence could be due to quantum processes, particularly during the early stages of the universe's expansion. V. Petrosian (USA) has shown that astronomical observations can only give an upper limit to the absolute value of  $\Lambda$ , viz.

$$|\Lambda| < 2 \times 10^{-56} \text{ cm}^{-2}. \quad (2.66)$$

Thus, the cosmological term is very small (if not zero), but it may be vital for the theory. Note that solutions other than Einstein's static model include a nonzero  $\Lambda$  for dust.

These include the models of De Sitter, Eddington-Lemaître, and Lemaître. At various stages of its evolution the universe may approach states quite close to these models.

Now, let us see what happens when  $\Lambda = 0$  and examine some details of the theory, deferring our concluding speculations until the following chapter. It is convenient to introduce the following notation:

(a) the Hubble parameter  $H(t) = \dot{R}(t)/R(t)$ ,

(b) the deceleration parameter  $q(t) = -R(t)\ddot{R}(t)/\dot{R}^2(t)$ ,  
and

(c) the density parameter  $\Omega(t) = \kappa\rho(t)/3H^2(t)$ .

They are normally evaluated for the recent epoch\*, given the subscript 0, e.g.  $H_0$ . An exact knowledge of these parameters helps to answer the question which may have occurred to the reader, viz. if the theory predicts three possible models of the world, which describes the reality? It follows from Einstein's field equations that there is a critical value of matter density,  $\rho_{cr} = 3H_0^2/8\pi G$  such that if it is compared with the observable average value  $\rho_0$  in the universe a decision can be made as to which model is true. If  $\rho_0 > \rho_{cr}$ , our world is closed, that is, the space in the comoving reference frame is Riemannian ( $K = +1$ ). If  $\rho_0 < \rho_{cr}$ , then the space in the frame is described by Lobachevski's geometry ( $K = -1$ ). Finally, for  $\rho_0 = \rho_{cr}$ , the frame is described by Euclidean geometry ( $K = 0$ ). At present, however, we still cannot get a definite answer.

Because of the difficulties in determining the distances in the universe, the value of the Hubble constant is only known approximately, i.e. we can bound it thus  $50 \text{ km/s.Mpc} \leq H_0 \leq 100 \text{ km/s.Mpc}$ . Since the lifetime of the Galaxy is from  $1.1 \times 10^{10}$  to  $1.8 \times 10^{10}$  years, the lowest bound of the Hubble constant is generally assumed to be the real value, but this is not a rigorous estimate. Hence, according to the formula for  $\rho_{cr}$ , the critical density is  $\rho_{cr} \sim 2 \times 10^{-29} \text{ g/cm}^3$ .

At present, the density of the observable part of the universe is believed to be  $\rho_0 \sim 5 \times 10^{-31} \text{ g/cm}^3$ . By extrapolating this value over the whole universe, we conclude that Friedmann's open model fits observational data best. Thus,

---

\* The term "recent epoch" in cosmology means the evolutionary stage of the universe in which we are living.

the spatial sections in the comoving reference frame are probably described by Lobachevski's geometry. However, it is possible that not all the matter in the universe has been accounted for, and that with the discovery of new objects and new types of matter, the average density may tend to the critical density  $\rho_{cr}$ . If we then assume that  $\Lambda$  may be positive, it becomes even more difficult at this stage to take a reliable choice of a best-fit model.

It is even more difficult to estimate the deceleration parameter. It could be determined by formally comparing the absolute and apparent magnitudes of extragalactic objects (we shall not give the relevant formula). Some researchers give a value  $q_0 = 1.6 \pm 0.4$ , but there is evidence that the errors in these measurements are much larger, so the probable bounds are  $-1 < q_0 < +2$ . Today, therefore, we only know for certain that our universe is expanding.

The "age" of our universe too cannot be fixed unequivocally—our estimate of the age depends on the model being used and the generalizations of the theory for the earliest stages of the universe. Recent estimates of the age are within some tens of thousands of millions of years, and this estimate will undoubtedly be made more precise in years to come.

Some of the other aspects of the universe's evolution are discussed elsewhere [54, 55, 56, 77, 78, 93, 96, 111]. Our cosmological studies of the universe have still to answer many questions, and in the following chapter we shall return to this problem in connection with the background microwave radiation and the early stages of evolution.

## **TOMORROW**

---

### **Modern Problems in the Theory of Gravitation. The Prospects for the Study of Space and Time**

When we were planning this book, we wanted at first to make each of the three chapters approximately equal in length. As our work proceeded, we found that Chapter Three, the one intended to cover the future of space-time physics, refused to fit this Procrustean bed. The future is always vaster in extent than the past, since the tiny bay of assimilated knowledge is a drop compared to the enormous ocean of the unknown. In addition, any prognosis of a future field of knowledge requires more words than a mature, elegant science, whose laconically stated principles unite many apparently isolated patterns and facts. We wrote this chapter in the knowledge that few of our readers will study every section in detail, simply because modern gravitational studies is so wide a subject area and our book is so short a presentation that we could not pay equal attention to every aspect of gravitation theory.

It is true that gravitation has maintained the interest of researchers at every stage in the history of physics, stimulating the formulation of new concepts in related fields. This is evident nowadays, but even in the past, Coulomb's law was modelled on Newton's gravitational potential formula. As new experimental techniques and the sensitivity of instruments have improved, theorists have started to compare the various theories of gravitation, beginning from those that were developed at the dawn of relativity theory. New theories of gravitation have now emerged which also need some comparative verification. Hence some researchers (Professor C. M. Will is one of the most active) have classified the theories of gravitation and formulated a deliberately phenomenological "parametrized post-Newtonian (PPN)

formalism" to analyze those theories which describe most realistic situations.

We cannot say that a phenomenological approach is satisfying, but it is used under the pressure of circumstances, when some of the steps in the investigation have to be done by feel. We believe it is worthwhile here, therefore, to consider how a physical theory is constructed, though we acknowledge our discussion is not exhaustive. In a very simplistic way, we can say that a theory is conceived by assimilating the information Man accumulates in the course of his practical activities. Naturally, this information is arranged in systematic patterns, and the idea that a theory is feasible is in itself a discovery that follows from an unexpectedly fruitful assimilation. Chaldean astronomy was first related to "earthly", pastoral needs, and, at the same time to "poetical" requirements that satisfied the Chaldeans religious aspirations (we would have said "heavenly" and not "poetical" had it not here been tautologous). The Euclidean geometry at this stage manifested the high level of human intellect and the maturity of Man's needs. But even now we are surprised that this very refined system of theoretical knowledge was in contrast to the contemporary low and primitive levels of technology and experimentation. The emergence of this theory was evidently stimulated by the earthly and spiritual needs of Man. On the one hand, postulates or axioms irreducible to more elementary ones were clearly defined. On the other, many new statements had been deduced from them by logical reasoning. These new statements or propositions had unprecedentedly rigorous derivations but were clearly related to recognizably real facts, and seemed to occur from necessity rather than by chance or coincidence. The most important objective of any theory, Einstein felt, was to have as few and as simple basic irreducible elements as were possible without contradicting experience. The practical application of such a theory solidifies our belief in its usefulness and encourages attempts to expand its boundaries so as to include new fields of knowledge and to establish relationships between them.

Another discovery was the establishment of affinities among apparently such different sciences as mechanics, electromagnetism and gravitation, and, later, quantum theory. In the past each new step in the development of science only came a few centuries after the previous one,

but nowadays the pulse of science may be called more tachycardiac than normal.

In modern times, the verification of a theory requires a wider variety of techniques: in addition to laboratory experiments, astronomical observations and practical application in production processes, a new theory must be checked for its relationship with previous theories in a given area, and with the theories in contiguous fields of physics, and with the structure of science in general. We call all this "human experience", and it undoubtedly includes the practice of constructing scientific theories. Man has now reached a stage again, at which the practical application of a theory is at a higher and more comprehensive level than was the case in the past.

The notion of the practical verification of a theory needs to be explained in a bit more detail, and we need to say more about how a theory is extracted from observations and experiment. No experiment can be interpreted unambiguously. On the one hand, the abstract presentation of real events with any given degree of precision (experiment always imposes this limitation!) is pluralistic. Nevertheless, a theory rarely has more than a few variants (remember matrix and wave quantum mechanics) because general attention is drawn to the first formulations. On the other hand, an experiment, as it is, does not produce either formulae or principles that make a theory (this is a fundamental truth). When a researcher passes to principles, his intuition becomes his most important creative tool. This and his imagination, which of course must be confined to real events, can derive basic principles that reach far beyond the limits of his epoch from the kaleidoscope of scattered facts and events. As the brilliant physicist and mathematician E. Wigner pointed out, the universal law of gravitation which Newton discovered without wishing so and which he could verify to only 4% has turned to be accurate to 0.0001%. It has become so closely associated with the notion of absolute precision that it was only recently that physicists dared to evaluate the accuracy limits. Other examples cited by Wigner are Maxwell's electrodynamics and quantum mechanics (Lamb's shift, for instance) [113]. These were not random deviations in the history of science but normal developments, the manifestations of the law of the interaction between conception and Nature. A scientific theory must be

governed by methodological regulators which include (we follow L. B. Bazhenov's terminology) verifiability, in principle, maximal generality, forecasting power, simplicity, and consistency. He also notes that a confirmation of theoretical consequences by experiment is valuable only if these consequences can equally be disproved by experiment. Any "confirmation" in an experiment of consequences which can in no way be disproved by the experiment is no confirmation at all [4]. The importance of this principle of the falsifiability of a scientific theory was established by K. Popper, but it is more fully developed in the theory of methodological regulators, though we can dwell no longer on this topic.

One such regulator reflects the quality of a physical theory we often call "beauty". Remember Dirac's dictum: "Physical laws should have mathematical beauty". The general theory of relativity is an impressive example in this respect: such that in contrast to the PPN formalism, it does not have a single extraneous parameter (if we do not consider the gravitational constant itself which has a dimension\*). Einstein's general relativity theory is an example of how from an extremely small number of facts (though very profound ones such as relativity principle) an elegant system of absolutely new and unexpected knowledge can be extracted, a system whose completeness and elegance challenges those of Euclidean geometry. Should the reader object that Euclidean geometry is mathematics and not physics, we have to disagree because Euclidean geometry purports to describe the world around us and exists as an objective system of knowledge. This means that it must be classified with the theories of physics (it is not phenomenological!).

### 3.1. GRAVITATIONAL WAVES

**3.1.1. The Elusive Radiation.** In the modern theory of gravitation the topic of gravitational radiation, its properties, its effect on measuring devices, and its possible sources has been given a great deal of attention. Einstein started the process in 1918, immediately after he had created the general theory of relativity. The very first estimates showed that

---

\* It is expected that this constant will not appear in future more general theories. In supergravity theories, for example, the constant arises when the conformal symmetry of the initial theory is spontaneously broken.



any real gravitational radiation must be exceptionally weak, and for a long time it was believed, quite rightly, that detecting it experimentally would be very difficult and could only be realized in the distant future.

In the early 1960s, however, Joseph Weber, an American physicist, announced an experimental program to detect gravitational waves of extraterrestrial origin. At the Third International Conference on General Relativity and Gravitation in Warsaw, in 1962, he described his detector. Three years later, at the next conference in London, Weber looked very tired. He had intensified the search for gravitational waves in order to present findings at the conference. His setup was very sensitive and responded to the noise produced by traffic and factories, so he had had to work at nights. No waves had been discovered by 1965, but Weber's experiments aroused the scientific community. A number of laboratories around the world began planning similar experiments. Finally in 1969, Weber announced his discovery of gravitational radiation. Scientists all over the world became excited. Weber, meanwhile, continued to present new findings. He reported that gravitational signals were being received with increasing regularity: once a month, twice a month, and then even more frequently. The direction from which the signals were coming (the centre of our Galaxy) and their polarization characteristics were now being determined. The signals had been received by several instruments at once and the readings processed by computers—the project seemed to be running quite satisfactorily. That was when the news hit the headlines.

The physics community responded to this in a variety of ways. Some of the participants of the Sixth International Conference on General Relativity and Gravitation in Copenhagen, in 1971, for example, were confident that Weber had discovered gravitational waves. Some people even shouted: "Hura! Bravo! Weber, the discoverer of gravitational waves!" Many theorists began working on models of the sources of the sort of radiation that could be detected by Weber's instruments. Quite a number of astrophysical mechanisms and phenomena were invented to fit the situation. There were, however, more careful researchers who remembered that any discovery must be verified in other laboratories before it can be claimed as valid. There were also outspoken sceptics.

Many laboratories rushed to repeat Weber's experiments. This was not, however, an easy endeavour since devices more sensitive than anything produced before had to be designed. The first group to reach the necessary sensitivity (in 1972) were the scientists at the Braginsky laboratory at Moscow State University. But the results of their experiments turned out to be negative. Then, reports from other laboratories (USA, UK, and Italy) also failed to confirm Weber's experiments.

More than ten years have passed and the experiments continue. Second-generation gravitational detectors have nearly been completed, still most physicists in the field do not think that gravitational waves have been discovered. We should not be critical of J. Weber for his failure, for every step forward in science is a complicated venture. Nevertheless, the broad scale of the gravitational experiments in progress now in the world is to his credit.\*

In parallel with the experimental work, theoretical research on gravitational waves intensified, and in particular on whether there are such things as gravitational waves because the theory also contains many unclear points. We shall briefly touch on some of the dubious aspects here.

Much of gravitation theory is analyzed by analogy with electromagnetic field. This is also true for gravitational radiation. A closer look, however, shows that the latter possesses a number of qualitatively different properties. It does not affect the measuring instruments in the same way, for example. When an electromagnetic wave interacts with a charged body, the motion of the body can be viewed relative to a neutral object that is not affected by the electromagnetic fields. With regards to gravitational waves, the very possibility of such a comparison is excluded because of the equivalence principle. We must therefore deal with the curvature of space-time which implies that all objects at a given place are equally disturbed (displaced). This universal kind of interaction requires a special method of description.

In the framework of the general theory of relativity, a profound unification of gravity and inertia has been attained. The equations of motion for a material point (the geo-

---

\* Weber continues to claim that his instruments are recording signals, but he is now more reticent to claim that they are due to gravitational radiation.

desic equation) contain quantities which describe both gravity and inertia rather than gravity alone. This situation is rather peculiar because in a curved space-time we cannot, in principle, distinguish which quantities describe inertia, and which "pure" gravitation. Therefore, it is more correct to speak about graviinertial rather than gravitational waves.

In descriptions of electromagnetic and other kinds of radiation, the notions of energy and momentum are widely used. For example, the terms "radiation power", "radiation flux density", "wave energy", etc. are only too frequent. Incidentally, in general relativity theory, these terms are, strictly speaking, incorrect when applied to a gravitational field. The reason is that the notions of energy and momentum are important only so long as their conservation laws are valid. In a flat space-time, the conservation of energy depends on the homogeneity of time; the conservation of momentum on the homogeneity of space, and the conservation of angular momentum on the isotropy of space. General relativity interprets gravitation in terms of the curvature of space-time, in which case the homogeneity and isotropy of space are thus violated. In other words, the basic principles of general relativity do not lead to the traditional laws for the conservation of energy, momentum, and angular momentum. Hence, they cannot be applied in their former sense to a gravitational field in general and to gravitational waves in particular. They may be valid in some special cases, but only in a limited and rather tenuous sense of the word. Although we often see in the literature estimates of the power or flux density of gravitational radiation, other authorities insist that these notions and formulae still need clarification. The topic is therefore very much a point of debate.

Electromagnetic waves in vacuum are as we know described by Maxwell's linear equations, whereas gravitational waves require Einstein's essentially nonlinear equations. This is one more source of difficulty. It is clear, for example, that in general relativity, the superposition principle could not hold for gravitational waves. It can be used when substantial simplifications are made, i.e. for a weak field, in which the nonlinear terms are ignored.

Calculations have shown that in all real situations gravitational waves are so weak that they cannot be detected, at least, by existing instruments. The experimental search

at present under way is for hypothetical superpowerful sources. As a result, theorists have been considering effects which have not yet been observed. This is a unique situation, unlike what happens in other sciences in which experimental detection of a phenomenon usually precedes its scientific description. However, almost all of the phenomena due to general relativity were first predicted theoretically.

In spite of the above, the prospects in this field of research are not hopeless. A number of important and interesting results have been obtained, and several new schools of thought about the theory of gravitational waves have formed. In the near future we would expect some crucial experimental discoveries.

Before considering the relatively simple case of a weak gravitational field, we shall briefly outline the state-of-the-art in the general theory. At present, quite a number of exact solutions to Einstein's equations have been derived, but only a few have any clear physical interpretation, and of these only a few can be classified as wave solutions. There are a number of criteria which, when applied to a solution, identify it as one pertaining to a space-time with gravitational waves. These criteria, however, do not exactly coincide, but they do outline (rather well) the contours of all the sets of cases in which a gravitational field can be regarded as a wave one. This part of the theory is rather complex and it uses a very powerful body of mathematics. Even so, we have to confine ourselves to approximate solutions of Einstein's field equations for weak gravitational fields, because in this case it is easier for us to reason by analogy with the theory of electromagnetic fields, which has been thoroughly proved experimentally.

**3.1.2. Weak Gravitational Waves.** Using the theory of electromagnetic fields as a basis, it was asserted that far from their source, gravitational waves must be described by a solution of a linear approximation of Einstein's field equations. This means that far from the source, the metric will be  $g_{\mu\nu} = g_{\mu\nu}^0 + h_{\mu\nu}$ , where  $g_{\mu\nu}^0$  is a flat space-time metric, and  $h_{\mu\nu}$  is a correction small in comparison with unity. Using auxiliary coordinate conditions, for example, harmonic coordinates (for  $h = h_{\mu\nu} g^{\mu\nu} = 0$ ), we get

$$\partial h^{\mu\nu} / \partial x^\nu = 0, \quad (3.1)$$

and it can then be shown that Einstein's field equations for a vacuum take the form

$$\left( \frac{\partial^2}{\partial x_0^2} - \frac{\partial^2}{\partial x_1^2} - \frac{\partial^2}{\partial x_2^2} - \frac{\partial^2}{\partial x_3^2} \right) h_{\mu\nu} = 0, \quad (3.2)$$

i.e. that of the wave equation.

Looking for the solution for the case of a plane wave, we obtain

$$h_{\mu\nu} = A_{\mu\nu} \sin(k_\alpha x^\alpha + \varphi_0), \quad (3.3)$$

where  $k_\alpha$  is the null wave vector. From the auxiliary conditions (3.1) it follows that only two of the ten components  $h_{\mu\nu}$  are independent. This means that if the plane wave propagates along  $x^1$  axis, the independent components will be the transversal-transversal components,  $b_1 = (1/2) \times (h_{22} - h_{33})$  and  $b_2 = h_{23}$ , that is, the gravitational wave is transversal and has two polarizations. This is like a plane electromagnetic wave with its two characteristic transversal (relative to  $x^1$  axis) components of the vector potential,  $A_2$  and  $A_3$ . Thus the metric of a plane gravitational wave is

$$ds^2 = dx_0^2 - dx_1^2 - (1 + b_1) dx_2^2 - (1 - b_1) dx_3^2 - 2b_2 dx_2 dx_3, \quad (3.4)$$

with

$$\begin{aligned} b_1 &= B_1 \sin \left[ \frac{\omega}{c} (x^0 - x^1)_t \right] \\ b_2 &= B_2 \sin \left[ \frac{\omega}{c} (x^0 - x^1) + \varphi \right], \end{aligned} \quad (3.5)$$

where  $B_1 \ll 1$  and  $B_2 \ll 1$  are the wave amplitudes.

Consider how a weak gravitational wave would act on free test particles (for example, particles of dust). We assume that the dust particles have point masses distributed before the arrival of the wave (i.e. when  $B_1 = B_2 = 0$ ) at a certain distance  $\delta\sigma_0$  from a given point 0. These particles lie on a sphere of radius  $\delta\sigma_0$ :  $\delta x_1^2 + \delta x_2^2 + \delta x_3^2 = \delta\sigma_0^2$ . When the wave arrives at point 0, the distances to the particles become\*

$$\begin{aligned} \delta\sigma &= \sqrt{h_{ik} \delta x^i \delta x^k} \\ &= \sqrt{\delta x_1^2 + \delta x_2^2 + \delta x_3^2 + b_1 (\delta x_2^2 - \delta x_3^2) + 2b_2 \delta x_2 \delta x_3}. \end{aligned}$$

\* The coordinates  $x^1$ ,  $x^2$  and  $x^3$  of the dust particles remain constant.

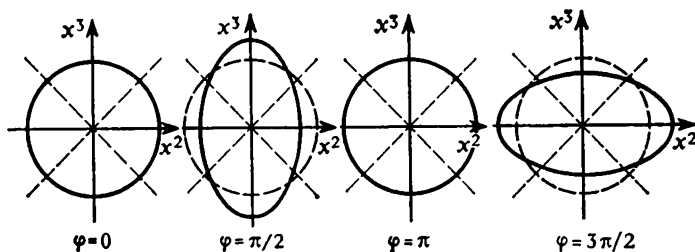


Fig. 13

For simplicity, let this wave be linearly polarized, that is,  $b_1 \neq 0$  and  $b_2 = 0$ . Then, in the coordinate space, the surface of the points equidistant from 0 (with  $\delta\sigma = \delta\sigma_0 = \text{const}$ ) will form an ellipsoid. A section  $\delta x_1 = 0$  (on the wave front) is an ellipse

$$\frac{\delta x_2^2 (1 + b_1)}{\delta \sigma_0^2} + \frac{\delta x_3^2 (1 - b_1)}{\delta \sigma_0^2} = 1.$$

As the wave, with a  $b_1$  determined by (3.5), propagates, the ellipse changes, a full cycle of changes during a wave period being shown in Fig. 13. From this we can see that during the first quarter cycle, the wave front expands vertically and contracts horizontally. During the second quarter cycle, the system gradually returns to the initial state. During the third quarter cycle, the wave plane expands horizontally and contracts vertically, and so on. Note that the points along the broken lines do not move.

If the wave has another polarization (that is,  $b_1 = 0$  and  $b_2 \neq 0$ ), the sphere of dust particles will evolve in a similar way, except that the expansion and contraction will occur along the broken lines (the figure should be rotated at  $45^\circ$ ), and the distances along the solid lines will remain constant.

If we put onto this plane a set of constrained points (in the simplest model these could be balls linked by springs), then cyclic strains should develop in the springs and hence it should be possible to detect gravitational waves. This simple model was the one used both in the first-generation gravitational antennae and is being incorporated in a number of the second- and third-generation devices now under development.

According to the theory of weak gravitational waves any system that had a changing quadrupole moment  $D_{ik}$ , i.e.

a system, in which the variations in the mass distribution are not spherically symmetric, can act as a source of gravitational radiation. If we consider weak gravitational waves to be analogues of electromagnetic ones, we can derive a formula for the energy loss rate due to the gravitational radiation, i.e.

$$-dE/dt = (G/45c^5) (\ddot{D}_{ik})^2, \quad (3.6)$$

where  $\ddot{D}_{ik}$  is the third time derivative of  $D_{ik}$ . In contrast to electromagnetic radiation, which basically has a dipole nature, gravitational radiation is quadrupole in nature.

To illustrate this, we give here a formula for the energy loss rate due to the gravitational radiation for a rod of length  $2l$  and mass  $M$  rotating with angular velocity  $\omega$ :

$$-dE/dt = \frac{128G}{45c^5} (M/l)^2 \omega^6 l^6.$$

In real situations, i.e. for  $\omega \sim 10^3 \text{ s}^{-1}$ ,  $l \sim 2 \times 10^3 \text{ cm}$ , and  $M \sim 6 \times 10^{10} \text{ g}$ , we have a negligibly small value of  $(-dE/dt) \sim 10^{-7} \text{ erg/s}$ .

The intensity of gravitational radiation is sometimes evaluated in terms of its energy flux density, which is measured in  $\text{erg/cm}^2 \cdot \text{s}$ . This is not necessary, however. If the space-time of metric is known, the system's behaviour can be described in the metric using general relativity theory alone. Gravitational waves can be most conveniently characterized by the relative deformation of a system, that is, by a dimensionless quantity  $h = \Delta l/l$ , where  $l$  is a linear dimension of the system in the absence of gravitational wave, and  $\Delta l$  is the change in the length due to the wave. Obviously, we would expect  $h$  to be very small ( $h \sim 10^{-20}$ , see Sec. 3.1.5).

**3.1.3. Sources of Gravitational Radiation.** It follows from the above estimates that we could hardly at present expect to construct a terrestrial source of gravitational waves, that is, to perform a Hertz-type gravitational experiment.\* For now, our main hope is cosmic sources. They can be grouped into two classes according to intensity and mode of the radiation, viz. (1) continuous radiation sources (e.g.

---

\* In 1888 Heinrich Hertz (1857-1894) experimentally generated electromagnetic waves and detected them under laboratory conditions.

binary stars, pulsars), and (2) burst sources (e.g. nonsymmetric star collapses, supernova, galactic nucleus collapses). The sources in the second category are on average 15 orders of magnitude more powerful than those in the first, but unfortunately it is impossible to predict when or where they might occur. We shall look briefly at these two categories.

**1. Continuous radiation sources.** These include planetary systems, rotating binary stars, and pulsating and rotating stars. About 50% of the stars in our Galaxy are probably binary (or multiple) systems, and so there should be an enormous number of sources of continuous gravitational radiation. The question is only how strong the flux of this radiation may be in the Earth's neighbourhood. Obviously, the flux  $F$  is related to the power  $L$  (erg/s) of the gravitation radiation, viz.  $F = L/4\pi R^2$ , where  $R$  is the distance from the source to the Earth. If the gravitational wave, with frequency  $\omega_g$ , is monochromatic, the dimensionless quantity  $h$  can be given as

$$h \sim [16\pi G F / c^3 \omega_g^2]^{1/2}.$$

The scientific literature now contains many papers devoted to the calculation of  $F$  and  $h$  for specific sources. Rather than listing them, we found it easier to summarize them in a chart (Fig. 14). The logarithm of the frequency (in hertz) is

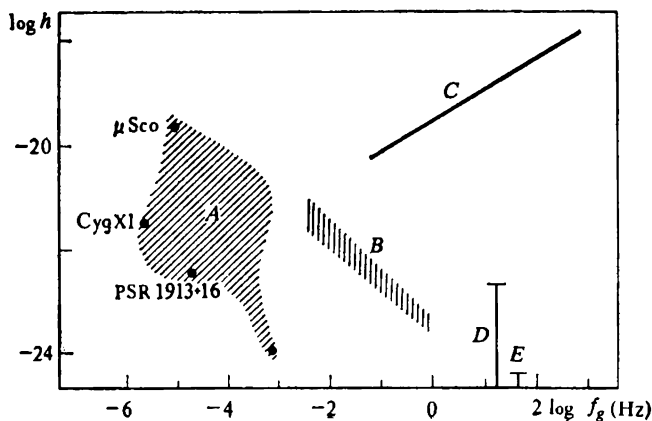


Fig. 14



plotted along the horizontal axis and the logarithm of  $h$  is plotted up the vertical axis. The sources fall into four domains as follows: (1) the classical binary stars (X-ray sources including) in domain  $A$ ; (2) domain  $B$  corresponds to white dwarfs after their emergence from novae ( $E \sim 10^{45}$  erg, distance from the Earth is  $R \sim 500$  pc); (3) domain  $C$  corresponds to compact multiple stellar systems (each component having masses of the order of the magnitude of the mass of the Sun, and  $R \sim 10^4$  pc); (4) while points  $D$  and  $E$  correspond to the Vela pulsar and the pulsar in the Crab Nebula.

More detailed estimates and references to the original publications can be found elsewhere [12, 69, 107, 114].

Typically, the sources of continuous gravitational radiation are extremely weak, and their frequency range lies well beyond the capabilities of the detectors presently being developed. Notwithstanding this we have strong indirect evidence of the existence of gravitational radiation. J. H. Taylor and J. M. Weisberg, for example, summarized the results of a seven-years' study of the binary pulsar PSR 1913 + 16 which they began in 1974. The pulsar has a mass of  $1.42 \pm 0.06 M_{\odot}$  and an invisible component (a neutron star or a black hole) of approximately the same mass. The orbital revolution cycle of the system is 7 hours and 45 minutes, and the pulsation period is 0.059 s. The observations were made at the Arecibo 305-m-dish radiotelescope in the 430 and 1410 MHz bands. It is reliably established that the system's orbital period is being shortened because of the lowering of the orbits due to the energy loss via gravitational radiation. The dimensionless relative rate of change of the period was found to be  $(-2.30 \pm 0.22) \times 10^{-12}$  whereas the linearized general relativity theory predicts a value  $-2.40 \times 10^{-12}$ .

**2. Burst sources.** To deal with this category we must introduce some additional parameters. Suppose a source has a mass  $M$  and an energy  $E_g$  is emitted as gravitational waves. Let  $E_g = \varepsilon M c^2$ , where  $\varepsilon$  is an efficiency factor, and we assume that  $\varepsilon < 0.5$ . The peak of the emission is characterized by time  $\tilde{\tau}$ , which determines characteristic frequency (radiation maximum for the whole spectrum)

$$f_g = (2\pi\tilde{\tau})^{-1} \simeq (4 \times 10^3 \text{ Hz}) (2M_{\odot}/M).$$

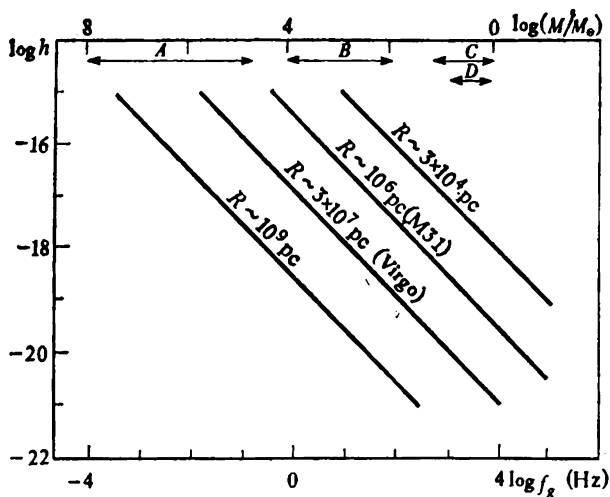


Fig. 15

The dimensionless parameter  $h$  can be written as

$$h \sim 4GM\epsilon^{1/2}/Rc^2.$$

As for continuous radiation sources, the radiation fluxes from burst sources in the neighbourhood of the Earth can be summarized in the form of a chart (Fig. 15). The same variables are plotted along the axes as in Fig. 14, and we assume  $\epsilon \sim 0.01$  throughout. Then, (1) domain  $A$  corresponds to galactic nuclei and quasars, (2) domain  $B$  represents star clusters and black holes, (3) domain  $C$  corresponds to black holes coming into being, and (4) domain  $D$  corresponds to neutron star formation.

In order to evaluate the strength of a pulsating source of gravitational radiation, we must take into account two more factors: (a) the frequency at which the radiation peaks of a given type occur (how many times may they be expected in one year), and (b) the location of the source of the radiation. Events of this sort are both extremely rare and open to various hypothetical interpretations. Even so the detectors now under development are being designed for just such sources.

**3.1.4. Gravitational Radiation Detectors.** The first gravitational-wave detector was built by J. Weber of Maryland in

the USA. It was a massive aluminum cylinder suspended by special means in a vacuum. Piezoelectric strain-transducers around the middle of the cylinder transformed any mechanical strain caused by oscillations into electric signals. The idea is very simple. When a gravitational wave hits the cylinder perpendicular to its axis, it must cause mechanical oscillations that are converted into electric signals which are then amplified.

The story goes that when Weber was specifying the parameters for his bar he was asked about the dimensions at workshop. Weber spread his arms and said about this long and this wide. As a result, a cylinder about 1.5 m long and 60 cm in diameter, weighing one and a half tonnes was made. Its natural oscillation was 1660 Hz, which was thus the frequency of the gravitational radiation that the setup could optimally detect. One has to begin with something, and this something must be technically simple. It was clear from the very start that the detector's frequency and the maximum attainable sensitivity fell far short of those necessary to detect continuous radiation from realistic sources. Weber had to rely on hypothetical burst sources which were suggested as abounding in the Universe by theorists.

Weber tuned his detector to register longitudinal oscillations with an amplitude of  $2 \times 10^{-14}$  cm. This corresponds to a dimensionless quantity  $h$  of about  $10^{-17}$ . About the same sensitivity was attained by other scientists on similar detectors. We called these instruments first generation detectors. The negative results from this first stage in the search for gravitational waves made the experimenters improve the sensitivities of their devices.

This engineering problem is now being attacked in several laboratories around the world, but any detailed discussion of the engineering problems goes beyond the scope of this book. We will only briefly indicate that the main obstacle is the background noise of the setup caused by: (a) Brownian movement in the cylinder itself, (b) noise from measuring and amplifying equipment, and (c) sharp pulses due to recrystallization of the aluminium in the cylinder. The Brownian noise in a cylinder of unit length is

$$(\Delta l/l) \sim (kT\omega_{\text{res}}\tau_{\text{mes}}/\pi M v_s Q)^{1/2}, \quad (3.7)$$

where  $\omega_{\text{res}}$  is the resonance frequency of the cylinder,  $\tau_{\text{mes}}$  is the period of measurement,  $v_s$  is the speed of sound

in the cylinder,  $T$  is the absolute temperature (in kelvin) of the bar, and  $Q$  is a mechanical quality factor. The quality factor indicates how many times the energy of the detector decreases per oscillation cycle.

It is obvious from the above formula that to reduce  $(\Delta l/l)$  we should increase the mass and decrease the temperature. This was the path taken by W. Fairbank with a team of scientists at Stanford University in the USA. Another approach would be to increase the mechanical quality factor  $Q$ . For this, aluminum, with  $Q \sim 10^5$ , should be replaced with some other material. The best materials would be monocrystals of quartz, silicon, niobium, or, the best of all, sapphire, whose quality is  $Q \sim 5 \times 10^9$ . This last approach was the one taken by V. B. Braginsky's group (Moscow State University) and D. Douglass (Rochester, USA). Of course, it would be desirable to improve the design in all above parameters, but this can only be done at the next stage of research [12].

There are some other technical problems. For example, other sorts of circuit for measuring and amplifying weak oscillations in the bar have to be developed because Weber's piezoelectric strain-transducers are no longer adequate. Braginsky used a capacitance readout system from the start of his research. It was shown that these transducers can in principle attain an  $h$  of around  $10^{-20}$ . The second-generation detectors now under development are expected to attain real sensitivities as good as  $h \sim 3 \times 10^{-19}$ .

A new type of gravitational-wave detector based on a laser interferometer is now being developed. This sort of device will be used to measure the relative oscillation of mirrors stationed at considerable distances apart. They are now under development in the USSR, the USA, UK, and West Germany. The enormity of the problem can be illustrated by pointing out that to detect radiation from supernovae in the Virgo cluster of galaxies requires a sensitivity of  $h \sim 10^{-21}$  which in turn requires a laser with some  $10^4$  W of power. Laser detectors have the other important advantage that they can be used over a wide frequency band.

The Doppler effect in combination with spacecraft tracking can be used to detect low frequency gravitational radiation (wavelength  $\lambda > 10^6$  km), but the radio receivers now available are not sensitive enough.

Second-generation gravitational-wave antennae should be

sensitive enough to detect gravitational waves from exploding supernovae in our Galaxy. These, however, occur very rarely: once every ten-thirty years.

All these considerations indicate that in order to obtain positive results, even with the second-generation detectors that are expected to be operational by the time this book is published, some good luck is needed. Great hopes have been pinned on the third-generation antennae, which should be able to detect displacements with  $h \geq 10^{-21}$ . At this accuracy, however, and, according to the Heisenberg uncertainty principle, longitudinal oscillations of a cylinder cannot be measured more accurately than with  $h$ 's of  $10^{-20}$ . As a result, the problem of developing a nondemolition measurement technique emerges. V. B. Braginsky and K. S. Thorne believe that the accuracy of the measurements can be as high as  $h \sim 10^{-21}$ , or even better.

We have therefore today some quite solid reasons for believing that gravitational radiation from cosmic sources may be detected in the near future. If no new factors enter the scene, we may expect that, with the discovery of gravitational waves, humanity will have a new information channel for studying the Universe; the science of gravitational wave astronomy will have been born.

## **3.2. BLACK HOLES AND RELATIVISTIC ASTROPHYSICS**

**3.2.1. Evolution of Stars.** The vacuum solution of Einstein's gravitational field equations which we discussed in Sec. 2.2, viz. the exterior Schwarzschild solution, can be interpreted as the gravitational field extending outwards from a star, if its rotation is negligibly small and the star is large enough for the part of the solution corresponding to the Schwarzschild radius and its neighbourhood to be dropped. The Schwarzschild field was therefore the first practical application of Einstein's equation in relativistic astrophysics. However, the interior solution, which Schwarzschild himself obtained, is also interesting. It holds when the right-hand sides of Einstein's equations reduce to a nonzero energy-momentum tensor. The tensor was chosen to be that of a perfect fluid and is similar to the one used in cosmology (see (2.60)). All of its components are functions of the radial coordinate only, and the 4-velocity vector of the medium is

for the medium at rest. This model of a star is also valid if its matter is a superdense state and Newton's standard law of gravity no longer holds. The more general case, dropping Schwarzschild's limitation to an incompressible fluid, can also be formulated. However, to be realistic when estimating the behaviour of stellar matter as the pressure rises, numerical or approximate solutions of the field equations and the equations of state are essential. It is important, though difficult, to select the equation of state that corresponds to reality since we scarcely know the properties of matter under the gigantic pressures that cause the density to approach or exceed nuclear density. This is where astrophysics converges with nuclear physics and the physics of elementary particles, and where cosmic objects have to be analyzed using quantum theory.

As we move in from the surface of a massive star, a number of distinctive domains are crossed. The density of the matter rises above that standard on the Earth, the potential barriers between individual nuclei become narrower, and the probability of particle exchange between the nuclei due to quantum tunnelling increases. At a certain density, the electrons become common to all the atoms, that is, practically free. This is when what is called the degenerate electron gas appears (along with atomic nuclei "gas"). This phase of matter is described by a special equation of state (the density of the matter in this domain is about  $10^7$  g/cm<sup>3</sup>).

Closer to the centre of star, at still higher densities, the protons swallow the electrons yielding neutrons. A neutron in its free state is unstable and tends to decay into a proton, an electron and a neutrino, but at great pressures, a neutron becomes stable, whereas the protons tend to merge with electrons to form neutrons and emit neutrinos. This is more favourable in energy terms. Now the whole system consists of three components: a relativistic degenerate electron gas (though less of it), an atomic nuclei (also less), and a degenerate neutron gas (the neutrinos "evaporate" from stars). As the star's density reaches  $10^{14}$  g/cm<sup>3</sup>, the neutron gas starts to prevail in the star, although some electrons and protons (but no nuclei) are left. The pressure-density equation of state (see Sec. 2.10) is not exactly known for this case since it should take into account the nuclear forces in the relativistic domain and the creation of new particles.

For this stratum of a star, the equation of state is normally written in the form of a table or complex phenomenological formulae, the high-density ( $\rho \geq 10^8 \text{ g/cm}^3$ ) equations of state drawn up by Harrison, Wheeler, and Sahakyan are examples. In practice, however, the simpler asymptotic equation of state used above in the cosmological section, i.e. the equation for ultrarelativistic ideal gas,  $p = \rho/3$ , can be used.

We assume that Einstein's field equations have been integrated for a given equation of state (the latter changing from stratum to stratum); this should not be difficult if a numerical technique is applied. The solution that will emerge will be static in structure, that is, it will be an equilibrium solution, but it may not be stable, and so it should be checked for stability. Typically, this is done by adding small disturbances dependent on time to the static solution, substituting the sum into Einstein's equations, and solving the resulting equations to reveal whether the disturbances tend to grow and so infer whether the equilibrium solution is stable. In practice all these calculations are made on a computer. The results for the Harrison-Wheeler-Sahakyan equation of state are shown graphically (Fig. 16) in terms of a

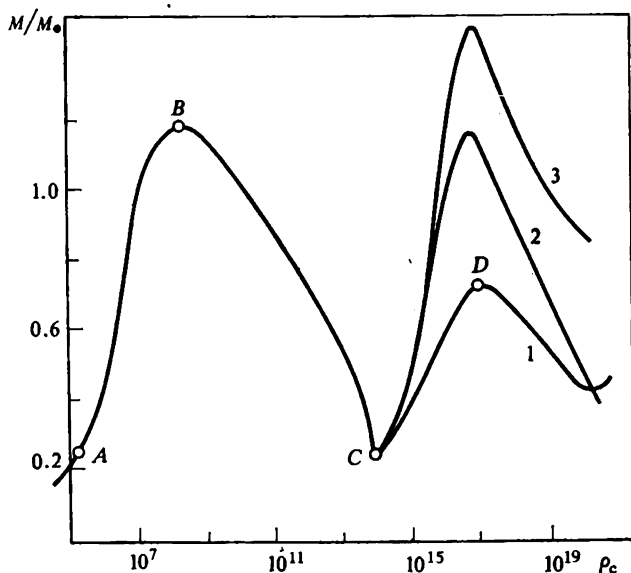


Fig. 16

given density  $\rho_c$  at the centre of the configuration as a function of total mass  $M$  (in Solar masses,  $M_\odot$ ). The relativistic domain lies to the right of point  $C$ . Three curves as drawn according to different theories are given in the diagram for comparison: (1) general relativity, (2) Newton's theory with an account for the gravitational mass defect, and (3) Newton's nonrelativistic theory. All the curves are almost the same qualitatively, but clearly the integral mass diminishes noticeably when the gravitation and relativistic effects are more accurately included. In the rising sections of the curve, the equilibrium configuration is stable, while in the descending sections, the total mass tends to "jump over" to the next minimum (the principle of energetically more favourable states). That is why points  $B$  and  $D$  represent the limits of the stability of the equilibrium configuration as the total mass grows. Point  $B$  is referred to as Chandrasekhar limit, and point  $D$  as the Oppenheimer-Volkoff-Landau limit.

The graph only covers supermassive stars, while the evolution of normal stars is described by a narrow strip just to the right of point  $A$ ; the states corresponding to white dwarfs are nearer to point  $B$ . These stars are still noticeably bright and superdense in earthbound terms, but they still have not reached the nuclear density. An example of a white dwarf is the famous satellite of Sirius, the brightest star in Earth's sky. We shall not go into the details of the evolution of normal stars as this is covered in a number of books both for the layman and the expert. We will only say that all stars are formed as a result of the gravitational compression of vast hydrogen clouds (these contain small quantities of helium and may contain other elements, which, unlike the hydrogen and helium, originated from the cores of stars that had already exploded). Under the compression, the gravitational potential energy is converted into the kinetic energy of thermal motion of the matter forming the new star, that is, the star gets hotter. At a certain temperature, the protons in the hydrogen nuclei start binding together to produce helium nuclei (neutrinos and positrons are also created): this is called the fusion reaction. Now the star really gets hot and lives for a long time using the energy of this reaction. The lifetime of a star depends on its initial mass: the larger the mass the shorter the helium synthesis phase lasts. Our Sun is at present in this helium synthesis



stage, which began several billion years ago and will end several billion years hence. The more massive stars (there are rather a lot of them in the Universe) pass the stage of helium synthesis in less than a billion years. It is only in the core of the star, where the temperature is especially high that the helium synthesis is possible, but when the hydrogen that feeds the reaction is exhausted in a massive star the internal pressure of the star falls, and the outgoing radiation weakens (and the reaction quenches). Since the mass of the star is no longer supported by the inner pressure the star begins again to contract under its own gravity. The gravitational potential energy is more than enough to warm up the star again and when the temperature in the core reaches a new critical level, a series of new fusion reactions begin. This time the helium nuclei merge to form the nuclei of carbon, oxygen, etc., up to iron, which is the most stable element in the Periodic Table. Of course, other elements that were not present in the initial cosmic gas clouds are also created in the star. But everything has its end, and these sources of energy are also exhausted (the end comes the sooner the greater the initial mass of star, because then the temperature rises faster due to the gravitational compression). Finally, if the star is smallish (the maximum is a little over the mass of the Sun), it enters a quiet "pension" age, contracting gradually as its thermonuclear furnaces shut down. This object is called a white dwarf.

In the range of densities corresponding to white dwarfs ( $AB$  on the graph in Fig. 16), the corrections introduced by general relativity into Newton's theory are small (of the order of 0.01%). But if the mass of the star is large it will "jump over" the maximum at  $B$ , and it will not become a white dwarf. What then happens is that the star explodes. The energy liberated by the gravitational contraction generates radiation with such a high pressure that star's external layers can no longer contain it. The outer layers (whose mass is, probably, much greater than that of the Sun) are then catapulted out with tremendous speeds into interstellar space. This is called a supernova. This catastrophic event transfers stars to the  $CD$  domain of curve 1. Here, they are "dead" objects, slowly cooling, with all the fusion reactions stopped. These stars consist of a degenerate neutron gas with the density of nuclei and are called neutron stars. The existence of such stars was predicted in the 1930s by Landau

and by Oppenheimer and Volkoff, but they were only discovered in 1967 as pulsars by radioastronomers in Cambridge. Pulsars emit a specific type of radio signals very much like the ticking of a clock (we mentioned pulsars in our discussion of the evidence for gravitational radiation in Sec. 3.1). The relativistic contributions to the description of stellar structures close to this second maximum—point *D*—is much more significant than those for point *B*. In general relativity, the instability of the star occurs for a central density of  $\rho_c \simeq 5 \times 10^{15} \text{ g/cm}^3$ , which is half the value predicted by Newton's theory. The critical mass of a neutron star is also lower at the instability point ( $1.2 M_\odot$  using Newton's theory compared to  $0.7 M_\odot$  using general relativity). It is true, however, that estimates of the critical mass are almost one hundred percent inaccurate because of the uncertainty in the equation of state for supernuclear densities of matter. Studies in this field may lead to a new criterion by which we can judge the validity of general relativity (along with the known criteria deduced from redshift effect, light deflection, and perihelion advance).

Pulsars characteristically rotate very quickly around their axes and have very strong magnetic fields (from  $10^{11}$  to  $10^{13}$  Gauss), but we cannot discuss these very interesting objects in any detail here. We will only mention that the exterior gravitational fields of rotating objects may be described by the Kerr metric, though with some ambiguity, whereas the exterior fields of spherically symmetric configurations are described by the Schwarzschild metric unambiguously. The rotation slightly increases the critical mass of a pulsar. But in any case, if the mass of a superdense configuration exceeds a certain critical mass (which is rather modest, measured in units of the mass of the Sun), the star passes over the second hunch of the curve in Fig. 16 (at point *D*), and after final contraction (collapse) becomes a black hole.

Before we start our discussion on black holes, we should note that in 1965, astrophysicists began to tackle another topic, i.e. the properties of relativistic star clusters. These objects were first studied by Ya. B. Zel'dovich and K. S. Thorne. A group of stars that held together by gravitational attraction but whose components rarely collide is called star cluster. It is unpractical and unnecessary to study the motion of each individual star in a cluster. There-

fore, a statistical description based on a distribution function for the stars in the cluster is used. The calculation is simplified if Boltzmann's distribution is used. As a result, star-cluster theory predicts equilibrium (stationary) configurations for relativistic star clusters; stability and instability conditions should be outlined on the next stage of research. Thus, astrophysicists hope to derive an explanation for explosions in the nuclei of some galaxies and of quasars. These are not so infrequent and are extremely catastrophic in nature.

**3.2.2. Black Holes.** After point *D* in Fig. 16, a star evolves towards an irreversible collapse and the formation of a black hole. In a black hole, which can no longer be called a star, all the information about the material of which it is made is lost. Incidentally, this loss of information means that a black hole can be described in terms of thermodynamic notions such as temperature and entropy. An understanding of what a black hole is can be obtained using Schwarzschild's field, which was derived in 1916. All the basic ideas we have about black holes could have been obtained at that time but due to a twist of fate the notion only appeared much later. The discovery and appreciation of the extraordinary structure awaited the studies of D. Finkelstein, M. D. Kruskal, I. D. Novikov, and P. Szekeres in the late 1950s and early 1960s. The terms "frozen" or "collapsed" star first adopted were succeeded by the new one "black hole". A substantial contribution to our understanding of the characteristics of the space-time of black holes and the physical fields around them was made by W. Israel and J. A. Wheeler (with their famous theorem that "black holes have no hair").

The history of the study of the Schwarzschild metric is an interesting example of the inertia of human thinking. We can now scarcely imagine how the researchers of the time (some of them were rather famous!), who knew the metric and the necessary mathematics, could not have noticed the special properties of Schwarzschild's space-time, properties characteristic for black holes in general. A similar situation occurred with the discovery of the laser (a quantum electromagnetic generator). For forty years the concepts of spontaneous and induced radiation, which Einstein also formulated, were known by everybody, but lay unutilized. In a new phenomenon we first tend to see the familiar things and to disregard the extremes, which may contain the essence of a

new theory or new principle (our logic and psychological attitude work in the same way when we try to identify familiar features in "riddle pictures").

Even when a "riddle picture" has been identified, we tend to lose it again. This is what happened when P. S. Laplace first predicted the existence of black holes late in the 18th century. He equated the second cosmic velocity of a super-massive body to the velocity of light. The reader can easily duplicate his calculations and find the minimum speed at which a particle can escape the surface of a body with a radius  $r_c$  and a mass  $m$  and move to infinity. If this velocity is assumed to be equal to that of light, we get

$$r_c = 2Gm/c^2, \quad (3.8)$$

which is the value of the Schwarzschild radius (cf. Sec. 2.2).

Of course, this prediction which depends on Newtonian theory could not have led to the particulars or even the core idea underlying black holes. According to Newton's theory an escaping particle being acted upon by gravity is decelerated. According to Einstein, a gravitational interaction occurs in the absence of any genuine forces (gravitational forces only appear relative to a reference frame), whereas light, according to general relativity, travels with the same fundamental velocity  $c$ . But in physics, there is a more powerful factor than any force, the causality principle. It imposes restrictions which can never be overcome without breaking causality itself (it is equally impossible to return to the past and change it). According to the relativistic causality principle, any consequence of a given local cause may only occur inside (or at least, on the surface) of the light cone of the future with its apex at the world point of the cause. You can ask: "What of it? It is so natural!" Yes, but in general relativity where the distribution and motion of matter cause space-time to curve, the light cone too may turn in a certain direction. If this turn is strong enough, nothing may leave a certain area of space. No matter what sort of engine you may have or how much power you possess, you will never leave an area which has been isolated from the rest of the world by the causality principle. You may get information or material objects from the outside, and you yourself have probably come from the outside, but nothing and nobody can ever leave the region. A black hole is an example of a region isolated from the rest

of the world due to causality. If this sounds strange, reflect that it arises from the same principle that makes travel to the past an impossibility, and that fact is somehow taken for granted.

**3.2.3. Horizon and Ergosphere.** Now we are ready to give a strict definition of the most important notion for a black hole: the horizon\*. Since by definition, the horizon can only be crossed one way, a light cone of the future with its apex at any event on the surface of the horizon may only point to one side (that is in the direction in which the motion is allowed). The horizon is the *first* of a family of surfaces that is not intersected by light cones whose apices lie on the surface. Clearly one of the generatrices of every cone of the future must touch the horizon while the rest of the points in the cone must lie to one side of the horizon. This part of surface is called a null (light-like) surface. The classification of hypersurfaces is quite easy: space-like surfaces have normal vectors that are time-like; time-like hypersurfaces have space-like normal vectors; and finally, the normal vectors of null hypersurfaces are themselves null vectors. It is interesting that a null vector, which is by definition orthogonal to itself ( $n^\mu n^\nu g_{\mu\nu} = 0$ ), may be orthogonal to another null vector if and only if both vectors are proportional (collinear). The reader might enjoy proving this simple theorem. The normal vector of null hypersurface is, at the same time, its own tangent vector!

The hypersurface of the horizon of a black hole is a null three-dimensional subspace of four-dimensional space-time; at the same time, it is, in terms of the three-dimensional physical space, a compact or closed surface. The equation of this hypersurface can be obtained if one condition is imposed on every point of space-time (then, one of the degrees of freedom is lost and only three degrees remain), viz.  $f(x^0, x^1, x^2, x^3) = 0$ . The normal vector to a hypersurface (or simply to a surface) has components which are the partial derivatives with respect to the relevant coordinates of the function  $f$ , i.e.

$$n_\mu = \partial f / \partial x^\mu. \quad (3.9)$$

---

\* We mean here the event horizon. In cosmology there is a more important notion of the particle horizon, which has somewhat different properties.

The condition that the normal vector is null is

$$n_\mu n_\nu g^{\mu\nu} = 0. \quad (3.10)$$

The horizon must be compact (otherwise, any light cone would satisfy this definition!), but how should this statement be interpreted? A rigorous definition of compactness can be found in a textbook on topological spaces, but here we shall confine ourselves to an intuitive formulation. A compact space is one in which it is impossible to move infinitely far away without leaving the space. A closed surface is, without doubt, compact. As to the horizon, a compact space, by definition, must be the intersection of the hypersurface of the horizon and the three-dimensional physical space, that is, a two-dimensional surface of the horizon. Only then can the horizon of a black hole act as a "one-way membrane". The horizon in some cases may be relative, that is, it exists, for one observer (one reference frame), while for another it does not (e.g. the Rindler horizon which exists for a uniformly accelerating observer). On the other hand, horizons can evolve, coming into being or disappearing (for example, during gravitational collapse of a mass greater than a certain critical mass, it may sink beneath its own gravitational radius, that is under the horizon that is formed in this process).

Let us determine the horizon in the Kerr space-time, using equation (3.10). Jumping ahead and simplifying the calculation, we assume that the equation of the horizon is  $f(\rho) = 0$  (that is, a closed surface whose coordinate  $x^1 = \rho$  is constant; it can be shown that this surface is an oblate spheroid of revolution). It follows from (3.9) that the normal vector is proportional to  $\delta_\mu^1$ , and from (2.47) we have

$$\delta_\mu^1 \delta_\nu^1 g^{\mu\nu} = g^{11} = 1/g_{11} = -\frac{\rho^2 c^2 - 2Gm\rho + a^2}{\rho^2 c^2 + a^2 \cos^2 \theta} = 0. \quad (3.11)$$

Hence the position of the horizon is determined by a quadratic equation, i.e.

$$c^2 \rho^2 - 2Gm\rho + a^2 = 0, \quad (3.12)$$

which, generally speaking, has two roots

$$\rho_\pm = Gm/c^2 \pm \sqrt{(Gm/c^2)^2 - a^2/c^2}. \quad (3.13)$$

These roots may only be effective if they are real. Thus, we have three different cases:

$$Gm/c^2 > |a/c|, \quad (3.14a)$$

$$Gm/c^2 = |a/c|, \quad (3.14b)$$

and

$$Gm/c^2 < |a/c|. \quad (3.14c)$$

In the first case we have two horizons: exterior and interior (Cauchy surface). If at  $a = 0$ , we use the Schwarzschild field, the exterior horizon transforms into the one we discussed in Sec. 2.2, whereas the interior horizon merges with the central singularity. As the Kerr parameter  $a$  (the central body's angular momentum divided by its mass) increases the two horizons tend to each other, and at (3.14b) they merge. This is the case for the extremal Kerr black hole. If the absolute value of parameter rises a little more, the horizons simply disappear; what remains is called a naked singularity since it is not hidden under a horizon. We have thus found the positions of the horizons of the Kerr field, and also established that the horizon in the Schwarzschild space-time satisfies (3.10).

However, before we reach the exterior horizon, moving from outside, we run into another interesting surface, this is a surface on which the  $00$ -component of the Kerr space-time metric tensor changes sign:

$$g_{00} \equiv 1 - \frac{2Gm \rho}{c^2 \rho^2 + a^2 \cos^2 \theta} = 0. \quad (3.15)$$

Here again we have a quadratic equation for  $\rho$ , but now it is dependent on  $\theta$ :

$$c^2 \rho^2 - 2Gm \rho + a^2 \cos^2 \theta = 0. \quad (3.16)$$

The solution of this equation is

$$\tilde{\rho}_{\pm} = Gm/c^2 \pm \sqrt{(Gm/c^2)^2 - (a/c)^2 \cos^2 \theta}. \quad (3.17)$$

Of the two surfaces determined by  $\tilde{\rho}_{\pm}$ , one lies above the exterior horizon and touches it at the poles (at  $\theta = 0$  or  $\pi$ ) since it is more oblate. Outside the exterior horizon but inside the new surface, the sign of  $g_{00}$  is negative and hence the coordinate  $x^0 = ct$  ceases to represent time. It can be shown, however, that in this domain (between  $\rho_+$  and  $\tilde{\rho}_+$ )

it is still possible at every individual point to combine the former coordinates  $t$  and  $\varphi$  (with constant coefficients specific for each point) so that a time-like coordinate  $x'^0$  that is "good" from the outside and up to the particular point emerges. By using this coordinate, we can thus get a  $g'_{00}$  component of the metric tensor that is positive at the point. Since the components of the metric tensor remain locally independent of the "new" time, we can say that the Kerr metric is stationary down to the exterior horizon, even though it is only *locally* stationary between  $\rho_+$  and  $\tilde{\rho}_+$ . Since the gravitational redshift of the radiation arriving from a fixed point (constant  $r$ ,  $\theta$ , and  $\varphi$ ) to a distant observer is determined by the 00-component of the initial metric tensor (see Sec. 2.1), the surface  $\rho = \tilde{\rho}_+$  is called the infinite redshift surface, or the static limit (see C. Misner, K. Thorne, and J. Wheeler: *Gravitation*, p. 880). On the surface no light source with a nonzero rest mass can ever be at rest, since otherwise it would have to move with the velocity of light! And this is the origin of the term "static limit", while the infinite redshift of a source moving with the velocity of light in this domain is due to an infinite Doppler effect as it is due to the properties of the gravitational field at  $\rho = \tilde{\rho}_+$ .

If a particle is orbiting around the central mass with a certain angular velocity, it can be considered at rest in the new system with coordinate  $x'^0$  (the combination of  $t$  and  $\varphi$ ). The light emitted by such a particle will have a finite redshift for a distant observer, but the particle must be outside the exterior horizon since no signal can be sent from within the horizon to infinity (or indeed beyond the horizon) as is the case for a Schwarzschild field (cf. Sec. 2.2).

The region between the horizon and the infinite redshift surface is called the ergosphere (it is only nonzero in volume for the Kerr space-time). A specific feature of this space is that in it physical processes are possible which extract energy from the central body, i.e. from rotating black hole. Particles with very great but negative total energies may be generated in any ergosphere; the particles need only be given a correctly directed momentum. If we launch a body consisting of two components from the outside into the ergosphere and it separates within the ergosphere, then one component will be given a high negative energy and the



other a corresponding positive energy. The first component will be absorbed through the horizon by the black hole, which gives up energy to do so, whereas the second component will escape from the ergosphere having gained energy. The main point, however, is that the black hole must lose a fraction of its angular momentum, i.e. its rotation must slow down. If this procedure is repeated, ultimately the angular momentum will fall to zero (as long as the energy of the hole is not exhausted). The ergosphere will then disappear together with the possibility of getting energy from the black hole.

Besides the Schwarzschild (nonrotating) and Kerr (rotating) black holes, which have no electric charge, we also have exact solutions of Einstein's equations when the source has an electric charge. If the system rotates, the electric field yields a dipole magnetic field. These solutions are referred to as the Reissner-Nordström and Kerr-Newman metrics. The first one looks very much like the Schwarzschild metric but it has two parameters: the mass  $m$  and charge  $q$ . In addition to these the Kerr-Newman metric has the Kerr parameter  $a$ . Normally, cosmic objects have no significant charge because it is quickly neutralized by diffusion into cosmic space, or by the attraction of opposite charges from interstellar gas or dust. Therefore, the Reissner-Nordström and Kerr-Newman black holes are mostly of academic interest, but they are important for the theory. Firstly, these four types of black holes cover all possible classes of collapsed massive objects, according to the Israel-Wheeler theorem (a black hole has no hair). Secondly, the Reissner-Nordström metric is simpler than the Kerr metric while both are very much alike in their physical and geometric properties. We shall see this later during our discussion of Penrose diagrams. This metric also can be used to show up an interesting property of gravitation, which has so far been avoiding our attention, namely the repulsion of electrically neutral particles.

The Reissner-Nordström metric can be "derived" using the same technique we used for the Schwarzschild metric. The only difference between the two is that the Newtonian potential for a point mass should be replaced by a solution of the Poisson equation for a distributed source,

$$\Delta\Phi_N = 4\pi G\rho_m. \quad (3.18)$$

The point mass remains at the origin and yields the same potential ( $-Gm/r$ , see Sec. 2.2), but the electrostatic source has a trick of its own\*. In Newton's theory, the  $\rho_m$  on the right-hand side of equation (3.18) is usually interpreted as the density of mass (or energy, since from special relativity mass and energy are equivalent). That was the case, however, only for nonrelativistic matter, whereas an electromagnetic, or even an electrostatic field is always relativistic though it might appear at rest. It can be rigorously shown that for such a field we have to take instead of  $\rho_m c^2$  *double* the energy density

$$\Delta\Phi_N = 8\pi Gw/c^2. \quad (3.19)$$

This doubling is generally characteristic for ultrarelativistic objects: remember, for example, that light rays are deflected by gravitational field twice as much according to general relativity than they are according to Newton's nonrelativistic dynamics. The density of the energy of a Coulomb electrostatic field (i.e. of  $E = q/r^2$ ) is

$$w = q^2/8\pi r^4, \quad (3.20)$$

and the Laplace operator  $\Delta$  in a spherically symmetric case (that is, when  $\Phi_N$  does not depend on angles) takes the form

$$\Delta\Phi_N = \frac{1}{r} \frac{d}{dr^2} (r\Phi_N). \quad (3.21)$$

Bearing this in mind, we have, as a complete solution of equation (3.19), the Newtonian potential

$$\Phi_N = -Gm/r + Gq^2/2c^2r^2, \quad (3.22)$$

which enters the 00-component of the metric tensor in the form:

$$g_{00} = 1 + 2\Phi_N/c^2 = 1 - 2Gm/c^2r + Gq^2/c^4r^2. \quad (3.23)$$

Hence, by doing exactly what we did in Sec. 2.2, we finally obtain the Reissner-Nordström field in the form

$$\begin{aligned} ds^2 = & (1 - 2Gm/c^2r + Gq^2/c^4r^2) c^2 dt^2 \\ & - (1 - 2Gm/c^2r + Gq^2/c^4r^2)^{-1} dr^2 - r^2 (d\theta^2 \\ & + \sin^2 \theta d\varphi^2). \end{aligned} \quad (3.24)$$

---

\* The form of equation (3.18) coincides with that from electrostatics (barring the interpretation of the functions involved), but with the opposite sign for the right-hand side which reflects the attraction of masses in Newton's theory.

It can be shown that this is an exact solution of the self-consistent system of Einstein's equations (with the electromagnetic field energy-momentum tensor on the right-hand side) and the vacuum (homogeneous) Maxwell equations with the electric field  $E$  given above.

If we consider this result in the same way as reasoned for the Kerr field, we can quickly conclude that the Reissner-Nordström space-time has two horizons (one exterior and one interior) with the infinite redshift surface coinciding with the exterior horizon, so that there is no ergosphere. If we express the Newtonian potential (3.22) as

$$\Phi_N = -\frac{G}{r}(m - q^2/2c^2r), \quad (3.25)$$

the term in the parentheses can be interpreted as a mass inside a sphere of radius  $r$  (the total mass of the system is obtained when  $r \rightarrow \infty$ ). This is because, for a spherically symmetric distribution of mass, the field at a point with a coordinate  $r$  is determined by the mass within the sphere with this radius, the field being the same if the mass were concentrated at the centre. This means that the mass inside a sphere of radius  $r$  is less than  $m$ , and the sphere of radius  $r_0 = q^2/2mc^2$  with its centre at the origin contains zero total mass. Spheres with shorter radii contain negative total masses which tend to negative infinity as  $r \rightarrow 0$ . The total positive mass of the whole system (that is,  $m$ ) is made up of an infinite negative mass concentrated at the centre, and an infinite, but positive, mass due to the electrostatic field that surrounds the centre. We would expect, therefore, that an electrically neutral body (say, a test particle) falling towards the centre of a Reissner-Nordström geometry would first be attracted and then (below  $r = r_0$ ) be repelled.

This solution is confirmed by an exact solution of the geodesic for a radial infall in a Reissner-Nordström field. A test particle falling inwards with an arbitrary initial velocity will eventually be stopped by the repulsion force of the infinitely great negative central mass, and then it will be repelled back again. If the motion begins from rest at a finite  $r_1$ , then the law of energy conservation dictates that the particle will return and stop at the same  $r_1$ , and then it will start falling inwards again. This oscillation will continue for an indefinite period of time. A simple calculation shows

that the period of the oscillation cannot be measured by a distant observer's clock but it can be measured by the particle's clock (assume a miniature clock is falling with the particle) if  $Gm \geq \sqrt{\bar{G}} |q|$  (that is, if the roots  $r_{\pm}$  that determine horizons are real). The clock that falls with the particle measures the proper time of the particle (the integral  $\int ds$  along its world line) and gives a finite period which depends on the initial energy of the particle. The clock of the distant observer (and of any observer in general at rest anywhere outside the exterior horizon) will record an infinitely long period\*. In other words, a distant observer's clock will show that it takes an infinitely long time for the particle to reach the horizon from the outside (as is generally the case with all black holes; and this is only a quarter of the period!). Only if  $Gm < \sqrt{\bar{G}} |q|$ , will both measurements record finite values for the period of the oscillation. But let us get back to the case of  $Gm \geq \sqrt{\bar{G}} |q|$ . The particle that is repulsed by the infinitely large negative mass that "lies in ambush" at the origin comes out again from under the exterior event horizon after an infinitely long period of time (which is indeed not the same as a "very long time"!) after it passed the same point on its inward journey. In other words, it will pass "onto another sheet" of space-time, as geometers put it, i.e. into *another universe* (since both the Reissner-Nordström and Kerr space-times represent infinitely many universes, we cannot consider that our Universe is unique anymore). This is how radically we change our conception of the world only after taking a passing glance at Einstein's general relativity!

**3.2.4. Penrose Diagrams.** In order to present the Schwarzschild, Kerr and Reissner-Nordström geometries more visually, we use Penrose diagrams, in which we picture an infinite space-time on a finite area of paper. To reduce this picture to two dimensions, we must suppress the angular coordinates  $\theta$  and  $\varphi$ . In the diagram we have the coordinates  $t$  and  $r$ , or a more usable combination of the two. R. Penrose, a British physicist, proposed a conformal mapping (multiplication of  $ds^2$  by a specially selected function of coordinates) to depict an infinite domain of natural coordinates ( $t$ , and

---

\* More exactly, the period is  $4\infty$  of time.

$r$ ) on a finite domain. The result is meaningful since we can always get back to the natural coordinate description of the world after such a mapping (no information is lost during a conformal transformation). This technique makes Minkowski's space-time in special relativity look like a triangle (see Fig. 17), the vertical side (left) representing points with an  $r$  coordinate equal to zero at all moments of time with the lower apex being the infinitely remote past of all spatial points (this is caused by the conformal compression of the picture) and the upper apex being the infinitely remote future of these points. The apex to the right is space-like (spatial) infinity at all finite moments of time (compressed to one point because of the conformal mapping). The slant sides of the triangle are "null infinities"\*) of the past (below) and the future (above), respectively. They are null because they depict infinitely remote domains of the space-time from which (past null infinity,  $\mathcal{I}^-$ ) and to where (future null infinity,  $\mathcal{I}^+$ ) light can go through all finite points of the space at all moments of time (all here-and-nows).

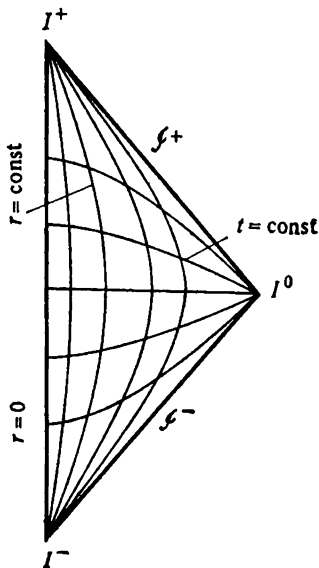


Fig. 17

There is a detailed theory covering this mapping technique. Our description has had to be brief but it should be sufficient to visualize the gravitational fields in question.

The conformal transformation has the important property that while it compresses or expands distances (depending on the position of the points in space-time), it does not change the properties of the null lines (e.g., the null geodesic of an initial metric remains a null geodesic for a conformally re-gauged metric). Therefore, the null lines on a conformally transformed graph of space-time are at  $45^\circ$  to the horizon-

\* The " $\mathcal{I}$ " on the diagram is pronounced "scri".

tal, while because of this transformation they all can be presented on a limited plane.

Let us now revisit the Schwarzschild field (Sec. 2.2). We know that the space-time path of a horizon point fixed in the space is described by a null line. This means that on a Penrose diagram it will be drawn at  $45^\circ$  (this is true for the horizon of any other gravitational field). The line may incline either way, as shown in Fig. 18. In Sec. 2.2, we found that on the line  $r = Gm/c^2$  the standard coordinates for the Schwarzschild field become inadmissible. It is not surprising, therefore, that as we pass over to admissible coordinates that are used on Penrose diagrams, the single  $r = 2Gm/c^2$  becomes two intersecting lines. A fixed point in space outside the horizon at all moments of time is represented by a line at an angle to the vertical of less than  $45^\circ$ ; this line is really time-like. Such a point inside the horizon, on the other hand (that is, at  $r < 2Gm/c^2$ ), will actually not be fixed in space and at best will only be fixed in time. The point is that the signs of  $g_{00}$  and  $g_{rr}$  under the horizon are opposite to those on the outside, and now the coordinate  $t$  becomes space-like, whereas  $r$  becomes time-like. Indeed, for  $r < 2Gm/c^2$ , the world line of a point with a constant  $r$  gives  $ds^2 < 0$  which holds for a space-like line. It is certainly true for the point  $r = 0$ , where there is a genuine singularity in Schwarzschild's space-time (compare this with Sec. 2.2). Therefore, the singularity at  $r = 0$  on a Pensore diagram (see Fig. 18) becomes horizontal. But there must be *two* such lines, one at the bottom, and one at the top. The origin, which is regular in a Minkowski world (the vertical side of the triangle in the diagram in Fig. 17) becomes a singularity line in a Schwarzschild world, breaking down into two sections, two horizontal lines. It is as if the world has decomposed into two worlds connected by horizons. To the left and to the right are the infinitely remote pasts and futures (it is strangely disturbing to use plurals for these subjects), two spatial infinities and two pairs of null infinities (for the past and the future).

Since a light cone with its apex at any given point is depicted on a Penrose diagram as an askew cross with a  $45^\circ$  slant intersecting at that point, it is obvious that a world line which does not intersect the horizon from the past (that is, from the bottom of the diagram) should lie in a narrow band between one of the generatrices of the cone and the horizon

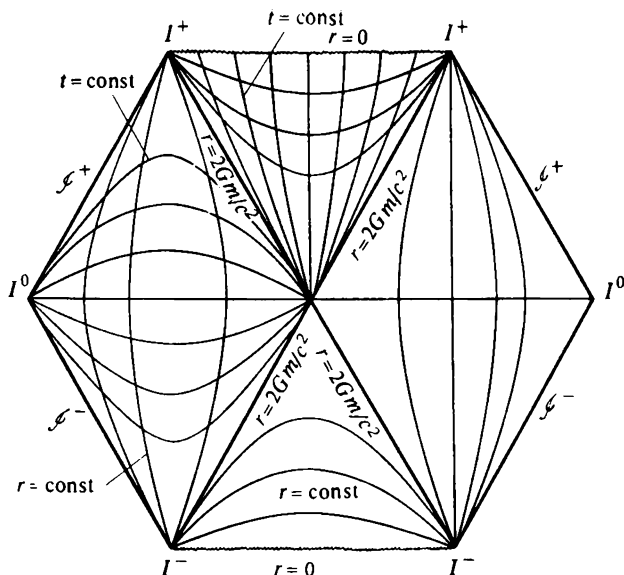


Fig. 18

line. This band contracts to zero as the apex of the cone approaches the horizon. Therefore, it becomes more difficult to avoid sinking beneath the horizon as the particle approaches it. Even if the particle can move away from the horizon then to escape to a significantly distant point from the horizon the particle will need a very long period of time (if the initial point is on the horizon, the escape will take an infinitely long time). All these considerations, together with the absolute impossibility of rising above the horizon, justify the name *black hole*. Do not, however, be puzzled by the fact that the domain corresponding to the black hole under the horizon is represented by the triangular region located *above* the horizon lines and below the singularity on the Penrose diagram for the Schwarzschild field (see Fig. 18). The lower triangle formed by the horizon lines and the lower singularity is called a *white hole*. This can only be left and not entered from an external universe (either left or right). A particle rising above the horizon of a white hole can enter either of these two universes and can (if it moves

properly) leave the universe to enter the black hole. Hence, the white hole is the place of departures and the black hole the place of arrivals. Since a departure differs from an arrival only by their orientation with respect to time, white and black holes change their places with the inversion of time.

Things change "from bad to worse" in the Reissner-Nordström and Kerr geometries. Let us have a look at the Reissner-Nordström world, which is slightly simpler. At the moment we will discuss the case of two horizons that do not coincide, that is,  $m\sqrt{G} > q$ . The singularity  $r = 0$  is no

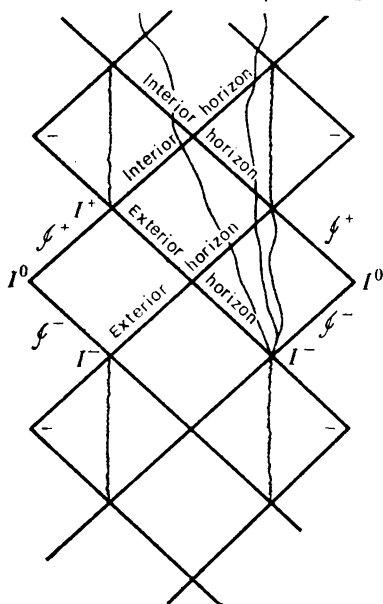


Fig. 19

longer space-like, as it was in the Schwarzschild geometry, but is time-like (like the nonsingular origin in Minkowski's universe). Then, all horizons, exterior and interior, are null and must be depicted by lines at  $45^\circ$  to the horizontal. As a result, we have the Penrose diagram in Fig. 19. We suggest the reader tries to explain why what we have here is not two universes with two sets of infinities, but two infinite series of universes, one set to the left and one to the right. We should note that the tradition of breaking up the line  $r = 0$  and changing its direction by a right angle, which was

started by the transition to the Schwarzschild geometry, is continued here. The fragment of the Penrose diagram in Fig. 19 shows the next change in the transition from the Schwarzschild to the Reissner-Nordström geometry. Finally, on the full diagram on the same figure we can easily see why the movement of the test particle we mentioned earlier should be oscillatory (when the repulsion action of the Reissner-Nordström centre was discussed). We can thus quite naturally arrive at the conclusion we gave earlier, that is



that the world line of a particle passes from one sheet of space-time to another, to a third, and so on.

Let us now consider the Kerr space-time. At  $\theta = \text{const}$ , it only differs from the Reissner-Nordström space-time by "a little". It has an ergosphere (which has little effect on the Penrose diagram), a ring-like singularity (see Sec. 2.6; we return to this subject), and it assumes the existence of a "negative space" (entry into the domain where  $\rho < 0$  is now allowed). The last two features arise because  $\rho = 0$  corresponds to a disc with a finite radius  $a/c$  and lying in the plane  $z = 0$  rather than to a point (as was the case for spherical coordinates and the corresponding fields). This happens because the coordinates for the Kerr metric are spheroidal and can be expressed in terms of the usual Cartesian coordinates  $(x, y, z)$  as

$$x^2 + y^2 = (\rho^2 + a^2/c^2) \sin^2 \theta, \quad z = \rho \cos \theta. \quad (3.26)$$

Starting from the Kerr metric (2.47), we see that a singularity (like the Schwarzschild one at  $r = 0$ ) is located at  $c^2\rho^2 + a^2\cos^2\theta = 0$ , so that if  $a \neq 0$ ,  $\rho$  vanishes, and  $\theta$  takes the value  $\pi/2$  (the equatorial plane). Singularities do not occur for zero  $\rho$  at other values of  $\theta$ , and the value  $\rho = 0$  can be freely passed and the domain of negative  $\rho$  ("negative space") can be entered. Here, however, repulsion acts on the test particle pushing it back. Then Penrose diagram of Fig. 19 describes this situation. The singularity at  $\rho = 0$  can only be reached at a single value of  $\theta$ , i.e.  $\pi/2$ . This is represented by a serated line (in the Reissner-Nordström field, this singularity occurs at  $r = 0$  approaching it from any direction and not just at "the equator", and beyond the serated line, space-time no longer exists for this field). If in the Kerr space-time we approach  $\rho = 0$  for a  $\theta \neq \pi/2$ , it is possible to enter a negative space (beyond the serated line) on the deepest boundaries of which negative infinities similar to standard ones occur. In all other aspects, the Penrose diagram for the Kerr space-time is the same as that for the Reissner-Nordström space-time.

**3.2.5. Evolution of Black Holes.** How do black holes form, how do they exist, and how do they die? We shall touch on each of these questions, but there is an extensive literature on the subject (both specialist and lay), much of it written by the scientists who have done most to construct the theory of black holes. In short, a black hole may be created by

the gravitational collapse of a very massive body when it contracts almost to its gravitational radius. This compression is opposed by the pressure of the matter and the radiation inside the body. The pressure is caused both by the thermal motion of the particles and by the repulsive action of the fields they create. The thermal motion slows as the body cools by emitting radiation into the surrounding space (the latter must be cool, that is, it must contain radiation of a lower temperature). The decisive aspect in the repulsive action of the particles at the critical stage is the Pauli quantum exclusion principle and not simply the force of their interaction. This principle controls the behaviour of particles with spin (particles that are described by Fermi-Dirac statistics) and it forbids more than one particle from occupying a single state. Hence, these particles repulse each other. Half-integer spin particles are electrons and nucleons (protons and neutrons), in particular. The last particles to surrender to gravitation (i.e., the causality principle) are neutrons. It is now believed that the most dense objects in cosmos are neutron stars (pulsars) which may contain at their cores even more massive nucleon states, called hyperons. A typical neutron star, with a mass of about one solar mass and a 10-km radius, will break down and collapse if it gains three more solar masses. The reader can easily calculate when the gravitational radius of such an object surpasses its geometrical radius. As a matter of fact, the gravitational radius is proportional to a body's mass, whereas its geometrical radius is proportional to the body's mass raised to the power of  $1/3$  (if, for a first approximation, we assume that the density remains constant, and the notion of radius does not change its sense in a strong gravitational field). Sooner or later, both curves will intersect signifying the gravitational collapse of the system.

This is the way in which a single black hole is formed, without an associated white hole, which should exist if we have a Schwarzschild or Kerr "primordial black hole". The collapse is completed within very short time, if it is timed by a clock travelling with the particles participating in the collapse. The collapse never ends if measured by the clock of a distant observer. However, even this observer will see (according to his clock) that the collapsing object has "frozen", and turned in practice into a black hole.

Consider a black hole of any origin and its behaviour.

Of course, it can swallow any object (light included) if it comes too close. As a result, it will gain mass. The reverse process is also possible, i.e. the "evaporation" of a black hole, even though its horizon may never let anything out (according to the classical theory). This was the conclusion reached by **S. W. Hawking**, an outstanding British physicist, in his fundamental work published in late 1974. He showed that this must be a sort of a quantum process, even though objects with macroscopic masses suffer from it. According to quantum theory, all black holes, whatever their mass, must radiate. However, since they are in a real Universe, they behave like objects immersed in a constantly heated medium with a nonzero mean temperature  $T$  (for this we have the background electromagnetic radiation with  $T = 2.7$  K). Depending on whether the temperature of the black hole radiation is higher or lower than that of the medium, the hole will "thaw" or grow (in energy, i.e. mass) as a result of the dynamic equilibrium between radiation and absorption. In fact, not every sort of radiation can be characterized by a temperature, but background radiation (see Sec. 3.2.6) and incidentally black hole radiation are those of black bodies, a type of radiation which is important in thermodynamics and which is determined by its temperature. The higher the frequency of the peak of the radiation spectrum (frequency distribution of radiation energy), the higher the radiation temperature normally ascribed to its source. Therefore, we can speak about the temperature of a black hole.

Hawking proved that a black hole can be described by its thermodynamic properties just like any other system with a complex internal structure. In this case, however, we mean the external (accessible to an outside observer) characteristics of a black hole and not its structural properties. These properties are independent of the objects that were involved in the collapse that led to its formation. Therefore, every black hole is a member of a large set of objects that can be individually characterized by mass, angular momentum and electric charge. In other words, black holes are essentially statistical entities. Hence, they naturally possess thermodynamic characteristics such as temperature, entropy, and black body radiation. Let us try to evaluate the temperature of a black hole which is characterized for simplicity only by its mass (i.e., a Schwarzschild black hole). To begin with we have to understand why black

hole can radiate though it has a horizon (so it seems to violate the causality principle). Note that black body radiation is absolutely chaotic and transmits no information (with the exception, perhaps, of the temperature of its source; but this sort of information is accessible for an observer outside a black hole in any case). Hence, the presence of black body radiation from a black hole does not contradict the causality principle in its general form. As to the mechanism underlying this emission, in strong fields (both gravitational and electromagnetic), particles such as electron-positron pairs are generated as a result of the transfer of energy from the field to the quantum vacuum, in which these particles are always present in a virtual state until energy is spent to create them. Let us assume that the field is strong enough to realize an electron-positron pair. Normally, the pair annihilates immediately to release at least one photon (since we consider here the process in the external field of a black hole). If this emission takes place outside the horizon, the photons may permanently escape the black hole, carrying away some of the energy of the gravitational field of the hole (that is, some of its mass). This is the radiation we are interested in.

When will it be more intensive, that is, which black holes have higher temperatures? Obviously, when the gravitational field is stronger outside the horizon from which photon can be radiated off to infinity. In Newton's theory, a gravitational field is characterized by its strength, the gradient of the potential taken with negative sign. After differentiating the Newtonian potential, we get an absolute value for gravitational field strength of  $Gm/r^2$ . At the horizon (where  $r = 2Gm/c^2$ ), the field becomes  $c^4 (4Gm)^{-1}$ , i.e. it is inversely proportional to the mass of the black hole. This is natural: the smaller the hole is in mass, the nearer the horizon is to the origin. The field at the horizon is stronger the smaller the black hole is. Its temperature will also be higher, i.e. a black hole's temperature is inversely proportional to its mass. The rigorous theory yields

$$T_{\text{b.h}} = \alpha/m, \quad \alpha = \hbar c^3/4kG, \quad (3.27)$$

where  $\hbar$  is Planck's constant,  $c$  is the velocity of light,  $k$  is Boltzmann's constant, and  $G$  is Newton's gravitational constant. If we now substitute in the values of the world constants, we will have  $\alpha = 0.77 \times 10^{27}$  (if the tempera-

ture is in kelvins and the mass in grams). This means that a black hole with a mass of one gram has a tremendously high temperature. The temperature of the background radiation, 2.7 K, corresponds to a black hole with a mass of  $2.85 \times 10^{26}$  g (this is somewhere between the mass of the Moon,  $7.35 \times 10^{25}$  g, and that of the Earth,  $6 \times 10^{27}$  g). However, lighter (consequently hotter and hence evaporating) black holes may live a very long time. For example, a black hole with a mass of  $10^{15}$  g may take about the lifetime of the Universe to evaporate away completely. In earlier epochs, of course, the background radiation was hotter, and at some stage, a black hole with this mass grew and did not evaporate. This improved the survival chances of those black holes that had been created at the earlier stages of the Universe's expansion (the chances are even greater for primordial black holes). The end of a black hole involves a catastrophically rapid evaporation at the highest possible temperature. The total energy at this last stage ( $mc^2$ ) is not very great because by that time most of the black hole's mass has already evaporated. The remaining mass will not therefore generate very many of the various particles that can be emitted at such high temperatures. Still it is enough to generate an X-ray flash.

At the moment, we cannot point to any observations that comprise unequivocal evidence of the existence of black holes, but candidates for black holes have long been discussed. For example, there is the X-ray source Cygnus X-1 in the Cygnus constellation and a mysterious object in the centre of our Galaxy. The X-ray and  $\gamma$ -ray radiations from these objects are due to relativistic processes occurring in their vicinity and not to the evaporation of the black holes. A black hole must constantly be drawing in matter from its neighbourhood, especially if it is a component of a binary or multiple star system. In this case the "hole" will "drag" substance from its neighbours via the tide mechanism. By observing the properties of such a system, we can estimate the mass of the objects suspected of being a black hole. If the object is sufficiently massive but it cannot be seen as a gigantic star emitting light independently, it is most probably a black hole, because objects with more than 3 to 5 solar masses should collapse at the end of their lifetimes to produce either a black hole or should explode as a supernova, thus losing most of their mass. The object at the centre of

our Galaxy probably has vast mass, millions of times greater than that of the Sun. At the same time, it is relatively very small, smaller than  $10^{18}$  cm (less than one light year). It may even be several orders of magnitude smaller. The figure of  $10^{18}$  cm was determined by multiplying the velocity of light by the time during which the radiation of the source varied significantly, which is the upper limit imposed by causality principle.

For the sake of balance it should be noted that doubts are often expressed with regard to the existence of black holes. Since no incontrovertible observational evidence has so far been produced for their existence, we cannot guarantee that there are not some as yet undiscovered general (and fundamental) physical principles that may prohibit their existence, in the same way as there are the principles of quantum physics that prohibit the electromagnetic collapse of electrons onto the nucleus in an atom. We should not think, however, that a science can dispose of phenomena it has predicted without radical changes. The structure of a scientific theory is rigid and cannot be voluntarily changed as it would not then be a science and its conclusions would be of no value. A genuine scientific theory cannot, when confronted with an incompatible fact (just one rigorously established fact will do), be corrected "a little" to make it agree. The theory must be radically changed. This does not mean, however, that the old theory is wrong, all it means is that we have found limits to its applicability, but it remains valid within the limits. If we discover such strict limits, we would understand the theory deeper because the emerging new theory allows us to view the old one from a new angle. Thus, the study of black hole phenomenon, like the study of any other extreme situation in science, is important and valuable. Extensive research on this phenomenon is being conducted both theoretically (at the boundary of general relativity and quantum theory) and practically in relativistic astrophysics through a large program of earthbound and space observations. We hope that within the next decade decisive evidence either for or against black holes will be obtained.

**3.2.6. Evolution of the Friedmann Universe.** From Sec. 2.9 we already know that our Universe is expanding and that it is described at this stage of its evolution by Friedmann's model, close to the intermediate one in which the 3-space moving with matter is flat. It is an interesting topic

to see how the Universe may have evolved in the past and may do so in the future. As to the past, the Friedmann model is valid from quite an early stage, but before that stage quantum effects must be accounted for. The latest trend in the theoretical investigations of these early stages is to take into account the intimate details of the theory of elementary particles, and particularly the Grand Unification Theory (GUT). According to the GUT, the Universe's evolution can be divided into stages corresponding to a series of spontaneous symmetry break-downs. This is probably the way in which the modern structure of space-time and of the intrinsic degrees of freedom of the elementary particles have evolved. Of course, the theory itself must evolve. Therefore, we start our discussion from the moment when the Friedmann model came into action (for a more detailed account see [13, 45, 77, 89, 96, 104]).

It can be shown that the density of matter and the density of radiation in this model depend on the scale factor, the radius, in different ways (this observation comes from the first law of thermodynamics and the equations describing the gravitational field and the equations of state). The density of matter behaves like  $R^{-3}$ , whereas the density of radiation behaves like  $R^{-4}$ . Hence, at the beginning (when  $R$  is small), radiation is dominant,  $\rho_r \gg \rho_m$ , whereas at the modern stage, its role is negligible. Simply speaking, at a certain stage during the Universe's expansion, the interaction between radiation and matter practically stops (the Universe becomes transparent) and thenceforth the radiation expands adiabatically like a gas. As a result, its energy (mass density) decreases monotonously. In 1948, **George Gamow** (1904-1968) analyzed this possibility and laid the foundations for the Big Bang theory of cosmology. In 1956, he used this theory to evaluate what the present density of radiation  $\rho_r$  should be; it turned out to correspond to a black body temperature of 6 K. This prediction was nearly forgotten until, in 1965, A. Penzias and R. Wilson, while doing some purely technical measurements at a wavelength of 7.4 cm, discovered an unaccountable noise at about 3 K. A group of theoretical physicists under R. H. Dicke immediately interpreted this result as the "hot Universe" effect, thus confirming Gamow's idea [21, 54]. A short time before the publication of the Penzias and Wilson result, Gamow's theory was refined by A. G. Doroshkevich and I. D. Novi-

kov (1964) who have theoretically predicted a temperature of about 3 K. According to the latest data, the temperature of this background cosmic radiation is 2.7 K. This radiation is very isotropic, and indeed the slight anisotropy which is present in the observational data can be ascribed to the motion of the Earth (together with the Solar system) relative to the averaged reference frame which is moving with the matter as described by Friedmann's cosmological model. Sometimes a layman claims that the experimental detection of this privileged reference frame contradicts the relativity principle, thus reaffirming the pre-Einstein concept of an ether. This is of course nonsense, and we suggest the reader formulates a reproof to these misinterpreters of cosmological theories.

Under fairly general assumptions it can be shown that in the radiation-dominated epoch

$$TR(t) = \text{const}, \quad (3.28)$$

hence it becomes obvious that during the initial stages of the Universe's expansion, the temperature was very high. Remembering the way radiation density depends on the scale factor, we find that as  $t \rightarrow 0$ , the density tends to infinity. Therefore, we must speak about the evolution of the Universe from some  $t_q > 0$ ; we will consider the state of the singularity later. Usually,  $t_q$  is assumed to be around  $10^{-43}$  s (Planckian time). At that moment, the density of matter was  $10^{90}$  kg/cm<sup>3</sup> and the temperature  $10^{31}$  K.

There is another significant moment of time, viz. at  $t_s \simeq 10^{-23}$  s. From this point on the notion of hadrons, an important class of elementary particles becomes meaningful; and their description becomes close to that now known in flat space-time. During the *hadron era* the strong interaction dominated and the Universe mainly consisted of baryons and antibaryons in a thermodynamic equilibrium.

The hadron era ended after a few milliseconds. Most hadrons were annihilated, the remaining ones being due to a low asymmetry of particles and antiparticles. At  $t \simeq 10^{-2}$  s ( $T \simeq 3 \times 10^{11}$  K), the *lepton era* began during which the weak and electromagnetic interactions became dominant. The Universe then consisted of the surviving heavy particles, and of photons and leptons (light particles, mainly electron-positron pairs and neutrinos and antineutrinos). The Universe became transparent to muon neutrinos whose



density thereafter decreases because of the adiabatic expansion of the neutrino gas together with the Universe as a whole. Somewhat later, the Universe became transparent to electron neutrinos, and then to photons (at  $t \simeq 0.3$  s). Earlier stages of the Universe's evolution cannot in principle be probed by optical or neutrino astronomy techniques (by looking into the Universe through a telescope, we look into its *past*).

It is believed that the average number of neutrinos per cubic centimeter approximately equals the number of photons (about 500 quanta per cubic centimeter). If the rest mass of a neutrino is zero, the mass density of neutrino radiation is very small (about  $10^{-34}$  g/cm<sup>3</sup>). If, however, the rest mass is above 10 eV, then neutrinos will be the major constituent of the average mass density of the Universe, which may thus be spatially closed. The closedness of the Universe is a question still under debate.

After a second, the density dropped to 10 kg/cm<sup>3</sup> and the temperature fell to  $10^{10}$  K. By this time, most electron-positron pairs had annihilated to produce photons. Gamma-radiation became the dominant ingredient of the Universe, one atomic nucleus per about 100 million photons.

Approximately a minute after the start of the Universe, the temperature dropped to  $10^9$  K. At this time, the Universe was like a gigantic hydrogen bomb: the thermonuclear fusion of nuclei of deuterium, then tritium and helium began. It is believed that the period from the first 10 to 100 seconds was crucial for the composition of young stars, which consist of approximately 70% hydrogen and 30% helium. This mixture must have been "cooked up" at these stages in the expansion. Note that the results of this calculation are very sensitive to changes in the parameters of the individual models. At these stages, the cosmological "kitchen" operated by nuclear reactions which rapidly ceased as the temperature dropped.

A few minutes from the start of the Universe, the temperature was so low that fusion stopped, and the *radiation era* began. In the ensuing three hundred thousand years, the Universe was an expanding fire ball made of matter and radiation. As it expanded, the radiation energy decreased both absolutely and as a proportion of the total mass of the Universe. At  $t = 100\,000$  years, the density of the radiation energy became less than the density of matter (e.g., electrons, protons, and light nuclei).

When the Universe was about one million years old, the temperature had dropped below 4000 K, and electrons and protons began combining into hydrogen atoms. At this period, the energy of the photons was too low to ionize hydrogen atoms. This is called the *era of the separation of matter and radiation* and was the end of the radiation-dominated epoch and the beginning of the *matter-dominated epoch*.

During this era, the dominant objects in the modern Universe, galaxies, began forming. The reader can find in the literature on astronomy a description of galaxies and their classification. However, their origins still remain unclear and are being studied. There are two main hypotheses, called the adiabatic and entropy hypotheses, according to the type of small initial disturbances they assume.

The hypothesis involving adiabatic disturbances has been the most thoroughly developed. In this theory it is assumed that before the combination of protons and electrons, the density disturbances in the Universe could be described thus:

$$\delta\rho/\rho = 10^{-4} e^{-(M_D/M)^{2/3}}, \quad M < M_D;$$

$$\delta\rho/\rho = 10^{-4}, \quad M_D < M < M_J;$$

$$\delta\rho/\rho = 10^{-4} (M_J/M)^{2/3}, \quad M_J < M,$$

where  $M_J$  is Jeans'\* mass and it limits the domain of gravitational instability,  $M_D$  is the mass that determines the damping domain, the damping being due to viscosity and heat conductivity. Both these masses depend on time, and the large-scale perturbations (for  $M > M_J$ ) grow with time as  $\delta\rho/\rho \sim t$ . At an early stage, the perturbations were still small and the Friedmann hot homogeneous model can be used throughout. By the time the combination starts, all perturbations with  $M < 10^{13} M_\odot$  must have smoothed out because of viscosity. Since in this epoch  $M_J \simeq 10^5 M_\odot$ , the role of pressure is unimportant for the evolution of growing perturbations with  $M > 10^{13} M_\odot$ . This leads to the formation of a new structure in the Universe, that is, to "pancakes" or very flat clusters of matter. It is believed that at a later stage

---

\* Sir James Hopwood Jeans (1877-1946), F.R.S., was a prominent English astrophysicist who worked for many years in the USA.

the pancakes degenerate into galaxies. It has moreover been determined how they could have started rotating. These details are on the frontline of modern cosmology.

The evolution of the small-mass primeval disturbances (cosmological black holes) and cosmological gravitational waves are also topics at the forefront of cosmology. Such waves are believed to have formed at approximately  $t \sim 10^{-43}$  s.

After galaxies, which represent the large-scale structure of the Universe, the stars, planets, Earth, life on Earth, etc., were formed. We would like to remind you that the present age of the Universe is estimated to be about 20 billion years; the Universe continues to expand, but probably less rapidly. The temperature of the primordial background electromagnetic radiation, which is a remnant of the early stages of the Universe's evolution, has now dropped to 3 K.

There are various scenarios for the evolution of the Universe from now on depending on whether it is open (infinite) or closed (finite). At the moment, the most probable model is believed to be the open one (there are, however, some doubts). According to this model, the Universe will continue to expand ad infinitum, and its fate is to become a limitless flat space, "frost and darkness, the future without further changes" [96]. Equally unappealing is the future of the other, closed, model of the world. The explosive expansion will be succeeded by an implosion and the Universe will burn in an "inferno" like the one from which it was born.

Notice that the evolution and future of the world described in this section are based on the known laws of the general theory of relativity and elementary particles physics. We are extrapolating these laws to vast extremes indeed to get these descriptions. This extrapolation must, however, be viewed dialectically. On the one hand, it is necessary to stretch the laws to cover the whole structure of the Universe and to test the limits of our concepts. On the other hand, however, we must realize that this is only a model that reflects the level of our understanding. As science progresses, we will inevitably be compelled both to correct some features of our picture and even to change it substantially. There is a good evidence which supports this viewpoint. In the course of the last 100 years, we saw how on many occasions the development of our knowledge shattered the "optimistic" claims that everything had been discovered about

the fundamentals of the world, and that all that remained were just one or two clouds in the sky. Then, we saw these clouds grow into relativity theory, quantum mechanics, microworld physics which, incidentally, happened all to be crucial to physics. Who of us can bet that there similar clouds overhead do not remain? We would prefer to say that the skies of our knowledge are instead full of storm-clouds.

We have also to bear in mind that the Friedmann solutions are only a very special class of cosmological models because of their homogeneity and isotropy. Of course, these properties are approximate and are realized as we consider larger and larger domains of the real Universe (these domains are becoming less and less well known). These are other arguments in favour of the more general, anisotropic and heterogeneous models of the Universe. One of the important reasons for their study is that the Friedmann models must have an initial (and, sometimes, final) singularity.

**3.2.7. Cosmological Singularity.** We touched on the topic of the cosmological singularity. We define a singularity to be a point in space-time at which the curvature is infinite\*. The Soviet physicists E. M. Lifshitz, I. M. Khalatnikov, and V. A. Belinsky succeeded in constructing a general solution of Einstein's equations near a singularity, and they found that it has nontrivial oscillatory properties. But they failed to attain the ultimate objective: to remove the cosmological singularity. This is, however, in accord with the Hawking-Penrose theorem, which states that a singularity is unavoidable (either a cosmological one or a black hole) if the sources of the gravitational field meet some quite reasonable energy conditions.

Perhaps, the cosmological singularity cannot be eliminated without new ideas, such as those certain to be met in a quantum theory of gravity and quantum cosmology. All other attempts in the framework of classical general relativity are doomed to be superficial. Even more, any modified (not just corrected) theory will sooner or later yield contradictions as the experiments become more sensitive since any

---

\* The problem of singularity, even of the meaning of the word in general relativity, is still not quite clear. There are also other kinds of singularity, e.g. a point with a finite curvature, but for which the ratio of the circumference of a circle with the centre at it to the circle's radius is not  $2\pi$  (the so-called conic point).

theory in the long run is simply an idealization of reality. We just cannot expect to have a "final" theory. It is important that modern cosmology agrees very well with a variety of independent and different observations, and the analysis of this agreement is one of the main trends in physical cosmology.

This is also true for the properties of singularities in black holes. The removal of singularities (note that any actual infinity in the local characteristics contradicts reality and demands a constructive removal) is possible only by violating the old physical concepts. Palliatives have been suggested, such as changes in the structure of the field equations, the introduction of auxiliary physical fields with the ad hoc properties which should prevent the formation of singularities or horizons. These may give a temporary cure but they cannot really resolve the problem.

Since we are discussing the properties of the real Universe (sometimes, the observable part of it is termed the Metagalaxy), it is worth mentioning the interesting fact that some types of elementary particles dominate over others in the Metagalaxy. We mean the prevalence of the particles proper (electrons and protons) over the antiparticles (positrons and antiprotons). Antiparticles are far fewer, although the theories of fields and matter indicate equal numbers to be produced under normal conditions. This actual asymmetry must be the reflection of some fundamental fact in the very nature of things. For a long time researchers have been seeking mechanisms that could separate particles from antiparticles and which could place antiparticles in some other parts of the Universe far from our Metagalaxy. If this idea is crazy there is a crazier idea now, namely that all elementary particles, including protons, are unstable. This means that they can be created in violation of the baryon number conservation law and this violation may be another spontaneous symmetry violation that may take place under the proper conditions (of energy, mass density, etc.). Using the Grand Unification Theory and the supergravity theory (see Sec. 3.4.2) some researchers have inferred that at a certain stage during the expansion of the Universe (at a very early stage) the symmetry that follows from the baryon number conservation law was broken, and that space-time (or whatever underlies it) created protons unbalanced by antiprotons, i.e. the matter in today's Universe. It looks now as if we are

returning to the majestic myths of antiquity in which the predictive power created images as grand as the ideas of philosophers who dreamt up the notion of the atom from just nothing.

Relevant ideas, now more than ten years old, have been very adequately put by Michael P. Ryan, Jr., and Lawrence C. Shepley (in their book *Homogeneous Relativistic Cosmologies*, Princeton, 1975, pp. 270-272), and we would like to recall some fragmentary citations here: "The most vague but vital program which will hopefully be carried through in coming years is the application of quantum principles to the universe. The quantum field theory of gravitation is still in dispute. Even the interpretation of basic quantum principles as applied to the structure of space and time is in doubt. These controversies must be resolved before any significant application to cosmological studies can be made. It is important to find quantum solutions which correspond to the known cosmological models, but it is the details of what goes on at the epoch where a classical model has a singularity that are the most significant and which are most affected by detailed interpretation of quantum principles... Finally, in this list of theoretical problems we mention further study of the nature and structure of singularities within general relativity. Although much has been accomplished in this study, too many people are forgetting that there is still much to accomplish. In particular, the relationship between mathematical and physical singularities as we have defined them here is only very poorly understood. It cannot be overemphasized that this question remains as vital as it has always been". Since that time the situation has changed, indeed, but the changes were not yet really decisive ones.

We shall now go on to a discussion of the possible generalizations of the theories which will lead us beyond the classical frame, but which will be based on rigorous logical reasoning and real experiments and observations. In these generalizations, the researcher shifts the emphasis and considers as central properties which were described (usually stated) by the classical theory, but which were considered as unimportant or self-evident. For example, the three-dimensionality of the physical space and one-dimensionality of time may be reconsidered.

### 3.3. GENERALIZATIONS OF EINSTEIN'S GRAVITATION THEORY

**3.3.1. Introductory Comments.** Attempts at generalizing Einstein's theory of gravitation started immediately after it appeared. Einstein's success in geometrizing gravitational interaction inspired others to try to geometrize the other type of interaction known at that time, viz. the electromagnetic one. This was quite logical, because long before general relativity had been conceived and the physical effects of non-Euclidean space were discussed, Clifford and other scientists had hypothesized connections between electromagnetic effects and the geometry of the real world.

Therefore, immediately after Einstein's fundamental papers of 1916-1917, **Hermann Weyl** (1885-1955) suggested a generalization of Riemannian geometry which allowed him to introduce an additional geometric field that could be interpreted as the electromagnetic one. Then, Sir Arthur Eddington found further generalizations of Weyl's geometry that served the same purpose. T. Kaluza began studying the Riemannian geometries of five-dimensional space and E. Cartan looked into the possibility of generalizing four-dimensional Riemannian space to a geometry with torsion. They were joined also by other physicists and mathematicians. Einstein himself spent the last 30 years of his life seeking a way of geometrizing gravitation and electromagnetism into a unified field theory. Looking back, we can say that in the 1920-1930s the development of a unified field theory seemed necessary and one of the most promising aspects of theoretical physics. Many hopes rested on it. It was expected that these studies would reveal new properties of space-time and would, sooner or later, lead to significant technological innovations.

With time, however, these hopes faded, but meanwhile important results were being obtained in quantum theory, and then the atomic nucleus was fissioned. As a result of these developments a variety of technological applications began rapidly to appear in electronics, atomic power engineering, etc.; the attention of physicists switched over to these problems. Thus, all interest in unified field theory was eclipsed in the 1940-1950s. Moreover, there is a psychological explanation for this in that if too much is expected from something and if the expectations are unfulfilled an aversion

to the object of the expectation sets in. The unified field theory was this "something". At the subject's lowest ebb, any investigation on this subject was considered a waste of time compared to the background of the then burgeoning scientific and technological revolution. This opinion also tainted to some extent on general relativity theory as a whole and these studies were not encouraged (to put it euphemistically) for some time.

But science is controlled by inexorable laws of dialectics. In the 1970s, the situation started to change rapidly and once again interest in a unified theory that would combine all known interactions and elementary particles reappeared. However, if in the 1920-1930s only two interactions were on the agenda (gravitational and electromagnetic), by the 1970s four basic types of interaction had to be united, since the weak and strong interactions had been discovered. In the interim a great deal of new information was accumulated about the Universe so we can say that the old idea was revitalized from a new basis.

In this new era, many old and forgotten papers were rediscovered and the ideas, techniques, laws, and regularities they contained were used when developing the more modern theories. It is therefore interesting for us to recapitulate the main trends and achievements of these former studies, to delineate their results and drawbacks, to bridge the gap between them and modern studies, and to show how these ripened seeds can be rationally used for the benefit of the future. There is no future without a past! Even though they are unrelated to the unified field theory, old papers on non-Einsteinian theories of gravitation revealed much that could be used to develop the unified theory.

The generalizations of Einstein's theory of gravitation can be divided into four groups according to the basic geometric concepts, such as connection, metrics, and dimensionality, that were used:

(1) theories based on four-dimensional manifolds describable by differential geometries more general than the Riemannian geometry used in the general theory of relativity;

(2) theories which use additional factors, such as a scalar field, different definitions of the metric, a second metric, etc.;

(3) theories based on different concepts of the physical picture of the Universe;



(4) theories in Riemannian spaces but with more than four dimensions.

These four classes can be supplemented by other groups of theories, as well as by those that are combinations of mentioned above. In this section we will briefly discuss the first three groups. We cover the last group, multidimensional geometries, in Sec. 3.5.

**3.2.2. Physical Theories Based on Four-Dimensional Geometries More General Than Riemannian.** 1. *The Weyl and Eddington unified theories of gravitation and electromagnetism* [25]. Remember that in the Riemannian geometry, which is the foundation of the general theory of relativity, the components of the vectors and tensors change when the vectors and tensors are parallelly transported from one point to another. This change is defined by the Christoffel symbols  $\Gamma_{\alpha\beta}^{\mu}$ . For the vector  $B^{\mu}(x)$ , the change has the form  $\delta B^{\mu} = \Gamma_{\alpha\beta}^{\mu} B^{\alpha}(x) \delta x^{\beta}$ , where  $\delta x^{\beta}$  is the difference between the coordinates of the two nearby points ( $\delta x^{\beta} = x'^{\beta} - x^{\beta}$ ) between which the vector is transported. The lengths of the transported values are however conserved.

In order to geometrize the electromagnetic field, Weyl proposed a more general geometry, in which a parallel transport of tensors would change both their components and their lengths. Mathematically this is done by replacing the Christoffel symbols by more general connection coefficients, i.e.  $\Gamma_{\alpha\beta}^{\mu} \rightarrow \tilde{\Gamma}_{\alpha\beta}^{\mu}$ . Equation (1.23), which connects  $\Gamma_{\alpha\beta}^{\mu}$  to the metric tensor  $g_{\alpha\beta}$ , is no longer valid for  $\tilde{\Gamma}_{\alpha\beta}^{\mu}$ . There is a physical hypothesis behind this: the change of length is due to the electromagnetic field and if there is no electromagnetic field, the Riemannian geometry and all the formulae of general relativity reappear. If any domain of space-time contains an electromagnetic field, the new geometry has to be used. Weyl specified the way in which the connection  $\tilde{\Gamma}_{\alpha\beta}^{\mu}$  is realized through electromagnetic field as

$$\tilde{\Gamma}_{\alpha\beta}^{\mu} = \Gamma_{\alpha\beta}^{\mu} - g_{\alpha}^{\mu} A_{\beta} - g_{\beta}^{\mu} A_{\alpha} + g_{\alpha\beta} A^{\mu},$$

where  $A_{\beta}$  is the electromagnetic vector potential times a dimensional constant.

Weyl went on to show that, besides the group of admissible coordinate transformations (1.8), the theory has a group

of scale transformations for all quantities, a group of conformal transformations  $\overset{*}{g}_{\mu\nu} = \varphi^2 g_{\mu\nu}$ , where  $\varphi$  is a scalar function of coordinates. Note that angles keep their values in the two Riemannian geometries whose metric tensors are connected in this way, but all the lengths are related in terms of the function  $\varphi$ .

Later Eddington showed that Weyl's results could be obtained by using an even more general relation between  $\tilde{\Gamma}_{\alpha\beta}^\mu$  and the electromagnetic vector potential  $A_\alpha$ .

The reader may naturally ask why these two theories were not generally accepted. There are two main reasons. Firstly, to obtain the standard theory of electromagnetism from these theories additional postulates which seem out of place have to be introduced. Secondly, the new effects predicted by these theories cannot possibly be tested experimentally.

Nonetheless, Weyl's theory has had a significant imprint in the theory of space and time. This is because it demonstrated a new type of differential geometry more general than Riemannian, and in this respect its effect was similar to that left by the Lobachevski geometry in relation to Euclidean. Secondly, his theory introduced conformal mappings and the notion of conformal invariance. These are now widely used in theoretical physics to describe zero-rest mass particle fields (see also Sec. 3.5).

2. *Theories of gravitation using torsion.* It must be emphasized that both in the Riemannian and in Weyl and Eddington's theories, the connection coefficients are symmetric with respect to the subscripts, that is  $\tilde{\Gamma}_{\alpha\beta}^\mu = \tilde{\Gamma}_{\beta\alpha}^\mu$ . What will happen if we take a nonsymmetric connection? This sort of geometry was first proposed and developed by Élie Cartan [14]. In this geometry the rule of parallelogram we learnt at school for adding vectors no longer holds. This effect can be explained as follows: if a small segment  $dx^\mu$  is parallelly transported along a small segment  $\delta x^\mu$ , and then the process is inverted by moving  $\delta x^\mu$  along  $dx^\mu$ , the ends of vectors  $dx^\mu + \delta x'^\mu$  and  $\delta x^\nu + dx'^\nu$  will not meet, the resultant gap being determined by the skew-symmetric part of the connection ( $\tilde{\Gamma}_{\alpha\beta}^\mu - \tilde{\Gamma}_{\beta\alpha}^\mu$ ). Here we shall point out where geometries with torsion are used:

(1) It is known that Einstein tested a variety of unified

theories of gravitation and electromagnetism. The last version he investigated used a geometry with torsion. His initial assumption was nonsymmetric metric, that is, a metric tensor in the form  $g_{\mu\nu} = g_{\mu\nu} + g_{\nu\mu}$ , where  $g_{\mu\nu} = g_{\nu\mu}$  is the symmetric part that is usually applied to determine the distance  $ds$ , whereas  $g_{\mu\nu} = -g_{\nu\mu}$  is the skew-symmetric part that does not affect distances (because  $g_{\mu\nu}dx^\mu dx^\nu = 0$ ).

The latter was compared to the field tensor of electromagnetic field  $F_{\mu\nu}$ . When the relationship between the metric and the connection coefficients is set up in this sort of theory the connection coefficients become nonsymmetric, that is, they contain a torsion tensor. Einstein studied several geometries with nonsymmetric metrics.

(2) R. Finkelstein was one of a number of physicists who used spaces with torsion tensors to introduce a matter with several types of geometric "charge". These charges were then compared against the various physical fields such as electromagnetic, meson, and others.

These two variants of unified field theory were also unsuccessful for the same reasons as was Weyl's.

(3) In the last decade, several investigators have reactivated studies of spaces with torsion for applications outside the unified field topic. What we mean here is the Einstein-Cartan theory, which has, in addition to Einstein's ten equations for the metric  $g_{\mu\nu}$ , a system of additional equations for the torsion tensor. In Einstein's equations, a source is represented by the energy-momentum tensor of the matter,  $T_{\mu\nu}$ , whereas in a torsion field, a source is represented by a tensor which is defined by the spin properties of the matter (its rotation). The physical possibilities of this sort of theory are under examination.

(4) It should also be noted that the modern theory of supergravity has, in a natural way, skew-symmetric connection, that is, torsion. This theory is being seriously developed and studied by many researchers.

3. *Schouten's differential geometries.* The cases we have discussed are not exhaustive of all the possible generalizations of Riemannian geometry. In the 1930s, J. A. Schouten (born 1883) analyzed the results of Weyl, Eddington, Cartan, and so on and formulated a number of general requirements for differential geometries. By doing so he found another type of geometry. To explain this category we must

first remind you that in the general theory of relativity all the tensors are either covariant or contravariant (with sub- or superscripts). In all the geometries we have so far discussed tensors can be parallelly transported using the same connection coefficients. In the new type of geometry, a parallel transport of covariant and contravariant tensors requires substantially different connection coefficients.

Schouten also showed that every possible geometry that meets the requirements he had formulated can be characterized by three and only three tensors of rank three. Each of these tensors may be zero, degenerate (that is, it can be decomposed into the product of a vector and a tensor) or nondegenerate. In the long run, if we consider all the possible combinations of these three tensors and their types, we find only that  $3^3 = 27$  types of differential geometry are possible. The Riemannian geometry of the general theory of relativity is the simplest case, for which all three Schouten tensors vanish.

**3.3.3. Gravitation Theories with Additional Factors.** These theories are developed not to combine gravitation and electromagnetism, but, mainly, to resolve other problems, such as whether Dirac's hypothesis that the gravitational constant is changing with time or Mach's principle that inertia is dependent on distant matter are theoretically justifiable. The additional factors introduced into the theory may include a new scalar, a vector or tensor field (second metric), or some other quantity or idea. We shall discuss three of the more hotly debated of these theories.

1. *Scalar-tensor theories of gravitation.* In this class of theory the tensor field of the metric tensor  $g_{\mu\nu}$  is considered along with a scalar field  $\varphi$ , which gives it this name. This approach was started by studies of the five-dimensional unified theories of gravitation and electromagnetism, into which an additional geometric scalar field was introduced (see Sec. 3.5). At a later stage, the five-dimensionality was dropped, and the scalar field began its independent existence. These studies were begun by P. Jordan (1948), Y. Thiry (1948), and W. Scherrer (1949). Then, a similar theory was developed by Brans and Dicke early in the 1960s [15]. The most complete variant is sometimes called the Jordan-Brans-Dicke theory. Interest to these theories flared up in the 1960s and the beginning of the 1970s.

The field  $\varphi$  is not directly connected to fundamental geo-

metric notions (metric, connection) and so it must be introduced by means of additional postulates. Without going into details, we shall briefly discuss here only main characteristics of such a theory:

(a) The scalar field  $\varphi$  is represented in Einstein's equations, which also emerge in this theory, by the coefficient in front of the  $T_{\mu\nu}$  on the right-hand side. Remember that in Einstein's own equations we have  $\kappa T_{\mu\nu}$ . Therefore, the gravitational constant  $\kappa$  now becomes a variable which depends on the scalar field  $\varphi$ . This is one way of implementing mathematically Dirac's hypothesis, which he only outlined qualitatively.

(b) In a scalar-tensor theory, the mass of a particle depends on the scalar field:  $m = m_0\varphi^{-1/2}$ , where  $m_0$  is a constant. Since  $\varphi$  is determined by the matter distribution in this theory, we can say that the inertial masses are dependent on the distribution of the surrounding matter. This means that Mach's principle is valid to a certain extent.

(c) In addition to Einstein's ten equations, the theory has an eleventh wave equation for a scalar field (with zero rest mass). In this way the field  $\varphi$  is an additional fundamental field, analogous, to a certain extent, with the electromagnetic field. We would therefore expect observational manifestations of this field.

2. *Finsler geometries.* This direction of research has a more profound geometric basis than scalar-tensor theories. Generally speaking, the metric in a geometry can be formulated independently of the connection. In Chapter One of *About the Hypotheses Lying at the Foundation of Geometry*, Riemann not only developed a new geometry on the basis of his method of length measurement, but also indicated that more general metrics were possible. We would like to remind the reader that in the Riemannian geometry an element of distance between two close points  $dx^\mu$  apart is given as the square root of the quadratic form,  $ds = \sqrt{g_{\mu\nu}(x) dx^\mu dx^\nu} \equiv F(x^\alpha, dx^\alpha)$ . Riemann wrote that an element of distance might be, for example, the fourth root of a quartic form, though this geometry would have to be very complicated. For simplicity, geometers have, in the sixty years since Riemann, concentrated on developing geometries with a quadratic form of distance measurement.

The development of a more general geometry based on a metric began with a thesis by D. Finsler in Göttingen,

Germany (1918). Finsler's metric  $F$  is a scalar function of the point  $x^\alpha$  of a manifold and the vector  $y^\mu$  at this point. Function  $F(x^\mu, y^\mu)$  is interpreted as the length of the vector  $y^\mu$  at point  $x^\mu$ . This function is restricted by one important condition, namely that it is assumed to be a homogeneous function of the first degree with respect to the vector parameter, that is,  $F(x^\mu, ky^\mu) = kF(x^\mu, y^\mu)$  for all  $k > 0$ . This means that for any collinear vectors,  $y_1^\mu$  and  $y_2^\mu$ , the ratios of their components are the same.

Many attempts have been made to reformulate the theory of gravitational field on the basis of Finsler's metric, and to study such a theory. But to do this the following question must be answered: How should Finsler's generalization be done for the gravitational field equations and the equations of motion of gravitational field sources?

Even though at the moment there is no hard experimental evidence to suggest that Riemannian geometry should be substituted by a Finslerian to describe the properties of physical space-time, the development of the Finslerian generalization seems to be an interesting program for the future. We expect that an application of this geometry to physical problems will be useful when the metric tensor in the neighbourhood of a point is anisotropic, that is, when it depends on the direction of the vector being measured.

3. *Bimetric theories of gravitation.* In this class of theories, the space-time manifold has two Riemannian metrics, viz.  $g_{\alpha\beta}(x)$  and  $g_{\alpha\beta}^*(x)$ , and not just one, as was the case in all earlier geometries. The extra metric may be introduced by any one of a number of possible ways. One of the best known is Rosen's bimetric theory (there are several versions). Its second metric is postulated. Everything that has been said about Riemannian spaces with one metric is duplicated in these theories but there are a number of additional aspects that are determined by the relation between the two metrics. An important feature of bimetric theories is that the difference between the Christoffel symbols of each metric is a tensor:  $\Gamma_{\alpha\beta}^\mu - \tilde{\Gamma}_{\alpha\beta}^\mu = F_{\alpha\beta}^\mu$ . This result can be proved starting from the transformation law of Christoffel symbols under coordinate transformations. Various other considerations are used to derive the equations for one metric against the background of the other. Gravitation in this geometry is described as the difference between the

two metrics, that is, the strength of the gravitational field is described by tensor quantity  $F_{\alpha\beta}^{\mu}$ . The advocates of this approach believe that this description has advantages over most general relativity. We should note that most of the work on the quantization of the gravitational field has in fact been done in the framework of bimetric theories.

The serious objections levelled against these theories are caused by the opacity of the physical sense and observability of the second metric.

4. We can mention *a number of other theories* which lead, for example, to fourth-order gravitational field equations rather than the second-order equations of Einstein, or use Lagrangians quadratic in the curvature tensor. The latter are much used in the conformal theory of gravity which also arises from the supergravitation (see Sec. 3.5.8) and the twistor theory of R. Penrose.

**3.3.4. Action-at-a-Distance Theories of Gravitation.** This theory is different from the other theories of physics because it is based on the concept of "action at a distance", which is opposite to the generally accepted concept of a field approach (short-range action). The notion of field no longer has significance in an action-at-a-distance theory. Particles affect each other "at a distance" symmetrically, that is both in a retarded and advanced way. The interaction between particles is defined by setting up an equation for the Green function against the background of a space-time manifold, which may either be flat (for convenience) or curved.

Remember that the battle between the advocates of the short-range and long-range interaction concepts has been going for several centuries, and no winner is yet in sight. Newton's theory of gravitation was a long-range one, as was the initial theory of electrical interactions. The emergence of Maxwell's theory of electromagnetic interaction seemed to signal the victory of the short-range concept. But this was not the case. It turned out that the concept of a long-range at-a-distance interaction needed some clarification by accounting for temporal retardation (or advance). The foundation for such a theory was laid by Gauss in the mid-nineteenth century. Then, in the 1920s, H. Tetrode and A. D. Fokker made another significant contribution. Fokker established principle of direct electromagnetic interaction and showed that the standard equations of charged particles' motion follow from it, while the second pair of Maxwell's

equations became identical. That retarded and advanced effects were equivalent in this theory is a drawback which has long delayed its acceptance. Everybody knows that in a real world cause precedes effect, and not vice versa. This drawback was removed by R. Feynman (b. 1918) and J. A. Wheeler (b. 1911) in the 1940s. They showed that all the difficulties had arisen because the way in which the rest of matter in the Universe affects the interacting particles had not been adequately allowed for. The introduction of an absolute absorber in the future gave a correct account of this influence of the Universe, thus eliminating all advanced interactions and bringing the theory into accord with the observed interactions, which are all retarded.

It has now become clear that the theory of direct interactions can describe the electromagnetic interaction as well as the generally accepted field theory does. It even has a number of advantages in that it naturally describes the connection between the local characteristics of matter (on a small scale) and the global properties of the whole Universe [13].

The above concerns direct electromagnetic interactions, but what about gravitational interactions? Remember that when Einstein was working on the general theory of relativity, he based his considerations on Mach's ideas and Mach's criticism of Newtonian mechanics. If we look closely at Mach's works, we see that this criticism was based on his belief in long-range interactions. But the theory Einstein created was a typical field theory, that is, it used the concept of short-range interactions. It is also known that Einstein was less enthusiastic about Mach's ideas after he had constructed his theory and even criticized them. How can we explain this? Obviously, he realized the difference between Mach's concept and general relativity. Further progress in physics however showed that the action-at-a-distance theory could be applied to gravitation as well. The theory has now been developed for the gravitational interaction and is even quite elegant. There are a number of versions, including a scalar-tensor one (the Hoyle-Narlikar theory) and one which coincides in terms of its practical results with standard general relativity.

We shall briefly discuss this last variant of gravitational interaction. We may use the following three main components to describe this theory:



(a) The postulate that gravity acts against the background of a fixed space-time; in the simplest case this background is flat. We can then show that the equations of motion for particles can be derived in the form of geodesics imbedded in an effective Riemannian metric. Difficulties arise because of the nonlinearity of gravitation (they are absent in electrodynamics) and to overcome them additional assumptions concerning the action principle have to be introduced. They are needed to account not only for pair interaction, but also for triple, quadruple, etc. interactions.

(b) It was shown that in the action-at-a-distance theory of gravitation, the relations that correspond to Einstein's standard equations are deduced from the rest of the theory.

(c) Finally, the third part of the theory involves the gravitational absorber. It was shown that by accounting for all matter in the Universe it is possible to exclude from the theory all advanced gravitational interactions. The equations of the geodesics acquire additional terms that are interpreted as the effect of gravitational radiation damping. As was the case for direct electromagnetic interaction, the theory contains questionable aspects with regard to the physics of the preferred direction of time (the time arrow).

At present, this theory is still under construction, and it faces a number of unanswered questions and unresolved problems, some of which pertain to general relativity while others are specific to the action-at-a-distance theory of direct interaction.

**3.3.5. Different Approaches to Constructing Physical Theories.** We have just described a variant of gravitational theory based on the unusual concept of action-at-a-distance interaction. As to whether there are more concepts the answer is closely related to two well-known philosophical attitudes towards space and time, the substantial and the relational.

The adherents of the substantial approach interpret space and time as a background substance, an independent container of all the observable types of matter. A similar concept was known in antiquity, during the time of Democritus, in particular. Newton and Galileo's classical concept of space (and time) belongs, strictly speaking, to the substantial approach.

The second, relational approach, rejects the idea that space and time are substantial independent entities. It postulates instead that space and time are only forms of

existence for material objects. The two entities can only describe relationships between material objects. Aristotle made statements in this spirit and later Leibniz and Mach were advocates. To which of the two concepts does general relativity belong? We have pointed out that it was formulated as a field theory, and therefore it follows the substantial approach. Einstein's equations, for instance, admit solutions for a vacuum, that is, in the absence on their right-hand sides of an energy-momentum tensor for any type of matter. The Minkowski metric, for example, is one such rigorous special solution.

At the same time we cannot claim that the theory of general relativity is an embodiment of the substantial concept of space-time in its extreme form. In fact, we have at present an intermediate concept in which geometrical characteristics (with which general relativity deals) partly describe space-time relationships, and partly a new type of geometrical or gravitational matter. There is a great deal of work aimed at classifying geometrical quantities explicitly according to these two functions.

Today's scientific literature includes a sizable number of papers (and generalizations) that present general relativity from the point of view of an extreme substantial approach to space-time. Clifford, whose ideas we discussed in Chapter One (see Sec. 1.4), is said to be a founder of this approach. Its advocates believe that their main objective is to define and obtain objects from geometrical descriptions of space-time (of metric or topological nature) which can then be identified with the observable types of matter. Attempts to obtain particle-like solutions of Einstein's equations (in a vacuum) are in this group. Another method of geometrizing particles uses more general topologies, the main idea being that the fabric of space is like a surface pierced by pairs of holes that are connected by tubes (topological handles). These wormholes in the surface are identified with pairs of particles. Electrical lines of force enter one end of a wormhole and leave the other, hence the particles in each pair must have opposite electrical charges. These investigations are still under way.

It follows from the above that the action-at-a-distance theory adopts the other, relational, concept of space and time. However if we use this approach, all the problems we were just discussing have no meaning. Particles are the

primary objects, whereas all the geometrical parameters and properties simply describe the relationships between the particles.

We must admit, however, that given the present state of the art the action-at-a-distance theory cannot fully reflect the relational concept of space-time. Indeed, to formulate the theory mathematically, both the presence of particles and the background space, which can be of any type, flat presumably, have to be postulated from the very beginning. A full realization of the relational concept, on the other hand, would require a formulation of the theoretical foundation which can stand by itself without a background space-time acting as a "support" [108]. There is no such theory at the moment, and meanwhile we only have some considerations and preliminary results from A. S. Eddington [26], D. van Dantzig, E. J. Zimmerman, among others. A more detailed discussion of these ideas lies beyond the scope of this book.

**3.3.6. Some Conclusions.** As we come to the close of this section we must regretfully conclude that we have fallen far short of covering in any adequate way the attempts that have been made to generalize Einstein's gravitation theory; we have even failed to mention all of the more ingenious proposals. Any detailed look at this field would require a book of its own. We should also say that the quantization of gravitation and multidimensional theories covered in the next sections are also attempts to generalize the standard theory of general relativity.

Hopefully we have created in the mind of the reader an image of a general relativity which has demonstrated a remarkable resistance to attempts to generalize or modify it. The above information has proved that the efforts of scientists over several generations have served only to build a compact outer shell of modifications around general relativity. But none of them has compelled us to replace the core—Einstein's gravitation theory—with any other theory. This is because Einstein's theory of general relativity has no match in elegance, beauty or laconicity in terms of its main assumptions. The one exception, perhaps, are the five-dimensional (or six-dimensional) theories, which naturally incorporate general relativity (see Sec. 3.5).

How then should we regard these generalizations of the theory? We can by no means agree with assertions of many

authors, some very prominent, that the evolution of science has shown up the unfeasibility of these ideas and approaches. The results these ideas have produced have contributed to the treasures of the science of physical space and time. They have allowed us to look deeper and broader into the nature of modern gravitation theory, have opened up many new angles of established regularities, and have indicated other possible generalizations of the theory. They must always be borne in mind and reviewed periodically to tie them in with new facts and data. Many "outworn" ideas have already been reestablished and incorporated in modern theoretical studies. Think how many other "old" ideas will reappear in future theories, which will naturally absorb the apparently dead ideas in a new light!

### **3.4. GRAVITATION AND QUANTUM PHYSICS**

**3.4.1. About the Need to Quantize Gravitation.** Theoretical physicists generally agree that quantizing gravitation, or as they often put it the development of a "quantum gravity theory", is a most important aspect of the modern theory of gravitation. It is in this field we expect the most substantial breakthroughs in the understanding of the nature of space-time and the whole of physics.

We should remind our readers that no quantum gravity theory yet exists. Moreover, at the moment we have no precise definition of what form such a theory should take or what exactly gravitational quantization would mean or involve. Scientists have a wide variety of ideas on this topic from the belief that existing quantum field theory must be accurately translated into gravitation, to the view that the solution of the problem requires an entirely new approach. Let us formulate this problem in its most general way. Modern theoretical physics has at the moment two fundamental divisions: quantum field theory and general relativity theory. So far, each of them has been developing independently, each with their own principles and notions. The problem now is to construct a theory that can combine the principles underlying the two theories.

The inquisitive eyes of our readers will now glisten in question: may be it is not all bad that there are two parallel theories? Or, perhaps, the two will merge sometimes in the future? We shall try to address these points in this section.

There are some very frequently debated topics underlying this problem.

1. Normally, gravitation is viewed as a physical field, equivalent in a number of ways with the other known fields, e.g. the electromagnetic and meson ones. It should therefore have the same general properties as all the other fields and, in particular, if all the other physical fields can be quantized, so should the gravitational field. It has been pointed out on several occasions that if the gravitational field is like the others in nature, then the coordinates and momentum of a particle could in principle be found more accurately using gravitational interactions than the respective uncertainty relations permit, and this would mean that the main principles of quantum theory could be violated.

2. We have now discovered in the Universe astrophysical objects with very curved space-times in their vicinities. These objects are the pulsars, which are neutron stars, and possibly quasars, while the black holes hypothesis is being widely discussed. The visible manifestations of these objects must substantially depend on the behaviour of the matter and elementary particles in their neighbourhoods. Hence, we must be able to describe quanta in significantly curved space-time. To do this, the laws of general relativity and quantum theory must be combined.

3. It has recently become clear that the general cosmological solution or at least approximate models of the Universe have, in the framework of Einstein's theory, tough singularities. In other words, the metric of space-time is only regular in a limited (from one side or from both) time period. These singularities must be interpreted as a proof that general relativity becomes invalid near these limits. Tremendous matter densities occur at these limits and these situations must be described by a new theory that substantially covers the laws of quantum and microphysics.

4. There are quite sound reasons for believing that the theory of gravitation could be used to construct a theory of elementary particles. In this field there are several directions in which to seek. One is to construct geometric models of the elementary particles, and to find particle-like solutions to Einstein's and Maxwell's equations. We discussed this in the previous section. There are other more sophisticated approaches which attempt to combine gravitation and the other fields.

5. Theoretical physicists are convinced that a new quantum gravitational theory will shed light on the basic difficulties of modern quantum electrodynamics and strong interactions theory. For example, it may allow them to eliminate in a correct way the divergences in these theories. That such seems possible is indicated by the appearance of characteristic distances, that is, the gravitational radii of particles and a new universal constant  $l_0 = \sqrt{\hbar G/c^3} \simeq 10^{-33}$  cm, Planck's distance, which includes the gravitational constant, the quantum constant  $\hbar$ , and the relativistic one  $c$ .

There are other arguments for the quantization of gravitation, and we shall discuss some of them later.

If we have convinced the reader that gravitation must be quantized, it is time for him to ask why it has not been quantized yet. There are several reasons for this and each can be discussed from different positions. The nature of these explanations depends on the formulation of the problem. If we straightforwardly set our task as the translation of quantum field theory into gravitation theory then two obstacles arise: (a) Einstein's equations are nonlinear, and (b) the theory is generally covariant in nature.

(a) Nonlinearity in equations usually means that sums of solutions are not themselves solution. Hence, the standard technique for quantizing free fields, that is representing solutions in the form of a sum of elementary quantum contributions of the field. Each elementary quantum being described by an exact solution is not possible for gravitation. This is because Einstein's equations are nonlinear.

(b) The covariance of gravitation theory leads, in the long run, to larger sets of variables (for example, the ten components of the metric tensor) to describe it than would be expected from the number of dynamic variables (two degrees of freedom). A similar situation occurs in the theory of electromagnetic fields, where the four components of the vector potential correspond to only two dynamic variables (two states of polarization). In electrodynamics, however, the equations are simple, and it is not difficult to eliminate the "redundant" variables. In general relativity, on the other hand, the equations are extremely complex, and nobody has yet succeeded in eliminating the "redundant" variables in the general case.

The lack of experimental evidence is another difficulty.

More often than not, physical theories evolve from known facts. A quantum gravity theory is being sought from purely logical arguments. At the present, even gravitational waves, the classical aspect of the even more hypothetical gravitons, have not yet been discovered.

**3.4.2. Preliminary Results.** 1. *Quantizing a weak gravitational field.* The difficulties we discussed concerning quantizing gravitation can be eliminated if we assume that for practical reasons a weak gravitational field is sufficient.

The metric tensor will then be  $g_{\mu\nu} = g_{\mu\nu}^0 + h_{\mu\nu}$ , where  $g_{\mu\nu}^0$  is the metric tensor of the Minkowski space-time, and  $|h_{\mu\nu}| \ll 1$  is a small addition. The gravitational field is thus described by  $h_{\mu\nu}$  which can be treated as a standard tensor field of rank 2 imbedded in a flat space-time. The usual technique for quantizing a linear field can now be applied and  $h_{\mu\nu}$  is quantized in exactly the same way the electromagnetic field is. The redundant components are excluded, leaving only transverse-transverse combinations  $h_{23}$  and  $1/2 (h_{22} - h_{33})$  to be quantized (if the field wave propagates along the  $x^1$  axis). In classical general relativity, these combinations would correspond to two polarizations of the gravitational waves, whereas in the quantum theory they describe two types of graviton.

It is next assumed that gravitons behave like photons in various ways, i.e. they can be created, they interact with each other and the quanta of other types of matter, and they can be absorbed. Incidentally, the nonlinear terms in Einstein's field equations are interpreted as the result of gravitons interacting among themselves. As in quantum electrodynamics, diagrams (like Feynman's but more complicated) are used for calculations.

The consequences of transmutating gravitons and other particles have been calculated in a number of papers. The effect of the gravitational annihilation of elementary particles is very interesting. For example, there is the annihilation of an electron and a positron to produce two gravitons, rather than the two photons as in electrodynamics. Roughly speaking, this process involves the transformation of conventional matter into a gravitational field. The calculations have shown that these effects must be extremely weak at elementary particle energies now obtainable. For example a head-on collision between an electron and a positron at a

speed of about one hundredth of the velocity of light would yield a cross-section for the gravitons creation of about  $10^{-110} \text{ cm}^2$ . This is fantastically small. It has been noted, however, that as the energy of the colliding particles increases, this cross-section grows in proportion to the square of the energy, whereas cross-section of the well-known electron-positron annihilation into two photons is inversely proportional to the square of the energy. At colossal energies of  $10^{21} mc^2$  (where  $m$  is the mass of electron), these two effects become equal, and as the energy rises further the two-graviton process dominates over the two-photon one. Theoretically the electron-positron pair may annihilate into one photon and one graviton. This is of higher probability than the annihilation into two-gravitons but it still remains too weak for practical study. If the particles moved at one percent the velocity of light the relevant cross-section is  $\sigma \simeq 10^{-75} \text{ cm}^2$ . Of course, the reverse processes are also possible, i.e. the conversion of two gravitons into electron-positron pairs or other particles. Other processes have also been discussed [70, 107], e.g. gravitational bremsstrahlung radiation of accelerated particles, graviton transmutation into photons and back. The most significant such effect (first order with respect to the gravitational constant) is the transformation of gravitons into photons and back in an external electromagnetic field. Simply this is the conversion of cosmic gravitons into photons within a high-capacity capacitor. This might be the basis of a graviton detector. Unfortunately, the effect is too weak to make an experiment feasible in the near future.\*

At first, the problem seems to have been solved in principle. However, the present situation cannot be accepted as satisfactory for several reasons:

(a) The quantum theory of a linearized gravitational field cannot be renormalized. This means that when the processes are recalculated using a higher approximation (in powers of the constant  $G$ ), divergences appear which cannot be eliminated by selecting a final number of counterterms (which is possible in electrodynamics). In other words, bringing back the terms that were neglected during the linearization, meaningless expressions are generated.

---

\* The conversion of gravitons into photons (and vice versa) must be more frequent in the tremendous magnetic fields surrounding pulsars (see some data concerning these fields in Sec. 2.7).



(b) It is very difficult to reconcile ourselves with the idea that a quantum theory of gravity might be a depressingly primitive duplicate of quantum electrodynamics. Below we bring out some evidence that makes us believe that the situation is in fact very different.

(c) In the weak-field limit which actually means the use of a flat background space-time, a number of additional assumptions foreign to general relativity implicitly appear which need a healthy foundation and rigorous analysis. So far, this has not been done.

2. *Particle pair generation in an expanding Universe.* Another important aspect of this investigation is the latest work on the effects of the generation of elementary particle pairs in nonstationary models of the Universe. We can say that these effects generalize, in some ways, the gravitational transmutation of gravitons and quanta of ordinary matter that we have discussed. The difference between nonstationary effects and the transmutation processes is that in the former case the pairs are generated by variations in the gravitational field which cannot in principle be described in terms of particles (as a superposition of solutions or a set of gravitons). Note that the situation is similar to that in a nonstationary classical electromagnetic field. The objective of this investigation is to combine the principles of quantum field theory with general relativity. The idea of the approach is that in a curved space-time, in a nonstationary Friedmann metric in particular, difficulties appear when stating an expression for the energy of secondarily quantized fields of ordinary matter. Since the components of metric tensor are present in all the expressions the field operators have factors which, in the general case, depend on time. As a result, the decomposition of field functions into positive (that correspond to the particle creation operators) and negative frequency parts (which correspond to the annihilation operators) becomes ambiguous and time dependent. Over time, these operators mix, and the number and density of the particles and antiparticles, however they were determined, are changed. This process is interpreted as the effect of generation (annihilation) of pairs in a nonstationary Universe.

This result was noticed for the first time by **Erwin Schrödinger** (1887-1961) in 1939, who found that when he was calculating particle scattering in an expanding Universe,

the operators for the particles and antiparticles got mixed. Much later, the same problem was studied by a number of other physicists, and the interpretation of the particle creation was widely accepted. **Wolfgang Rindler** (b. 1924) showed that even in the flat Minkowski universe, there is a horizon for an accelerating observer. **Leonard Parker** of the University of Wisconsin and a group of Leningrad physicists (**A.A. Grib**, **S.G. Mamaev**†, **V.M. Mostepanenko**) looked further into the topic and they found that a vacuum assumes an entirely new property in the presence of a horizon, viz. temperature, and this makes particle creation possible. However, this type of particle creation in a flat universe (as well as in some similar cosmological situations) depends on the act of measurement. If we were to accelerate uniformly in the Minkowski universe, for example, we would detect photons and charged particle pairs, which would be created in the detector itself as it accelerates. These kinds of particles must in fact exist and they could be detected in an inertial frame as by-products of a foreign object accelerating in this frame of reference. In cosmology, there is another situation in which the same process may occur, since a non-stationary gravitational field (like an analogous electromagnetic field) can transform its energy, momentum, and angular momentum into the same characteristics of particles, thus creating them. It has been shown that in the framework of the nonquantum model of the universe under consideration, most pairs should have been created during the early stages of the universe's expansion, or more exactly, close to the time  $\tau_0 \sim \hbar/mc^2$ . (Note that  $\tau_0$  is considerably greater than Planck's time  $\tau_{Pl} \simeq 10^{-43}$  s, where the quantum gravitation theory is most essential.) Later, at  $\tau \gg \tau_0$ , negligibly few pairs are created.

Some investigators have shown that more particles would be expected in models in which the universe is anisotropic. In these models, zero-restmass particles (photons, neutrinos, etc.) can be created. In Friedmann models, only particles with a nonzero restmass can be created because of conformally flat nature of these models (see the Weyl conformal transformations in the previous section). Studies of particle pair creation during periods of gravitational collapse and in the neighbourhoods of black holes belong to this category of investigation.

These studies have undoubtedly made some of the theo-

retical consequences of combining the principles of quantum field theory and general relativity, as adopted in certain approaches, clearer. These results must however be viewed as preliminary, and there are several drawbacks:

(a) The formulation of the problem is not self-consistent because the creation of matter is determined by external factors. It would therefore be more helpful to have a consistent quantum theory of general relativity.

(b) The approach is based on classical (nonquantum) concepts of space-time in the neighbourhood of singularity, but it is doubtful that they are valid at the time and in the space domain where the calculated effects should be most intensive.

(c) Calculations show that the density of the matter created in this way in the Universe is smaller than the observed density\*.

3. *Quantization of closed cosmological models.* These investigations counter completely different problems in that they try to explain several of the underlying features of the quantum theory of the world as a whole (quantized Universe). Of course, everyone realized that creating such a complete theory is extremely difficult and that it must be attacked in several stages, beginning from a very simple approach that is gradually made more sophisticated. In other words, to begin with, a symmetric model of the Universe is used and an appropriate quantization scheme is adjusted to it.

It was natural that Friedmann's homogeneous, isotropic, closed model was selected as the first step. Bryce S. DeWitt of Texas University at Austin studied one possible quantization procedure. He took one of Einstein's equations (the 00-component or what is called the Hamiltonian constraint) in the form corresponding to the expression for energy (Hamiltonian). The geometric (left-hand) part of this can easily be represented as a sum of two parts, one of which can be interpreted as kinetic energy and the other as potential energy. This cosmological model is very simple because it has only one parameter: the radius of the universe,  $R$ . The kinetic part is represented by the square of the time derivat-

---

\* Recently, scientists have succeeded in elucidating why matter dominates over antimatter. The Grand Unification Theory was used and accounted for the spontaneous breaking of the baryon number conservation law during the early stages of the universe's expansion (see below).

ive of  $R$ , that is, we can say that it has the form of squared momentum. Then, the model is quantized using a procedure similar to that for quantizing a particle, i.e. the relation between the energy and the momentum of the particle is replaced by an operator expression which results in a Schrödinger or Klein-Fock-Gordon wave equation. Finally, we arrive at the DeWitt equation, which is reminiscent of a unidimensional Klein-Fock-Gordon equation. This wave equation is solved to yield the evolution of the universe (an expansion followed by a contraction) as that of a unidimensional wave packet in a special type of a superspace.

The next step was made by **Charles W. Misner** of the University of Maryland. He operated on a more complex homogeneous anisotropic model. Evolution in his model came out as the alternating contraction and expansion of the universe along different coordinate axes. In this case the universe is described by three parameters: one corresponds to the general scale of the universe, and the other two determine the extent of its anisotropy. Like DeWitt's procedure this approach leads to a three-dimensional equation (the Misner equation) which is reminiscent of a Klein-Fock-Gordon equation. The scale factor in this model plays the role of time, and the anisotropy parameters act as spatial coordinates. Here again, the evolution of the universe is described by a wave packet.

These two models are just first steps towards a more complete theory and they do not show how models with fewer symmetries can be obtained. The experience of constructing these two models has revealed a number of difficulties as to the selection of an appropriate time factor, and the establishment of the initial and boundary conditions when solving the wave equations. It has been even more difficult to ascribe physical meanings to the solutions that have been obtained. For example, what is the meaning of a state functional of a quantum universe, if there cannot exist a measuring device outside of it? How should we treat a superposition of the harmonics of a universe, each describing a single universe, some contracting and some expanding, whose number is infinite? Thus, even in the initial stages there are many questions which do not have unambiguous answers.

4. *Gravitational modelling of elementary particles.* In this section we shall explain how various authors are trying to construct models of elementary particles as objects with

a nearly closed intrinsic metric [68]. This is the same as saying that elementary particles are made up of collapsed objects, Friedmann or other types of universes, that is, matter which has shrunk within its gravitational radius. Externally these particles manifest themselves as small objects of mass  $m \sim \sqrt{\hbar c/\kappa} \simeq 10^{-5}$  g. Different writers have given them different names, e.g. maximons, friedmons, and planckeons [68, 99, 105, 112]. The trouble is that this hypothesis poses too many conceptual and philosophical problems, many more, in our opinion, than there are yet solutions from the theory.

5. *Supersymmetry and supergravity.* There are now a growing number of theorists who are getting fascinated with the Grand Unification idea, viz. a theory to unify the strong, electromagnetic and weak (and possibly gravitational) interactions on the basis of local symmetry. The principle of local symmetry allows investigators to look differently at physical fields. It is assumed that the known group properties (symmetries) such as the Lorentz (or Poincaré) transformations, gauge transformation in electrodynamics, etc., only hold locally, that is, at each point of space-time separately. Directions, gauges, and measurement standards are strictly individual for each point. When moving from one point to another they have to be compared. The notions that are needed for this comparison and which describe it in the language of each group are called gauge fields. For example, a gravitational field is viewed as gauge one with respect to a Lorentz (or Poincaré) group, the electromagnetic field as gauge one with respect to electromagnetic gauge transformations, etc. The  $SU(3)$  and  $SU(2)$  groups of symmetries are introduced into the theory that describes the strong and weak interactions and corresponding gauge fields are determined as the carriers of these interactions. The Grand Unification Theory they are seeking will be a local unification of all the known group symmetries.

The initial variants of the Grand Unification Theory that have been suggested all lead to the conclusion that at low energies the three types of interaction (weak, electromagnetic and strong) are manifested in different ways as seen experimentally. But at energies of about  $10^{15}$ - $10^{16}$  GeV, however, no distinction can be made between them, i.e. they are the same interaction. At about  $10^{19}$  GeV gravitational interaction is expected to follow suit. This means that

in about this range, a theory that unifies all four types of interaction should become valid.

This theory is based on a newly discovered supersymmetry which unifies the descriptions of bosons and fermions (that is, particles with integral and half-integral spins). Supersymmetry means that the contributions of the boson and fermion fields to all processes are symmetric, hence they can be mixed. Local supersymmetry means that this can be done differently at different points in space; thus gravitational interaction can be included in a natural way into the theory.

The first results from the supersymmetry theories have shown that it may help avoid the difficulties of divergence, the problems we mentioned above concerning the unrenormalizability of quantum gravity. Any discussion of these questions at the moment is still premature (especially in a book like this). We can only hope that the work on this topic will yield a better understanding both of the relations between the theories describing the various interactions and of the structure of space-time.

**3.4.3. Conceptual Problems of a Quantum Theory of Gravity.** The many years already spent attempting to quantize gravitation have shown that it has much deeper roots than was first expected. Once resolved however, it will probably shake the entire foundation of our concepts of space and time.

Here we shall show the reader why the parallel existence of quantum theory and general relativity poses limits on our notions of spatial distance and time intervals. Consider a particle with which we want to measure a very small distance  $\Delta x$ . It follows from the quantum-mechanical uncertainty principle that we must have  $\Delta x \geq \hbar/\Delta p$ , where  $\Delta p$  is the uncertainty in the particle's momentum. At very small  $\Delta x$  we have a very large  $\Delta p \simeq \Delta E/c \sim c\Delta m$ , where  $E$  is energy, and  $m$  is the particle's mass. Hence

$$\Delta x \geq \hbar/c\Delta m. \quad (3.29)$$

According to general relativity, the metric near a point mass is Schwarzschildian, that is,  $g_{00} = 1 - 2Gm/c^2\Delta x$ . Distance only remains distance and time time whilst the metric retains its standard signature, that is whilst  $g_{00} > 0$ . Hence

$$r \sim \Delta x > mG/c^2 \simeq G\Delta m/c^2. \quad (3.30)$$

By multiplying equations (3.29) and (3.30) together, we find a lower limit for  $\Delta x$ , which is a combination of fundamental constants, viz.

$$\Delta x > \sqrt{\hbar G/c^3} \equiv l_0 \simeq 1.6 \times 10^{-33} \text{ cm.}$$

This is called Planck's distance.

The same result can be obtained in another way, i.e. by considering the distance that can be measured by a light signal, by weighing, etc. Whatever the method the conclusion is the same, namely distances shorter than Planck's distance are impossible. Similarly, we have to conclude that time intervals shorter than  $10^{-43}$  s are equally impossible. There are also similar limitations for all the other geometric values, the Christoffel symbols and the metric tensor, for instance. While these considerations are of course somewhat approximate, they nevertheless raise doubts about any claim that we may really deal with infinitesimally small distances in the mathematical sense.

We should recall, in this connection, a problem Riemann formulated more than a hundred years ago, i.e. are our geometric assumptions valid for infinitesimal distances? He pointed out that the question is intimately related to the problem of intrinsic causation of metric relations in space and that it was basic to the science of space. When it is considered, we must, Riemann believed, take into account that if we have a discrete manifold, then the principle that metric relations are valid is embedded in the notion of manifold, but if we have a continuous manifold, then the principle has to be sought elsewhere. Riemann then pointed out that it follows that either the reality yielding the idea of space forms a discrete manifold, or else we have to try to explain the origin of metric relations by something external: constraints imposed upon this reality. Riemann went on to assert that we stand, at this point in our thoughts, on the threshold of another science, physics, and the state-of-the-art as it then was precluded further progress [85].

Why have we still made no progress over the last hundred years, a period during which science has penetrated the microcosm and the quantum theory was developed and remarkably proved experimentally? Some people believe that quantum mechanics is directly concerned with the creation of the classical notion of metric relations. Perhaps, the construction of a quantum theory of gravity should begin with

an analysis of the notions and postulates which we use in the classical\* model of space-time. This analysis would seek the reasons underlying the classical notions, and find a more fundamental description of the microcosm than we use now (i.e. more fundamental than the theory of elementary particles and quantum theory). Space-time relationships might arise as a result of the description of micro-particle combinations. This approach is sometimes called macroscopic or statistical view of the nature of classical space-time. These ideas may, if omitted from consideration, turn into dangerous underwater rocks capable of sinking the most sophisticated variants of a quantum theory of gravity.

How should we assess from what we have said the present state-of-the-art and the prospects for the quantization of gravity? There is no single generally accepted opinion among theoretical physicists. In the literature and in personal conversations with various researchers we have found a very diverse range of assessments: from those that claim that the problem will soon be resolved after a few minor details have been cleared up to a conviction that we are still at the beginning of a very long and difficult road. These assessments reflect similarly the variety of opinions as to how profound the problem is that is being faced and as to how it should be solved. We, ourselves, believe that there is no need to take matters to extremes, and that we should neither oversimplify the situation nor be entirely pessimistic. Many years of work have undoubtedly taught theorists quite a bit and that based on a sober account of what has been achieved and the finding from the various lines of study on gravitational theory, a decisive breakthrough should be made in the near future. It is possible that everything is now ready for this step, but what is needed is a catalyst.

### **3.5. DIMENSIONALITY OF PHYSICAL SPACE-TIME**

**3.5.1. Formulation of the Problem of Dimensionality of Space-Time.** Have you, dear reader, ever thought about how surprising and mysterious it is that the physical space of the universe is three-dimensional or the space-time manifold four-dimensional. Let us do so now. Consider the uni-

---

\* Classical here means nonquantum.



verse from the following standpoint: suppose it has  $N$  particles so that we can expect  $N(N - 1)/2$  symmetric relationships (distances) to exist between them, and that generally speaking, all these relationships are independent of each other. If we say however that the particles are in an  $n$ -dimensional space, we mean that the position of each particle must be described by  $n$  coordinates, that is, only  $nN$  of the distances among the particles are independent. In a physical space,  $n$  is three but why is our universe described by a finite number of dimensions, and why is it three?

We all realize that such straightforward questions are the most difficult to answer in physics, and they have to be formulated in a more correct form. The ancient Greeks realized that nothing can be made from nothingness, that, for example, to build a geometry there must be a set of axioms from which by certain rules (the laws of logic) all the rest can be derived. To return to dimensionality, we can look at the problem from two points of view, i.e. either the three-dimensionality of space (the four-dimensionality of space-time) is taken as an axiom (this is what is actually done in modern physics), or it must be explained (derived as a theorem) from some simpler or more general physical notion.

The developments in theoretical physics over the last hundred years have shaken the tolerability of three-dimensionality as an axiom. Indeed, we have come to understand that space is not a priori a container of matter and that its basic properties are in fact either affected by or indeed caused by the physical presence of matter. For example, the general theory of relativity involves a rejection of the idea of absolute ordering (simple ordering) of events in time and replaces it with one of partial ordering, the theory being based on the velocity of light as a fundamental constant. General relativity, as we showed in Chapter One, is based on a generalization of the metric properties of space-time, the metric then accounting for the existence of the gravitational interaction. Quantum theory has united in a new way the space-time relationships of matter and its dynamic characteristics. In the last few years, a great deal of effort has been spent to analyze the whole set of ideas and axioms that make up the mathematical model of (classical) space-time. Also attempts have been made to generalize or to revise almost all these notions and axioms and to substantiate the

dimensionality of space-time using various physical considerations [76].

The problem of dimensionality has in fact been debated for a long time. **Immanuel Kant** (1724-1804) noticed (he seems to have been the first to do so) that the inverse square law for gravitational and electrostatic forces was related to the three-dimensionality of our space. Multi-dimensional spaces (with dimensions more than three) appeared in mathematics for the first time, and this development must be considered an important milestone in the study of the structure of physical space (and time).

It is difficult to say who started this generalization but the ideas concerning multi-dimensionality are clearest as expressed by **Hermann Grassman** (1809-1877) and **Arthur Cayley** (1821-1885). Profound ideas about the dimensionality of physical space can be found in the memoirs of **Riemann** and **Mach**. **Riemann** wrote that he had set himself the task of creating a general multiply extended quantity. Having done so he concluded that different measures of distance are possible in such a case and that space is nothing but a particular case of triply extended quantity [85]. Later, **Riemann** went on to discuss  $n$ -extended quantities, without specifying the value of  $n$ . **Mach** drew a number of examples from physics where the idea of multi-dimensionality was fruitful and he began to ponder the nature of physical space, formulating the problem explicitly: viz. Why is space three-dimensional? [66]. However long before these developments, **Joseph Louis Lagrange** (1736-1813) had already analyzed four-dimensional configuration spaces in mechanics.

This line of reasoning has made several important contributions to physics during its history. We shall discuss three contributions. We showed in Chapter One how **Riemann** arrived at his idea of curved spaces by generalizing **Gauss's** theory of two-dimensional curved surfaces (this was a generalization from two dimensions to three). Then special relativity emerged after a generalization of three-dimensional space and unidimensional time to a unified four-dimensional space-time manifold (from three dimensions to four). And, thirdly, as we will discuss below, the physical laws of this universe can be more naturally described by specially curved five-dimensional Riemannian manifold (a generalization from four dimensions to five).

This last idea was first tried out by **Theodor Kaluza** (1885-1954) [53].

Now theoretical physics must solve another task, that is, it must create a physical picture of the universe that is based on some sort of abstract regularity that stems from the physics of the microworld. In this sort of model, the classical concepts of space-time for macroscale phenomena will not be assumed from the very beginning but will instead be derived with all their consequences at some later stage. This must be true for the notion of dimensionality. We would find in such a theory (or in consistent fragments of it) a physical reason for the mathematical postulate of the three-dimensionality of space.

In order to bring about this kind of theory, it is very important that the basic rules that must lie at its foundation be discovered, and in particular that the origin of the observable three-dimensionality of space be identified. Several approaches are possible.

**3.5.2. Physical Features of Universes in Four-Dimensional Space-Time in Comparison with Universes in Manifolds with Other Dimensionalities.** The first approach, which was followed by Eddington, Ehrenfest, and Einstein among others, was to analyze the physical notions and laws in manifolds of  $1 + n$  dimensions ( $n$  being the number of spatial dimensions). The research was based on selecting a physical law considered to be fundamental and investigating whether it depended on the dimension of the manifold. The only laws that were examined are those which take place in a four-dimensional space, or for which the limiting dimensionality from which the law begins to be valid or at which the law ceases to be valid is four. We refer the reader to more detailed reviews and original papers [70, 76, 109] for a formal treat and we shall just discuss a few points.

1. Einstein's theory for a vacuum is only pithy in terms of manifolds with four or more dimensions. For Einstein's vacuum equations  $R_{\mu\nu} = 0$  to be meaningful the  $(1 + n)$ -dimensional manifold must be curved, this corresponding, you may remember, to a nonzero Riemann-Christoffel tensor. In a three-dimensional manifold the tensor  $R^\alpha_{\beta\mu\nu}$  is algebraically expressed in terms of  $R_{\mu\nu}$ , in particular, they both vanish simultaneously.

2. Four-dimensional manifolds are the only ones in which Maxwell's equations for vacuum are conformally invariant.

As we indicated in Sec. 3.3, conformal invariance means that these equations are scale independent.

3. Circular orbits of test bodies in a central gravitational field (for example, the planetary orbits in the Sun's field) are unstable in spaces with  $n \geq 4$  and stable in spaces with  $n \leq 3$  irrespective of whether the universe is flat or curved. This means that planetary systems around stars cannot survive for long in spaces with more than three dimensions.

4. Stable atoms are only possible in space-times with four or less dimensions. In manifolds with more dimensions, the solutions of the Schrödinger-type equations may either have no negative energy levels (i.e. no bound states), or the set of these energy levels includes states with infinitely large negative energies. In the latter case an electron in an atom will, for any energy level, tend to jump to one which is lower, that is, the electrons will for ever be radiating, hence there can be no stable states for matter.

5. Huygens' principle is only valid for spaces with an odd number of dimensions, that is, for  $n = 3, 5, 7, \dots$ . In simple terms, this means that a signal carried by a wave reaches a receiver (observer) at a single point in time and that it cannot be observed once the wave has passed (there must not be a tail or aftereffect). If we impose a nondistortion condition on the signal, then the space can only be three-dimensional.

6. Quantum electrodynamics is only renormalizable in spaces with  $n \leq 3$  dimensions. It is only in such spaces that the infinities that arise in the theory can be eliminated by conventional means. When there are more spatial dimensions the theory becomes considerably less satisfactory.

7. The "stiffness" of the field equations for gravitational, electromagnetic and two-component spinor fields is the same only for  $n = 3$ . We refer the interested reader to Einstein's work [29] for detailed information about this interesting notion or to a review of four-dimensionality characteristics [109].

This list of restrictions on the dimensions of a physical space can be continued.

An interesting picture emerges from an analysis of these characteristics of a four-dimensional universe. We should admit, however, that there is no real reason for believing that any of the restrictions we have mentioned is more fundamental than the simple postulate that  $1 + n = 4$ . At

the moment, however, it is not clear how the classical view of space-time can be developed from any of these restrictions.

In the following sections we will discuss some of the other approaches to the problem of the dimensionality of physical space. Since we will not return to the approach we have just discussed, we should emphasize here its three most characteristic features:

(a) all  $n$  spatial dimensions are equivalent, as has been postulated;

(b) there is only one time-like dimension;

(c) once an electromagnetic, scalar or other field is mentioned, it is assumed that it is introduced in a nongeometric way and is external.

**3.5.3. Five-Dimensionality: Preliminary Results.** In the 1920s, another attack on the problem of dimensionality began. T. Kaluza, H. Mandel, L. de Broglie, and O. Klein initiated work on a unified five-dimensional theory of gravity and electromagnetism. This line of thought faced new problems, which can be illustrated by the questions: How would a five-dimensional physical space-time manifold manifest itself? Why should a fifth dimension be different? Should we limit our consideration to five dimensions?

Here we shall briefly discuss the essence of unified five-dimensional theories. The theories postulate that at the foundations of every mathematical description of the universe there is a curved five-dimensional world with one time and four spatial coordinates. In such a manifold, the metric tensor  ${}^5G_{AB}$  (from now on each capital Latin index  $A, B, C$ , etc., can take the values 0, 1, 2, 3, 5\*) has 15 components. They correspond to the ten components of the four-dimensional metric tensor  $g_{\mu\nu}$ , the four components of the electromagnetic vector potential  $A_\mu$ , and one more as yet unidentified component. The main field equations that we have in the standard theory must now be interpreted as four-dimensional projections of Einsteinian five-dimensional equations.

We now enumerate the indubitable successes of the old five-dimensional unified theories:

(1) the fifteen five-dimensional "Einstein equations" in a vacuum decompose into the standard system of the ten

---

\* Later  $x^4$  will be used to refer to an extra time-like coordinate.

four-dimensional equations of Einstein (of electrovacuum type), and a system of the four standard equations (the second pair of Maxwell's equations without sources on the right-hand sides), and, in general, one extra scalar equation;

(2) if the fifth coordinate is space-like, then we get a standard tensor for the energy-momentum of the electromagnetic field (with an appropriate sign) on the right-hand side of the four-dimensional Einstein equations;

(3) four of the five geodesic equations are standard equations of motion for an electrically charged particle in a gravitational and an electromagnetic field;

(4) given that all geometrical values are independent of the fifth coordinate, transformations of the latter, which are admissible in the theory, correspond to gradient (gauge) transformations in standard electrodynamics.

The older variants of the five-dimensional theory did not gain the trust of physicists for some fairly serious reasons. As Einstein and others pointed out the main drawbacks to five-dimensional theories are:

1. The physical sense of the fifth coordinate was unclear.

2. That all the geometrical values are independent of the fifth coordinate seemed artificial. Einstein pointed out that one consideration that made him doubt the theory was that it did not seem reasonable to substitute a four-dimensional continuum with a five-dimensional one and then to impose in an artificial way a restriction on the fifth dimension just to explain why it does not manifest itself physically [31].

3. The extra fifteenth component of the metric tensor  ${}^5G_{55}$  could not be interpreted physically [34].

4. If an additional restriction is imposed, viz.  ${}^5G_{55} = -1$ , the theory relied on only 14 of the 15 equations, so one equation was redundant.

5. Einstein himself was not really satisfied by the voluntary introduction into the right-hand side of his equations of a nongeometrical quantity (the energy-momentum tensor of matter,  $T_{\mu\nu}$ ). He felt it violated the consistency and integrity of the geometrical approach. Therefore, he tried wherever possible to avoid the  $T_{\mu\nu}$  of external matter. He pointed out that in the context of an exclusively five-dimensional theory the equations do not admit nonzero electric charge or current densities and that the last of Max-

well's equations, which sets to zero the divergence of the electromagnetic field tensor, would seem to exclude the existence of a charge density, and therefore that of electrically charged particles [34].

6. The theory only formally unified gravitation and electromagnetism. Einstein believed that though the objective of Kaluza was to reassess gravitation and electricity by introducing a unified structure to space, he had not been successful [35].

7. The theory led through a unified approach to equations of gravitational and electromagnetic fields, but it added nothing new to our understanding either of particles or of the results obtained in quantum mechanics [34].

8. There are alternatives to the theory (see Sec. 3.3), but it was unclear which is better. Einstein noted that until his time two rather simple and natural attempts had been made to unify gravitation and electricity into a unified field theory: one by Weyl and the other by Kaluza [29].

In the following section we shall show that most of these objections can be removed by modifying the five-dimensional theory.

By the end of the 1920s, almost all these difficulties had been well understood. However, the merits offered by five-dimensional theories continue to attract researchers. In trying to overcome the difficulties they have suggested many interesting ideas. Even though none of them separately removes the objections to five-dimensional theory, they were not useless and some important results arose from this line of research.

(a) Five-dimensionality was the background against which the relativistic Klein-Fock-Gordon wave equation was first derived.

(b) The technique of manifold splitting in a fifth coordinate and in local four-dimensional space-time orthogonal to the fifth axis was developed in the framework of five-dimensional theory. This technique, now called the monad method (see Secs. 2.8 and 2.9), or its special case (the chronometric invariants technique) is used to describe reference frames in the framework of four-dimensional general relativity theory.

(c) The scalar-tensor theory of gravitation (the Jordan and Brans-Dicke theories, see Sec. 3.3), which was widely discussed in the literature during the 1960s and 1970s, emerged

from the five-dimensional theories. We might say that these later theories were due to the 15th component of the five-dimensional metric tensor which had till then seemed redundant.

We believe that in general those people who developed the five-dimensional theories were on the right track and that the grains of truth they obtained may help progress in the same direction, as we shall see below.

**3.5.4. The Pro's and Con's of the Unified Five-Dimensional Theories of Gravitation, Electromagnetism, and Electrically Charged Matter.** The building blocks for a five-dimensional theory are its geometrical quantities, viz. the co- and contravariant components of its metric tensor ( ${}^5G_{AB}$  and  ${}^5G^{AB}$ ) as well as their first and higher order derivatives with respect to the coordinates. As to how to make use of this vast array of material, the study of several variants of the theory showed the task to be more difficult than it might seem. Typical problems that arose involve the selection of which values and in what combinations to be ascribed to which physical quantities. It was equally difficult to understand how to move from a theory dependent on a fifth coordinate to the standard theory, in which there was, as most people believed, no such dependence. It happens that these questions can all be answered unambiguously and the solutions arise using modifications of the techniques of (1 + 4)-splitting and conformal mapping which first appeared in the early unified field theories. Other ideas from these older versions have also been used. Without going into too much detail, we shall summarize these techniques and procedures.

1. *The conformal regauging procedure.* It was found that the identification of geometrical values with physical ones should not, as early workers did, start directly from the  ${}^5G_{AB}$  metric of an initial five-dimensional curved manifold, but it should begin with a conformal regauging. In the general case, this involves a transformation from one metric to another, the latter differing from the former by a scalar factor which depends on coordinates (see Sec. 3.3.2, Weyl's theory). This means that the angles between given directions remain the same in both geometries but the distances between given points change. In this particular procedure, the metric component  ${}^5G_{55} = -\varphi^2 ({}^5G_{AB} = \varphi^2 {}^5G_{AB}^*)$  turns out to be a scalar function. In the new metric  $\hat{G}_{55} = -1$ , as



happens in the first versions of five-dimensional theory, but this time the other components differ from those in the older theories by a scalar multiplicative function  $\varphi$ . Only after this procedure are the geometrical values identified with physical values.

2. The  $(1 + 4)$ -splitting technique was first developed in the 1930s and 1940s because of the need to move from the five-dimensional manifold into four-dimensional space-time. In the meantime, this technique has been substantially improved in the framework of four-dimensional space-time in order to refine reference frame description methods. It is now being reimported in five-dimensional theory.

3. The theory uses the condition of *quasi-cylindricity along the fifth coordinate*. In other words, it is assumed that the fifth coordinate only influences the scalar field  $\varphi$ , whereas the physically meaningful four-dimensional quantities  $g_{\mu\nu}$  and  $A_\mu$ , which are obtained from  $\hat{G}_{AB}^*$ , are independent of it.

4. What should we do with the scalar field  $\varphi$ ? This is a matter of principle. The physical meaning of this field can be interpreted in several ways:

(a) The field can be thought of as a yet undiscovered, supplementary neutral and massless fundamental field. This would make this five-dimensional theory a special case of Brans-Dicke scalar-tensor theories but with an electromagnetic field. This is the interpretation now cited most frequently in the literature. The representation of nongeometrized matter in the equations is now necessary, and there is no need for us to discuss the dependence of  $\varphi$  on the fifth coordinate.

(b) The  $\varphi$  field can also be viewed as a rough representation of charged matter (electrons, nucleons, etc.) in the same way as the scalar field of the Klein-Fock-Gordon equation corresponds to the types of matter more properly described by the Dirac equation. This is an interesting approach which one of us has probed [109], and it means that five-dimensional theory is an approximation to a unified geometrical field theory.

A concrete form of such a theory has been found, but it was necessary to introduce a special cyclic dependence of  $\varphi$  on the fifth coordinate, viz.

$$\varphi^{3/2} = 1 + b\Psi \exp(-i\alpha x^5) - b\Psi^* \exp(i\alpha x^5), \quad (3.31)$$

where  $\Psi$  is a complex scalar function that depends only on the four standard space-time coordinates ( $\Psi^*$  is the complex conjugate function),  $\alpha$  and  $b$  are constants obtained by bringing this theory into line with the four-dimensional theory. They take the values

$$\alpha = \frac{ec}{2\sqrt{G}\hbar}, \quad b^2 = \frac{3\kappa\hbar^2}{32m}, \quad (3.32)$$

where  $e$  is the electrical charge on an electron,  $G$  and  $\kappa$  are the Newtonian and Einsteinian gravitational constants,  $\hbar$  is Planck's constant, and  $m$  is the particle's mass.

It follows from (3.32) that the period  $T$  of the dependence of  $\phi$  on  $x^5$  is extremely short in comparison with the distances for which standard equations hold ( $T = 2\pi/\alpha \simeq 10^{-33}$  cm). It is natural, therefore, to assume that when these equations are averaged over  $x^5$ , the dependence on the fifth coordinate drops out, leaving the fifth component of the momentum to become (when multiplied by a constant factor) the electrical charge of the particle. In a way, this reasoning reflects the hypothesis that Einstein and Bergmann made in the 1930s that the universe is closed with respect to the fifth coordinate. As a result, we have second pair of Maxwell's equations with current formed of the  $\Psi$  field on the right-hand side; Einstein's equations with the energy-momentum tensor of the  $\Psi$  field, which has a geometrical origin, and a Klein-Fock-Gordon-type equation for the  $\Psi$  field with an electrical charge and rest mass. We might say that this approach realizes Einstein's desire to geometrize matter completely. However, significant difficulties have arisen.

One complication is that the rest mass of the scalar field  $\Psi$  is rigidly related to the charge:  $m = e/2\sqrt{G}$ , and if  $e$  is the charge of an electron, the mass would be of the order of  $10^{-6}$  g. This is an extremely large mass in comparison with the masses of the elementary particles, and so this theory clearly cannot yield the masses (or mass spectra) for real particles. To obtain these, it must account for the whole variety of interactions and the symmetries of the microcosm. Remember that in the standard theory, masses are normally introduced either phenomenologically or if their field origin is assumed they become infinitely large. In this theory, the masses are large but finite, and this result alone is an achievement for the theory.

It is therefore necessary to make mass renormalizable, that is, to reduce it to standard values. This can be done using a cosmological constant, but which particle should be chosen to have the standard mass? Since  $\phi$  is geometric in nature, we would expect that any particle it produces would play a fundamental role in the universe, but which particle of the elementary particles known should it be related to? The objective of five-dimensional theory is only limited to a unification of the gravitational and electromagnetic interactions, and hence this question becomes unanswerable.

Another difficulty is that the theory is entirely geometrical and so cannot produce spinor particles such as leptons and nucleons, i.e. the main types of matter. They must be specially introduced into the geometry additionally.

(c) In the light of the above, a third approach to the scalar field may be reasonable. This is to assume that the scalar field describes only some and not all types of matter, and that it may even not be an approximation to their description. Moreover, the validity of Einstein's program to geometrize completely the right-hand sides of his equations must in general be put in question.

It would now be helpful to return to the topic of Sec. 3.3.5 and remind the reader that there are several physical pictures of this universe. We can say that our modern view is based on three basic physical categories: (a) space-time notions, (b) fields as carriers of interactions, and (c) particles (at quantum level they are, probably, fermions). In the framework of the generally accepted approach, all these classes are independent and cannot be reduced to one another. Let us call this the first approach. The literature contains a variety of other approaches which attempt to bring together two of these three classes. Let us call the approach proposed by Clifford and Wheeler the second one, and in it the two categories brought together are the space-time notions and the fields. The third class, particles, are then derived from the first two. Einstein's geometrization program should therefore be classed as a second approach. There is a third approach (the Mach-Feynman one), which is realizable as the action-at-a-distance theory (for either electromagnetic, gravitational, or scalar interaction, see Sec. 3.3.4). In this approach, the basic categories are space-time and particles, with the interaction carrying fields

being derived from these two. A fourth approach would thus involve the remaining two categories as the basic ones.

We may thus reason that Einstein's program can only claim to geometrize all the interaction fields, i.e. the gravitational, electromagnetic, weak (vector and scalar bosons), and strong (gluons) fields, but it can do no more. This would again suggest that  $\varphi$  should be interpreted as one of the interaction carriers. If we take this position we can see that in the framework of five-dimensional theory, which unifies the gravitational and electromagnetic interactions, it is impossible to interpret the  $\varphi$  field unambiguously.

To sum up the work on the five-dimensional theories, we can say that some of the drawbacks of the older variants have been removed, that is points two to five of the list we discussed. The eighth drawback is also removed because Kaluza's approach is combined naturally with Weyl's. The notion of conformal correspondence is an integral part of the theory. The physical meaning of the fifth coordinate is related to electrical charge because charge appears in the theory when the functions are differentiated with respect to the fifth coordinate. The charges of particles in  $e$  units are determined by integer coefficients in front of the argument  $i\alpha x^5$  of exponents in (3.31).

**3.5.5. A Unified Six-Dimensional Theory of Gravitational and Electroweak Interactions.** From the viewpoint of modern physics, a program to unify only two types of interactions, the gravitational and electromagnetic ones, appears outdated. There is now an accepted unified theory for the electromagnetic and weak interactions (the Weinberg-Salam model) and so we have two substantial fragments of a unified theory for all interactions, namely, the Kaluza-Klein unified five-dimensional theory for gravitation and electromagnetism and the Weinberg-Salam theory of electroweak interactions. Electromagnetic interactions lie at the intersection of these two theories. Naturally, we want to know how to unify these two unified theories (three types of interaction) into a single theory.

One obstacle is the qualitatively different basis of the two theories. The Kaluza-Klein theory is five-dimensional and geometrical. The transformations of the fifth coordinate form, in this theory, an Abelian group which is supplementary to the group of the four-dimensional coordinate

transformations in general relativity. The Weinberg-Salam model, on the other hand, is not based on geometrical notions, it uses local symmetries of the Abelian  $U(1)$  group and the non-Abelian  $SU(2)$  group. In order to combine the two theories, a richer structure capable of including both theories is needed. In fact, a six-dimensional geometry with torsion and an external spinor field has this sort of structure.

Before discussing this theory, we shall try to outline the Weinberg-Salam model of electroweak interactions. It is called gauge theory, in that it is based on symmetries with respect to transformations of the  $U(1)$  and  $SU(2)$  groups. The  $U(1)$  group can be interpreted as a group of rotations in two-dimensional space, hence it depends on one parameter, the angle of rotation. The  $SU(2)$  group is related to the rotations in three-dimensional space and so it is determined by three parameters. Generally speaking, these spaces are not related to the four-dimensional space-time, they are instead interpreted as the spaces of the internal symmetries of elementary particles attached to every point of classical space-time. The idea of local symmetry is thus introduced, that is, of the dependence of the group parameters on the coordinates of space-time. In order to maintain the symmetry of the theory, the number of vector fields to be introduced must therefore be equal to the number of parameters possessed by the groups in question. In our case, we have to introduce four vector fields: one, the  $B_\alpha$  field, corresponds to the  $U(1)$  group, and the remaining three fields, the  $A_\alpha(k)$  fields, where  $k = 1, 2, 3$ , correspond to the  $SU(2)$  group. The fields are manifested as combinations but not in their pure form. One linear combination of the  $B_\alpha$  and  $A_\alpha(3)$  fields forms the vector potential of the electromagnetic field  $A_\alpha$  and another forms the massive neutral field of the  $Z$  vector boson. Combinations of the remaining two fields, i.e.  $A_\alpha(1)$  and  $A_\alpha(2)$  make up the fields of the massive charged  $W^\pm$  vector bosons. All these vector bosons have recently been observed experimentally and they are the carriers of the weak interactions.

Another essential element of the Weinberg-Salam model is the Higgs mechanism of spontaneous symmetry breaking. Without going into details, we will just say that this technique allows us to introduce into the theory the rest masses

of the leptons and vector bosons. Initially, these particles were all massless, but the model includes the doublets of the scalar particles: the Higgs bosons. Interactions with  $Z$  and  $W$  vector bosons and electrons obtain their masses by interacting with the Higgs bosons (it is also possible in principle to ascribe rest masses to neutrinos, if experimentally it is proved that they are nonzero). It is interesting that Higgs's mechanism leaves photons massless.

It is important to note that the spinor wave-functions of the leptons (the electron and neutrino) are split in the model into left- and right-handed components which interact differently with the four vector fields. The left-handed components of the leptons are said to form a doublet in internal isotopic space, whereas the right-handed components form singlets.

Returning to the problem of a unified multi-dimensional theory of several interactions, we must notice, in the first place, that the gauge approach used in the Weinberg-Salam model is very close to a geometrical description. Before this model of the electroweak interactions was constructed, a gauge theory of electromagnetic interactions had been developed which was almost equivalent to the Kaluza-Klein five-dimensional theory. The literature is rather rich with discussions of gauge formulation of the general theory of relativity (localizations of the six-parametric Lorentz group or of the ten-parametric Poincaré group had been made). In some ways the gauge and geometrical approaches are two languages for explaining the same phenomena. One of us feels, however, that the geometrical approach is more fundamental because geometrical notions are at the basis of classical physics.

A future unified multi-dimensional theory ought to include all the notions and fields of the Weinberg-Salam model, and they should be in particular the four vector fields. We know from the Kaluza-Klein theory, that incrementing the dimensions of the space-time manifold by one, a new geometrical vector field, the electromagnetic one, comes into being. It is easy to see that each increment in the dimensionality will bring in another vector field. Simple arithmetic indicates that in order to describe four vector fields we must have an eight-dimensional manifold with four additional spatial coordinates. Some work has been done on this subject, but one question that arises is how do we know

that the additional fields should arise out of the metric? We know that geometry has properties other than metric ones. It is also determined by the connection (which is equivalent to the definition of the parallel transport, see Sec. 3.3.2). The connection, in turn, is built of three tensors (Schoutens) as well as a residual nontensor quantity. The best known Schouten is the torsion tensor. It turns out that a unified theory for the gravitational and electroweak interactions can be built using only six dimensions but a geometry with torsion.

This six-dimensional theory has several features.

1. The neutral vector fields in the Weinberg-Salam model correspond to two vector fields of metric origin (a diad) and thus to two extra spatial coordinates. The vector potential of the electromagnetic field  $A_\alpha$  and the field of neutral  $Z$  boson are their linear combinations. This is a way of generalizing the Kaluza-Klein theory.

2. The charged vector bosons (or the  $A_\alpha$  (1) and  $A_\alpha$  (2) vector potentials) are constructed from the components of six-dimensional torsion tensor  ${}^6S_{MNP}$ . It can be shown that by a choice of the spinors' components in six-dimensional space-time, we can get the same interactions with four vector fields as there are in the standard Weinberg-Salam model. Thus, the many years of discussion (since the 1920s) as to the possible physical manifestations of space-times with torsion (see Sec. 3.3.2) have come to a close.

3. The doublet of the Higgs scalar fields  $\varphi^0$  and  $\varphi^+$  in the Weinberg-Salam model occurs from the conformal factor of the theory, in the same way  $\Psi$  arose, as discussed in the previous section. The difference is that in six-dimensional theory, the two extra coordinates  $x^5$  and  $x^6$  must be cyclic. Then, equation (3.31) becomes

$$\begin{aligned} \varphi = & 1 + b_1 \varphi^0 \exp(i\alpha x^5 - i\beta x^6) - b_1^* \varphi^{0*} \exp(-i\alpha x^5 \\ & + i\beta x^6) + b_2 \varphi^+ \exp(i\alpha x^5 + i\beta x^6) \\ & - b_2^* \varphi^{+*} \exp(i\alpha x^5 - i\beta x^6), \end{aligned} \quad (3.33)$$

where  $\alpha$ ,  $\beta$ ,  $b_1$  and  $b_2$  are constants that are found by a comparison with the standard theory,  $\varphi^0$  and  $\varphi^+$  are complex scalar fields which correspond to the neutral and charged components of the Higgs doublet. This interpretation gives

physical sense to the  $\varphi$  field which emerged in five-dimensional theory.

4. In the standard Weinberg-Salam model, the interaction between particles and the vector fields is described by two special numbers: a hypercharge  $Y$  and the projection  $T_3$  of the isotopic spin of the elementary particles. It happens that in six-dimensional theory these quantum numbers are due to the way the wave-functions depend on the extra coordinates. The integer coefficient  $\varepsilon_5$  in the exponent with  $i\beta x^6$  corresponds to the hypercharge  $Y$ , whereas the integer coefficient  $\varepsilon_6$  of  $i\beta x^6$  corresponds to double the projection of the isospin  $T_3$  in the Weinberg-Salam model. Now, we can rewrite the formula for the electrical charge  $Q$  of a particle (in  $e$  units) in terms of the harmonics  $\varepsilon_5$  and  $\varepsilon_6$ , i.e.

$$Q = \frac{Y}{2} + T_3 = \frac{1}{2} (\varepsilon_5 + \varepsilon_6). \quad (3.34)$$

Note that this rule is valid both for leptons and for the scalar bosons in (3.33) and for vector bosons. The torsion tensor is a function of  $x^6$ .

5. Spinors in a six-dimensional manifold are described by eight-component wave-functions. These are split by a special technique into two four-component spinors which correspond to the electron and electron neutrino.

6. In six-dimensional theory, all particles are originally massless. Their rest masses arise due to interactions with the scalar fields (from the conformal factor) via a mechanism like that of Higgs's.

**3.5.6. Constructing a Multi-Dimensional Theory Which Unites All Four Types of Interaction.** The focus in elementary particle physics has now shifted towards the search for a theory that will unite the electroweak interactions with strong interactions, the latter theory being called chromodynamics. In the light of our previous discussions it is logical to ask whether the techniques used in six-dimensional geometry can be applied to this problem? An analysis one of us (Y.S.V.) has completed showed that this can be done, and that the geometry is so rich that this may be accomplished in several ways. In chromodynamics, all interactions are mediated by gluons, which are neutral bosons and of which there are eight altogether. If we follow the Kaluza-Klein idea, that is describe them using a metric, the number of dimensions will have to be increased by eight. But our



experience of six-dimensional theory has shown that a tensor of torsion can be used instead and the number of dimensions needed can be substantially decreased. However, when we come to consider which variant is optimal, we find ourselves at a frontier. We will limit our explanation by saying that at the moment there is one model that realistically describes many of the properties of the theories under synthesis in terms of seven dimensions.

**3.5.7. Multi-Dimensional Theories with Two Time-Like Coordinates.** So far we have discussed multi-dimensional theories with respect to the unification of the various interactions. This sort of problem requires extra space-like coordinates. However, there are a number of facts in theoretical physics which appear to be manifestations of an additional time-like coordinate.

(a) *The group of conformal transformations.* The Lorentz transformations were first found as a group of transformations with respect to which Maxwell's equations are invariant. It was only some time later that it was established that space and time form a single four-dimensional manifold in which the Lorentz transformations represent a rotation group. It seems as if similar situation has evolved again but with a larger number of dimensions. It was noticed some years ago that the basic equations in physics (for massless fields) in a flat four-dimensional space-time are invariant with respect to a broader (15-parametric) group of conformal transformations, of which the six-parametric Lorentz group is a subgroup. But the most interesting fact is that the group of conformal transformations is as a whole nonlinear and can be represented as a group of linear transformations (rotations) of a flat six-dimensional manifold with two time axes and four space axes. The generalization from this sort of flat six-dimensional manifold to the curved one, like the generalization from four-dimensional Minkowski space-time to the space-time of general relativity, has been studied.

How should this six-dimensional manifold be interpreted in the light of the unification of the various interactions? Obviously, one of the extra space-like coordinates must contain, in this particular case, all the other extra space-like dimensions in a condensed form. Something analogous can be seen in five-dimensional theory. The electromagnetic field turned out to consist of two fields the description of which

requires two space-like coordinates. In exactly the same way we should interpret one space-like coordinate in the examples that follow.

(b) *Multi-dimensional optics*. According to six-dimensional unified theory of gravitational and electroweak interactions, all particles originally have no rest mass, that is, they are light-like. The theory can therefore be regarded as a sort of multi-dimensional optics. This takes us back to the ideas of F. Klein in the last century and the work of Yu. B. Rumer in the 1950s on "5-optics" [61, 87]. There is no doubt that they contain a grain of truth. The failure of Rumer's five-dimensional theory was mainly due to the assumption that the five-dimensional interval along the particle's path was assumed to vanish, i.e.  $dI^2 = ds^2 - (dx_5)^2 = 0$ . This meant that the fifth component of five-velocity  $dx^5/ds$  was unity and hence the electrical charge of the particles could not be brought into the equations of motion. Rumer attempted to find a way out of this situation by identifying extra components of the five-metric with physical quantities, i.e.  ${}^5G_{5\mu} \sim A_\mu e/mc^2$ , where  $e$  is the particle's electrical charge. However, space then became dependent upon the properties of a concrete particle, that is, it became configurational (each particle had its own space). Moreover, a universal space-time had also to be postulated beside the configurational one. How to combine the two spaces was a question that defeated him (and it may even be unanswerable in the framework of a five-dimensional theory).

If, however, Rumer's work is generalized to a theory with an extra time-like dimension, these difficulties do not occur. For example, in a six-dimensional theory with a signature  $(+ - - - + -)$ , the optics condition means that

$$\begin{aligned} d\Sigma^2 &= ds^2 - (dx^5)^2 + (dx^4)^2 \\ &= 0 \rightarrow 1 - \left(\frac{dx^5}{ds}\right)^2 + \left(\frac{dx^4}{ds}\right)^2 = 0, \end{aligned} \quad (3.35)$$

where  $x^4$  is the extra time-like coordinate. The component  $dx^5/ds$  can be identified with  $e/2m \sqrt{G}$ , as is needed in order to describe charged particles in an electromagnetic field. In this case, the relation (3.35) still holds due to an extra degree of freedom in  $dx^4/ds$ .

(c) *Renormalizing mass using an extra time-like dimension.* When we discussed five-dimensional theories, we spoke about difficulties associated with the values of the rest masses of the particles related to their electrical charge. If, however, we assume that there is another time-like coordinate (for which an exponential expression like (3.31) or (3.33) is postulated), the contributions of all the extra coordinates to the rest mass become opposite in sign. For example, for a six-dimensional theory with a  $(+ - - - + -)$  signature, the rest mass is  $\sqrt{e^2/4G - (\hbar\gamma/c^2)}$ , where  $\gamma$  is the constant in the term  $\exp(i\gamma x^4)$ . We can obtain any observable value for the rest mass by setting the value of  $\gamma$ . Notice that the same assertion can in fact be directly derived from (3.35).

(d)  *$O(4)$  and  $O(4,2)$ -symmetries in the problem of a hydrogen-like atom.* In the 1930s, Fock showed that in the quantum-mechanical description of the hydrogen atom by Schrödinger's equation, there was a symmetry with respect to rotations in a four-dimensional space ( $O(4)$ -symmetry) rather than one with respect to a three-dimensional space, as had been expected. This result was surprising, but now it fits naturally into five-dimensional theory with four space-like coordinates. In the last few years, however, it has been found that there is an even higher dynamic  $O(4,2)$ -symmetry with respect to rotations in the six-dimensional momentum manifold.

We could have presented several more sophisticated indications but they made much stronger mathematical and theoretical demands on our readers than we wished. Obviously, increasing the number of space-like dimensions is more easily appreciable than increasing the number of time-like dimensions. In the latter case, more conceptual questions arise but their resolution will uncover new and interesting aspects of our universe.

We must admit that for the time being, all these constructions may only be formally significant and reduce to modified ways of explaining known facts. But think how important it is to recognize that there may be a parallel existence of universe descriptions both in four-dimensional and multi-dimensional space-time manifolds! Which way will the balance tilt in favour of one description? This question is especially important because the analysis of the theoretical and experimental power of the multi-dimensional approach has not yet been exhausted.

**3.5.8. Supergravity and the Extra Dimensions.** We mentioned earlier that the supergravity theory accounts for the supersymmetry that exists between bosons and fermions. This theory is based on multi-dimensional superspace, which in contrast with the usual space-time and even its multi-dimensional generalizations, consists of inhomogeneous dimensions, i.e. a combination of (a) usual variables, the coordinates of four-dimensional space-time and other commuting numbers, and (b) anti-commuting quantities, Grassmanian variables. The latter may be thought of as unusual extra dimensions.

The construction of a supersymmetric theory is difficult because it requires to unify a great a number of interacting fields and to incorporate in it supersymmetry properly (i.e., bosons and fermions must be treated on equal footing). The task can be made easier by using additional dimensions and geometrical techniques such as dimensional reduction. This is exactly the way in which the promising maximally expanded supergravity theory is being developed. It has been shown that the maximum number of dimensions for which we can still have a reasonable (from the viewpoint of phenomenology) theory after the dimensional reduction (transition to four dimensions) is eleven.

On the other hand, it has also been shown that eleven is a good number for the Kaluza-Klein  $(4 + r)$ -dimensional approach too. In fact, the minimum gauge group needed for a unified theory is  $SU(3) \times SU(2) \times U(1)$  and to describe this group the minimum number of extra dimensions must be seven. Hence, if we add in the four dimensions of space-time, we have  $7 + 4 = 11$ . How important this coincidence is, we shall see in the future. So far it is clear that multi-dimensional manifolds have conquered the physics of the microworld. A further study of the problems we have discussed would bring us closer to a solution of the more global problem, namely why our observable macroscopic space-time is four-dimensional.

## CONCLUSION

In Chapter Three, we tried to project the future of space-time physics through the prism of the lines of investigations which have been mostly developed recently. Of course, it is impossible to say now where and when a decisive breakthrough to a new physics will occur. The laws of the physics operate all around us in the real world at their individual levels, and we try to decipher the laws' deeper structures from their heterogeneous external manifestations. This requires researchers both individually and in teams to employ rigorous logic and critical assessments of the available solutions. It is important that we do not invent generalizations for the sake of generalization (there have been many such, and the "activity" becomes fashionable from time to time). We believe a problem oriented search is possible in this context, and that it supports the creativity of every scientist and our experience solidifies this support. At the beginning of this book, we tried to present the various lines of argument in this science in the past. Then we passed on to later developments to show how scientific ideas evolve dynamically, fed by both separate discoveries and general reviews of accumulated knowledge.

We have found that in the past, new ideas took hundreds of years to formulate and now this process has quickened immeasurably, so that now the situation changes before our very eyes. At the same time, we have seen many examples of how ideas or even fundamental results arose prematurely and as a consequence they could not be understood or assimilated. The speculations of the ancient philosophers and the engineering projects of **Leonardo da Vinci** (1452-1519) can be put besides the Lorentz transformations that **V. Voigt** (1850-1919) published almost 20 years before spe-

cial relativity, or the technical prerequisites for the discovery of the laser which could have occurred in 1916 when Einstein introduced the ideas of induced and spontaneous radiation. We saw that sometimes a premature result remained unknown to physicists (for example, Dirac rediscovered spinors, which had previously been introduced by **Élie Cartan** (1869-1951) long before). As soon as the conditions for solving a problem are ready, a breakthrough occurs and more often than not it is realized by several people at once (recall such a situation for quantum mechanics). These are lessons humanity has learned from its experience of constructing scientific theories.

It is possible that even now, the conditions are all ready in space-time physics for a new step forward from the jungle of established fact, ideas, and logical constructions. Only we have to overcome a sort of psychological inertia. Every specialist acquires a certain number of clichés and professional prejudices in addition to his positive skills and knowledge. A new attitude towards reality takes time to grow, the largest proportion of which is spent on overcoming the older ideas. Perhaps, this is why great discoveries are made by comparatively young men, by those who while they have studied their field in depth have still retained for some reason the ability to be amazed by the harmony of the universe and have not been fenced in by cliché and dogma.

We concluded our survey by describing the progress recently made in this science and outlined some of the important frontiers. Of course, we have not exhausted this topic and it is quite possible that the long-awaited breakthrough will arise out of a line of thought we have not covered. We simply could not cover them all in a small chapter in a small book because there are too many, and even our own investigations have not been described here fully. For example, one of us places great hopes on a new concept of space and time that he feels is an alternative to field theory and the theory of direct interparticle interaction, being close to a statistical (macroscopic) interpretation of space-time [108, 110]. Another of us divides his sympathies between a topologically nontrivial (multiply connected) two-dimensional structure which, when averaged, has the properties of a four-dimensional curved space-time (this space-time "foam" is different to that proposed by Hawking for four-dimensional one), and a technique involving complex and

conformal manifolds that are related to Penrose's twistors and their extension to supergravity. The third author is developing a space-time philosophy as related to relativistic astrophysics and cosmology. He aims to account for the characteristic symmetries in the configurations of cosmic objects. Thus, the differences among us, as scientists, do not just concern differences in our favoured notations (we used a sort of average in our book), but they involve completely different lines of investigation and estimates of the prospects for the development of our theory.

The evolution of science however is a complex process which leads both to new concepts and also to reassessments of older ones (see, for example, the discussion by D. Bohm [7]). Therefore, a success in one direction does not cancel efforts made in other directions because they each have their own domain and separate existence, intersecting non-trivially to contain known facts (otherwise they would not have been discussed at all). For example, five- and six-dimensional symmetries undoubtedly reflect a substantive aspect of reality; they are not just coincidences. The black hole phenomenon, although an extreme aspect of modern physics may only exist in theory (which is what some scientists would like), but even so it is a valuable concept, because it is a concentration of the structural characteristics of a theory which objectively reflects many of the properties of this universe. An unremitting struggle against all singularities in general (something many researchers in fact do) is like jousting at windmills. Every theory is only a step forward and is always limited. As we approach the limits, singularities (of some kind) appear. Thus a struggle (interpreted as a search for new paths) against concrete singularities is useful and helps science, but it is futile when directed against singularities in general. Some ideas leap forward at certain stages in the evolution of science and eclipse others which temporarily retreat (but their constructive nature is reflected by the debate they generate). A genuine and final replacement of ideas consists of their thorough reassessment and merger so that a more universal understanding of natural phenomena emerges. Our discussion in the last chapter may therefore have seemed motley, but an apparent incongruity was unavoidable.

When we looked back at what we had written for this book and thought about the variety and riches of modern

investigations, we felt like ending it with a bit of poetry. This is a "gravitational" adaptation of a short poem by Friedrich Schiller. It was first written in German by one of us (N.V.M.) for when in 1980 he was awarded the Schiller Bronze Medal of the Jena University, it was published together with an English version in the [89], and the following is a further revision by the author:

#### A THEORETICAL PHYSICIST'S SCHILLER-DREAM

From the depths of gravitation  
Which makes sober thought so dense,  
I would rush to sweet temptation,  
A search for final evidence!  
There I see its sheen all-winning,  
There I hear its singing strings!  
I should be at world's beginning,  
If I had a pair of wings.

Congruences bring News Functions  
With derivatives of Lie,  
And I am so very anxious  
To trace the tunes of theory.  
'Tween the thicket of equations  
Exact solutions bloom and grow  
Which in deep intoxication  
Physicists attempt to know.  
O it ever must be summer  
There in Newman's Heavenly Space;  
All is full of light cones' glamour,  
Their null lines my dreams embrace.  
But I fear the Killing-vector  
Which grows light-like sans remorse.  
Black holes spring up, thirst and hector,  
With a crushing tidal force.

New ideas I feel in motion—  
Still, the apparatus fails...  
Forward, friends! Forget your caution!  
Our souls will fill our sails!  
With our faith and firm persistence  
(Else our pledge the gods refuse)  
Can we cross the rocky distance  
To the Holy Land of Truth.



1. Akhundov, M. D., *Space and Time Concepts: Sources, Evolution, and Prospects*, Nauka, Moscow, 1982 (in Russian).
2. Arifov, L. Y., *General Relativity and Gravitation*, Fan, Tashkent, 1983 (in Russian).
3. Aristotle, *Physics*, ed. by W. D. Ross, Clarendon Press, Oxford, 1950.
4. Bazhenov, L. B., *The Structure and Functions of Natural Sciences*, Nauka, Moscow, 1978 (in Russian).
5. Bergmann, P. G., *Introduction to the Theory of Relativity*, Prentice-Hall, Englewood Cliffs, 1958.
6. Bergmann, P. G., *The Riddle of Gravitation*, Charles Scriber's Sons, New York, 1968.
7. Bohm, D., *The Special Theory of Relativity*, W. A. Benjamin, New York-Amsterdam, 1965.
8. Bondi, H., *Relativity and Common Sense*, Anchor Books, Doubleday & Co., Garden City, New York, 1964.
9. Born, M., *Einstein's Theory of Relativity*, Dover Publications, New York, 1962.
10. Bowler, M. G., *Gravitation and Relativity*, Pergamon Press, Oxford, 1976.
11. Braginsky, V. B. and Manukin, A. B., *Measurement of Weak Forces in Physical Experiment*, Nauka, Moscow, 1974 (in Russian).
12. Braginsky, V. B. and Thorne, K. S., "Present Status of Gravitational-Wave Experiments", in: *Proceedings of the 9th International Conference on General Relativity and Gravitation* (Jena, July 14-19, 1980), Cambridge University Press, Cambridge, 1983.
13. Carelli, A. (ed.), *Astrofisica e Cosmologia, Gravitazione, Quanti e Relatività. Negli sviluppi del pensiero scientifico di Albert Einstein*, Giunti Barbèra, Firenze, 1979.
14. Cartan, É., *Compt. Rend. (Paris)*, 174, 593 (1922).
15. Chiu, Hong-Yee and Hoffmann, W. F. (eds.), *Gravitation and Relativity*, W. A. Benjamin, New York-Amsterdam, 1964.
16. Clifford, W. K., *Lectures and Essays*, ed. by L. Stephen and F. Pollock, Macmillan, London, 1879.
17. Clifford, W. K., *Mathematical Papers*, ed. by R. Tucker, Macmillan, London, 1882.
18. Copernicus, N., in: Dobson, J. F. and Brodetsky, S., *Nicolaus*

- Copernicus's De Revolutionibus*, Royal Astronomical Society, London, 1947.
19. Delokarov, K. Kh., *Philosophical Problems of the Theory of Relativity*, Nauka, Moscow, 1973 (in Russian).
  20. Deloné, B. N., *A Short Exposition of the Proof of the Lobachevski Planimetry Consistency*, Nauka, Moscow, 1953 (in Russian).
  21. Dicke, R. H., *Gravitation and the Universe*, American Philosophical Society, Philadelphia, 1970.
  22. Dirac, P. A. M., *General Theory of Relativity*, Wiley, New York, 1975.
  23. Durell, C. V., *Readable Relativity*, London, 1962.
  24. Dyson, F. K., Eddington, A. S., and Davidson, C., *Phil. Transactions of the Roy. Soc.*, A 220, 291 (1920).
  25. Eddington, A. S., *The Mathematical Theory of Relativity*, Cambridge University Press, Cambridge, 1924.
  26. Eddington, A. S., *Fundamental Theory*, Cambridge University Press, Cambridge, 1946.
  27. Einstein, A., "Über das Relativitätsprinzip und die aus demselben gezogenen Folgerungen", *Jahrbuch d. Radioaktivität u. Elektronik*, 4, 411 (1907).
  28. Einstein, A., "Ernst Mach", *Phys. Zs.*, 17, 7, 101 (1916).
  29. Einstein, A., "Die Grundlage der allgemeinen Relativitätstheorie", *Ann. Phys. (Leipzig)*, 49, 769 (1916); *The Meaning of Relativity*, 4th ed., Princeton University Press, Princeton, 1953.
  30. Einstein, A., "Nichteuklidische Geometrie in der Physik", *Neue Rundschau*, Januar 1925, p. 16.
  31. Einstein, A., "Gravitational and Electrical Fields", *Science*, 74, 438 (1930).
  32. Einstein, A., "Lens-Like Action of a Star by the Deviation of Light in the Gravitational Field", *Science*, 84, 506 (1936).
  33. Einstein, A. and Grossmann, M., "Entwurf einer verallgemeinerten Relativitätstheorie und der Theorie der Gravitation", *Zs. Math. und Phys.*, 62, 225 (1913).
  34. Einstein, A. and Mayer, W., *Einheitliche Theorie von Gravitation und Elektrizität*, Part I, Sitzungsber. Preuß. Akad. Wiss., phys.-math. Kl., 1931, p. 541; Part II, 1932, p. 130.
  35. Einstein, A. and Bergmann, P. G., "Generalization of Kaluza's Theory of Electricity", *Ann. Math.*, 39, 683 (1938).
  36. Einstein, A. and Infeld, L., *The Evolution of Physics*, New York, Simon and Schuster, 1961.
  37. Fock, V. A., *The Theory of Space, Time, and Gravitation*, Pergamon Press, Oxford, 1959.
  38. Fock, V. A., *Einstein's Theory and Relativity in Physics*, Znanie, Moscow, 1967 (in Russian).
  39. Fomalont, E. B. and Sramek, R. A., *Phys. Rev. Lett.*, 36, 1475 (1976).
  40. Friedmann, A. A., *The World as Space and Time*, 1st ed., Petrograd, 1923; 2nd ed., Nauka, Moscow, 1965 (in Russian).
  41. Gliozzi, M., *Storia della fisica*, Torino, 1965.
  42. Gray J., *Ideas of Space: Euclidean, Non-Euclidean, and Relativistic*, Clarendon Press, Oxford, 1979.
  43. Gribanov, D. P., *Philosophical Foundations of the Theory of Relativity*, Nauka, Moscow, 1982 (in Russian).

44. Gurevich, L. E. and Chernin, A. D., *The General Theory of Relativity in Physical Picture of the Universe*, Znanie, Moscow, 1970 (in Russian).
45. Hawking, S. W. and Israel, W. (eds.), *General Relativity. An Einstein Centenary Survey*, Cambridge University Press, Cambridge, 1979.
46. Herneck, F., *Albert Einstein. Ein Leben für Wahrheit. Menschlichkeit und Frieden*, Der Morgen, Berlin, 1963.
47. Hilbert, D., "Die Grundlagen der Physik", *Nachr. Königl. Gesellsch. Wiss. Göttingen, math.-phys. Kl.* 1915, Heft 3, p. 395.
48. Hoffmann, B., *Albert Einstein: Creator and Rebel*, The Viking Press, New York, 1972.
49. Horský, Jan, *Úvod do teorie relativity*, Praha, SNTL, 1975.
50. Ivanenko, D. D. and Sokolov, A. A., *Classical Field Theory*, GITTL, Moscow-Leningrad, 1952 (in Russian).
51. Ivanitskaya, O. S., *Lorentz's Basis and Gravitational Effects in Einstein's Theory of Gravity*, Nauka i Tekhnika, Minsk, 1979 (in Russian).
52. Kagan, V. F., *Essays on Geometry*, Moscow University Press, Moscow, 1963 (in Russian).
53. Kaluza, Th., *Zum Unitätsproblem der Physik*, Sitzungsber. Preuß. Akad. Wiss., 1921, p. 966.
54. Kaufmann, W. J., *The Cosmic Frontiers of General Relativity*, Little, Brown & Co., Boston-Toronto, 1977.
55. Kaziutinsky, V. V. et al. (eds.), *Philosophical Problems of the 20th Century Astronomy*, Nauka, Moscow, 1976 (in Russian).
56. Kaziutinsky, V. V. et al. (eds.), *Astronomy, Methodology, Weltanschauung*, Nauka, Moscow, 1979 (in Russian).
57. Kerr, R. P., *Phys. Rev. Lett.*, 11, 237 (1963).
58. Kompaneyetz, A. S., *Gravitation, Quanta, and Shock Waves*, Znanie, Moscow, 1968 (in Russian).
59. Kuznetsov, B. G., *Einstein*, USSR Academy of Sciences Press, Moscow, 1962 (in Russian).
60. Kuznetsov, B. G., *Evolution of Physical Ideas from Galilei to Einstein*, USSR Academy of Sciences Press, Moscow, 1963 (in Russian).
61. Kuznetsov, I. V. (ed.), *Space, Time, Motion*, Nauka, Moscow, 1971 (in Russian).
62. Landau, L. D. and Lifshitz, E. M., *The Classical Theory of Fields*, 3d ed., Pergamon Press, Oxford, 1971.
63. Landau, L. D. and Rumer, Yu. B., *What Is the Theory of Relativity?* Mir Publishers, Moscow, 1981.
64. Livanova, A., *Three Fates. Exploration of the Universe*, Znanie, Moscow, 1969 (in Russian).
65. Mach, E., *The Science of Mechanics*, Open Court, La Salle, Ill., 1960.
66. Mach, E., *Erkenntnis und Irrtum*, 1968.
67. Marder, L., *Time and the Space-Traveller*, George Allen and Unwin, London, 1971.
68. Markov, M. A., *On the Nature of Matter*, Nauka, Moscow, 1977, p. 199 (in Russian).
69. Misner, C. W., Thorne, K. S., and Wheeler, J. A., *Gravitation*, W. H. Freeman and Co., San Francisco, 1973.

70. Mitskiévich, N. V., *Physical Fields in General Relativity*, Nauka, Moscow, 1969 (in Russian).
71. Mitskiévich, N. V., "On the Difference Between the Notions of Reference Frame and System of Coordinates", in: *Physical Science and Philosophy*, Nauka, Moscow, 1973 (in Russian).
72. Mitskiévich, N. V., "Cosmology, Relativistic Astrophysics, and the Physics of Elementary Particles", in: *Philosophical Problems of the 20th Century Astronomy*, Nauka, Moscow, 1976 (in Russian).
73. Mitskiévich, N. V., "Paradoxes of Space-Time in Modern Cosmology", in: *Astronomy, Methodology, Weltanschauung*, Nauka, Moscow, 1968 (in Russian).
74. Mitskiévich, N. V., Yefremov, A. P., and Nesterov, A. I., *Dynamics of Fields in General Relativity*, Energoatomizdat, Moscow, 1985 (in Russian).
75. Moriyasu, K., *An Elementary Primer for Gauge Theory*, World Scientific, Singapore, 1983.
76. Mostepanenko, A. M. and Mostepanenko, M. V., *The Four-Dimensionality of Space-Time*, Nauka, Moscow-Leningrad, 1966 (in Russian).
77. Novikov, I. D., *The Evolution of the Universe*, Nauka, Moscow, 1979 (in Russian).
78. Novikov, I. D., *Black Holes and the Universe*, Molodaya Gvardia, Moscow, 1985 (in Russian).
79. Omelianovsky, M. E. et al. (eds.), *Physical Science and Philosophy*, Nauka, Moscow, 1973 (in Russian).
80. Pauli, W., *Theory of Relativity*, Pergamon Press, Oxford, 1958.
81. Petrov, A. Z., *Einstein Spaces*, Pergamon Press, Oxford, 1954.
82. Petrov, A. Z., *Modern Methods in General Relativity*, Nauka, Moscow, 1966 (in Russian).
83. Poincaré, H., "Sur la dynamique de l'électron", *Rendiconti del Circolo Matematico di Palermo*, 21, 129 (1906).
84. Rashevsky, P. K., *Riemannian Geometry and Tensor Calculus*, Nauka, Moscow, 1967 (in Russian).
85. Riemann, B., *Gesammelte mathematische Werke*, Dover, New York, 1953 (see also English translation (by W. K. Clifford) of the *Habilitationsvorlesung* in: *Nature*, 8, 14 (1873)).
86. Rosenfeld, L., "Newton and the Law of Gravitation", in: *Archives for the History of Exact Sciences*, Vol. 2, 1962-1965, p. 365.
87. Rumer, Yu. B., *Studies of 5-Optics*, Gostekhizdat, Moscow, 1956 (in Russian).
88. Schmutzer, E., *Relativitätstheorie—aktuell. Ein Beitrag zur Einheit der Physik*, Teubner, Leipzig, 1979.
89. Schmutzer, E. (ed.), *Proceedings of the 9th International Conference on General Relativity and Gravitation (Jena, July 14-19, 1980)*, Cambridge University Press, Cambridge, 1983.
90. Schwartz, J. T., *Relativity in Illustration*, New York University Press, New York, 1965.
91. Schwarzschild, K., *Sitzungsber. Preuß. Akad. Wiss.*, 1916, p. 189.
92. Sciama, D. W., *The Physical Foundations of General Relativity*, Doubleday & Co., Carden City, New York, 1969.

93. Sciama, D. W., *Modern Cosmology*, Cambridge University Press, Cambridge, 1971.
94. Seelig, C., *Albert Einstein. Eine dokumentarische Biographie*, Zürich-Stuttgart-Wien, 1954.
95. Shapiro, I. I., et al., *Journal of Geophysical Research*, June 1977 (preprint).
96. Silk, J., *The Big Bang. The Creation and Evolution of the Universe*, W. H. Freeman & Co., San Francisco, 1980.
97. Sokolov, A. A. and Ivanenko, D. D., *Quantum Field Theory*, GITTL, Moscow-Leningrad, 1953 (in Russian).
98. Sommerfeld A., *Elektrodynamik*, Akad. Verlagsges., Geest & Portig K.-G., Leipzig, 1949.
99. Staniukowicz, K. P. and Melnikov, V. N., *Hydrodynamics, Fields, and Constants in the Theory of Gravitation*, Energoatomizdat, Moscow, 1983 (in Russian).
100. Synge, J. L., *Relativity: The Special Theory*, North-Holland, Amsterdam, 1965.
101. Synge, J. L., *Relativity: The General Theory*, North-Holland, Amsterdam, 1960.
102. Terletskii, Ya. P., *Paradoxes in the Theory of Relativity*, Plenum Press, New York, 1968.
103. Treder, H.-J., *Gravitationstheorie und Äquivalenzprinzip*, Akademie-Verlag, Berlin, 1971.
104. Treder, H.-J. (ed.), *Einstein-Centenary*, Akademie-Verlag, Berlin, 1979.
105. Vassiliev, M., Klimontovich, N., and Staniukowicz, K. P., *The Force Which Rules the Worlds*, Atomizdat, Moscow, 1978 (in Russian).
106. Vizgin, V. P., *The Relativistic Theory of Gravitation*, Nauka, Moscow, 1981 (in Russian).
107. Vladimirov, Yu. S., "Quantum Theory of Gravitation", in: *Einstein's Miscellany*, 1972, Nauka, Moscow, 1974, p. 280 (in Russian).
108. Vladimirov, Yu. S., "On the Development of the Notions of Space and Time", *History and Methodology of Natural Sciences*, Issue 26 (Physics), Moscow University Press, Moscow, p. 76, 1981 (in Russian).
109. Vladimirov, Yu. S., *Reference Frames in the Theory of Gravitation*, Energoizdat, Moscow, 1982 (in Russian).
110. Vladimirov, Yu. S. and Turygin, A. Yu., *The Action-at-a-Distance Theory*, Energoatomizdat, Moscow, 1985 (in Russian).
111. Weinberg, S., *The First Three Minutes*, Vincent Torre, 1977.
112. Wheeler, J. A., *Einstein's Vision*, Springer-Verlag, New York, 1968.
113. Wigner, E. P., *Commun. Pure and Appl. Math.*, 13, 1 (1960).
114. Zakharov, V. D., *Gravitational Waves in Einstein's Theory*, Halsted Press, New York, 1973; *Gravitational Waves in Einstein's Theory of Gravitation*, Statistical Publishing Soc., Calcutta, 1977.
115. Zeldovich, Ya. B. and Novikov, I. D., *Relativistic Astrophysics*, University of Chicago Press, Chicago, Ill. Vol. I—*Stars and Relativity*, 1971, Vol. II—*The Universe and Relativity*, 1974.

1. The first part of the document is a list of names and titles, including "The Hon. Mr. Justice" and "The Hon. Mr. Justice".



# Space Time Gravitation

---

An historical survey of our ideas about space, time and gravitation.

In three chapters, the authors consider how space and time were perceived from ancient times to the present (as marked by the publication of Einstein's theory of relativity)

how they are now perceived,  
and finally how this field of mathematics  
and physics might develop in the future.

Although there are many presentations of our current understanding of this subject on the market, the topics covered in the second chapter of the book are not those usually examined.

However, the main difference between this and other popularizations is the third chapter,

in which the authors try get to grips with subjects such as cosmological singularities, generalizations of Einstein's gravitational theory, the quantization of gravitational fields, and the dimensionality of space-time.

Written for students and professionals specializing in physics or related sciences.

