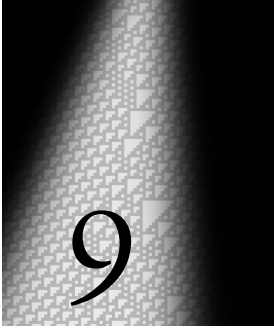


EXCERPTED FROM

STEPHEN
WOLFRAM
A NEW
KIND OF
SCIENCE

CHAPTER 9

Fundamental Physics



Fundamental Physics

The Problems of Physics

In the previous chapter, we saw that many important aspects of a wide variety of everyday systems can be understood by thinking in terms of simple programs. But what about fundamental physics? Can ideas derived from studying simple programs also be applied there?

Fundamental physics is the area in which traditional mathematical approaches to science have had their greatest success. But despite this success, there are still many central issues that remain quite unresolved. And in this chapter my purpose is to consider some of these issues in the light of what we have learned from studying simple programs.

It might at first not seem sensible to try to use simple programs as a basis for understanding fundamental physics. For some of the best established features of physical systems—such as conservation of energy or equivalence of directions in space—seem to have no obvious analogs in most of the programs we have discussed so far in this book.

As we will see, it is in fact possible for simple programs to show these kinds of features. But it turns out that some of the most important unresolved issues in physics concern phenomena that are in a sense more general—and do not depend much on such features.

And indeed what we will see in this chapter is that remarkably simple programs are often able to capture the essence of what is going on—even though traditional efforts have been quite unsuccessful.

Thus, for example, in the early part of this chapter I will discuss the so-called Second Law of Thermodynamics or Principle of Entropy Increase: the observation that many physical systems tend to become irreversibly more random as time progresses. And I will show that the essence of such behavior can readily be seen in simple programs.

More than a century has gone by since the Second Law was first formulated. Yet despite many detailed results in traditional physics, its origins have remained quite mysterious. But what we will see in this chapter is that by studying the Second Law in the context of simple programs, we will finally be able to get a clear understanding of why it so often holds—as well as of when it may not.

My approach in investigating issues like the Second Law is in effect to use simple programs as metaphors for physical systems. But can such programs in fact be more than that? And for example is it conceivable that at some level physical systems actually operate directly according to the rules of a simple program?

Looking at the laws of physics as we know them today, this might seem absurd. For at first the laws might seem much too complicated to correspond to any simple program. But one of the crucial discoveries of this book is that even programs with very simple underlying rules can yield great complexity.

And so it could be with fundamental physics. Underneath the laws of physics as we know them today it could be that there lies a very simple program from which all the known laws—and ultimately all the complexity we see in the universe—emerges.

To suppose that our universe is in essence just a simple program is certainly a bold hypothesis. But in the second part of this chapter I will describe some significant progress that I have made in investigating this hypothesis, and in working out the details of what kinds of simple programs might be involved.

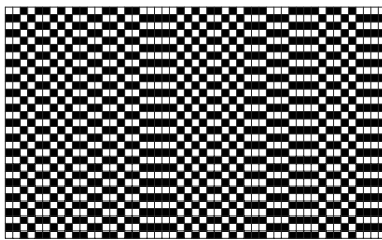
There is still some distance to go. But from what I have found so far I am extremely optimistic that by using the ideas of this book the most fundamental problem of physics—and one of the ultimate problems of all of science—may finally be within sight of being solved.

The Notion of Reversibility

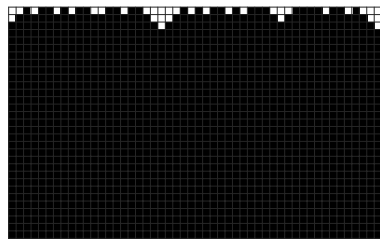
At any particular step in the evolution of a system like a cellular automaton the underlying rule for the system tells one how to proceed to the next step. But what if one wants to go backwards? Can one deduce from the arrangement of black and white cells at a particular step what the arrangement of cells must have been on previous steps?

All current evidence suggests that the underlying laws of physics have this kind of reversibility. So this means that given a sufficiently precise knowledge of the state of a physical system at the present time, it is therefore possible to deduce not only what the system will do in the future, but also what it did in the past.

In the first cellular automaton shown below it is also straightforward to do this. For any cell that has one color at a particular step must always have had the opposite color on the step before.



rule 51



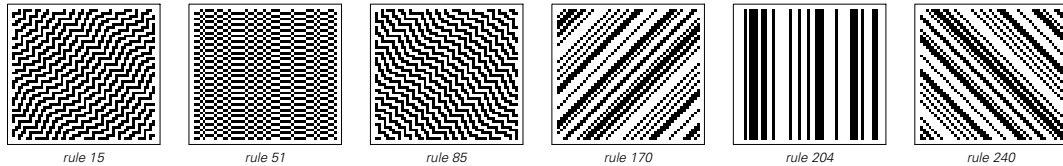
rule 254

Examples of cellular automata that are and are not reversible. Rule 51 is reversible, so that it preserves enough information to allow one to go backwards from any particular step as well as forwards. Rule 254 is not reversible, since it always evolves to uniform black and preserves no information about the arrangement of cells on earlier steps.

But the second cellular automaton works differently, and does not allow one to go backwards. For after just a few steps, it makes every cell black, regardless of what it was before—with the result that there is no way to tell what color might have occurred on previous steps.

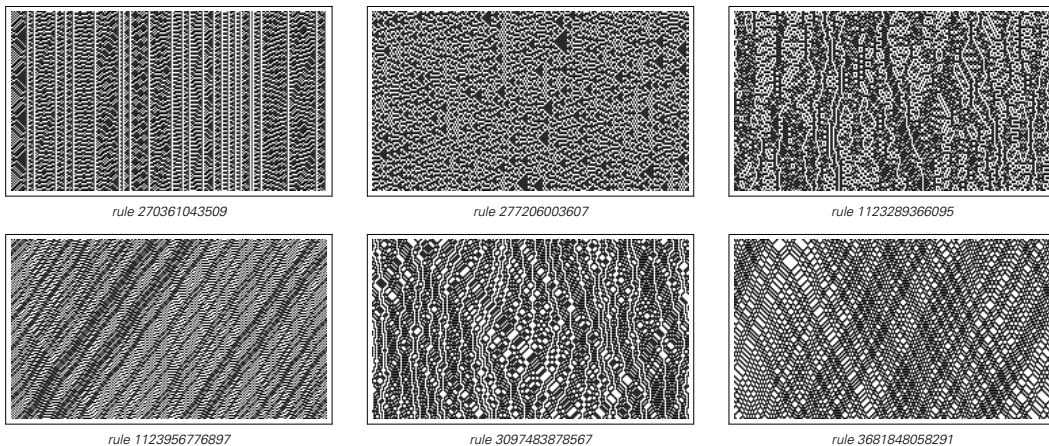
There are many examples of systems in nature which seem to organize themselves a little like the second case above. And indeed the conflict between this and the known reversibility of underlying laws of physics is related to the subject of the next section in this chapter.

But my purpose here is to explore what kinds of systems can be reversible. And of the 256 elementary cellular automata with two colors and nearest-neighbor rules, only the six shown below turn out to be reversible. And as the pictures demonstrate, all of these exhibit fairly trivial behavior, in which only rather simple transformations are ever made to the initial configuration of cells.



Examples of the behavior of the six elementary cellular automata that are reversible. In all cases the transformations made to the initial conditions are simple enough that it is straightforward to go backwards as well as forwards in the evolution.

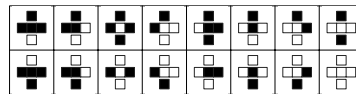
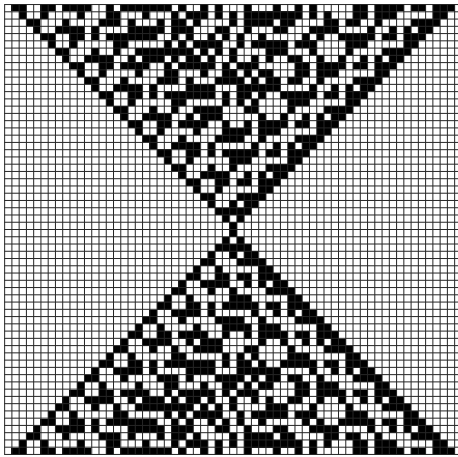
So is it possible to get more complex behavior while maintaining reversibility? There are a total of 7,625,597,484,987 cellular automata with three colors and nearest-neighbor rules, and searching through these one finds just 1800 that are reversible. Of these 1800, many again exhibit simple behavior, much like the pictures above. But some exhibit more complex behavior, as in the pictures below.



Examples of some of the 1800 reversible cellular automata with three colors and nearest-neighbor rules. Even though these systems exhibit complex behavior that scrambles the initial conditions, all of them are still reversible, so that starting from the configuration of cells at the bottom of each picture, it is always possible to deduce the configurations on all previous steps.

How can one now tell that such systems are reversible? It is no longer true that their evolution leads only to simple transformations of the initial conditions. But one can still check that starting with the specific configuration of cells at the bottom of each picture, one can evolve backwards to get to the top of the picture. And given a particular rule it turns out to be fairly straightforward to do a detailed analysis that allows one to prove or disprove its reversibility.

But in trying to understand the range of behavior that can occur in reversible systems it is often convenient to consider classes of cellular automata with rules that are specifically constructed to be reversible. One such class is illustrated below. The idea is to have rules that explicitly remain the same even if they are turned upside-down, thereby interchanging the roles of past and future.

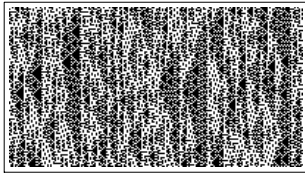


An example of a cellular automaton that is explicitly set up to be reversible. The rule for the system remains unchanged if all its elements are turned upside-down—effectively interchanging the roles of past and future. Patterns produced by the rule must exhibit the same time reversal symmetry, as shown on the left. The specific rule used here is based on taking elementary rule 214, then adding the specification that the new color of a cell should be inverted whenever the cell was black two steps back. Note that by allowing a total of four rather than two colors, a version of the rule that depends only on the immediately preceding step can be constructed.

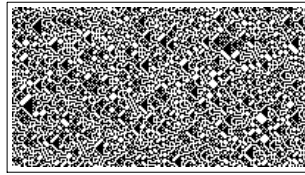
Such rules can be constructed by taking ordinary cellular automata and adding dependence on colors two steps back.

The resulting rules can be run both forwards and backwards. In each case they require knowledge of the colors of cells on not one but two successive steps. Given this knowledge, however, the rules can be used to determine the configuration of cells on either future or past steps.

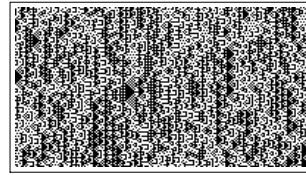
The next two pages show examples of the behavior of such cellular automata with both random and simple initial conditions.



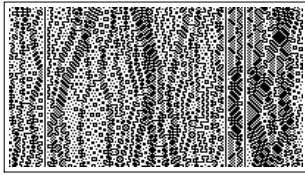
rule 13R



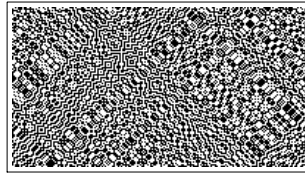
rule 30R



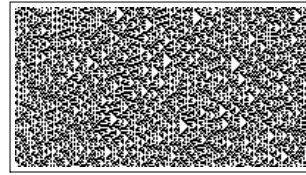
rule 67R



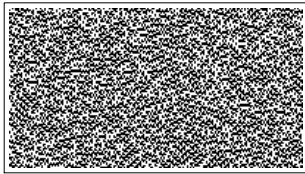
rule 173R



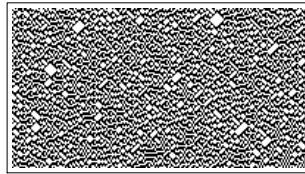
rule 90R



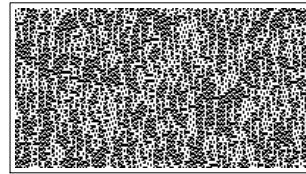
rule 142R



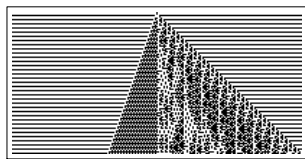
rule 173R



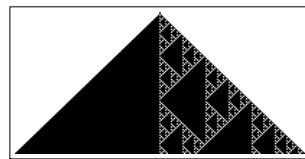
rule 190R



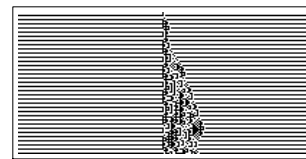
rule 197R



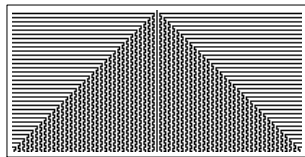
rule 13R



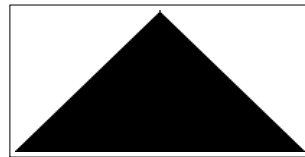
rule 30R



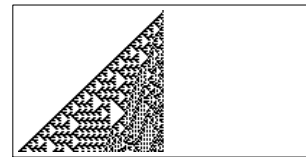
rule 67R



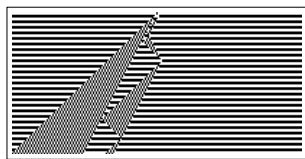
rule 173R



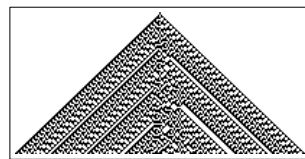
rule 90R



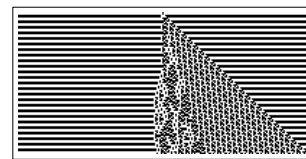
rule 142R



rule 173R

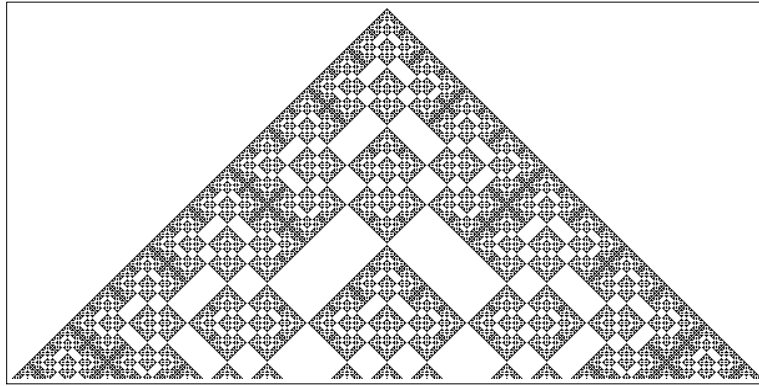


rule 190R

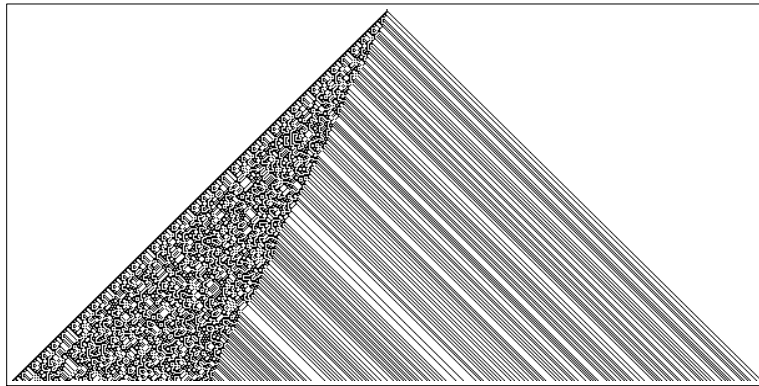


rule 197R

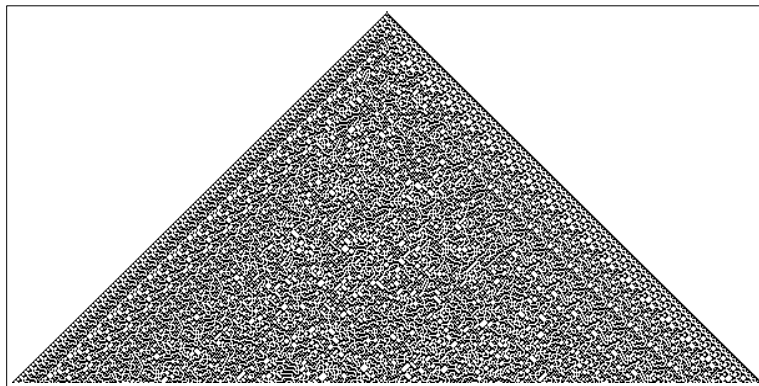
Examples of reversible cellular automata starting from random and from simple initial conditions. In the upper block of pictures, every cell is chosen to be black or white with equal probability on the two successive first steps. In the lower block of pictures, only the center cell is taken to be black on these steps.



rule 150R

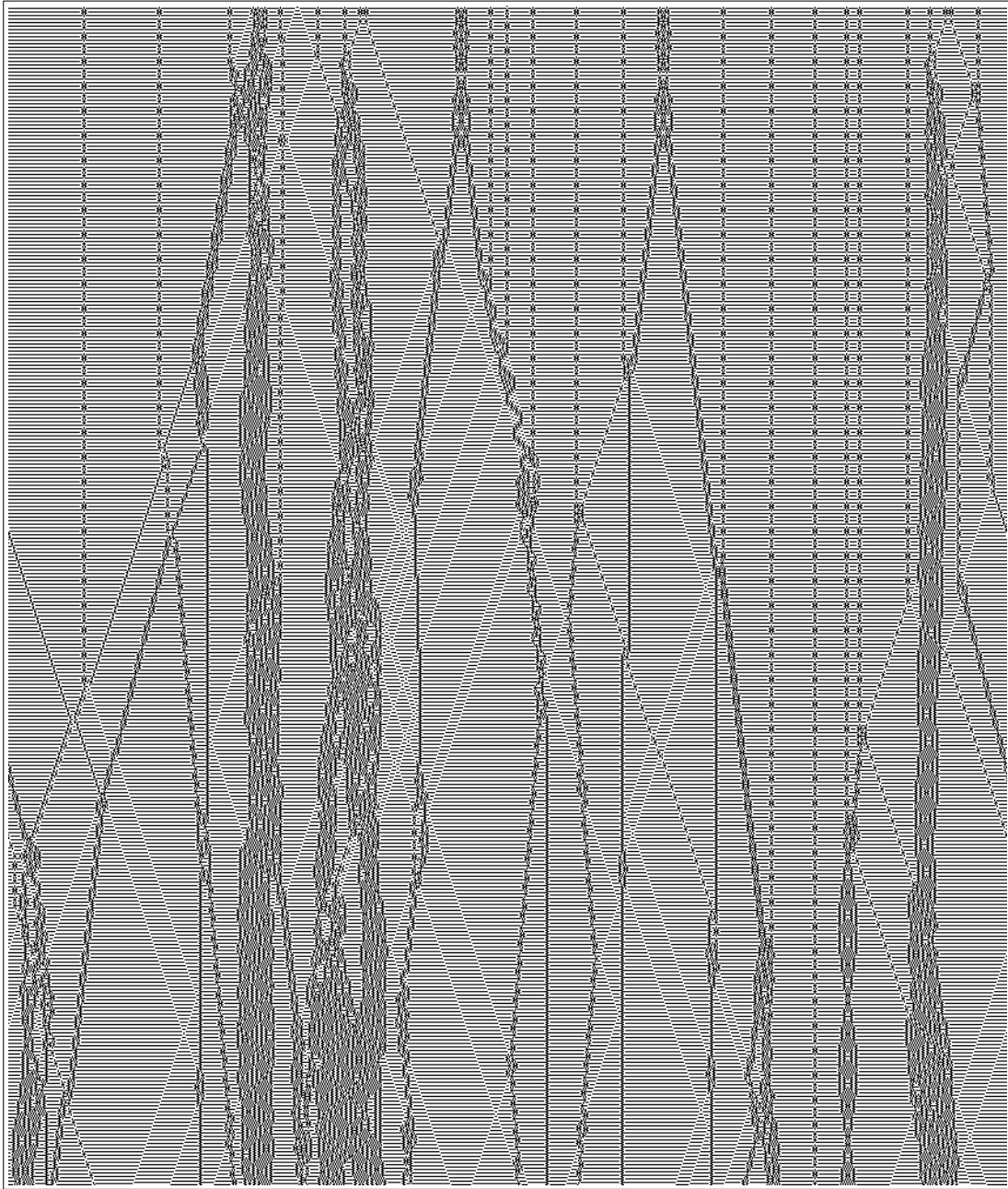


rule 154R



rule 214R

The evolution of three reversible cellular automata for 300 steps. In the first case, a regular nested pattern is obtained. In the other cases, the patterns show many features of randomness.



rule 37R

An example of a reversible cellular automaton whose evolution supports localized structures. Because of the reversibility of the underlying rule, every collision must be able to occur equally well when its initial and final states are interchanged.

In some cases, the behavior is fairly simple, and the patterns obtained have simple repetitive or nested structures. But in many cases, even with simple initial conditions, the patterns produced are highly complex, and seem in many respects random.

The reversibility of the underlying rules has some obvious consequences, such as the presence of triangles pointing sideways but not down. But despite their reversibility, the rules still manage to produce the kinds of complex behavior that we have seen in cellular automata and many other systems throughout this book.

So what about localized structures?

The picture on the facing page demonstrates that these can also occur in reversible systems. There are some constraints on the details of the kinds of collisions that are possible, but reversible rules typically tend to work very much like ordinary ones.

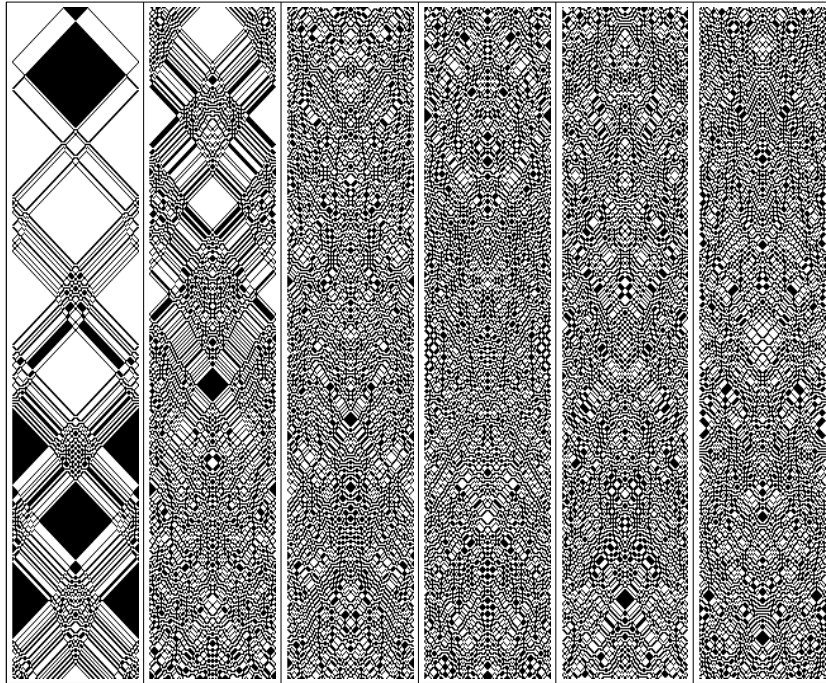
So in the end it seems that even though only a very small fraction of possible systems have the property of being reversible, such systems can still exhibit behavior just as complex as one sees anywhere else.

Irreversibility and the Second Law of Thermodynamics

All the evidence we have from particle physics and elsewhere suggests that at a fundamental level the laws of physics are precisely reversible. Yet our everyday experience is full of examples of seemingly irreversible phenomena. Most often, what happens is that a system which starts in a fairly regular or organized state becomes progressively more and more random and disorganized. And it turns out that this phenomenon can already be seen in many simple programs.

The picture at the top of the next page shows an example based on a reversible cellular automaton of the type discussed in the previous section. The black cells in this system act a little like particles which bounce around inside a box and interact with each other when they collide.

At the beginning the particles are placed in a simple arrangement at the center of the box. But over the course of time the picture shows that the arrangement of particles becomes progressively more random.

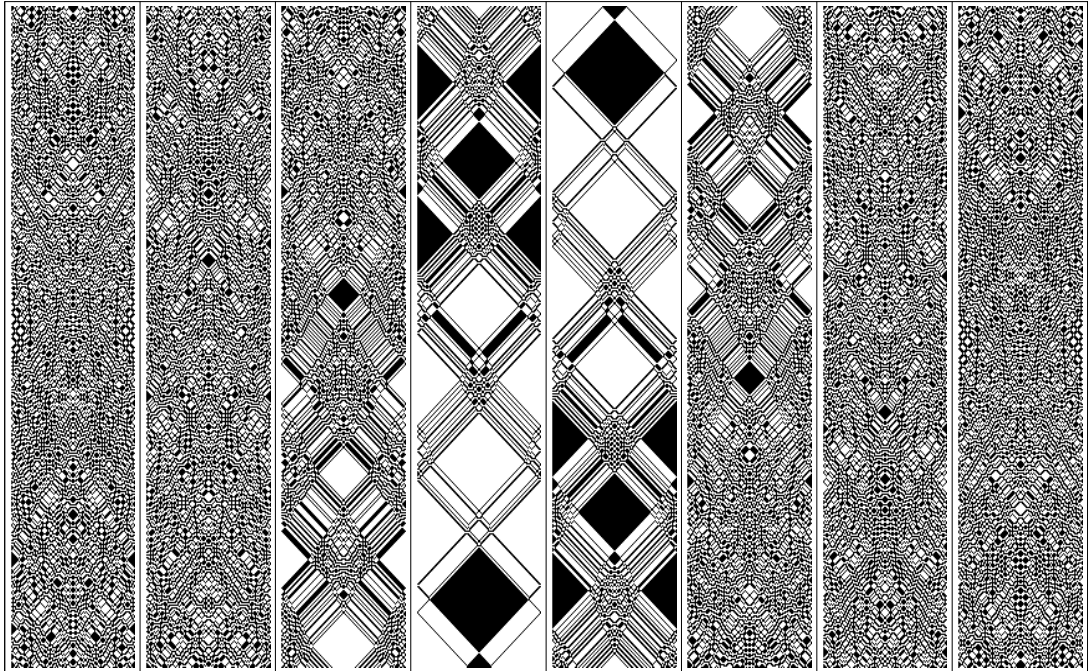


A reversible cellular automaton that exhibits seemingly irreversible behavior. Starting from an initial condition in which all black cells or particles lie at the center of a box, the distribution becomes progressively more random. Such behavior appears to be the central phenomenon responsible for the Second Law of Thermodynamics. The specific cellular automaton used here is rule 122R. The system is restricted to a region of size 100 cells.

Typical intuition from traditional science makes it difficult to understand how such randomness could possibly arise. But the discovery in this book that a wide range of systems can generate randomness even with very simple initial conditions makes it seem considerably less surprising.

But what about reversibility? The underlying rules for the cellular automaton used in the picture above are precisely reversible. Yet the picture itself does not at first appear to be at all reversible. For there appears to be an irreversible increase in randomness as one goes down successive panels on the page.

The resolution of this apparent conflict is however fairly straightforward. For as the picture on the facing page demonstrates, if the



An extended version of the picture on the facing page, in which the reversibility of the underlying cellular automaton is more clearly manifest. An initial condition is carefully constructed so that halfway through the evolution shown a simple arrangement of particles will be produced. If one starts with this arrangement, then the randomness of the system will effectively increase whether one goes forwards or backwards in time from that point.

simple arrangement of particles occurs in the middle of the evolution, then one can readily see that randomness increases in exactly the same way—whether one goes forwards or backwards from that point.

Yet there is still something of a mystery. For our everyday experience is full of examples in which randomness increases much as in the second half of the picture above. But we essentially never see the kind of systematic decrease in randomness that occurs in the first half.

By setting up the precise initial conditions that exist at the beginning of the whole picture it would certainly in principle be possible to get such behavior. But somehow it seems that initial conditions like these essentially never actually occur in practice.

There has in the past been considerable confusion about why this might be the case. But the key to understanding what is going on is simply to realize that one has to think not only about the systems one is studying, but also about the types of experiments and observations that one uses in the process of studying them.

The crucial point then turns out to be that practical experiments almost inevitably end up involving only initial conditions that are fairly simple for us to describe and construct. And with these types of initial conditions, systems like the one on the previous page always tend to exhibit increasing randomness.

But what exactly is it that determines the types of initial conditions that one can use in an experiment? It seems reasonable to suppose that in any meaningful experiment the process of setting up the experiment should somehow be simpler than the process that the experiment is intended to observe.

But how can one compare such processes? The answer that I will develop in considerable detail later in this book is to view all such processes as computations. The conclusion is then that the computation involved in setting up an experiment should be simpler than the computation involved in the evolution of the system that is to be studied by the experiment.

It is clear that by starting with a simple state and then tracing backwards through the actual evolution of a reversible system one can find initial conditions that will lead to decreasing randomness. But if one looks for example at the pictures on the last couple of pages the complexity of the behavior seems to preclude any less arduous way of finding such initial conditions. And indeed I will argue in Chapter 12 that the Principle of Computational Equivalence suggests that in general no such reduced procedure should exist.

The consequence of this is that no reasonable experiment can ever involve setting up the kind of initial conditions that will lead to decreases in randomness, and that therefore all practical experiments will tend to show only increases in randomness.

It is this basic argument that I believe explains the observed validity of what in physics is known as the Second Law of Thermodynamics. The law was first formulated more than a century

ago, but despite many related technical results, the basic reasons for its validity have until now remained rather mysterious.

The field of thermodynamics is generally concerned with issues of heat and energy in physical systems. A fundamental fact known since the mid-1800s is that heat is a form of energy associated with the random microscopic motions of large numbers of atoms or other particles.

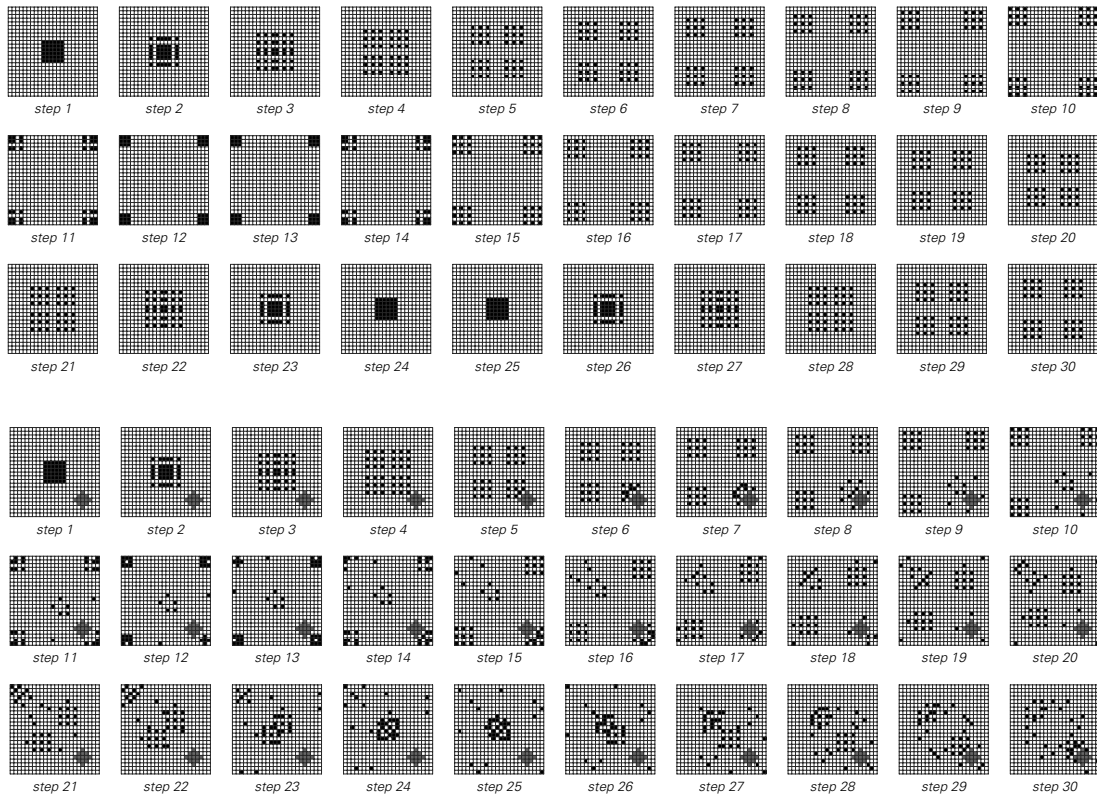
One formulation of the Second Law then states that any energy associated with organized motions of such particles tends to degrade irreversibly into heat. And the pictures at the beginning of this section show essentially just such a phenomenon. Initially there are particles which move in a fairly regular and organized way. But as time goes on, the motion that occurs becomes progressively more random.

There are several details of the cellular automaton used above that differ from actual physical systems of the kind usually studied in thermodynamics. But at the cost of some additional technical complication, it is fairly straightforward to set up a more realistic system.

The pictures on the next two pages show a particular two-dimensional cellular automaton in which black squares representing particles move around and collide with each other, essentially like particles in an ideal gas. This cellular automaton shares with the cellular automaton at the beginning of the section the property of being reversible. But it also has the additional feature that in every collision the total number of particles in it remains unchanged. And since each particle can be thought of as having a certain energy, it follows that the total energy of the system is therefore conserved.

In the first case shown, the particles are taken to bounce around in an empty square box. And it turns out that in this particular case only very simple repetitive behavior is ever obtained. But almost any change destroys this simplicity.

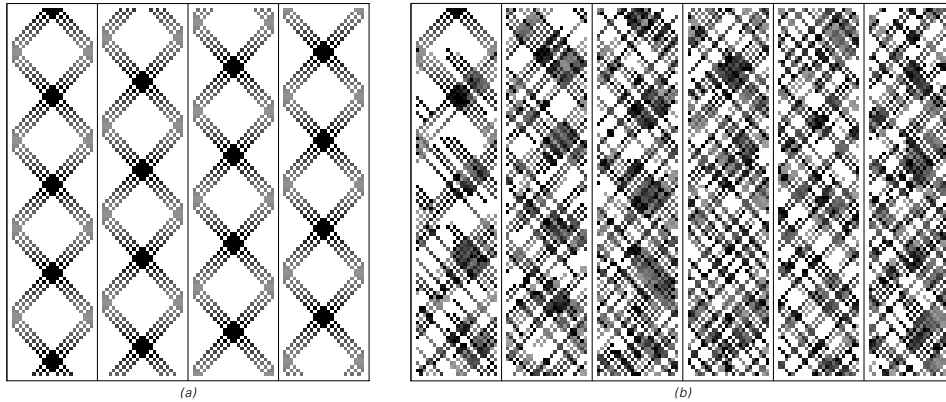
And in the second case, for example, the presence of a small fixed obstacle leads to rapid randomization in the arrangement of particles—very much like the randomization we saw in the one-dimensional cellular automaton that we discussed earlier in this section.



The behavior of a simple two-dimensional cellular automaton that emulates an ideal gas of particles. In the top group of pictures, the particles bounce around in an empty square box. In the bottom group of pictures, the box contains a small fixed obstacle. In the top group of pictures, the arrangement of particles shows simple repetitive behavior. In the bottom group, however, it becomes progressively more random with time. The underlying rules for the cellular automaton used here are reversible, and conserve the total number of particles. The specific rules are based on 2×2 blocks—a two-dimensional generalization of the block cellular automata to be discussed in the next section. For each 2×2 block the configuration of particles is taken to remain the same at a particular step unless there are exactly two particles arranged diagonally within the block, in which case the particles move to the opposite diagonal.

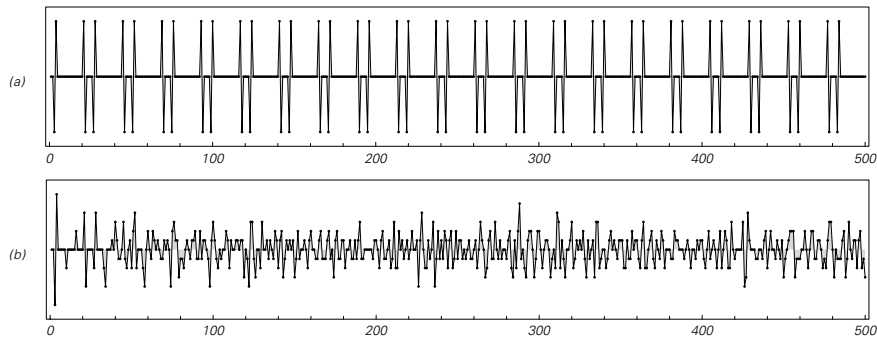
So even though the total of the energy of all particles remains the same, the distribution of this energy becomes progressively more random, just as the usual Second Law implies.

An important practical consequence of this is that it becomes increasingly difficult to extract energy from the system in the form of systematic mechanical work. At an idealized level one might imagine trying to do this by inserting into the system some kind of paddle which would experience force as a result of impacts from particles.



Time histories of the cellular automata from the facing page. In each case a slice is taken through the midline of the box. Black cells that are further from the midline are shown in progressively lighter shades of gray. Case (a) corresponds to an empty square box, and shows simple repetitive behavior. Case (b) corresponds to a box containing a fixed obstacle, and in this case rapid randomization is seen. Each panel corresponds to 100 steps in the evolution of the system; the box is 24 cells across.

The pictures below show how such force might vary with time in cases (a) and (b) above. In case (a), where no randomization occurs, the force can readily be predicted, and it is easy to imagine harnessing it to produce systematic mechanical work. But in case (b), the force quickly randomizes, and there is no obvious way to obtain systematic mechanical work from it.



The force on an idealized paddle placed on the midline of the systems shown above. The force reflects an imbalance in the number of particles at each step arriving at the midline from above and below. In case (a) this imbalance is readily predictable. In case (b), however, it rapidly becomes for most practical purposes random. This randomness is essentially what makes it impossible to build a physical perpetual motion machine which continually turns heat into mechanical work.

One might nevertheless imagine that it would be possible to devise a complicated machine, perhaps with an elaborate arrangement of paddles, that would still be able to extract systematic mechanical work even from an apparently random distribution of particles. But it turns out that in order to do this the machine would effectively have to be able to predict where every particle would be at every step in time.

And as we shall discuss in Chapter 12, this would mean that the machine would have to perform computations that are as sophisticated as those that correspond to the actual evolution of the system itself. The result is that in practice it is never possible to build perpetual motion machines that continually take energy in the form of heat—or randomized particle motions—and convert it into useful mechanical work.

The impossibility of such perpetual motion machines is one common statement of the Second Law of Thermodynamics. Another is that a quantity known as entropy tends to increase with time.

Entropy is defined as the amount of information about a system that is still unknown after one has made a certain set of measurements on the system. The specific value of the entropy will depend on what measurements one makes, but the content of the Second Law is that if one repeats the same measurements at different times, then the entropy deduced from them will tend to increase with time.

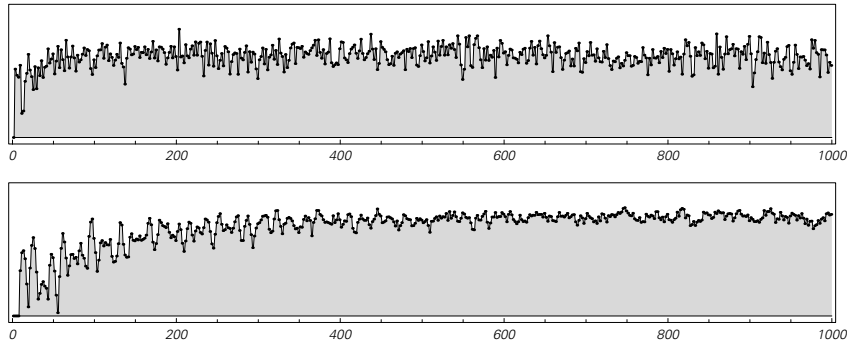
If one managed to find the positions and properties of all the particles in the system, then no information about the system would remain unknown, and the entropy of the system would just be zero. But in a practical experiment, one cannot expect to be able to make anything like such complete measurements.

And more realistically, the measurements one makes might for example give the total numbers of particles in certain regions inside the box. There are then a large number of possible detailed arrangements of particles that are all consistent with the results of such measurements. The entropy is defined as the amount of additional information that would be needed in order to pick out the specific arrangement that actually occurs.

We will discuss in more detail in Chapter 10 the notion of amount of information. But here we can imagine numbering all the possible arrangements of particles that are consistent with the results of our

measurements, so that the amount of information needed to pick out a single arrangement is essentially the length in digits of one such number.

The pictures below show the behavior of the entropy calculated in this way for systems like the one discussed above. And what we see is that the entropy does indeed tend to increase, just as the Second Law implies.



The entropy as a function of time for systems of the type shown in case (b) from page 447. The top plot is exactly for case (b); the bottom one is for a system three times larger in size. The entropy is found in each case by working out how many possible configurations of particles are consistent with measurements of the total numbers of particles in a 6×6 grid of regions within the system. Just as the Second Law of Thermodynamics suggests, the entropy tends to increase with time. Note that the plots above would be exactly symmetrical if they were continued to the left: the entropy would increase in the same way going both forwards and backwards from the simple initial conditions used.

In effect what is going on is that the measurements we make represent an attempt to determine the state of the system. But as the arrangement of particles in the system becomes more random, this attempt becomes less and less successful.

One might imagine that there could be a more elaborate set of measurements that would somehow avoid these problems, and would not lead to increasing entropy. But as we shall discuss in Chapter 12, it again turns out that setting up such measurements would have to involve the same level of computational effort as the actual evolution of the system itself. And as a result, one concludes that the entropy associated with measurements done in practical experiments will always tend to increase, as the Second Law suggests.

In Chapter 12 we will discuss in more detail some of the key ideas involved in coming to this conclusion. But the basic point is that the phenomenon of entropy increase implied by the Second Law is a more or less direct consequence of the phenomenon discovered in this book that even with simple initial conditions many systems can produce complex and seemingly random behavior.

One aspect of the generation of randomness that we have noted several times in earlier chapters is that once significant randomness has been produced in a system, the overall properties of that system tend to become largely independent of the details of its initial conditions.

In any system that is reversible it must always be the case that different initial conditions lead to at least slightly different states—otherwise there would be no unique way of going backwards. But the point is that even though the outcomes from different initial conditions differ in detail, their overall properties can still be very much the same.

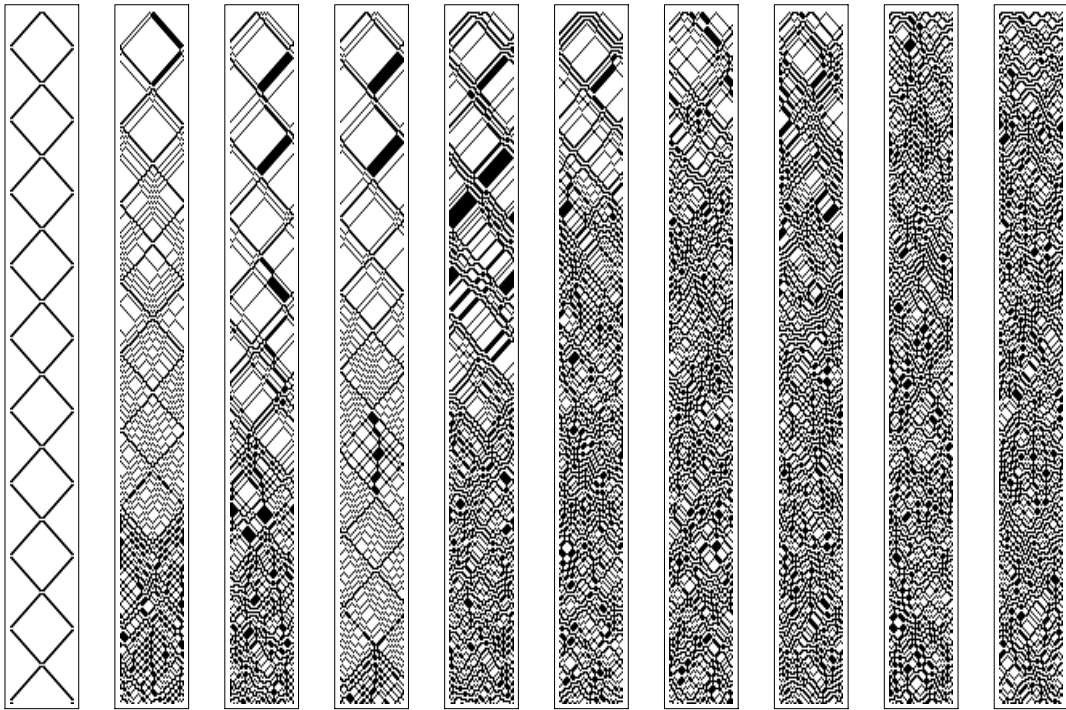
The pictures on the facing page show an example of what can happen. Every individual picture has different initial conditions. But whenever randomness is produced the overall patterns that are obtained look in the end almost indistinguishable.

The reversibility of the underlying rules implies that at some level it must be possible to recognize outcomes from different kinds of initial conditions. But the point is that to do so would require a computation far more sophisticated than any that could meaningfully be done as part of a practical measurement process.

So this means that if a system generates sufficient randomness, one can think of it as evolving towards a unique equilibrium whose properties are for practical purposes independent of its initial conditions.

This fact turns out in a sense to be implicit in many everyday applications of physics. For it is what allows us to characterize all sorts of physical systems by just specifying a few parameters such as temperature and chemical composition—and avoids us always having to know the details of the initial conditions and history of each system.

The existence of a unique equilibrium to which any particular system tends to evolve is also a common statement of the Second Law of



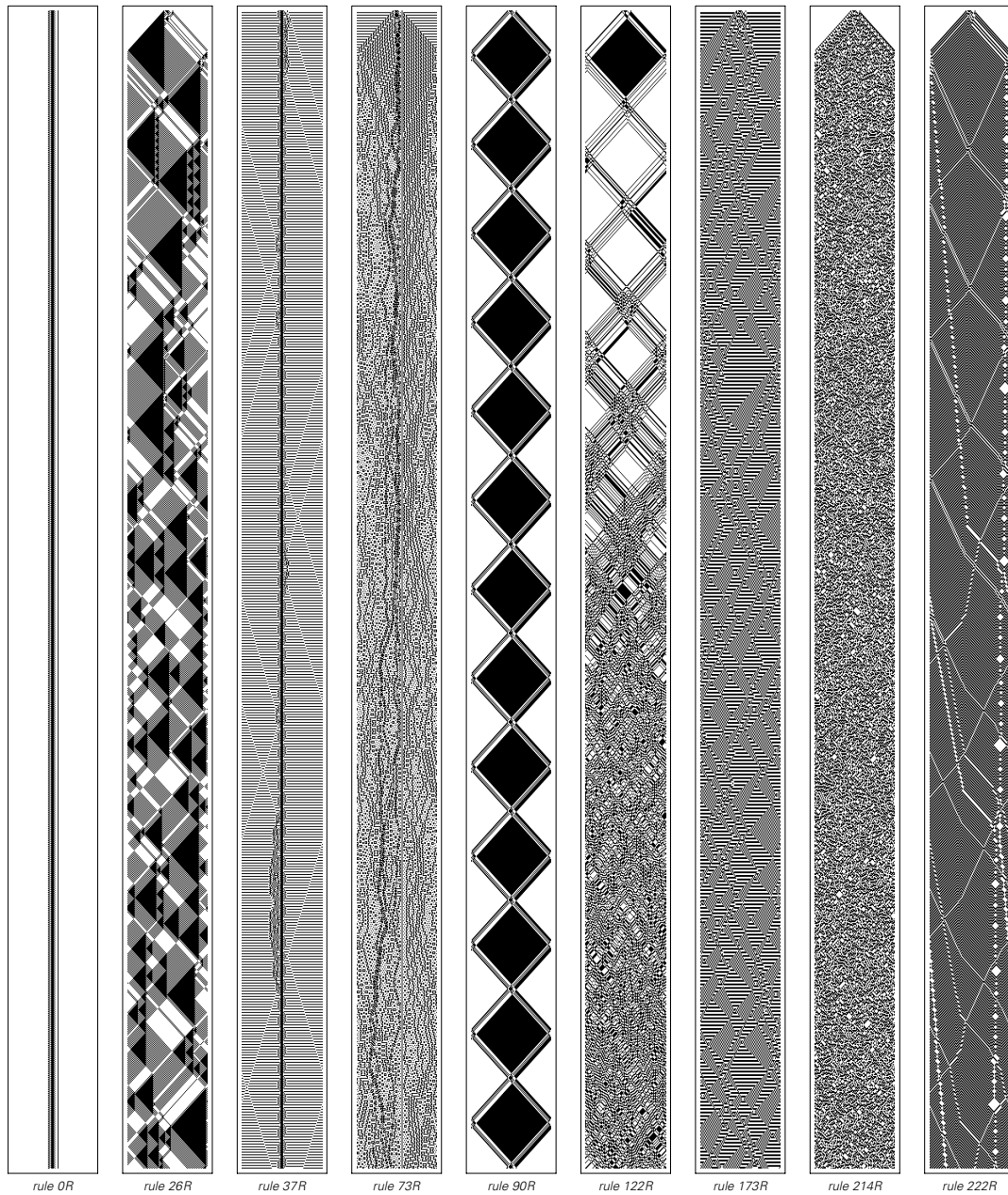
The approach to equilibrium in a reversible cellular automaton with a variety of different initial conditions. Apart from exceptional cases where no randomization occurs, the behavior obtained with different initial conditions is eventually quite indistinguishable in its overall properties. Because the underlying rule is reversible, however, the details with different initial conditions are always at least slightly different—otherwise it would not be possible to go backwards in a unique way. The rule used here is 122R. Successive pairs of pictures have initial conditions that differ only in the color of a single cell at the center.

Thermodynamics. And once again, therefore, we find that the Second Law is associated with basic phenomena that we already saw early in this book.

But just how general is the Second Law? And does it really apply to all of the various kinds of systems that we see in nature?

Starting nearly a century ago it came to be widely believed that the Second Law is an almost universal principle. But in reality there is surprisingly little evidence for this.

Indeed, almost all of the detailed applications ever made of the full Second Law have been concerned with just one specific area: the behavior of gases. By now there is therefore good evidence that gases obey the Second Law—just as the idealized model earlier in this section suggests. But what about other kinds of systems?



Examples of reversible cellular automata with various rules. Some quickly randomize, as the Second Law of Thermodynamics would suggest. But others do not—and thus in effect do not obey the Second Law of Thermodynamics.

The pictures on the facing page show examples of various reversible cellular automata. And what we see immediately from these pictures is that while some systems exhibit exactly the kind of randomization implied by the Second Law, others do not.

The most obvious exceptions are cases like rule 0R and rule 90R, where the behavior that is produced has only a very simple fixed or repetitive form. And existing mathematical studies have indeed identified these simple exceptions to the Second Law. But they have somehow implicitly assumed that no other kinds of exceptions can exist.

The picture on the next page, however, shows the behavior of rule 37R over the course of many steps. And in looking at this picture, we see a remarkable phenomenon: there is neither a systematic trend towards increasing randomness, nor any form of simple predictable behavior. Indeed, it seems that the system just never settles down, but rather continues to fluctuate forever, sometimes becoming less orderly, and sometimes more so.

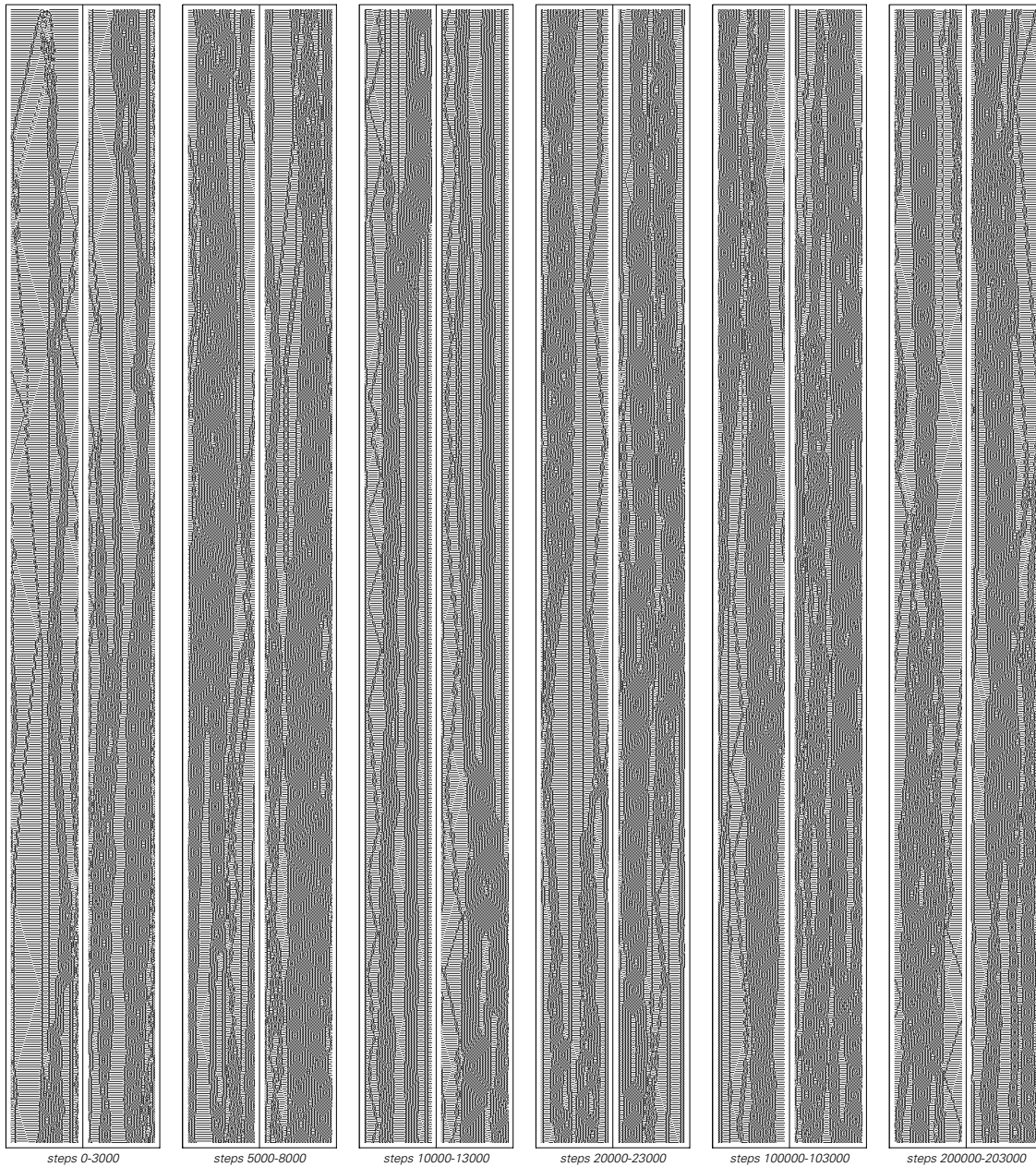
So how can such behavior be understood in the context of the Second Law? There is, I believe, no choice but to conclude that for practical purposes rule 37R simply does not obey the Second Law.

And as it turns out, what happens in rule 37R is not so different from what seems to happen in many systems in nature. If the Second Law was always obeyed, then one might expect that by now every part of our universe would have evolved to completely random equilibrium.

Yet it is quite obvious that this has not happened. And indeed there are many kinds of systems, notably biological ones, that seem to show, at least temporarily, a trend towards increasing order rather than increasing randomness.

How do such systems work? A common feature appears to be the presence of some kind of partitioning: the systems effectively break up into parts that evolve at least somewhat independently for long periods of time.

The picture on page 456 shows what happens if one starts rule 37R with a single small region of randomness. And for a while what one sees is that the randomness that has been inserted persists. But eventually the system instead seems to organize itself to yield just a small number of simple repetitive structures.



More steps in the evolution of the reversible cellular automaton with rule 37R. This system is an example of one that does not in any meaningful way obey the Second Law of Thermodynamics. Instead of exhibiting progressively more random behavior, it appears to fluctuate between quite ordered and quite disordered states.

This kind of self-organization is quite opposite to what one would expect from the Second Law. And at first it also seems inconsistent with the reversibility of the system. For if all that is left at the end are a few simple structures, how can there be enough information to go backwards and reconstruct the initial conditions?

The answer is that one has to consider not only the stationary structures that stay in the middle of the system, but also all various small structures that were emitted in the course of the evolution. To go backwards one would need to set things up so that one absorbs exactly the sequence of structures that were emitted going forwards.

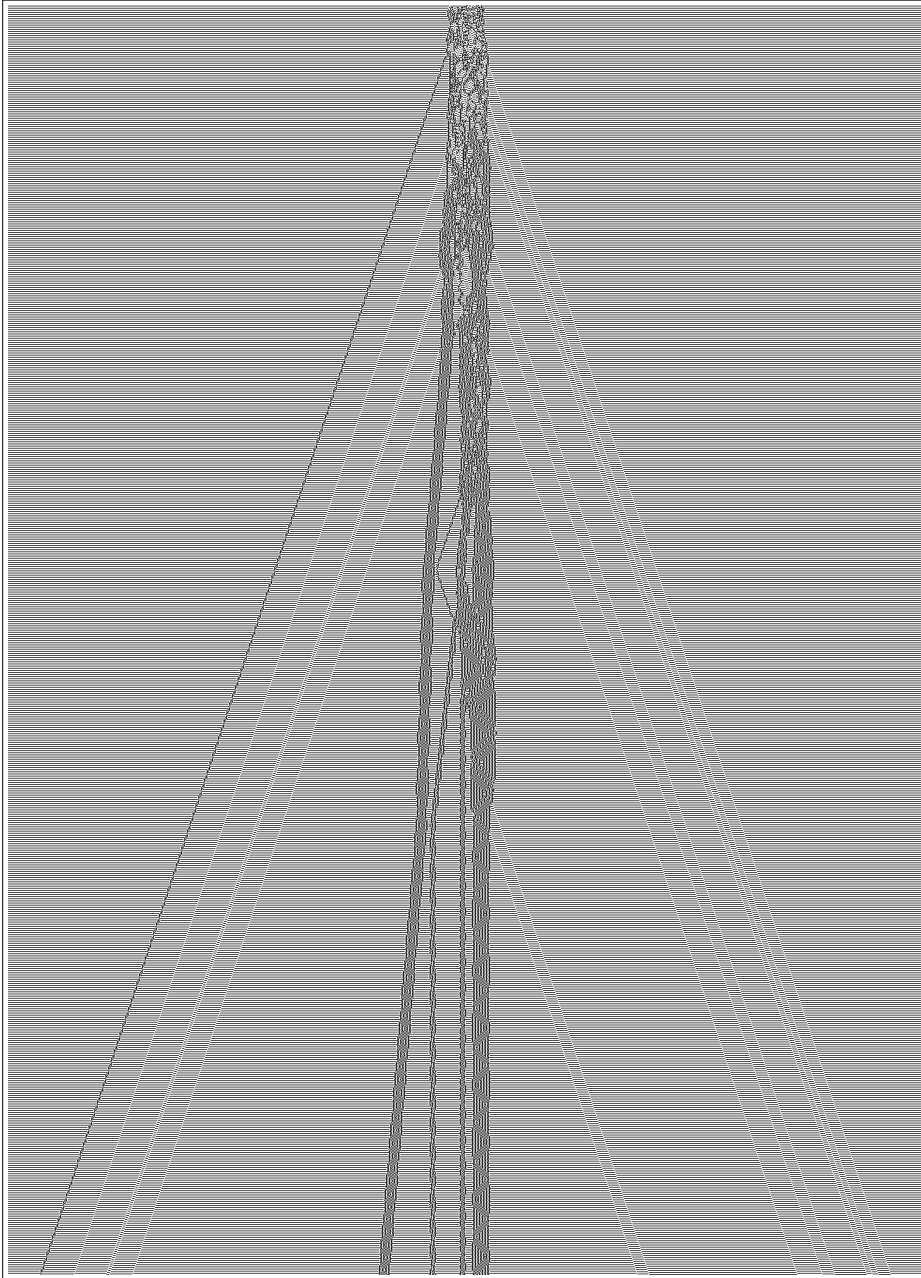
If, however, one just lets the emitted structures escape, and never absorbs any other structures, then one is effectively losing information. The result is that the evolution one sees can be intrinsically not reversible, so that all of the various forms of self-organization that we saw earlier in this book in cellular automata that do not have reversible rules can potentially occur.

If we look at the universe on a large scale, then it turns out that in a certain sense there is more radiation emitted than absorbed. Indeed, this is related to the fact that the night sky appears dark, rather than having bright starlight coming from every direction. But ultimately the asymmetry between emission and absorption is a consequence of the fact that the universe is expanding, rather than contracting, with time.

The result is that it is possible for regions of the universe to become progressively more organized, despite the Second Law, and despite the reversibility of their underlying rules. And this is a large part of the reason that organized galaxies, stars and planets can form.

Allowing information to escape is a rather straightforward way to evade the Second Law. But what the pictures on the facing page demonstrate is that even in a completely closed system, where no information at all is allowed to escape, a system like rule 37R still does not follow the uniform trend towards increasing randomness that is suggested by the Second Law.

What instead happens is that kinds of membranes form between different regions of the system, and within each region orderly behavior can then occur, at least while the membrane survives.



An example of evolution according to rule 37R from an initial condition containing a fairly random region. Even though the system is reversible, this region tends to organize itself so as to take on a much simpler form. Information on the initial conditions ends up being carried by localized structures which radiate outwards.

This basic mechanism may well be the main one at work in many biological systems: each cell or each organism becomes separated from others, and while it survives, it can exhibit organized behavior.

But looking at the pictures of rule 37R on page 454 one may ask whether perhaps the effects we see are just transients, and that if we waited long enough something different would happen.

It is an inevitable feature of having a closed system of limited size that in the end the behavior one gets must repeat itself. And in rules like 0R and 90R shown on page 452 the period of repetition is always very short. But for rule 37R it usually turns out to be rather long. Indeed, for the specific example shown on page 454, the period is 293,216,266.

In general, however, the maximum possible period for a system containing a certain number of cells can be achieved only if the evolution of the system from any initial condition eventually visits all the possible states of the system, as discussed on page 258. And if this in fact happens, then at least eventually the system will inevitably spend most of its time in states that seem quite random.

But in rule 37R there is no such ergodicity. And instead, starting from any particular initial condition, the system will only ever visit a tiny fraction of all possible states. Yet since the total number of states is astronomically large—about 10^{60} for size 100—the number of states visited by rule 37R, and therefore the repetition period, can still be extremely long.

There are various subtleties involved in making a formal study of the limiting behavior of rule 37R after a very long time. But irrespective of these subtleties, the basic fact remains that so far as I can tell, rule 37R simply does not follow the predictions of the Second Law.

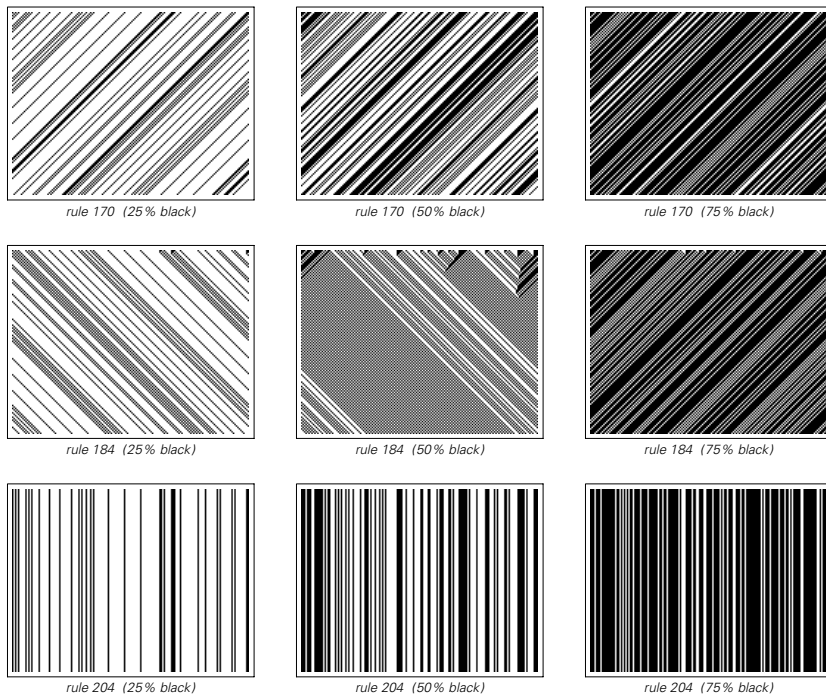
And indeed I strongly suspect that there are many systems in nature which behave in more or less the same way. The Second Law is an important and quite general principle—but it is not universally valid. And by thinking in terms of simple programs we have thus been able in this section not only to understand why the Second Law is often true, but also to see some of its limitations.

Conserved Quantities and Continuum Phenomena

Reversibility is one general feature that appears to exist in the basic laws of physics. Another is conservation of various quantities—so that for example in the evolution of any closed physical system, total values of quantities like energy and electric charge appear always to stay the same.

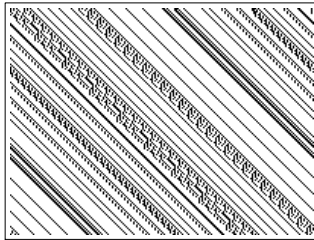
With most rules, systems like cellular automata do not usually exhibit such conservation laws. But just as with reversibility, it turns out to be possible to find rules that for example conserve the total number of black cells appearing on each step.

Among elementary cellular automata with just two colors and nearest-neighbor rules, the only types of examples are the fairly trivial ones shown in the pictures below.

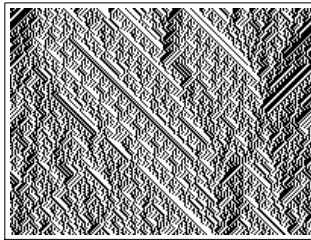


Elementary cellular automata whose evolution conserves the total number of black cells. The behavior of the rules shown here is simple enough that in each case it is fairly obvious how the number of black cells manages to stay the same on every step.

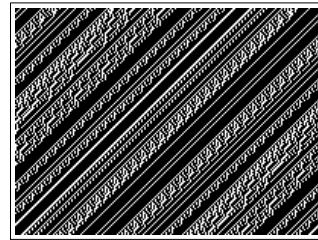
But with next-nearest-neighbor rules, more complicated examples become possible, as the pictures below demonstrate.



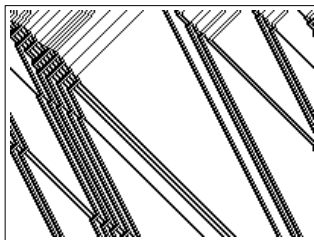
rule 3450663328 (25% black)



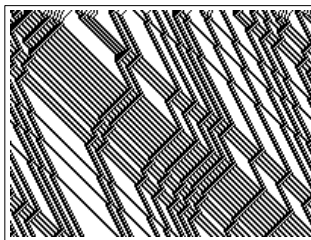
rule 3450663328 (50% black)



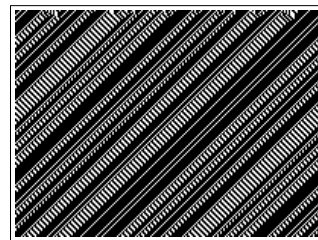
rule 3450663328 (75% black)



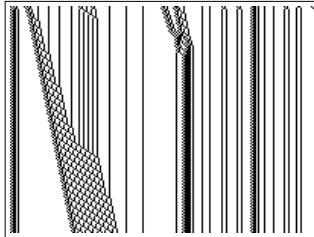
rule 3484741764 (25% black)



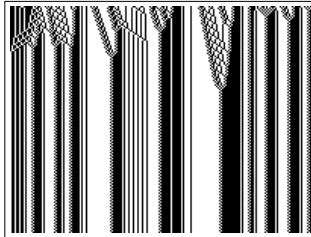
rule 3484741764 (50% black)



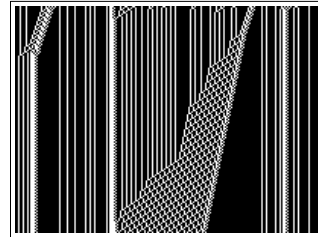
rule 3484741764 (75% black)



rule 3822644248 (25% black)



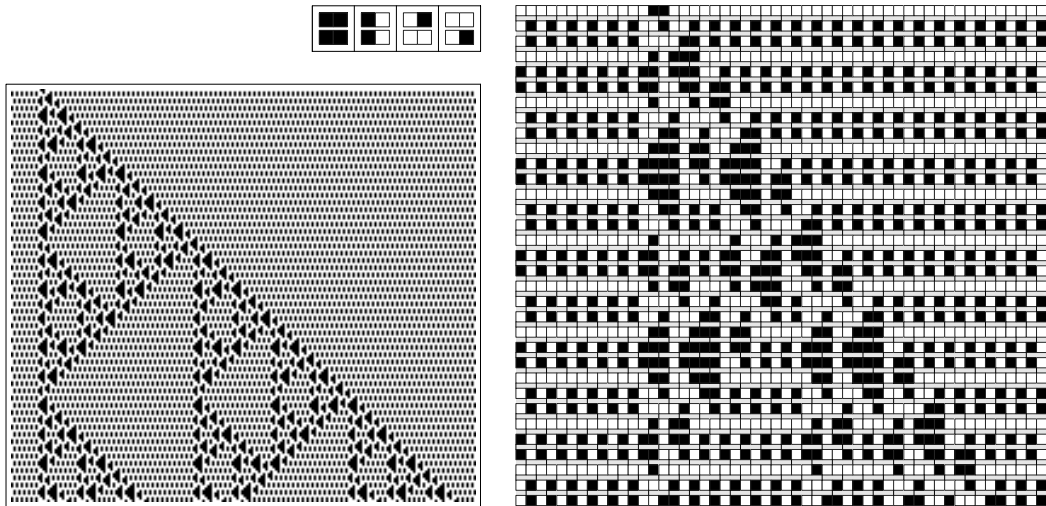
rule 3822644248 (50% black)



rule 3822644248 (75% black)

Examples of cellular automata with next-nearest-neighbor rules whose evolution conserves the total number of black cells. Even though it is not immediately obvious by eye, the total number of black cells stays exactly the same on each successive step in each picture. Among the 4,294,967,296 possible next-neighbor rules, only 428 exhibit the kind of conservation property shown here.

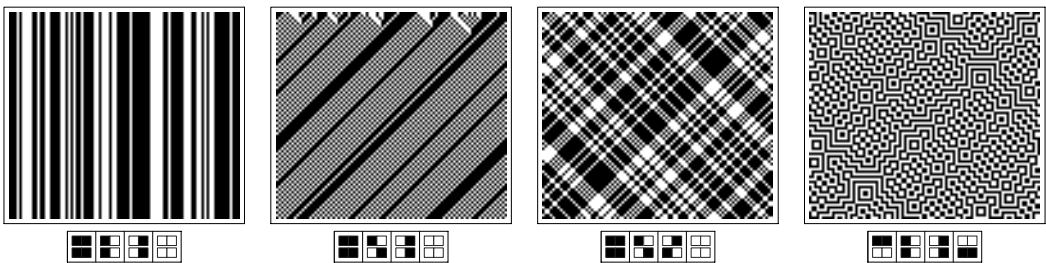
One straightforward way to generate collections of systems that will inevitably exhibit conserved quantities is to work not with ordinary cellular automata but instead with block cellular automata. The basic idea of a block cellular automaton is illustrated at the top of the next page. At each step what happens is that blocks of adjacent cells are replaced by other blocks of the same size according to some definite rule. And then on successive steps the alignment of these blocks shifts by one cell.



An example of a block cellular automaton. The system works by partitioning the sequence of cells that exists at each step into pairs, then replacing these pairs by other pairs according to the rule shown. The choice of whether to pair a cell with its left or right neighbor alternates on successive steps. Like many block cellular automata, the system shown is reversible, since in the rule each pair has a unique predecessor. It does not, however, conserve the total number of black cells.

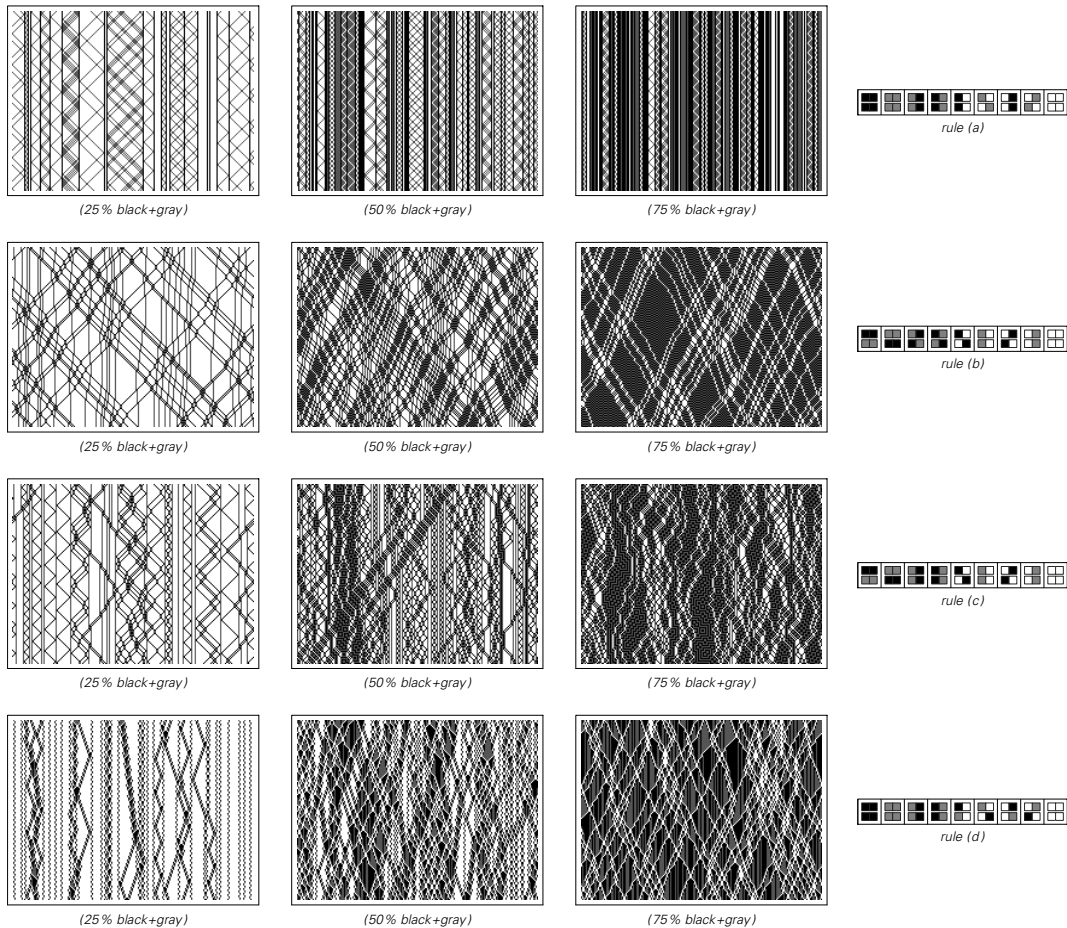
And with this setup, if the underlying rules replace each block by one that contains the same number of black cells, it is inevitable that the system as a whole will conserve the total number of black cells.

With two possible colors and blocks of size two the only kinds of block cellular automata that conserve the total number of black cells are the ones shown below—and all of these exhibit rather trivial behavior.



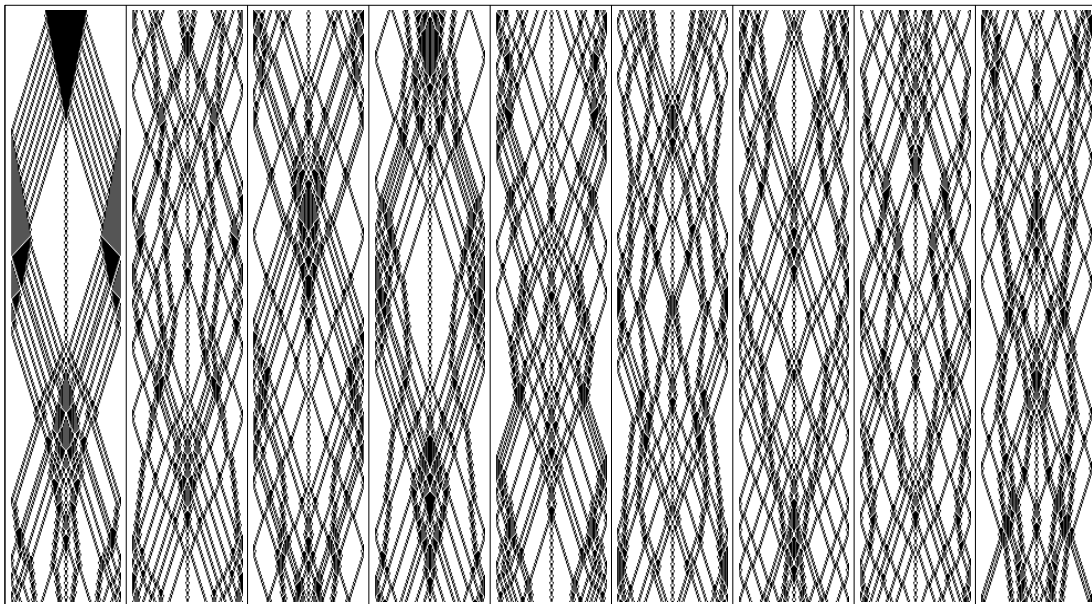
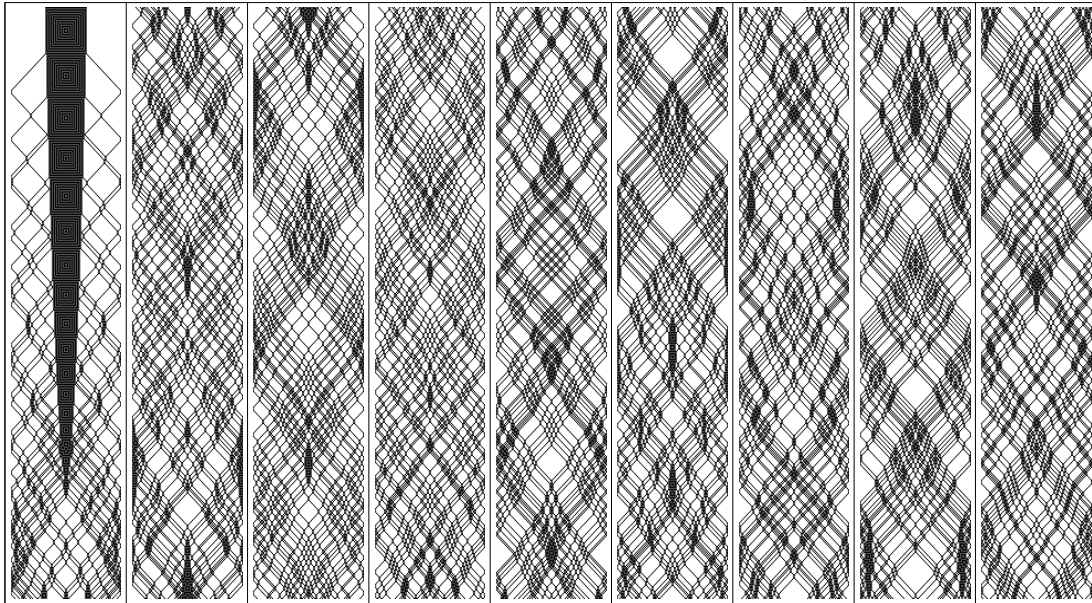
Block cellular automata with two possible colors and blocks of size two that conserve the total number of black cells (the last example has this property only on alternate steps). It so happens that all but the second of the rules shown here not only conserve the total number of black cells but also turn out to be reversible.

But if one allows three possible colors, and requires, say, that the total number of black and gray cells together be conserved, then more complicated behavior can occur, as in the pictures below.



Block cellular automata with three possible colors which conserve the combined number of black and gray cells. In rule (a), black and gray cells remain in localized regions. In rule (b), they move in fairly simple ways, and in rules (c) and (d), they move in a seemingly somewhat random way. The rules shown here are reversible, although their behavior is similar to that of non-reversible rules, at least after a few steps.

Indeed, as the pictures on the next page demonstrate, such systems can produce considerable randomness even when starting from very simple initial conditions.

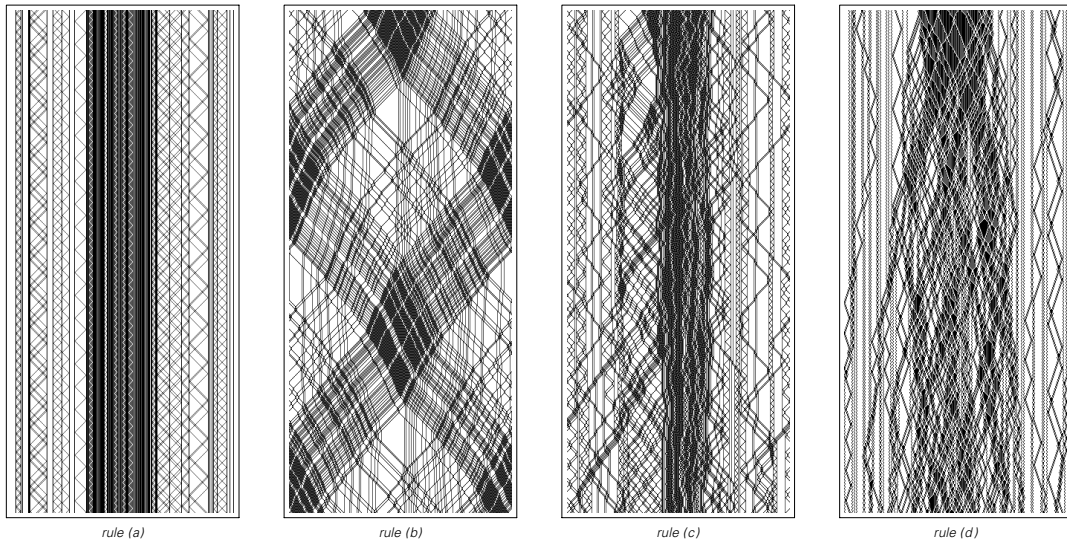


The behavior of rules (c) and (d) from the previous page, starting with very simple initial conditions. Each panel shows 500 steps of evolution, and rapid randomization is evident. The black and gray cells behave much like physical particles: their total number is conserved, and with the particular rules used here, their interactions are reversible. Note that the presence of boundaries is crucial; for without them there would in a sense be no collisions between particles, and the behavior of both systems would be rather trivial.

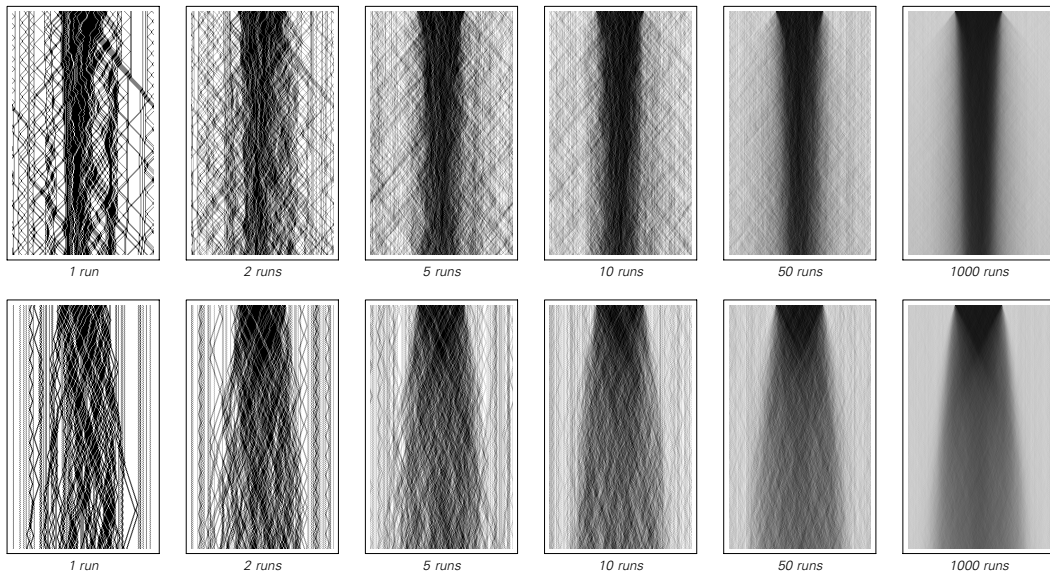
But there is still an important constraint on the behavior: even though black and gray cells may in effect move around randomly, their total number must always be conserved. And this means that if one looks at the total average density of colored cells throughout the system, it must always remain the same. But local densities in different parts of the system need not—and in general they will change as colored cells flow in and out.

The pictures below show what happens with four different rules, starting with higher density in the middle and lower density on the sides. With rules (a) and (b), each different region effectively remains separated forever. But with rules (c) and (d) the regions gradually mix.

As in many kinds of systems, the details of the initial arrangement of cells will normally have an effect on the details of the behavior that occurs. But what the pictures below suggest is that if one looks only at the overall distribution of density, then these details will become largely irrelevant—so that a given initial distribution of density will always tend to evolve in the same overall way, regardless of what particular arrangement of cells happened to make up that distribution.



The block cellular automata from previous pages started from initial conditions containing regions of different density. In rules (a) and (b) the regions remain separated forever, but in rules (c) and (d) they gradually diffuse into each other.



The evolution of overall density for block cellular automata (c) and (d) from the previous page. Even though at an underlying level these systems consist of discrete cells, their overall behavior seems smooth and continuous. The results shown here are obtained by averaging over progressively larger numbers of runs with initial conditions that differ in detail, but have the same overall density distribution. In the limit of an infinite number of runs (or infinite number of cells), the behavior in the second case approaches the form implied by the continuum diffusion equation. (In the first case correlations in effect last too long to yield exactly such behavior.)

The pictures above then show how the average density evolves in systems (c) and (d). And what is striking is that even though at the lowest level both of these systems consist of discrete cells, the overall distribution of density that emerges in both cases shows smooth continuous behavior.

And much as in physical systems like fluids, what ultimately leads to this is the presence of small-scale apparent randomness that washes out details of individual cells or molecules—as well as of conserved quantities that force certain overall features not to change too quickly. And in fact, given just these properties it turns out that essentially the same overall continuum behavior always tends to be obtained.

One might have thought that continuum behavior would somehow rely on special features of actual systems in physics. But in fact what we have seen here is that once again the fundamental mechanisms responsible already occur in a much more minimal way in programs that have some remarkably simple underlying rules.

Ultimate Models for the Universe

The history of physics has seen the development of a sequence of progressively more accurate models for the universe—from classical mechanics, through quantum mechanics, to quantum field theory, and beyond. And one may wonder whether this process will go on forever, or whether at some point it will come to an end, and one will reach a final ultimate model for the universe.

Experience with actual results in physics would probably not make one think so. For it has seemed that whenever one tries to get to another level of accuracy, one encounters more complex phenomena. And at least with traditional scientific intuition, this fact suggests that models of progressively greater complexity will be needed.

But one of the crucial points discovered in this book is that more complex phenomena do not always require more complex models. And indeed I have shown that even models based on remarkably simple programs can produce behavior that is in a sense arbitrarily complex.

So could this be what happens in the universe? And could it even be that underneath all the complex phenomena we see in physics there lies some simple program which, if run for long enough, would reproduce our universe in every detail?

The discovery of such a program would certainly be an exciting event—as well as a dramatic endorsement for the new kind of science that I have developed in this book.

For among other things, with such a program one would finally have a model of nature that was not in any sense an approximation or idealization. Instead, it would be a complete and precise representation of the actual operation of the universe—but all reduced to readily stated rules.

In a sense, the existence of such a program would be the ultimate validation of the idea that human thought can comprehend the construction of the universe. But just knowing the underlying program does not mean that one can immediately deduce every aspect of how the universe will behave. For as we have seen many times in this book, there is often a great distance between underlying rules and overall

behavior. And in fact, this is precisely why it is conceivable that a simple program could reproduce all the complexity we see in physics.

Given a particular underlying program, it is always in principle possible to work out what it will do just by running it. But for the whole universe, doing this kind of explicit simulation is almost by definition out of the question. So how then can one even expect to tell whether a particular program is a correct model for the universe? Small-scale simulation will certainly be possible. And I expect that by combining this with a certain amount of perhaps fairly sophisticated mathematical and logical deduction, it will be possible to get at least as far as reproducing the known laws of physics—and thus of determining whether a particular model has the potential to be correct.

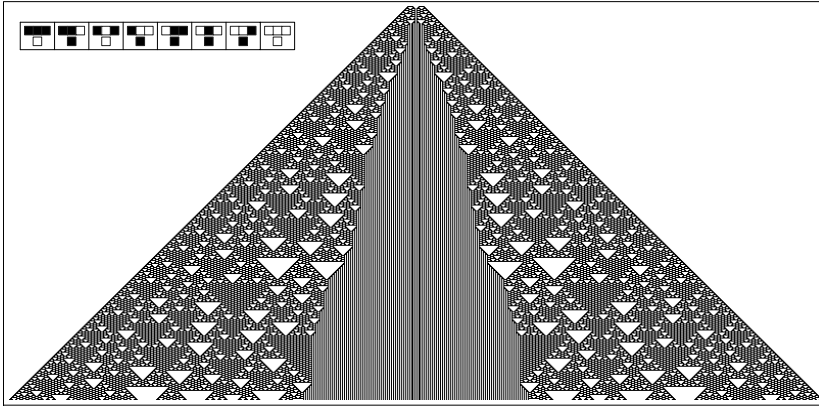
So if there is indeed a definite ultimate model for the universe, how might one set about finding it? For those familiar with existing science, there is at first a tremendous tendency to try to work backwards from the known laws of physics, and in essence to try to “engineer” a universe that will have particular features that we observe.

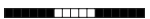
But if there is in fact an ultimate model that is quite simple, then from what we have seen in this book, I strongly believe that such an approach will never realistically be successful. For human thinking—even supplemented by the most sophisticated ideas of current mathematics and logic—is far from being able to do what is needed.

Imagine for example trying to work backwards from a knowledge of the overall features of the picture on the facing page to construct a rule that would reproduce it. With great effort one might perhaps come up with some immensely complex rule that would work in most cases. But there is no serious possibility that starting from overall features one would ever arrive at the extremely simple rule that was actually used.

It is already difficult enough to work out from an underlying rule what behavior it will produce. But to invert this in any systematic way is probably even in principle beyond what any realistic computation can do.

So how then could one ever expect to find the underlying rule in such a case? Almost always, it seems that the best strategy is a simple one: to come up with an appropriate general class of rules, and then just



A typical example of a situation where it would be very difficult to deduce the underlying rule from a description of the overall behavior that it produces. There is in a sense too great a distance between the simple rule shown and the behavior that emerges from it. I suspect that the same will be true of the basic rule for the universe. The particular rule shown here is the elementary cellular automaton with rule number 94, and with initial condition .

to search through these rules, trying each one in turn, and looking to see if it produces the behavior one wants.

But what about the rules for the universe? Surely we cannot simply search through possible rules of certain kinds, looking for one whose behavior happens to fit what we see in physics?

With the intuition of traditional science, such an approach seems absurd. But the point is that if the rule for the universe is sufficiently simple—and the results of this book suggest that it might be—then it becomes not so unreasonable to imagine systematically searching for it.

To start performing such a search, however, one first needs to work out what kinds of rules to consider. And my suspicion is that none of the specific types of rules that we have discussed so far in this book will turn out to be adequate. For I believe that all these types of rules in some sense probably already have too much structure built in.

Thus, for example, cellular automata probably already have too rigid a built-in notion of space. For a defining feature of cellular automata is that their cells are always arranged in a rigid array in space. Yet I strongly suspect that in the underlying rule for our universe there will be no such built-in structure. Rather, as I discuss in the sections

that follow, my guess is that at the lowest level there will just be certain patterns of connectivity that tend to exist, and that space as we know it will then emerge from these patterns as a kind of large-scale limit.

And indeed in general what I expect is that remarkably few familiar features of our universe will actually be reflected in any direct way in its ultimate underlying rule. For if all these features were somehow explicitly and separately included, the rule would necessarily have to be very complicated to fit them all in.

So if the rule is indeed simple, it almost inevitably follows that we will not be able to recognize directly in it most features of the universe as we normally perceive them. And this means that the rule—or at least its behavior—will necessarily seem to us unfamiliar and abstract.

Most likely for example there will be no easy way to visualize what the rule does by looking at a collection of elements laid out in space. Nor will there probably be any immediate trace of even such basic phenomena as motion.

But despite the lack of these familiar features, I still expect that the actual rule itself will not be too difficult for us to represent. For I am fairly certain that the kinds of logical and computational constructs that we have discussed in this book will be general enough to cover what is needed. And indeed my guess is that in terms of the kinds of pictures—or *Mathematica* programs—that we have used in this book, the ultimate rule for the universe will turn out to look quite simple.

No doubt there will be many different possible formulations—some quite unrecognizably different from others. And no doubt a formulation will eventually be found in which the rule somehow comes to seem quite obvious and inevitable.

But I believe that it will be essentially impossible to find such a formulation without already knowing the rule. And as a result, my guess is that the only realistic way to find the rule in the first place will be to start from some very straightforward representation, and then just to search through large numbers of possible rules in this representation.

Presumably the vast majority of rules will lead to utterly unworkable universes, in which there is for example no reasonable notion of space or no reasonable notion of time.

But my guess is that among appropriate classes of rules there will actually be quite a large number that lead to universes which share at least some features with our own. Much as the same laws of continuum fluid mechanics can emerge in systems with different underlying rules for molecular interactions, so also I suspect that properties such as the existence of seemingly continuous space, as well as certain features of gravitation and quantum mechanics, will emerge with many different possible underlying rules for the universe.

But my guess is that when it comes to something like the spectrum of masses of elementary particles—or perhaps even the overall dimensionality of space—such properties will be quite specific to particular underlying rules.

In traditional approaches to modelling, one usually tries first to reproduce some features of a system, then goes on to reproduce others. But if the ultimate rule for the universe is at all simple, then it follows that every part of this rule must in a sense be responsible for a great many different features of the universe. And as a result, it is not likely to be possible to adjust individual parts of the rule without having an effect on a whole collection of disparate features of the universe.

So this means that one cannot reasonably expect to use some kind of incremental procedure to find the ultimate rule for the universe. But it also means that if one once discovers a rule that reproduces sufficiently many features of the universe, then it becomes extremely likely that this rule is indeed the final and correct one for the whole universe.

And I strongly suspect that even in many of the most basic everyday physical processes, every element of the underlying rule for the universe will be very extensively exercised. And as a result, if these basic processes are reproduced correctly, then I believe that one can have considerable confidence that one in fact has the complete rule for the universe.

Looking at the history of physics, one might think that it would be completely inadequate just to reproduce everyday physical processes. For one might expect that there would always be some other esoteric phenomenon, say in particle physics, that would be discovered and would show that whatever rule one has found is somehow incomplete.

But I do not think so. For if the rule for our universe is at all simple, then I expect that to introduce a new phenomenon, however esoteric, will involve modifying some basic part of the rule, which will also affect even common everyday phenomena.

But why should we believe that the rule for our universe is in fact simple? Certainly among all possible rules of a particular kind only a limited number can ever be considered simple, and these rules are by definition somehow special. Yet looking at the history of science, one might expect that in the end there would turn out to be nothing special about the rule for our universe—just as there has turned out to be nothing special about our position in the solar system or the galaxy.

Indeed, one might assume that there are in fact an infinite number of universes, each with a different rule, and that we simply live in a particular—and essentially arbitrary—one of them.

It is unlikely to be possible to show for certain that such a theory is not correct. But one of its consequences is that it gives us no reason to think that the rule for our particular universe should be in any way simple. For among all possible rules, the overwhelming majority will not be simple; in fact, they will instead tend to be almost infinitely complex.

Yet we know, I think, that the rule for our universe is not too complex. For if the number of different parts of the rule were, for example, comparable to the number of different situations that have ever arisen in the history of the universe, then we would not expect ever to be able to describe the behavior of the universe using only a limited number of physical laws.

And in fact if one looks at present-day physics, there are not only a limited number of physical laws, but also the individual laws often seem to have the simplest forms out of various alternatives. And knowing this, one might be led to believe that for some reason the universe is set up to have the simplest rules throughout.

But, unfortunately perhaps, I do not think that this conclusion necessarily follows. For as I have discussed above, I strongly suspect that the vast majority of physical laws discovered so far are not truly fundamental, but are instead merely emergent features of the large-scale behavior of some ultimate underlying rule. And what this

means is that any simplicity observed in known physical laws may have little connection with simplicity in the underlying rule.

Indeed, it turns out that simple overall laws can emerge almost regardless of underlying rules. And thus, for example, essentially as a consequence of randomness generation, a wide range of cellular automata show the simple density diffusion law on page 464—whether or not their underlying rules happen to be simple.

So it could be that the laws that we have formulated in existing physics are simple not because of simplicity in an ultimate underlying rule, but rather because of some general property of emergent behavior for the kinds of overall features of the universe that we readily perceive.

Indeed, with this kind of argument, one could be led to think that there might be no single ultimate rule for the universe at all, but that instead there might somehow be an infinite sequence of levels of rules, with each level having a certain simplicity that becomes increasingly independent of the details of the levels below it.

But one should not imagine that such a setup would make it unnecessary to ask why our universe is the way it is: for even though certain features might be inevitable from the general properties of emergent behavior, there will, I believe, still be many seemingly arbitrary choices that have to be made in arriving at the universe in which we live. And once again, therefore, one will have to ask why it was these choices, and not others, that were made.

So perhaps in the end there is the least to explain if I am correct that the universe just follows a single, simple, underlying rule.

There will certainly be questions about why it is this particular rule, and not another one. And I am doubtful that such questions will ever have meaningful answers.

But to find the ultimate rule will be a major triumph for science, and a clear demonstration that at least in some direction, human thought has reached the edge of what is possible.

The Nature of Space

In the effort to develop an ultimate model for the universe, a crucial first step is to think about the nature of space—for inevitably it is in space that the processes in our universe occur.

Present-day physics almost always assumes that space is a perfect continuum, in which objects can be placed at absolutely any position. But one can certainly imagine that space could work very differently. And for example in a cellular automaton, space is not a continuum but instead consists just of discrete cells.

In our everyday experience space nevertheless appears to be continuous. But then so, for example, do fluids like air and water. And yet in the case of these fluids we know that at an underlying level they are composed of discrete molecules. And in fact over the course of the past century a great many aspects of the physical world that at first seemed continuous have in the end been discovered to be built up from discrete elements. And I very strongly suspect that this will also be true of space.

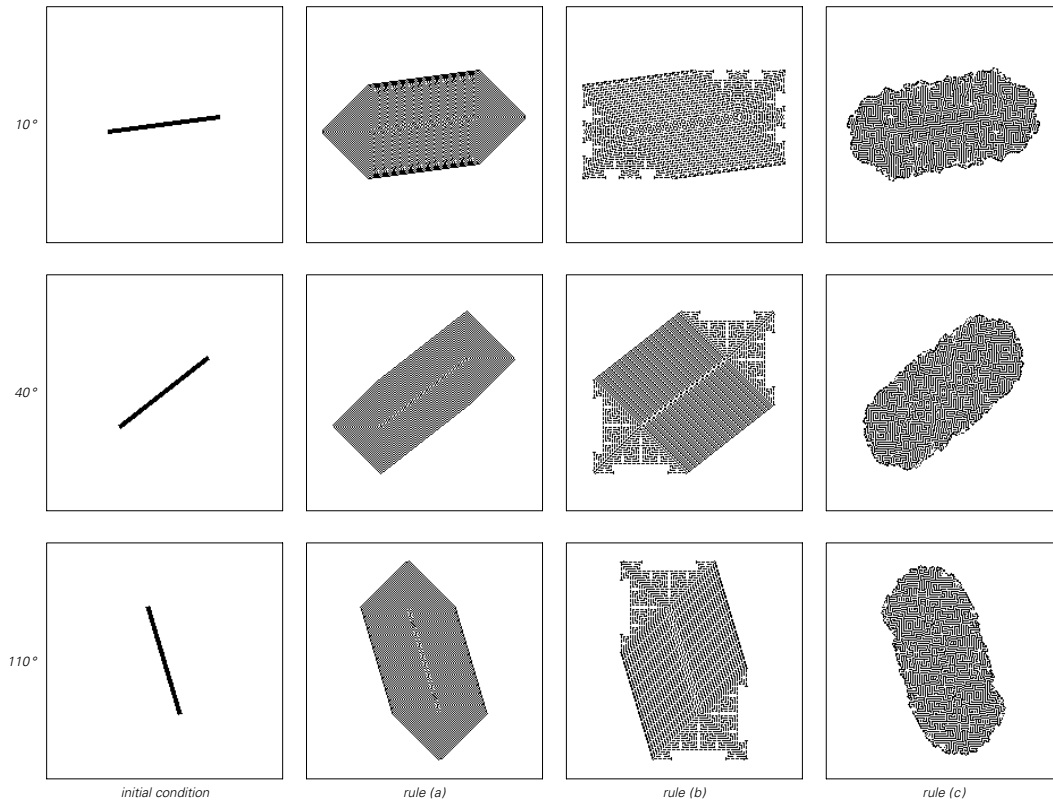
Particle physics experiments have shown that space acts as a continuum down to distances of around 10^{-20} meters—or a hundred thousandth the radius of a proton. But there is absolutely no reason to think that discrete elements will not be found at still smaller distances.

And indeed, in the past one of the main reasons that space has been assumed to be a perfect continuum is that this makes it easier to handle in the context of traditional mathematics. But when one thinks in terms of programs and the kinds of systems I have discussed in this book, it no longer seems nearly as attractive to assume that space is a perfect continuum.

So if space is not in fact a continuum, what might it be? Could it, for example, be a regular array of cells like in a cellular automaton?

At first, one might think that this would be completely inconsistent with everyday observations. For even though the individual cells in the array might be extremely small, one might still imagine that one would for example see all sorts of signs of the overall orientation of the array.

The pictures below show three different cellular automata, all set up on the same two-dimensional grid. And to see the effect of the grid, I show what happens when each of these cellular automata is started from blocks of black cells arranged at three different angles.



Examples of orientation dependence in the behavior of two-dimensional cellular automata on a fixed grid. Three different initial conditions, consisting of blocks at three different angles, are shown. For rules (a) and (b) the patterns produced always exhibit features that remain aligned with directions in the underlying grid. But with rule (c) essentially the same rounded pattern is obtained regardless of orientation. The rules shown here are outer totalistic: (a) 4-neighbor code 468, (b) 4-neighbor code 686 and (c) 8-neighbor code 746. In cases (a) and (b) 40 steps of evolution are used; in case (c) 100 steps are used.

In all cases the patterns produced follow at least to some extent the orientation of the initial block. But in cases (a) and (b) the effects of the underlying grid remain quite obvious—for the patterns produced always have facets aligned with the directions in this grid. But in case (c) the situation is different, and now the patterns produced turn out

always to have the same overall rounded form, essentially independent of their orientation with respect to the underlying grid.

And indeed what happens is similar to what we have seen many times in this book: the evolution of the cellular automaton generates enough randomness that the effects of the underlying grid tend to be washed out, with the result that the overall behavior produced ends up showing essentially no distinction between different directions in space.

So should one conclude from this that the universe is in fact a giant cellular automaton with rules like those of case (c)?

It is perhaps not impossible, but I very much doubt it.

For there are immediately simple issues like what one imagines happens at the edges of the cellular automaton array. But much more important is the fact that I do not believe in the distinction between space and its contents implied by the basic construction of a cellular automaton.

For when one builds a cellular automaton one is in a sense always first setting up an array of cells to represent space itself, and then only subsequently considering the contents of space, as represented by the arrangement of colors assigned to the cells in this array.

But if the ultimate model for the universe is to be as simple as possible, then it seems much more plausible that both space and its contents should somehow be made of the same stuff—so that in a sense space becomes the only thing in the universe.

Several times in the past ideas like this have been explored. And indeed the standard theory for gravity introduced in 1915 is precisely based on the notion that gravity can be viewed merely as a feature of space. But despite various attempts in the 1930s and more recently it has never seemed possible to extend this to cover the whole elaborate collection of forces and particles that we actually see in our universe.

Yet my suspicion is that a large part of the reason for this is just the assumption that space is a perfect continuum—described by traditional mathematics. For as we have seen many times in this book, if one looks at systems like programs with discrete elements then it immediately becomes much easier for highly complex behavior to emerge. And this is fundamentally what I believe is happening at the lowest level in space throughout our universe.

Space as a Network

In the last section I argued that if the ultimate model of physics is to be as simple as possible, then one should expect that all the features of our universe must at some level emerge purely from properties of space. But what should space be like if this is going to be the case?

The discussion in the section before last suggests that for the richest properties to emerge there should in a sense be as little rigid underlying structure built in as possible. And with this in mind I believe that what is by far the most likely is that at the lowest level space is in effect a giant network of nodes.

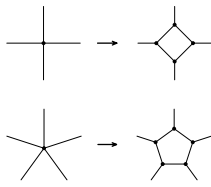
In an array of cells like in a cellular automaton each cell is always assigned some definite position. But in a network of nodes, the nodes are not intrinsically assigned any position. And indeed, the only thing that is defined about each node is what other nodes it is connected to.

Yet despite this rather abstract setup, we will see that with a sufficiently large number of nodes it is possible for the familiar properties of space to emerge—together with other phenomena seen in physics.

I already introduced in Chapter 5 a particular type of network in which each node has exactly two outgoing connections to other nodes, together with any number of incoming connections. The reason I chose this kind of network in Chapter 5 is that there happens to be a fairly easy way to set up evolution rules for such networks. But in trying to find an ultimate model of space, it seems best to start by considering networks that are somehow as simple as possible in basic structure—and it turns out that the networks of Chapter 5 are somewhat more complicated than is necessary.

For one thing, there is no need to distinguish between incoming and outgoing connections, or indeed to associate any direction with each connection. And in addition, nothing fundamental is lost by requiring that all the nodes in a network have exactly the same total number of connections to other nodes.

With two connections, only very trivial networks can ever be made. But if one uses three connections, a vast range of networks immediately become possible. One might think that one could get a

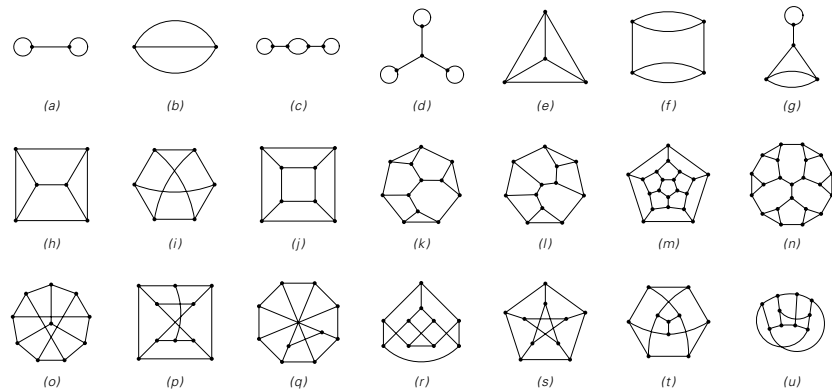


Examples of how nodes with more than three connections can be decomposed into collections of nodes with exactly three connections.

fundamentally larger range if one allowed, say, four or five connections rather than just three. But in fact one cannot, since any node with more than three connections can in effect always be broken into a collection of nodes with exactly three connections, as in the pictures on the left.

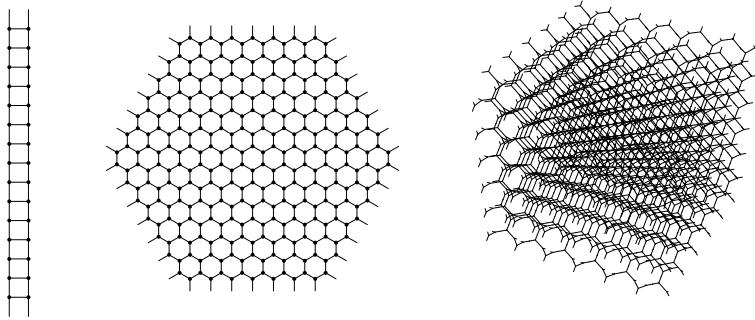
So what this means is that it is in a sense always sufficient to consider networks with exactly three connections at each node. And it is therefore these networks that I will use here in discussing fundamental models of space.

The pictures below show a few small examples of such networks. And already considerable diversity is evident. But none of the networks shown seem to have many properties familiar from ordinary space.



Examples of small networks with exactly three connections at each node. The first line shows all possible networks with up to four nodes. In what follows I consider only non-degenerate networks, in which there is at most one connection between any two nodes. Example (i) is the smallest network that cannot be drawn in two dimensions without lines crossing. Examples (k) and (l) are the smallest networks that have no symmetries between different nodes. Example (e) corresponds to the net of a tetrahedron, (j) to the net of a cube, and (m) to the net of a dodecahedron. Examples (o) through (u) show seven ways of drawing the same network, in this case the so-called Petersen network.

So how then can one get networks that correspond to ordinary space? The first step is to consider networks that have much larger numbers of nodes. And as examples of these, the pictures at the top of the facing page show networks that are specifically constructed to correspond to ordinary one-, two- and three-dimensional space.

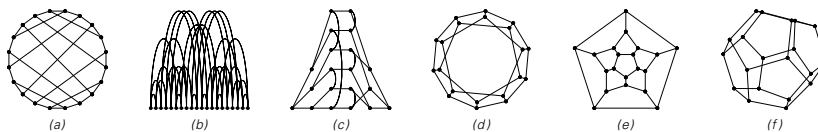


Examples of networks with three connections at each node that are effectively one, two and three-dimensional. These networks can be continued forever, and all have the property of being homogeneous, in the sense that every node has an environment identical to every other node.

Each of these networks is at the lowest level just a collection of nodes with certain connections. But the point is that the overall pattern of these connections is such that on a large scale there emerges a clear correspondence to ordinary space of a particular dimension.

The pictures above are drawn so as to make this correspondence obvious. But what if one was just presented with the raw pattern of connections for some network? How could one see whether the network could correspond to ordinary space of a particular dimension?

The pictures below illustrate the main difficulty: given only its pattern of connections, a particular network can be laid out in many completely different ways, most of which tell one very little about its potential correspondence with ordinary space.



Six different ways of laying out the same network. (a) nodes arranged around a circle; (b) nodes arranged along a line; (c) nodes arranged across the page according to distance from a particular node; (d) 2D layout with network and spatial distances as close as possible; (e) planar layout; (f) 3D layout.

So how then can one proceed? The fundamental idea is to look at properties of networks that can both readily be deduced from their pattern of connections and can also be identified, at least in some

large-scale limit, with properties of ordinary space. And the notion of distance is perhaps the most fundamental of such properties.

A simple way to define the distance between two points is to say that it is the length of the shortest path between them. And in ordinary space, this is normally calculated by subtracting the numerical coordinates of the positions of the points. But on a network things become more direct, and the distance between two nodes can be taken to be simply the minimum number of connections that one has to follow in order to get from one node to the other.

But can one tell just by looking at such distances whether a particular network corresponds to ordinary space of a certain dimension?

To a large extent one can. And a test is to see whether there is a way to lay out the nodes in the network in ordinary space so that the distances between nodes computed from their positions in space agree—at least in some approximation—with the distances computed directly by following connections in the network.

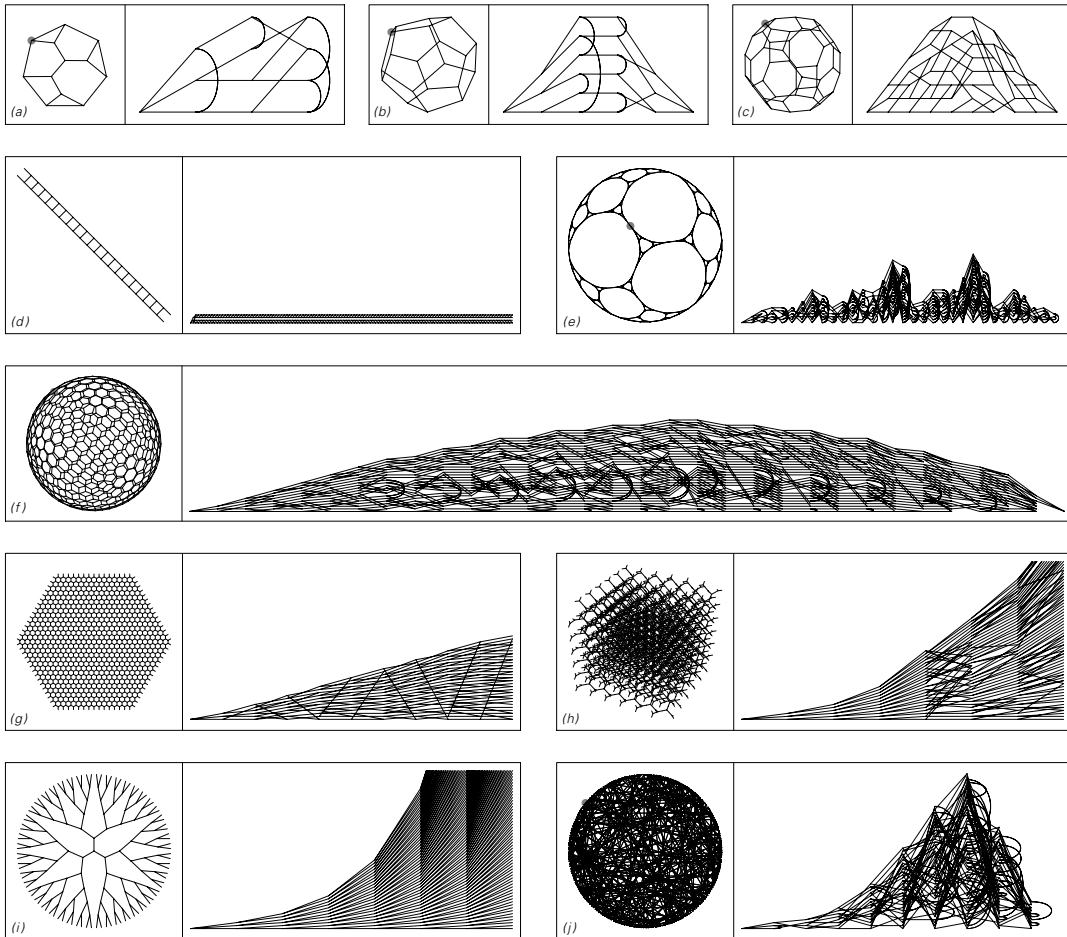
The three networks at the top of the previous page were laid out precisely so as to make this the case respectively for one, two and three-dimensional space. But why for example can the second network not be laid out equally well in one-dimensional rather than two-dimensional space? One way to see this is to count the number of nodes that appear at a given distance from a particular node in the network.

And for this specific network, the answer for this is very simple: at distance r there are exactly $3r$ nodes—so that the total number of nodes out to distance r grows like r^2 . But now if one tried to lay out all these nodes in one dimension it is inevitable that the network would have to bulge out in order to fit in all the nodes. And it turns out that it is uniquely in two dimensions that this particular network can be laid out in a regular way so that distances based on following connections in it agree with ordinary distances in space.

For the other two networks at the top of the previous page similar arguments can be given. And in fact in general the condition for a network to correspond to ordinary d -dimensional space is precisely that the total number of nodes that appear in it out to distance r grows in some limiting sense like r^d —a result analogous to the standard

mathematical fact that the area of a two-dimensional circle is πr^2 , while the volume of a three-dimensional sphere is $4/3 \pi r^3$, the volume of a four-dimensional hypersphere is $1/2 \pi^2 r^4$, and so on.

Below I show pictures of various networks. In each case the first picture is drawn to emphasize obvious regularities in the network. But the second picture is drawn in a more systematic way—by picking a specific starting node, and then laying out other nodes so that those at



Examples of various networks, shown first to emphasize their regularities, and second to illustrate the number of nodes reached by going successively more steps from a given node. For networks that in a limiting sense correspond to ordinary d -dimensional space, this number grows like r^{d-1} . All the larger networks shown are approximately uniform, in the sense that similar results are obtained starting from any node. Network (e) effectively has limiting dimension $\text{Log}[2, 3] \approx 1.58$.

successively greater network distances appear in successive columns across the page. And this setup has the feature that the height of column r gives the number of nodes that are at network distance r .

So by looking at how these heights grow across the page, one can see whether there is a correspondence with the r^{d-1} form that one expects for ordinary d -dimensional space. And indeed in case (g), for example, one sees exactly r^1 linear growth, reflecting dimension 2.

Similarly, in case (d) one sees r^0 growth, reflecting dimension 1, while in case (h) one sees r^2 growth, reflecting dimension 3.

Case (f) illustrates slightly more complicated behavior. The basic network in this case locally has an essentially two-dimensional form—but at large scales it is curved by being wrapped around a sphere. And what therefore happens is that for fairly small r one sees r^1 growth—reflecting the local two-dimensional form—but then for larger r there is slower growth, reflecting the presence of curvature.

Later in this chapter we will see how such curvature is related to the phenomenon of gravity. But for now the point is just that network (f) again behaves very much like ordinary space with a definite dimension.

So do all sufficiently large networks somehow correspond to ordinary space in a certain number of dimensions? The answer is definitely no. And as an example, network (i) from the previous page has a tree-like structure with 3^r nodes at distance r . But this number grows faster than r^d for any d —implying that the network has no correspondence to ordinary space in any finite number of dimensions.

If the connections in a network are chosen at random—as in case (j)—then again there will almost never be the kind of locality that is needed to get something that corresponds to ordinary finite-dimensional space.

So what might an actual network for space in our universe be like?

It will certainly not be as simple and regular as most of the networks on the previous page. For within its pattern of connections must be encoded everything we see in our universe.

And so at the level of individual connections, the network will most likely at first look quite random. But on a larger scale, it must be arranged so as to correspond to ordinary three-dimensional space. And somehow whatever rules update the network must preserve this feature.

The Relationship of Space and Time

To make an ultimate theory of physics one needs to understand the true nature not only of space but also of time. And I believe that here again the idea of thinking in terms of programs provides some crucial insights.

In our everyday experience space and time seem very different. For example, we can move from one point in space to another in more or less any way we choose. But we seem to be forced to progress through time in a very specific way. Yet despite such obvious apparent differences, almost all models in present-day fundamental physics have been built on the idea that space and time somehow work fundamentally the same.

But for most of the systems based on programs that I have discussed in this book this is certainly not true. And thus for example in a cellular automaton moving from one point in space to another just corresponds to shifting from one cell to another. But moving from one point in time to another involves actually applying the cellular automaton rule.

When we make a picture of the behavior of a cellular automaton, however, we do nevertheless tend to represent space and time in the same visual kind of way—with space going across the page and time going down. And in fact the basic notion of extending the idea of position in space to an idea of position in time has been common in scientific thought for more than five centuries.

But in the past century what has happened is that space and time have come to be thought of as being much more fundamentally similar. As we will discuss later in this chapter, the main origin of this is that in relativity theory certain aspects of space and time seem to become interchangeable. And from this there emerged the idea of thinking in terms of a spacetime continuum in which time appears merely as a fourth dimension just like the three ordinary dimensions of space.

So while in a system like a cellular automaton one typically imagines that a new and separate state of the system is somehow produced at each step in time, present-day physics more tends to think of the complete history of the universe throughout time as being just a single structure laid out in the four dimensions of spacetime.

So what then might determine the form of this structure?

The laws of physics in effect provide a collection of constraints on the structure. And while these laws are traditionally stated in terms of sophisticated mathematical equations, their basic character is similar to the simple constraints on arrays of black and white cells that I discussed at the end of Chapter 5. But now instead of defining constraints just in space, the laws of physics can be thought of as defining constraints on what can happen in both space and time.

Just as for space, it is my strong belief that time is fundamentally discrete. And from the discussion of networks for space in the previous section, one might imagine that perhaps the whole history of the universe in spacetime could be represented by a giant four-dimensional network.

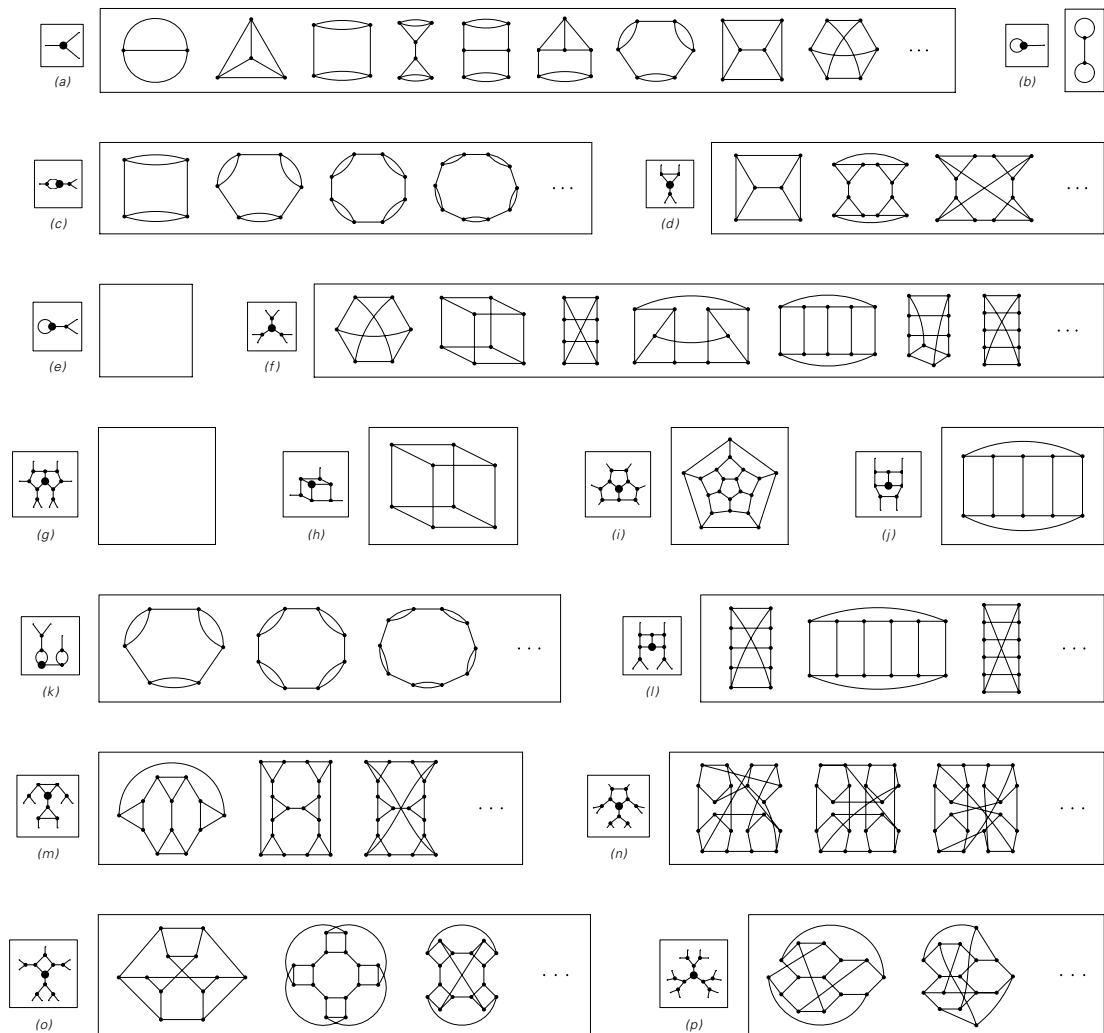
By analogy with the systems at the end of Chapter 5 a simple model would then be that this network is determined by the constraint that around every one of its nodes the overall arrangement of other nodes must match some particular template or set of templates.

Yet much as in Chapter 5 it turns out often not to be especially easy to find out which networks, if any, satisfy specific constraints of this kind. The pictures on the facing page nevertheless show results for quite a few choices of templates—where in each case the dangling connections in a template are taken to go to nodes that are not part of the template itself.

Pictures (a) and (b) show what happens with the two very simplest possible templates—involving just a single node. In case (a), all networks are allowed except for ones in which a node is connected directly to itself. In case (b), only the single network shown is allowed.

With templates that involve nodes out to distance one there are a total of 11 distinct non-trivial cases. And of these, 8 allow no complete networks to be formed, as in picture (e). But there turn out to be three cases—shown as pictures (c), (d) and (f)—in which complete networks can be formed, and in each of these one discovers that a fairly simple infinite set of networks are actually allowed.

In order to have a meaningful model for the universe, however, what must presumably happen is that essentially just one network can satisfy whatever constraints there are, and this one network must then represent all of the complex spacetime history of our universe.



Examples of networks determined by constraints. In each case the networks shown are required to satisfy the constraint that around every node their form must correspond to the template shown, in such a way that no dangling connections in the template are joined to each other. The pictures include all 14 templates that involve nodes out to distance at most two for which complete networks can be formed. In most cases where any such network can be formed, an infinite sequence of networks is allowed. But in cases (b), (h), (i) and (j) just a single network turns out to be allowed. The network constraint systems shown here are analogs of the two-dimensional systems based on constraints discussed at the end of Chapter 5.

So what does one find if one allows templates that include nodes out to distance two? There are a total of 690 distinct non-trivial such templates—and of these, 681 allow no complete networks to be formed, as in case (g). Six of the remaining templates then again allow an infinite sequence of networks. But there are three templates—shown as cases (h), (i) and (j)—that turn out to allow just single networks. These networks are however rather simple, and indeed the most complicated of them—case (i)—has just 20 nodes, and corresponds to a dodecahedron.

So are there in fact reasonably simple sets of constraints that in the end allow just one highly complex network, or perhaps a family of similar networks? I tend to doubt it. For our experience in Chapter 5 was that even in the much more rigid case of arrays of black and white squares, it was rather difficult to find constraints that would succeed in forcing anything but very simple patterns to occur.

So what does this mean for getting the kind of complexity that we see in our universe? We have not had difficulty in getting remarkable complexity from systems like cellular automata that we have discussed in this book. But such systems work not by being required to satisfy constraints, but instead by just repeatedly applying explicit rules.

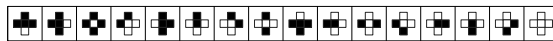
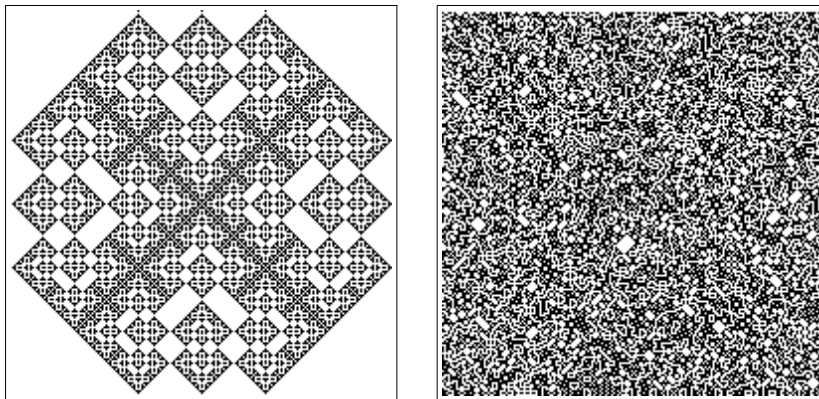
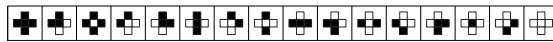
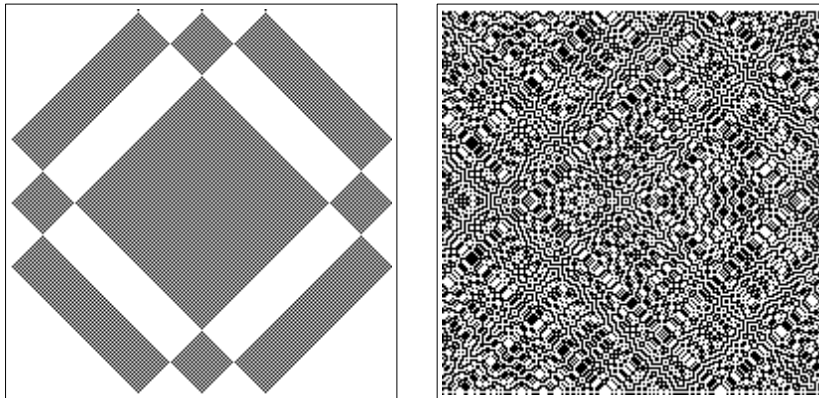
So is it in the end sensible to think of the universe as a single structure in spacetime whose form is determined by a set of constraints? Should we really imagine that the complete spacetime history of the universe somehow always exists, and that as time progresses, we are merely exploring different parts of it? Or should we instead think that the universe—more like systems such as cellular automata—explicitly evolves in time, so that at each moment a new state of the universe is in effect created, and the old one is lost?

Models based on traditional mathematical equations—in which space and time appear just as abstract symbolic variables—have never had to make much distinction between these two views. But in trying to understand the ultimate underlying mechanisms of the universe, I believe that one must inevitably distinguish between these views.

And I strongly believe that the second view is the one most likely to provide a meaningful underlying model for our universe. But while this view is closer to our everyday perception of time, it seems to

contradict the correspondence between space and time that is built into most of present-day physics. So one might wonder how then it could be consistent with experiments that have been done in physics?

One possibility, illustrated in the pictures below, is to have a system that evolves in time according to explicit rules, but for these rules to have built into them a symmetry between space and time.



Examples of one-dimensional cellular automata which exhibit a symmetry between space and time. Each picture can be generated by starting from initial conditions at the top, and then just evolving down the page repeatedly applying the cellular automaton rule. The particular rules shown are reversible second-order ones with numbers 90R and 150R.

But I very much doubt that any such obvious symmetry between space and time exists in the fundamental rules for our universe. And instead what I expect is much like we have seen many times before in this book: that even though at the lowest level there is no direct correspondence between space and time, such a correspondence nevertheless emerges when one looks in the appropriate way at larger scales of the kind probed by practical experiments.

As I will discuss in the next several sections, I suspect that for many purposes the history of the universe can in fact be represented by a certain kind of spacetime network. But the way this network is formed in effect treats space and time rather differently. And in particular—just as in a system like a cellular automaton—the network can be built up incrementally by starting with certain initial conditions and then applying appropriate underlying rules over and over again.

Any such rules can in principle be thought of as providing a set of constraints for the spacetime network. But the important point is that there is no need to do a separate search to find networks that satisfy such constraints—for the rules themselves instead immediately define a procedure for building up the necessary network.

Time and Causal Networks

I argued in the last section that the progress of time should be viewed at a fundamental level much like the evolution of a system like a cellular automaton. But one of the features of a cellular automaton is that it is set up to update all of its cells together, as if at each tick of some global clock. Yet just as it seems unreasonable to imagine that the universe consists of a rigid grid of cells in space, so also it seems unreasonable to imagine that there is a global clock which defines the updating of every element in the universe synchronized in time.

But what is the alternative? At first it may seem bizarre, but one possibility that I believe is ultimately not too far from correct is that the universe might work not like a cellular automaton in which all cells get updated at once, but instead like a mobile automaton or Turing machine, in which just a single cell gets updated at each step.

As discussed in Chapter 3—and illustrated in the picture on the right—a mobile automaton has just a single active cell which moves around from one step to the next. And because this active cell is the only one that ever gets updated, there is never any issue about synchronizing behavior of different elements at a given step.

Yet at first it might seem absurd to think that our universe could work like a mobile automaton. For certainly we do not notice any kind of active cell visiting different places in the universe in sequence. And indeed, to the contrary, our perception is that different parts of the universe seem to evolve in parallel and progress through time together.

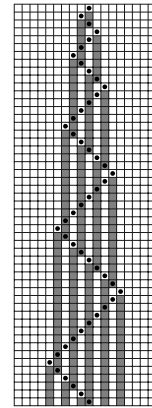
But it turns out that what one perceives as happening in a system like a mobile automaton can depend greatly on whether one is looking at the system from outside, or whether one is oneself somehow part of the system. For from the outside, one can readily see each individual step in the evolution of a mobile automaton, and one can tell that there is just a single active cell that visits different parts of the system in sequence. But to an observer who is actually part of the mobile automaton, the perception can be quite different.

For in order to recognize that time has passed, or indeed that anything has happened, the state of the observer must somehow change. But if the observer itself just consists of a collection of cells inside a mobile automaton, then no such change can occur except on steps when the active cell in the mobile automaton visits this collection of cells.

And what this means is that between any two successive moments of time as perceived by an observer inside the mobile automaton, there can be a great many steps of underlying mobile automaton evolution.

If an observer could tell what was happening on every step, then it would be easy to recognize the sequential way in which cells are updated. But because an observer who is part of a mobile automaton can in effect only occasionally tell what has happened, then as far as such an observer is concerned, many cells can appear to have been updated in parallel between successive moments of time.

To see in more detail how this works it could be that it would be necessary to make a specific model for the observer. But in fact, it turns out that it is sufficient just to look at the evolution of the mobile



A mobile automaton in which only the single active cell indicated by a dot is updated at each step, thereby avoiding the issue of global synchronization.

automaton not in terms of individual steps, but rather in terms of updating events and the causal relationships between them.

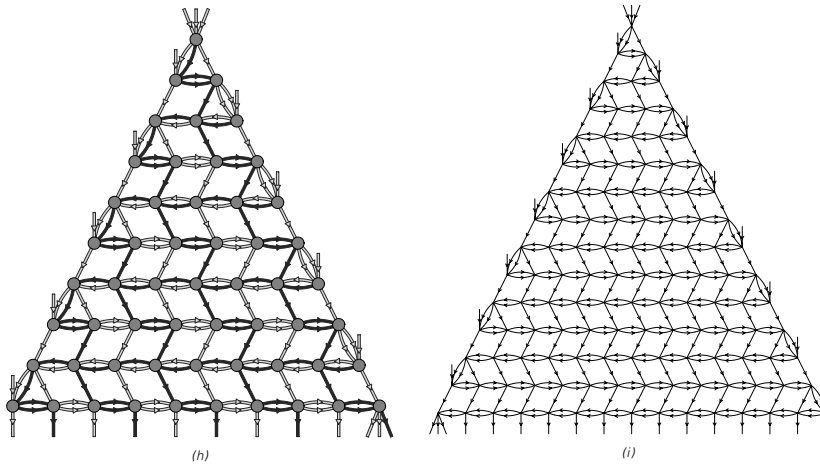
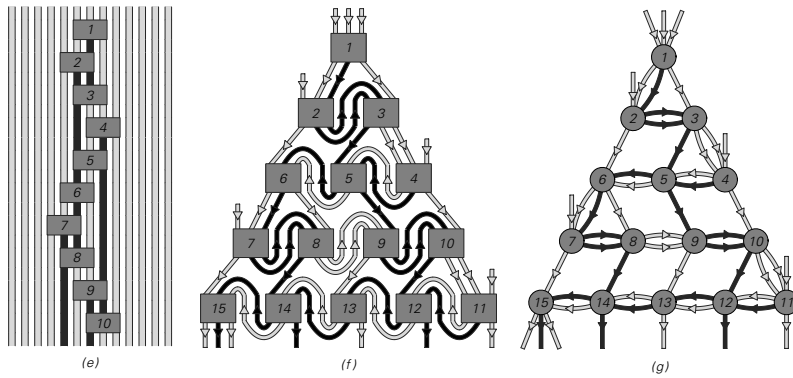
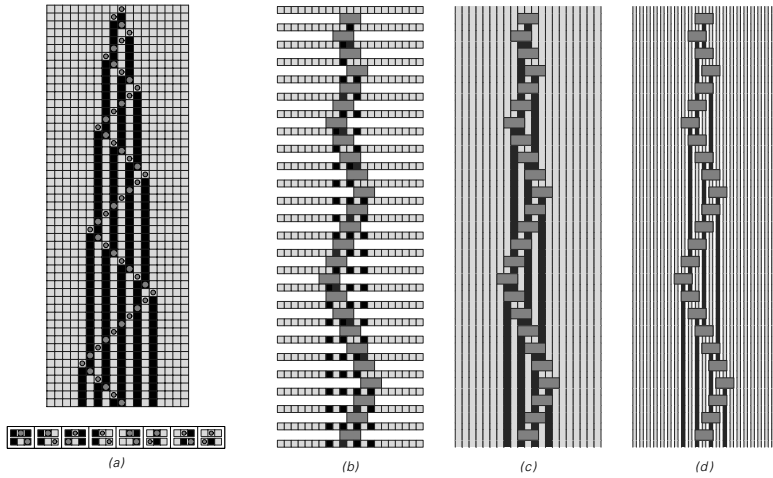
The pictures on the facing page show an example of how this works. Picture (a) is a version of the standard representation that I have used for mobile automaton evolution elsewhere in the book—in which successive lines give the colors of cells on successive steps, and the position of the active cell is indicated at each step by a gray dot. The subsequent pictures on the facing page all ultimately give essentially the same information, but gradually present it to emphasize more a representation in terms of updating events and causal relationships.

Picture (b) is very similar to (a), but shows successive steps of mobile automaton evolution separated, with gray blobs in between indicating “updating events” corresponding to each application of the underlying mobile automaton rule. Picture (b) still has a definite row of cells for each individual step of mobile automaton evolution. But in picture (c) cells not updated on a given step are merged together, yielding vertical stripes of color that extend from one updating event to another.

So what is the significance of these stripes? In essence they serve to carry the information needed to determine what the next updating event will be. And as picture (d) begins to emphasize, one can think of these stripes as indicating what causal relationships or connections exist between updating events.

And this notion then suggests a quite different representation for the whole evolution of the mobile automaton. For rather than having a picture based on successive individual steps of evolution, one can instead form a network of the various causal relationships between updating events, with each updating event being a node in this network, and each stripe being a connection from one node to another.

A sequence of views of the evolution of a mobile automaton, showing how a network of causal relationships between updating events can be created. This network provides a very simple model for spacetime in the universe. Picture (a) is essentially the standard representation of mobile automaton evolution that I have used in this book. Picture (b) includes gray blobs to indicate updating events. Picture (c) merges cells that are not being updated. Picture (d) emphasizes the role of vertical stripes as connections between updating events. Pictures (e) through (g) show how a network can be formed with nodes corresponding to updating events. Pictures (h) and (i) demonstrate that with the particular underlying rule used here, a highly regular network is produced. ▶



Picture (e) shows the updating events and stripes from the top of picture (d), with the updating events now explicitly numbered. Pictures (f) and (g) then show how one can take the pattern of connectivity from picture (e) and lay out the updating events as nodes so as to produce an orderly network. And for the particular mobile automaton rule used here, the network one gets ends up being highly regular, as illustrated in pictures (h) and (i).

So what is the significance of this network? It turns out that it can be thought of as defining a structure for spacetime as perceived by an observer inside the mobile automaton—in much the same way as the networks we discussed two sections ago could be thought of as defining a structure for space. Each updating event, corresponding to each node in the network, can be imagined to take place at some point in spacetime. And the connections between nodes in the network can then be thought of as defining the pattern of neighbors for points in spacetime.

But unlike in the space networks that we discussed two sections ago, the connections in the causal networks we consider here always go only one way: each connection corresponds to a causal relationship in which one event leads to another, but not the other way around.

This kind of directionality, however, is exactly what is needed if a meaningful notion of time is to emerge. For the progress of time can be defined by saying that only those events that occur later in time than a particular event can be affected by that event.

And indeed the networks in pictures (g) through (i) on the previous page were specifically laid out so that successive rows of nodes going down the page would correspond, at least roughly, to events occurring at successively later times.

As the numbering in pictures (e) through (g) illustrates, there is no direct correspondence between this notion of time and the sequence of updating events that occur in the underlying evolution of the mobile automaton. For the point is that an observer who is part of the mobile automaton will never see all the individual steps in this evolution. The most they will be able to tell is that a certain network of causal relationships exists—and their perception of time must therefore derive purely from the properties of this network.

So does the notion of time that emerges actually have the familiar features of time as we know it? One might think for example that in a network there could be loops that would lead to a deviation from the linear progression of time that we appear to experience. But in fact, with a causal network constructed from an underlying evolution process in the way we have done it here no such loops can ever occur.

So what about traces of the sequential character of evolution in the original mobile automaton? One might imagine that with only a single active cell being updated at each step different parts of the system would inevitably be perceived to progress through time one after another. But what the pictures on page 489 demonstrate is that this need not be the case. Indeed, in the networks shown there all the nodes on each row are in effect connected in parallel to the nodes on the row below. So even though the underlying rules for the mobile automaton involve no global synchronization, it is nevertheless possible for an observer inside the mobile automaton to perceive time as progressing in a synchronized way.

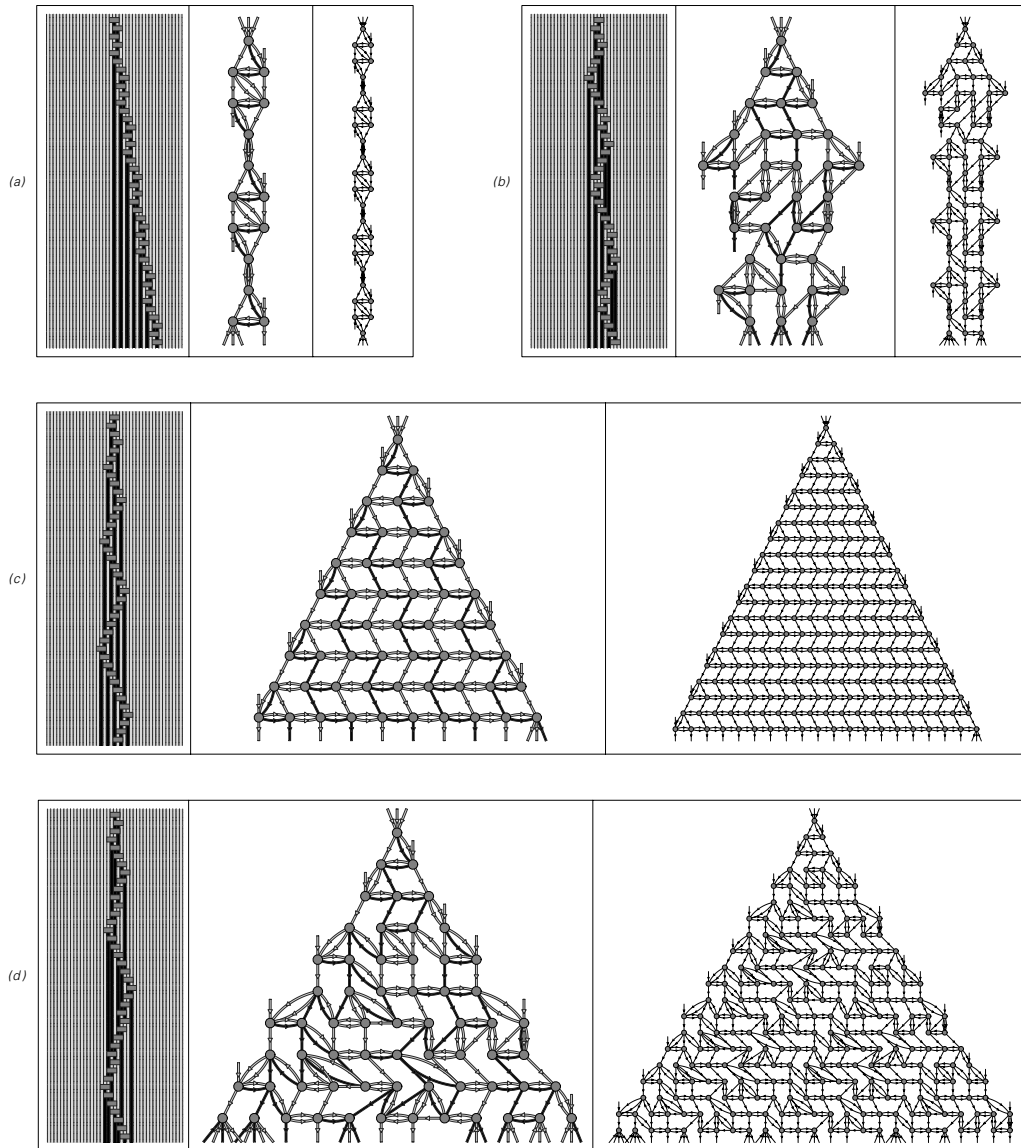
Later in this chapter I will discuss how space works in the context of causal networks—and how ideas of relativity theory emerge. But for now one can just think of networks like those on page 489 as being laid out so that time goes down the page and space goes across. And one can then see that if one follows connections in the network, one is always forced to go progressively down the page, even though one is able to move both backwards and forwards across the page—thus agreeing with our everyday experience of being able to move in more or less any direction in space, but always being forced to move onward in time.

So what happens with other mobile automata?

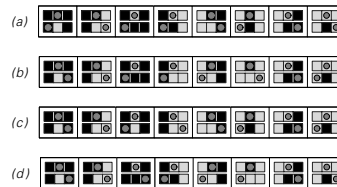
The pictures on the next two pages show a few examples.

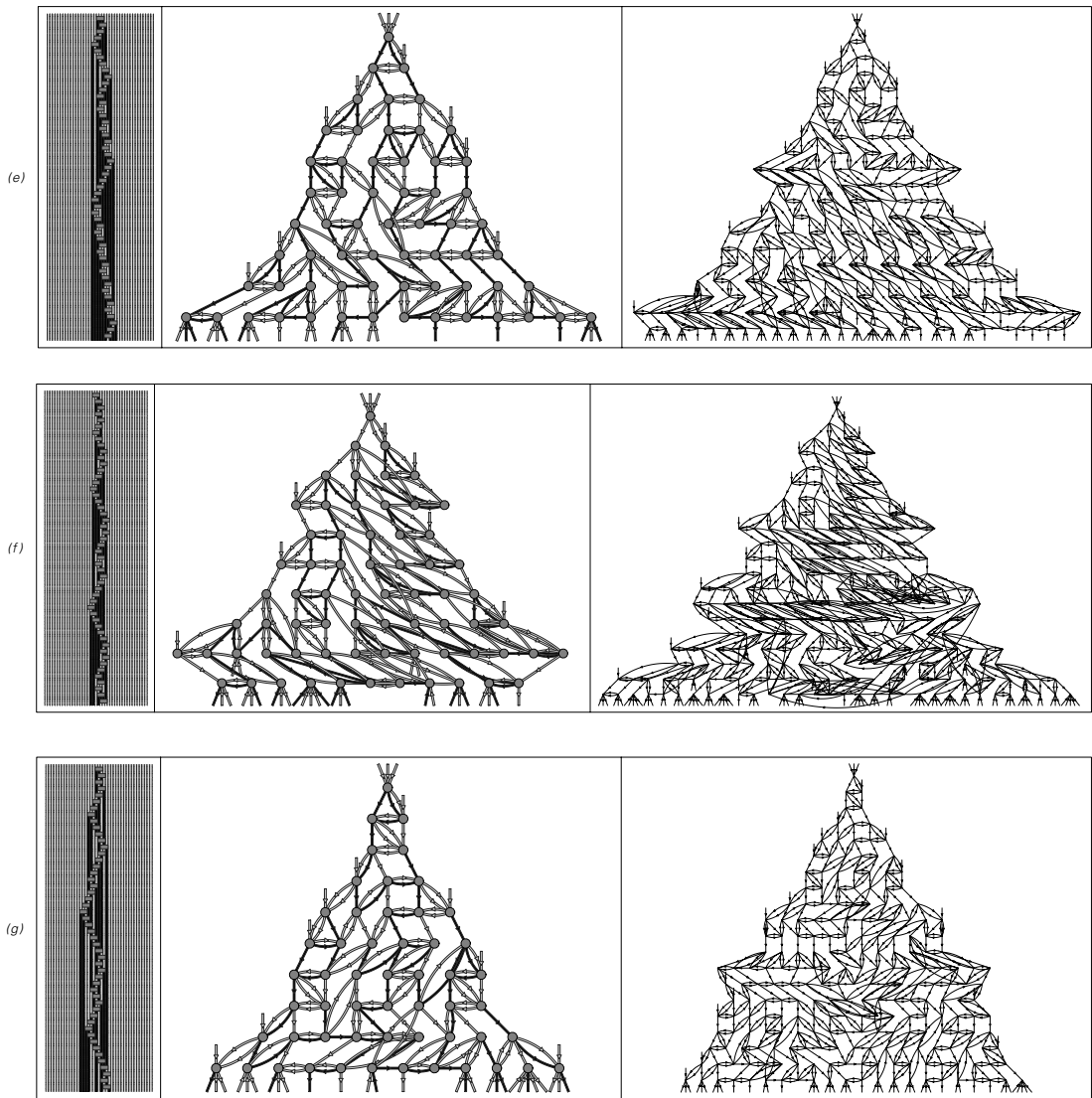
Rules (a) and (b) yield very simple repetitive networks in which there is in effect a notion of time but not of space. The underlying way any mobile automaton works forces time to continue forever. But with rules (a) and (b) only a limited number of points in space can ever be reached.

The other rules shown do not, however, suffer from this problem: in all of them progressively more points are reached in space as time goes on. Rules (c) and (d) yield networks that can be laid out in a quite

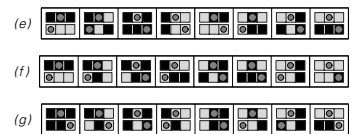


Examples of mobile automata from Chapter 3 and the causal networks they generate. In each case the picture on the left is essentially the standard representation of mobile automaton evolution used in Chapter 3. The pictures on the right are then causal network representations of the same evolution. The networks are laid out in analogy to the space networks on page 479, with nodes being placed on successive rows if they take progressively more connections to reach from the top node.





Note that a single connection can join events that occur at very different steps in the evolution of the underlying mobile automaton. And indeed to construct even a small part of the causal network can require an arbitrarily long computation in the underlying mobile automaton. Thus for example to make the causal networks in pictures (e), (f) and (g) requires looking respectively at 2447, 731 and 322 steps of mobile automaton evolution. And indeed in some cases there can be connections that are in effect never resolved. And thus for example in picture (a) there are downward connections that never reach any other node—reflecting the presence of positions on the left in the mobile automata evolution to which the active cell never returns.



regular manner. But with rules (e), (f) and (g) the networks are more complicated, and begin to seem somewhat random.

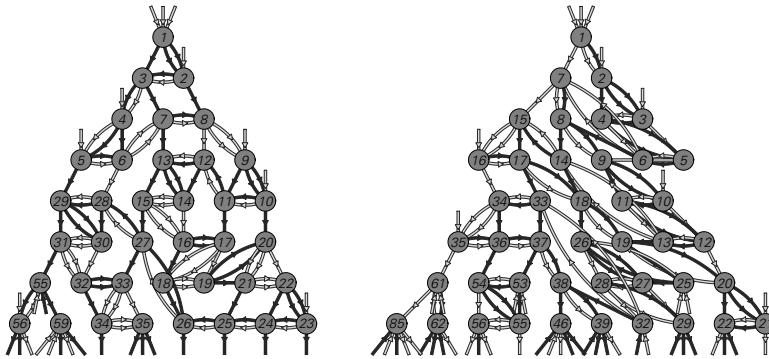
The procedure that is used to lay out the networks on the previous two pages is a direct analog of the procedure used for space networks on page 479: the row in which a particular node will be placed is determined by the minimum number of connections that have to be followed in order to reach that node starting from the node at the top.

In cases (a) and (c) the networks obtained in this way have the property that all connections between nodes go either across or down the page. But in every other case shown, at least some connections also go up the page. So what does this mean for our notion of time? As mentioned earlier, there can never be a loop in any causal network that comes from an evolution process. But if one identifies time with position down the page, the presence of connections that go up as well as down the page implies that in some sense time does not always progress in the same direction. Yet at least in the cases shown here there is still a strong average flow down the page—agreeing with our everyday perception that time progresses only in one direction.

Like in so many other systems that we have studied in this book, the randomness that we find in causal networks will inevitably tend to wash out details of how the networks are constructed. And thus, for example, even though the underlying rules for a mobile automaton always treat space and time very differently, the causal networks that emerge nevertheless often exhibit a kind of uniform randomness in which space and time somehow work in many respects the same.

But despite this uniformity at the level of causal networks, the transformation from mobile automaton evolution to causal network is often far from uniform. And for example the pictures at the top of the facing page show the causal networks for rules (e) and (f) from the previous page—but now with each node numbered to specify the step of mobile automaton evolution from which it was derived.

And what we see is that even nodes that are close to the top of the causal network can correspond to events which occur after a large number of steps of mobile automaton evolution. Indeed, to fill in just twenty rows



Causal networks corresponding to rules (e) and (f) from page 493, with each node explicitly labelled to specify from which step of mobile automaton evolution it is derived. Even to fill in the first few rows of such causal networks, many steps of underlying mobile automaton evolution must be traced.

of the causal networks for rules (e) and (f) requires following the underlying mobile automaton evolution for 2447 and 731 steps respectively.

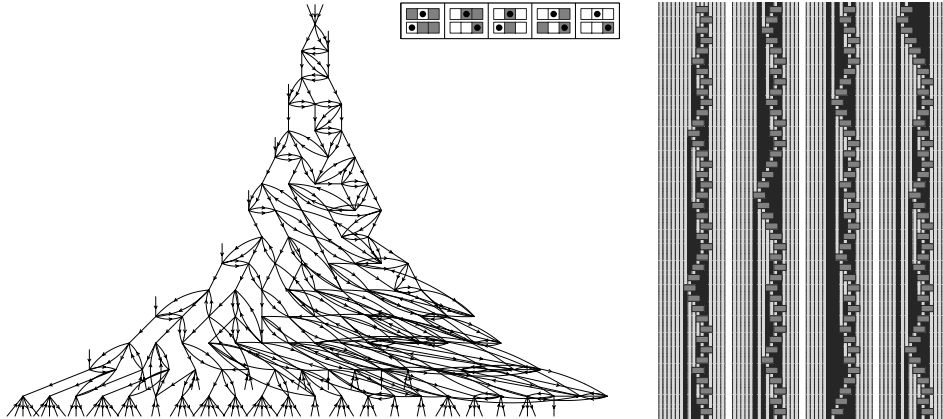
One feature of causal networks is that they tell one not only what the consequences of a particular event will be, but also in a sense what its causes were. Thus, for example, if one starts, say, with event 17 in the first causal network above, then to find out that its causes were events 11 and 16 one simply has to trace backwards along the connections which lead to it.

With the specific type of underlying mobile automaton used here, every node has exactly three incoming and three outgoing connections. And at least when there is overall apparent randomness, the networks that one gets by going forwards and backwards from a particular node will look very similar. In most cases there will still be small differences; but the causal network on the right above is specifically constructed to be exactly reversible—much like the cellular automata we discussed near the beginning of this chapter.

Looking at the causal networks we have seen so far, one may wonder to what extent their form depends on the particular properties of the underlying mobile automata that were used to produce them.

For example, one might think that the fact that all the networks we have seen so far grow at most linearly with time must be an inevitable consequence of the one-dimensional character of the mobile

automaton rules we have used. But the picture below demonstrates that even with such one-dimensional rules, it is actually possible to get causal networks that grow more rapidly. And in fact in the case shown below there are roughly a factor 1.22 more nodes on each successive row—corresponding to overall approximate exponential growth.



A one-dimensional mobile automaton which yields a causal network that in effect grows exponentially with time. The underlying mobile automaton acts like a binary counter, yielding a pattern whose width grows logarithmically with the number of steps. The three cases not shown in the rule are never used with the initial conditions given here.

The causal network for a system is always in some sense dual to the underlying evolution of the system. And in the case shown here the slow growth of the region visited by the active cell in the underlying evolution is reflected in rapid growth of the corresponding causal network.

As we will see later in this chapter there are in the end some limitations on the kinds of causal networks that one-dimensional mobile automata and systems like them can produce. But with different mobile automaton rules one can still already get tremendous diversity.

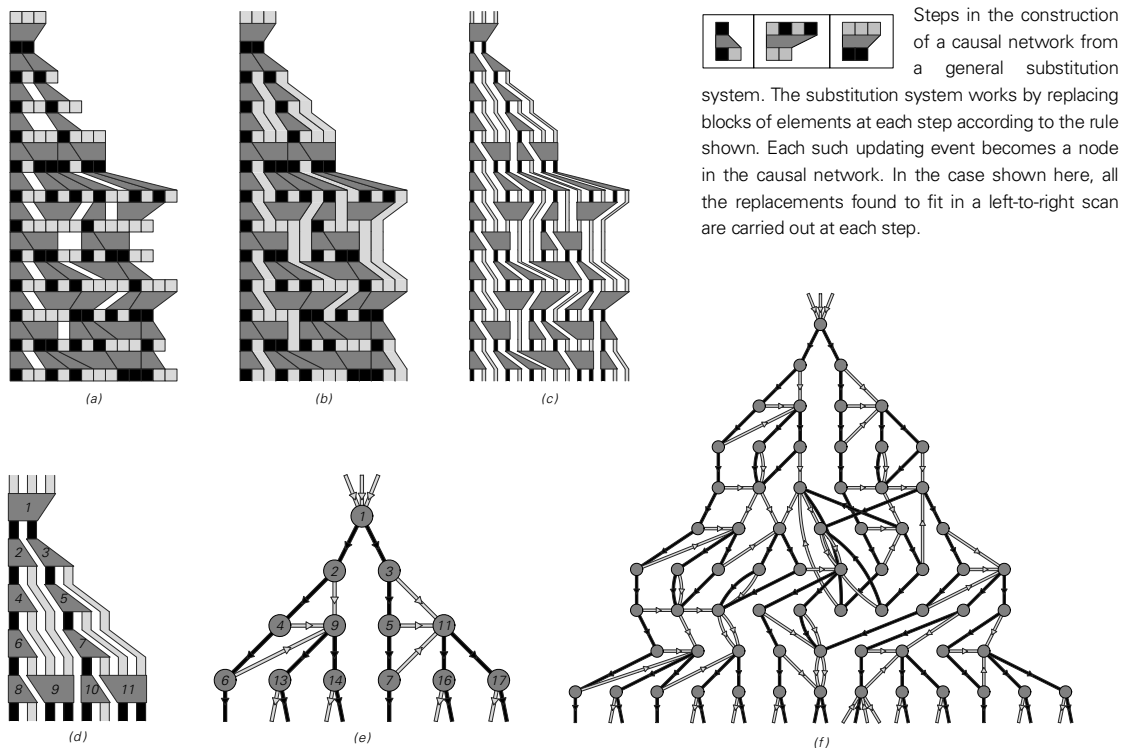
And even though when viewed from outside, systems like mobile automata might seem to have almost none of the familiar features of our universe, what we see is that if we as observers are in a sense part of such systems then immediately some major features quite similar to those of our universe can emerge.

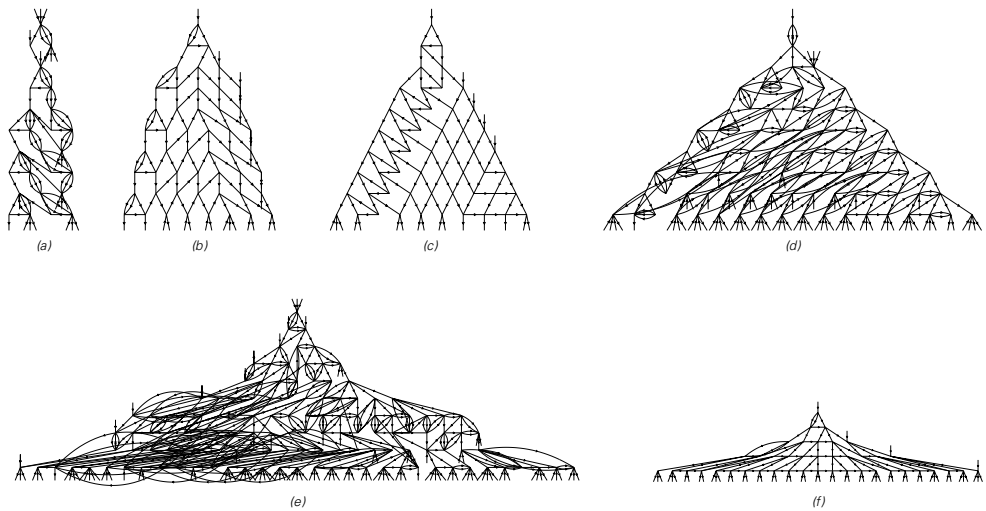
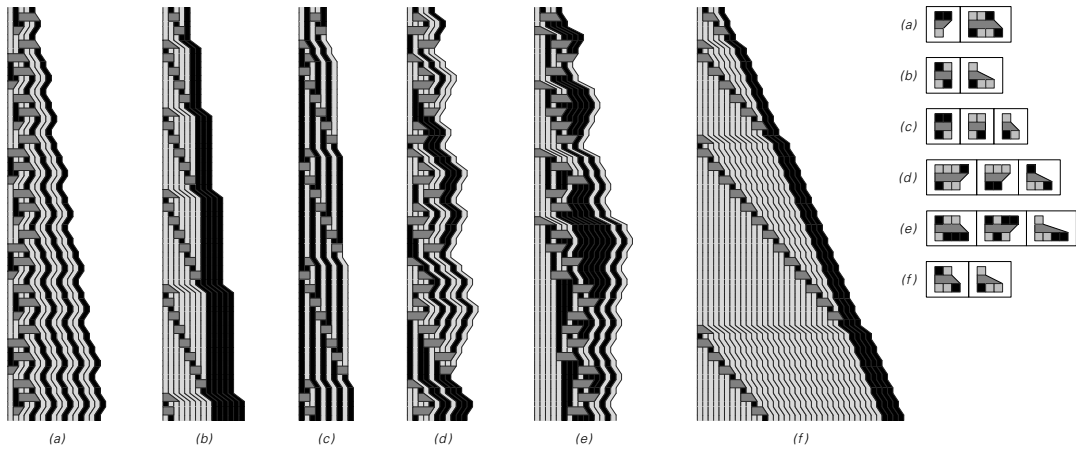
The Sequencing of Events in the Universe

In the last section I discussed one type of model in which familiar notions of time can emerge without any kind of built-in global clock. The particular models I used were based on mobile automata—in which the presence of a single active cell forces only one event ever to occur in the universe at once. But as we will see in this section, there is actually no need for the setup to be so rigid, or indeed for there to be any kind of construct like an active cell.

One can think of mobile automata as being special cases of substitution systems of the type I introduced in Chapter 3. Such systems in general take a string of elements and at each step replace blocks of elements with other elements according to some definite rule.

The picture below shows an example of one such system, and illustrates how—just like in a mobile automaton—relations between updating events can be represented by a causal network.





Examples of sequential substitution systems of the type discussed on page 88, and the causal networks that emerge from them. In a sequential substitution system only the first replacement that is found to apply in a left-to-right scan is ever performed at any step. Rule (a) above yields a causal network that is purely repetitive and thus yields no meaningful notion of space. Rules (b), (c) and (d) yield causal networks that in effect grow roughly linearly with time. In rule (f) the causal network grows exponentially, while in rule (e) the causal network also grows quite rapidly, though its overall growth properties are not clear. Note that to obtain the 10 levels shown here in the causal network for rule (e), it was necessary to follow the evolution of the underlying substitution system for a total of 258 steps.

Substitution systems that correspond to mobile automata can be thought of as having rules and initial conditions that are specially set up so that only one updating event can ever occur on any particular step. But with most rules—including the one shown on the previous page—there are usually several possible replacements that can be made at each step.

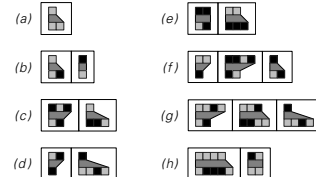
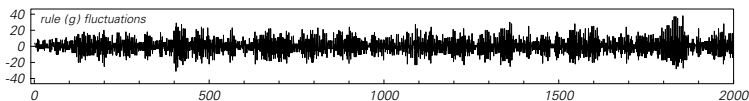
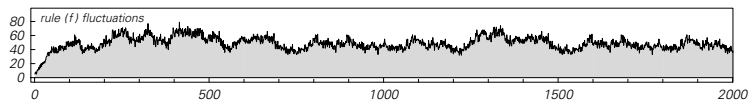
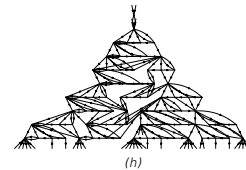
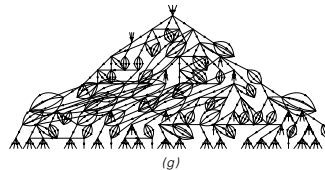
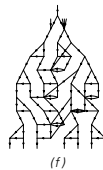
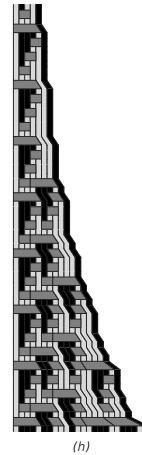
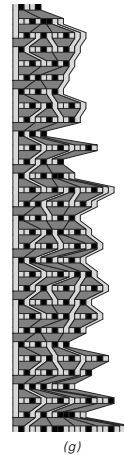
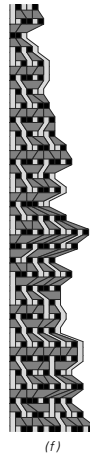
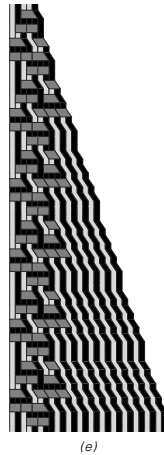
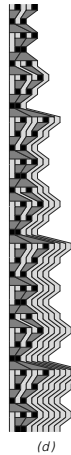
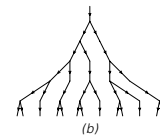
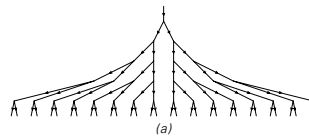
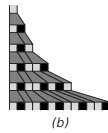
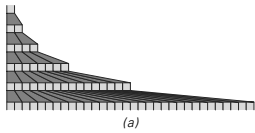
One scheme for deciding which replacement to make is just to scan the string from left to right and then pick the first replacement that applies. This scheme corresponds exactly to the sequential substitution systems we discussed in Chapter 3.

The pictures on the facing page show a few examples of what can happen. The behavior one gets is often fairly simple, but in some cases it can end up being highly complex. And just as in mobile automata, the causal networks that emerge typically in effect grow linearly with time. But, again as in mobile automata, there are rules such as (a) in which there is no growth—and effectively no notion of space. And there are also rules such as (f)—which turn out to be much more common in general substitution systems than in mobile automata—in which the causal network in effect grows exponentially with time.

But why do only one replacement at each step? The pictures on the next page show what happens if one again scans from left to right, but now one performs all replacements that fit, rather than just the first one.

In the case of rules (a) and (b) the result is to update every single element at every step. But since the replacements in these particular rules involve only one element at a time, one in effect has a neighbor-independent substitution system of the kind we discussed on page 82. And as we discovered there, such systems can only ever produce rather simple behavior: each element repeatedly branches into several others, yielding a causal network that has the form of a regular tree.

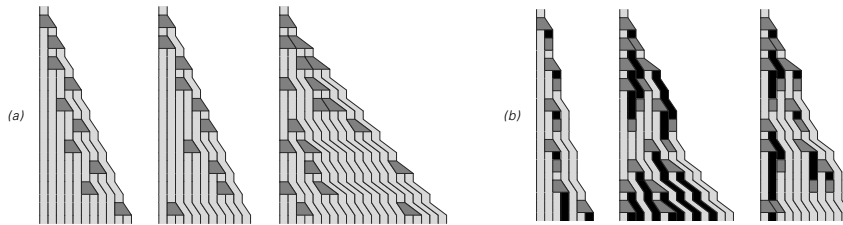
So what happens with replacements that involve more than just one element? In many cases, the behavior is still quite simple. But as several of the pictures on the next page demonstrate, fairly simple rules are sufficient—as in so many other systems that we have discussed in this book—to obtain highly complex behavior.



Examples of general substitution systems and the causal networks that emerge from them. In the pictures shown here, every replacement that is found to fit in a left-to-right scan is performed at each step. Rules (a) and (b) act like neighbor-independent substitution systems of the type discussed on page 84, and yield exponentially growing tree-like causal networks. The plots at the bottom show the growth rates of the patterns produced by rules (f) and (g). In the case of rule (f) the pattern turns out to be repetitive, with a period of 796 steps.

One may wonder, however, to what extent the behavior one sees depends on the exact scheme that one uses to pick which replacements to apply at each step. The answer is that for the vast majority of rules—including rules (c) through (g) in the picture on the facing page—using different schemes yields quite different behavior—and a quite different causal network.

But remarkably enough there do exist rules for which exactly the same causal network is obtained regardless of what scheme is used. And as it turns out, rules (a) and (b) from the picture on the facing page provide simple examples of this phenomenon, as illustrated in the pictures below.

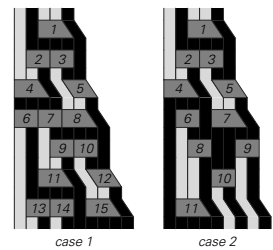


The behavior of rules (a) and (b) from the facing page when replacements are performed at random. Even though the detailed patterns obtained are different, the causal networks in these particular rules that represent relationships between replacement events are always exactly the same.

For each rule, the three different pictures shown above correspond to three different ways that replacements can be made. And while the positions of particular updating events are different in every picture, the point is that the network of causal connections between these events is always exactly the same.

This is certainly not true for every substitution system. Indeed, the pictures on the right show how it can fail, for example, for rule (e) from the facing page. What one sees in these pictures is that after event 4, different choices of replacements are made in the two cases, and the causal relationships implied by these replacements are different.

So what could ensure that no such situation would ever arise in a particular substitution system? Essentially what needs to be true is that the sequence of elements alone must always uniquely determine what replacements can be made in every part of the system. One still has a

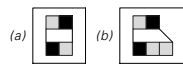
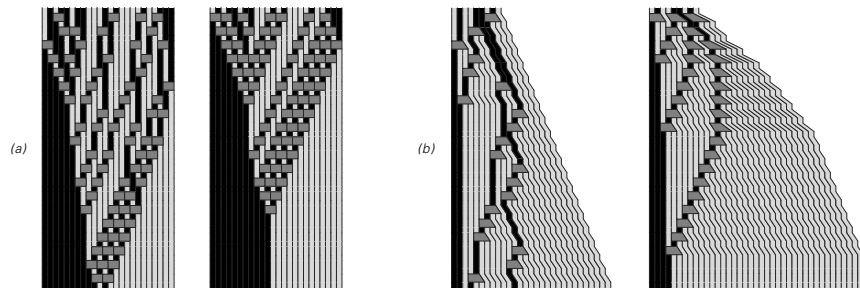


Examples of two different ways of performing replacements in rule (e) from the facing page, yielding two different causal networks.

choice of whether actually to perform a given replacement at a particular step, or whether to delay that replacement until a subsequent step. But what must be true is that there can never be any ambiguity about what replacement will eventually be made in any given part of the system.

In rules like the ones at the top of page 500 where each replacement involves just a single element this is inevitably how things must work. But what about rules that have replacements involving blocks of more than one element? Can such rules still have the necessary properties?

The pictures below show two examples of rules that do. In the first picture for each rule, replacements are made at randomly chosen steps, while in the second picture, they are in a sense always made at the earliest possible step. But the point is that in no case is there any ambiguity about what replacement will eventually be made at any particular place in the system. And as a result, the causal network that represents the relationships between different updating events is always exactly the same.



Examples of substitution systems in which the same causal networks are obtained regardless of the way in which replacements are performed. In the first picture for each rule, the replacements are performed essentially at random. In the second picture they are performed on the earliest possible step. Note that rule (a) effectively sorts the elements in its initial conditions, always placing black before white.

So what underlying property must the rules for a substitution system have in order to make the system as a whole operate in this way? The basic answer is that somehow different replacements must never be able to interfere with each other. And one way to guarantee this is if the blocks involved in replacements can never overlap.

In both the rules shown on the facing page, the only replacement specified is for the block \blacksquare . And it is inevitably the case that in any sequence of \square 's and \blacksquare 's different blocks of the form \blacksquare do not overlap. If one had replacements for blocks such as \blacksquare , \square or \blacksquare then these could overlap. But there is an infinite sequence of blocks such as \blacksquare , \blacksquare or \blacksquare for which no overlap is possible, and thus for which different replacements can never interfere.

If a rule involves replacements for several distinct blocks, then to avoid the possibility of interference one must require that these blocks can never overlap either themselves or each other. The simplest non-trivial pair of blocks that has this property is \blacksquare , \blacksquare , while the simplest triple is \blacksquare , \blacksquare , \blacksquare . And any substitution system whose rules specify replacements only for blocks such as these is guaranteed to yield the same causal network regardless of the order in which replacements are performed.

In general the condition is in fact somewhat weaker. For it is not necessary that no overlaps exist at all in the replacements—only that no overlaps occur in whatever sequences of elements can actually be generated by the evolution of the substitution systems.

And in the end there are then all sorts of substitution systems which have the property that the causal networks they generate are always independent of the order in which their rules are applied.

So what does this mean for models of the universe?

In a system like a cellular automaton, the same underlying rule is in a sense always applied in exact synchrony to every cell at every step. But what we have seen in this section is that there also exist systems in which rules can in effect be applied whenever and wherever one wants—but the same definite causal network always emerges.

So what this means is that there is no need for any built-in global clock, or even for any mechanism like an active cell. Simply by choosing the appropriate underlying rules it is possible to ensure that any sequence of events consistent with these rules will yield the same causal network and thus in effect the same perceived history for the universe.

Uniqueness and Branching in Time

If our universe has no built-in global clock and no construct like an active cell, then it is almost inevitable that at the lowest level there will be at least some arbitrariness in how its rules can be applied.

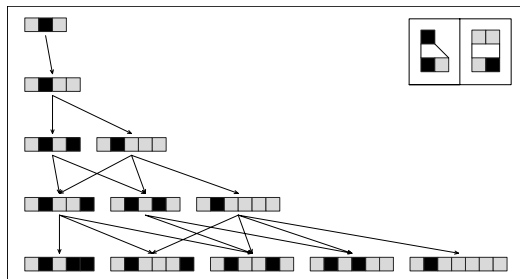
Yet in the previous section we discovered the rather remarkable fact that there exist rules with the property that essentially regardless of how they are applied, the same causal network—and thus the same perceived history for the universe—will always emerge.

But must it in the end actually be true that the underlying rules for our universe force there to be a unique perceived history? Near the end of Chapter 5 I introduced multiway systems as examples of systems that allow multiple histories. And it turns out that multiway systems are actually extremely similar in basic structure to the substitution systems that I discussed in the previous section.

Both types of systems perform the same type of replacements on strings of elements. But while in a substitution system one always carries out just a single set of replacements at each step, getting a single new string, in a multiway system one instead carries out every possible replacement, thereby typically generating many new strings.

The picture below shows a simple example of how this works. On the first step in this particular picture, there happens to be only one replacement that can be performed consistent with the rules, so only a single string is produced. But on subsequent steps several different replacements are possible, so several strings are produced. And in general every path through a picture like this corresponds to a possible history that exists in the evolution of the multiway system.

A simple example of a multiway system in which replacements are applied in all possible ways to each string at each step.



So is it conceivable that the ultimate model for our universe could be based on a multiway system? At first one might not think so. For our everyday impression is that our universe has just one definite history, not some kind of whole collection of different histories. And assuming that one is able to look at a multiway system from the outside, one will immediately see that different paths exist corresponding to different histories.

But the crucial point is that if the complete state of our universe is in effect like a single string in a multiway system, then there is no way for us ever to look at the multiway system from the outside. And as entities inside the multiway system, our perception will inevitably be that just a single path was followed, corresponding to a single history.

If one were able to look at the multiway system from the outside, this path would seem quite arbitrary. But for us inside the multiway system it is the unique path that represents the thread of experience we have had.

Up until a few centuries ago, it was widely believed that the Earth had some kind of fundamentally unique position in space. But gradually it became clear that this was not so, and that in a sense it was merely our own presence that made our particular location in space seem in any way unique. Yet for time the belief still exists that we—and our universe—somehow have a unique history. But if in fact our universe is part of a multiway system, then this will not be true. And indeed the only thing that will be unique about the particular history that our universe has had will be that it is the one we have experienced.

At a purely human level I find it rather disappointing to think that essentially none of the details of our existence are in any way unique, and that there might be other paths in the multiway system on which everything would be different. And scientifically it is also unsatisfying to have to say that there are features of our universe which are not determined by any finite set of underlying rules, but are instead in a sense just pure accidents of history associated with the particular path that we have happened to follow in a multiway system.

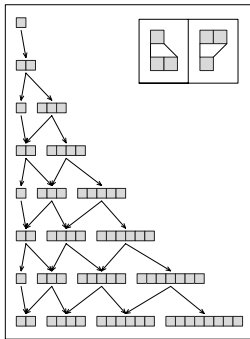
In the early parts of Chapter 7 we discussed various possible origins for the apparent randomness that we see in many natural systems. And if the universe is described by a multiway system, then

there will be an additional source of randomness: the arbitrariness of the path corresponding to the history that we have experienced.

In many respects this randomness is similar to the randomness from the environment that we discussed at the beginning of Chapter 7. But an important difference is that it would occur even if one could in effect perfectly isolate a system from the rest of the universe. If in the past one had seen apparent randomness in such a system there might have seemed to be no choice but to assume something like an underlying multiway system. But one of the discoveries of this book is that it is actually quite possible to generate what appears to be almost perfect randomness just by following definite underlying rules.

And indeed I would not expect that observations of randomness could ever reasonably be used to show that our universe is part of a multiway system. And in fact my guess is that the only way to show this with any certainty would be actually to find a specific set of multiway system rules with the property that regardless of the path that gets followed these rules would always yield behavior that agrees with the various observed features of our universe.

At some level it might seem surprising that a multiway system could ever consistently exhibit any particular form of behavior. For one might imagine that with so many different paths to choose from it would often be the case that almost any behavior would be able to occur on some path or another. And indeed, as the picture on the left shows, it is not difficult to construct multiway systems in which all possible strings of a particular kind are produced.



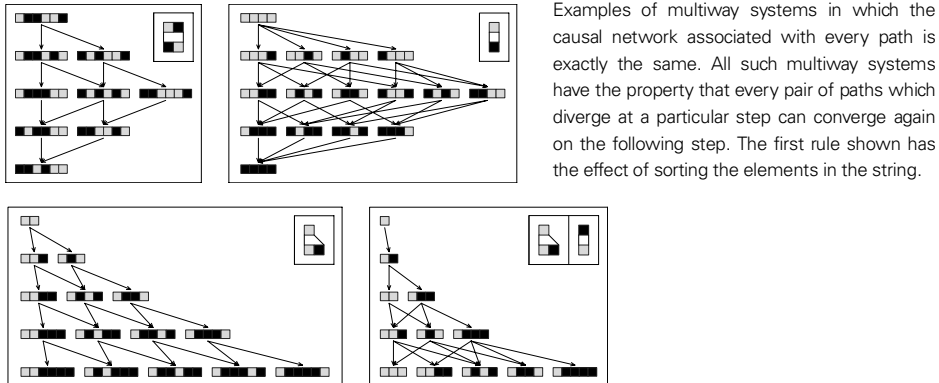
A multiway system in which strings of any length can be generated—but in which only specific sequences of lengths actually occur on any path.

But if one looks not just at individual strings but rather at the sequences of strings that exist along paths in the multiway system, then one finds that these can no longer be so arbitrary. And indeed, in any multiway system with a limited set of rules, such sequences must necessarily be subject to all sorts of constraints.

In general, each path in a multiway system can be thought of as being defined by a possible sequence of ways in which the replacements specified by a multiway system rule can be applied. And each such path in turn then defines a causal network of the kind we discussed in the previous section. But as we saw there, certain underlying rules have the

property that the form of this causal network ends up being the same regardless of the order in which replacements are applied—and thus regardless of the path that is followed in the multiway system.

The pictures below show some simple examples of rules with this property. And as it turns out, it is fairly easy to recognize the presence of the property from the overall pattern of multiway system paths that occur.



If one starts from a given initial string, then typically one will generate different strings by applying different replacements. But if one is going to get the same causal network, then it must always be the case that there are replacements one can apply to the strings one has generated that yield the same final string. So what this means is that any pair of paths in the multiway system that diverge must be able to converge again within just one step—so that all the arrows in pictures like the ones above must lie on the edges of quadrilaterals.

Most multiway systems, however, do not have exactly this property, and as a result the causal networks that are obtained by following different paths in them will not be absolutely identical. But it still turns out that whenever paths can always eventually converge—even if not in a fixed number of steps—there will necessarily be similarities on a sufficiently large scale in the causal networks that are obtained.

At the level of individual events, the structure of the causal networks will typically vary greatly. But if one looks at large enough collections of events, these details will tend to be washed out, and

regardless of the path one chooses, the overall form of causal network will be essentially the same. And what this means is that on a sufficiently large scale, the universe will appear to have a unique history, even though at the level of individual events there will be considerable arbitrariness.

If there is not enough convergence in the multiway system it will still be possible to get stuck with different types of strings that never lead to each other. And if this happens, then it means that the history of the universe can in effect follow many truly separate branches. But whenever there is significant randomness produced by the evolution of the multiway system, this does not typically appear to occur.

So this suggests that in fact it is at some level not too difficult for multiway systems to reproduce our everyday perception that more or less definite things happen in the universe. But while this means that it might be possible for there to be arbitrariness in the causal network for the universe, it still tends to be my suspicion that there is not—and that in fact the particular rules followed by the universe do in the end have the property that they always yield the same causal network.

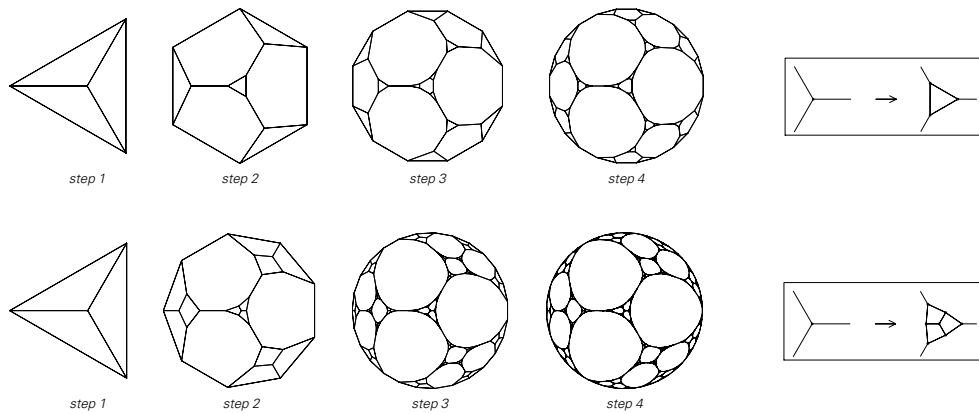
Evolution of Networks

Earlier in this chapter, I suggested that at the lowest level space might consist of a giant network of nodes. But how might such a network evolve?

The most straightforward possibility is that it could work much like the substitution systems that we have discussed in the past few sections—and that at each step some piece or pieces of the network could be replaced by others according to some fixed rule.

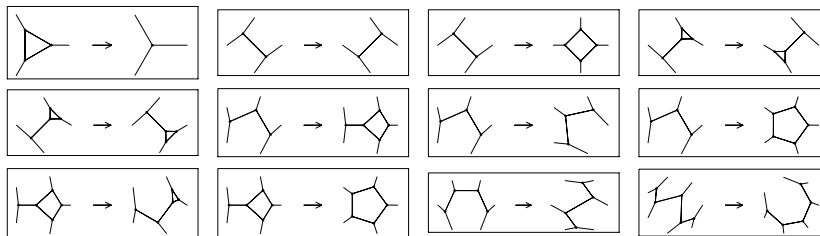
The pictures at the top of the facing page show two very simple examples. Starting with a network whose connections are like the edges of a tetrahedron, both the rules shown work by replacing each node at each step by a certain fixed cluster of nodes.

This setup is very much similar to the neighbor-independent substitution systems that we discussed on pages 83 and 187. And just as in these systems, it is possible for intricate structures to be produced, but the structures always turn out to have a highly regular nested form.



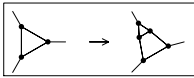
Network evolution in which each node is replaced at each step by a fixed cluster of nodes. The resulting networks have a regular nested form. The dimensions of the limiting networks are respectively $\text{Log}[2, 3] \approx 1.58$ and $\text{Log}[3, 7] \approx 1.77$.

So what about more general substitution systems? Are there analogs of these for networks? The answer is that there are, and they are based on making replacements not just for individual nodes, but rather for clusters of nodes, as shown in the pictures below.



Examples of rules that involve replacing clusters of nodes in a network by other clusters of nodes. All these rules preserve the planarity of a network. Notice that some of them cannot be reversed since their right-hand sides are too symmetrical to determine which orientation of the left-hand side should be used.

In the substitution systems for strings discussed in previous sections, the rules that are given can involve replacing any block of elements by any other. But in networks there are inevitably some restrictions. For example, if a cluster of nodes has a certain number of connections to the rest of the network, then it cannot be replaced by a cluster which has a different number of connections. And in addition, one cannot have replacements



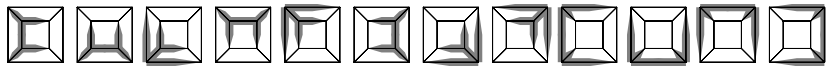
A replacement whose outcome orientation cannot be determined.

like the one on the left that go from a symmetrical cluster to one for which a particular orientation has to be chosen.

But despite these restrictions a fairly large number of replacements are still possible; for example, there are a total of 419 distinct ones that exist involving clusters with no more than five nodes.

So given a replacement for a cluster of a particular form, how should such a replacement actually be applied to a network? At first one might think that one could set up some kind of analog of a cellular automaton and just replace all relevant clusters of nodes at once.

But in general this will not work. For as the picture below illustrates, a particular form of cluster can in general appear in many overlapping ways within a given network.

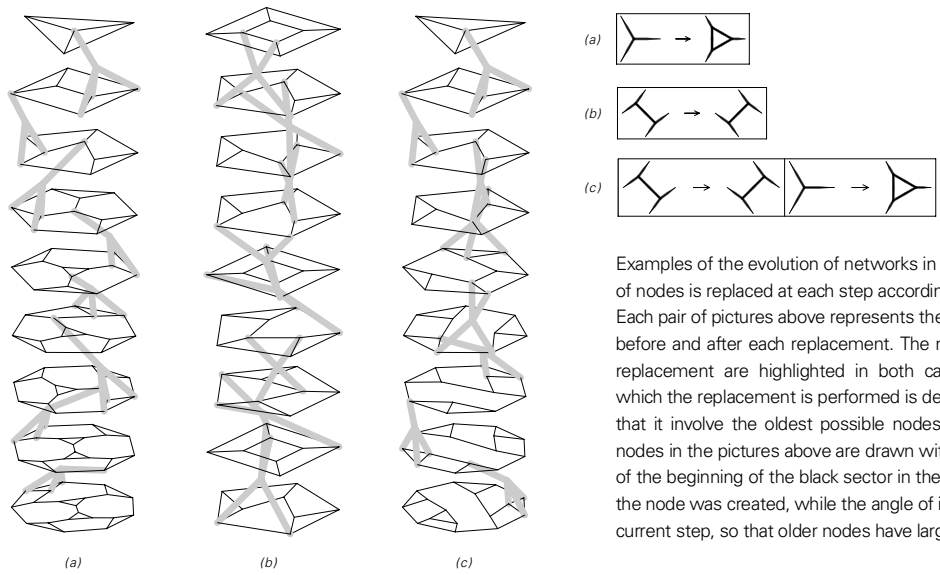
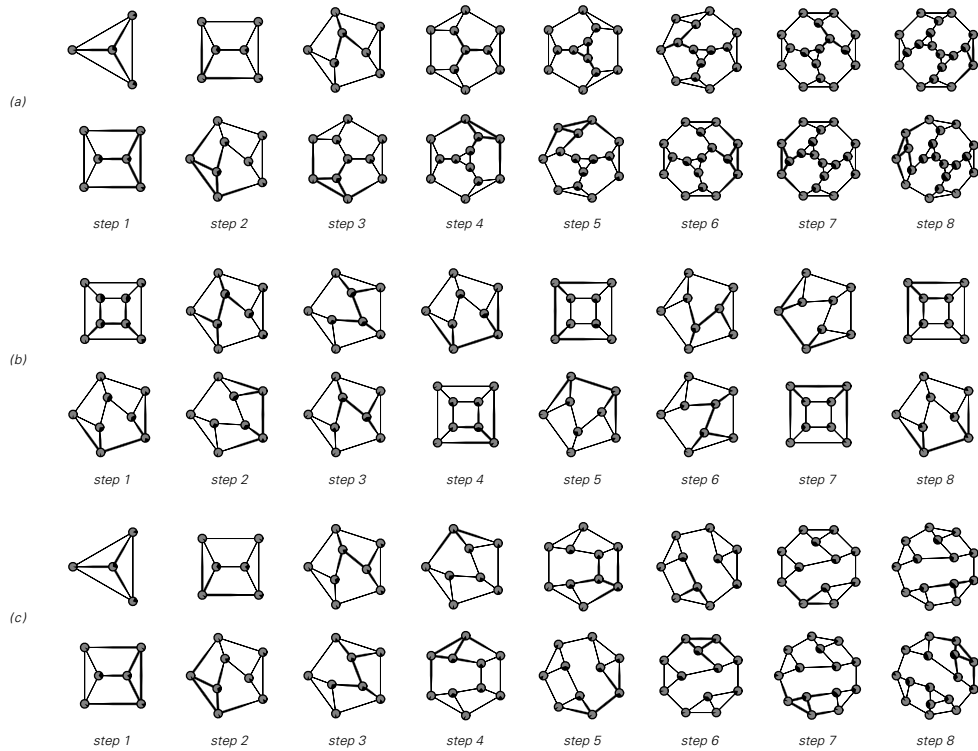


The 12 ways in which the cluster of nodes on the left occurs in a particular network. In the particular case shown, each way turns out to overlap with nodes in exactly four others.

The issue is essentially no different from the one that we encountered in previous sections for blocks of elements in substitution systems on strings. But an additional complication is that in networks, unlike strings, there is no immediately obvious ordering of elements.

Nevertheless, it is still possible to devise schemes for deciding where in a network replacements should be carried out. One fairly simple scheme, illustrated on the facing page, allows only a single replacement to be performed at each step, and picks the location of this replacement so as to affect the least recently updated nodes.

In each pair of pictures in the upper part of the page, the top one shows the form of the network before the replacement, and the bottom one shows the result after doing the replacement—with the cluster of nodes involved in the replacement being highlighted in both cases. In the 3D pictures in the lower part of the page, networks that arise on successive steps are shown stacked one on top of the other, with the nodes involved in each replacement joined by gray lines.



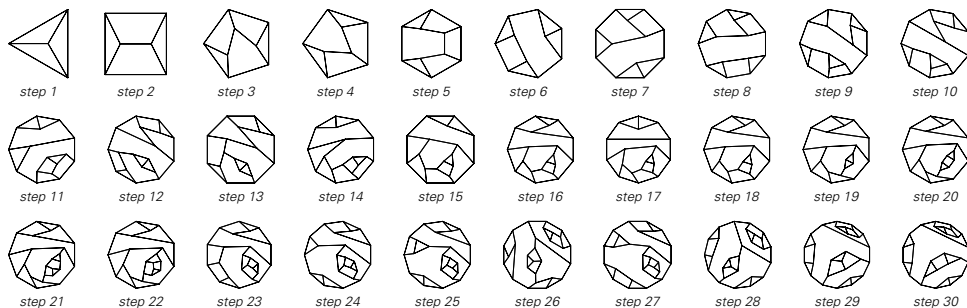
Examples of the evolution of networks in which a single cluster of nodes is replaced at each step according to the rules shown. Each pair of pictures above represents the state of the network before and after each replacement. The nodes affected by the replacement are highlighted in both cases. The location at which the replacement is performed is determined by requiring that it involve the oldest possible nodes in the network. The nodes in the pictures above are drawn with a "clock." The angle of the beginning of the black sector in the clock indicates when the node was created, while the angle of its end represents the current step, so that older nodes have larger black sectors.

Inevitably there is a certain arbitrariness in the way these pictures are drawn. For the underlying rules specify only what the pattern of connections in a network should be—not how its nodes should be laid out on the page. And in the effort to make clear the relationship between networks obtained on different steps, even identical networks can potentially be drawn somewhat differently.

With rule (a), however, it is fairly easy to see that a simple nested structure is produced, directly analogous to the one shown on page 509. And with rule (b), obvious repetitive behavior is obtained.

So what about more complicated behavior? It turns out that even with rule (c), which is essentially just a combination of rules (a) and (b), significantly more complicated behavior can already occur.

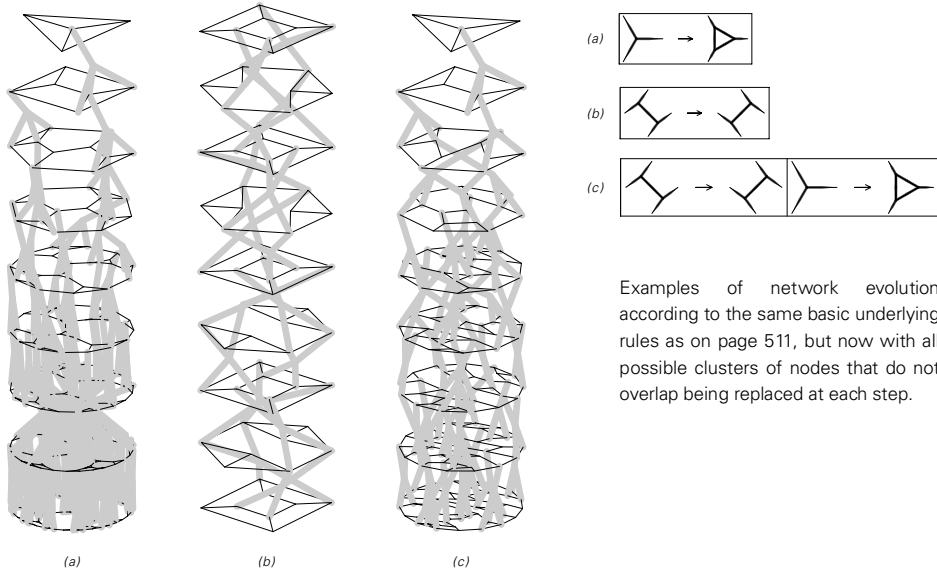
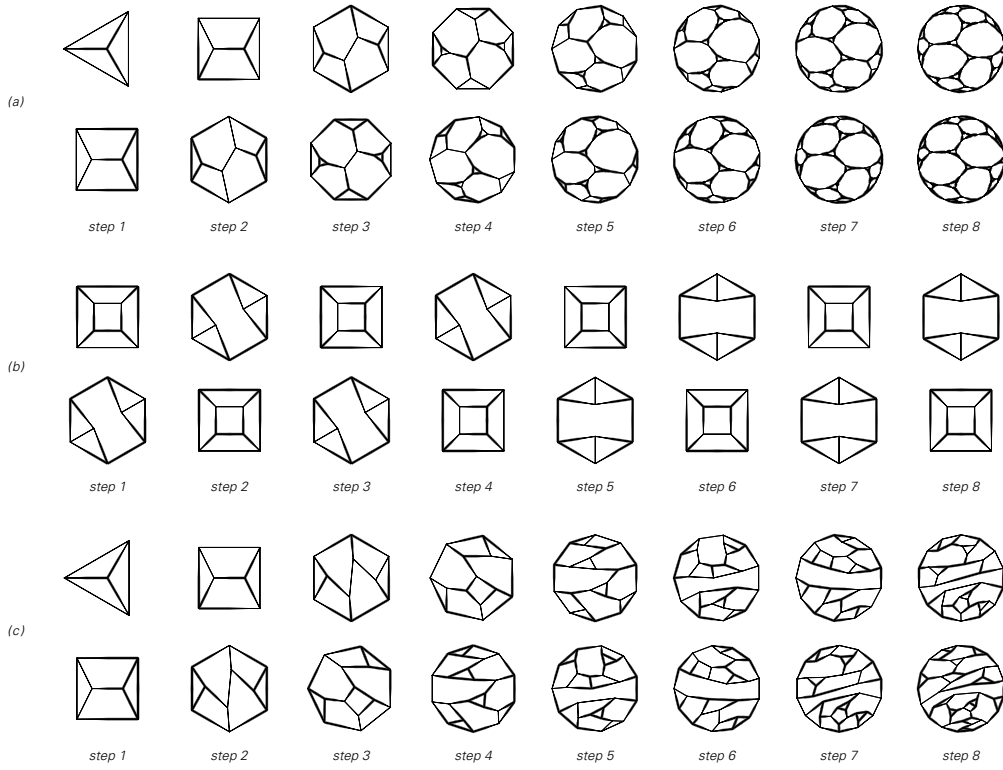
The picture below shows a few more steps in the evolution of this rule. And the behavior obtained never seems to repeat, nor do the networks produced exhibit any kind of obvious nested form.



More steps in the evolution of rule (c) from the previous page. The number of nodes increases irregularly (though roughly linearly) with successive steps.

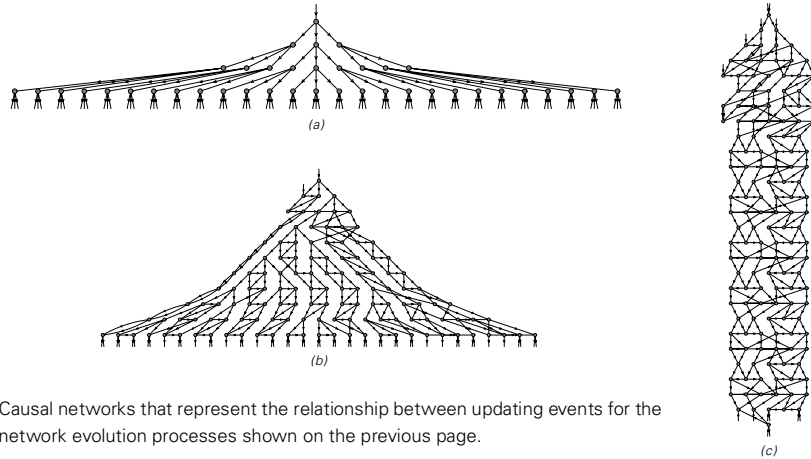
What about other schemes for applying replacements? The pictures on the facing page show what happens if at each step one allows not just a single replacement, but all replacements that do not overlap.

It takes fewer steps for networks to be built up, but the results are qualitatively similar to those on the previous page: rule (a) yields a nested structure, rule (b) gives repetitive behavior, while rule (c) produces behavior that seems complicated and in some respects random.



Examples of network evolution according to the same basic underlying rules as on page 511, but now with all possible clusters of nodes that do not overlap being replaced at each step.

Just as for substitution systems on strings, one can find causal networks that represent the causal connections between different updating events on networks. And as an example the pictures below show such causal networks for the evolution processes on the previous page.

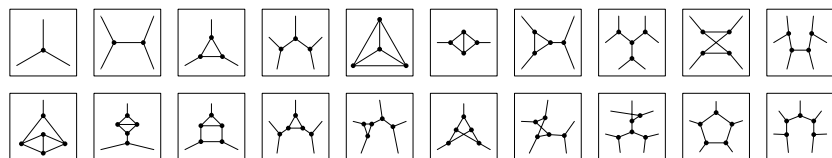


Causal networks that represent the relationship between updating events for the network evolution processes shown on the previous page.

In the rather simple case of rule (a) the results turn out to be independent of the updating scheme that was used. But for rules (b) and (c), different schemes in general yield different causal networks.

So what kinds of underlying replacement rules lead to causal networks that are independent of how the rules are applied? The situation is much the same as for strings—with the basic criterion just being that all replacements that appear in the rules should be for clusters of nodes that can never overlap themselves or each other.

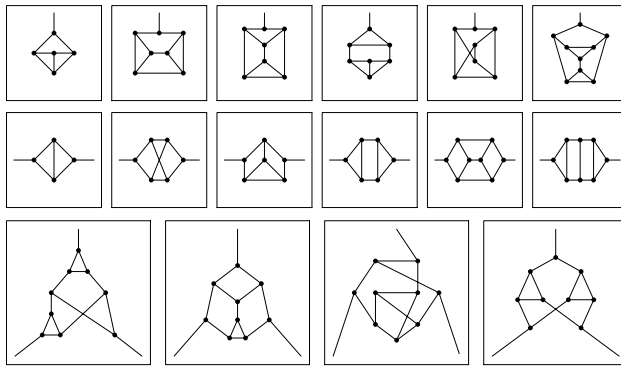
The pictures below show all possible distinct clusters with up to five nodes—and all but three of these already can overlap themselves.



All possible distinct clusters containing up to five nodes, with planarity not required.

But among slightly larger clusters there turn out to be many that do not overlap themselves—and indeed this becomes common as soon as there are at least two connections between each dangling one.

The first few examples are shown below. And in almost all of these, there is no overlap not only within a single cluster, but also between different clusters. And this means that rules based on replacements for collections of these clusters will have the property that the causal networks they produce are independent of the updating scheme used.



The simplest clusters that have no overlaps with themselves—and mostly have no overlaps with each other. Replacements for sets of clusters that do not overlap have the property of causal invariance.

One feature of the various rules I showed earlier is that they all maintain planarity of networks—so that if one starts with a network that can be laid out in the plane without any lines crossing, then every subsequent network one gets will also have this property.

Yet in our everyday experience space certainly does not seem to have this property. But beyond the practical problem of displaying what happens, there is actually no fundamental difficulty in setting up rules that can generate non-planarity—and indeed many rules based on the clusters above will for example do this.

So in the end, if one manages to find the ultimate rules for the universe, my expectation is that they will give rise to networks that on a small scale look largely random. But this very randomness will most likely be what for example allows a definite and robust value of 3 to emerge for the dimensionality of space—even though all of the many complicated phenomena in our universe must also somehow be represented within the structure of the same network.

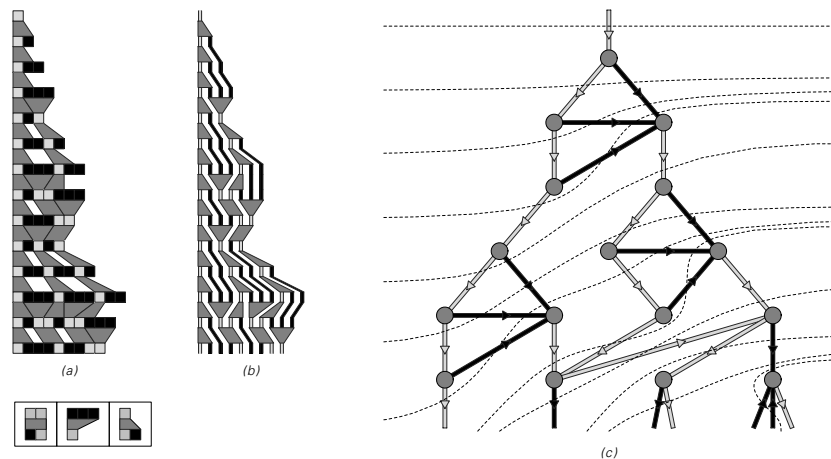
Space, Time and Relativity

Several sections ago I argued that as observers within the universe everything we can observe must at some level be associated purely with the network of causal connections between events in the universe. And in the past few sections I have outlined a series of types of models for how such a causal network might actually get built up.

But how do the properties of causal networks relate to our normal notions of space and time? There turn out to be some slight subtleties—but these seem to be exactly what end up yielding the theory of relativity.

As we saw in earlier sections, if one has an explicit evolution history for a system it is straightforward to deduce a causal network from it. But given only a causal network, what can one say about the evolution history?

The picture below shows an example of how successive steps in a particular evolution history can be recovered from a particular set of slices through the causal network derived from it. But what if one were to choose a different set of slices? In general, the sequence of strings that one would get would not correspond to anything that could arise from the same underlying substitution system.



An example of how the succession of states in an evolution history can be recovered by taking appropriate slices through a causal network. Any consistent choice of such slices will correspond to a possible evolution history—with the same underlying rules, but potentially a different scheme for determining the order in which to apply replacements.

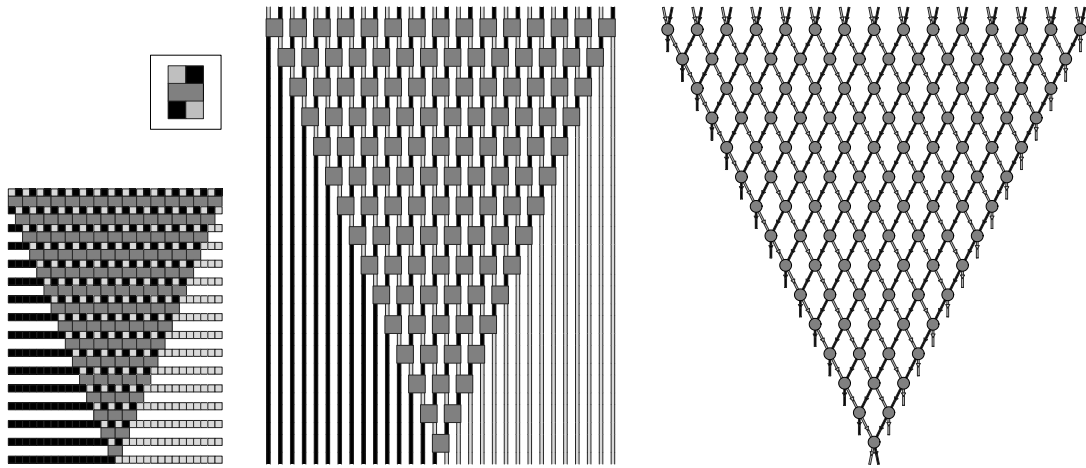
But if one has a system that yields the same causal network independent of the scheme used to apply its underlying rules, then the situation is different. And in this case any slice that consistently divides the causal network into a past and a future must correspond to a possible state of the underlying system—and any non-overlapping sequence of such slices must represent a possible evolution history for the system.

If we could explicitly see the particular underlying evolution history for the system that corresponds to our universe then this would in a sense immediately provide absolute information about space and time in the universe. But if we can observe only the causal network for the universe then our information about space and time must inevitably be deduced indirectly from looking at slices of causal networks.

And indeed only some causal networks even yield a reasonable notion of space at all. For one can think of successive slices through a causal network as corresponding to states at successive moments in time. But for there to be something one can reasonably think of as space one has to be able to identify some background features that stay more or less the same—which means that the causal network must yield consistent similarities between states it generates at successive moments in time.

One might have thought that if one just had an underlying system which did not change on successive steps then this would immediately yield a fixed structure for space. But in fact, without updating events, no causal network at all gets built up. And so a system like the one at the top of the next page is about the simplest that can yield something even vaguely reminiscent of ordinary space.

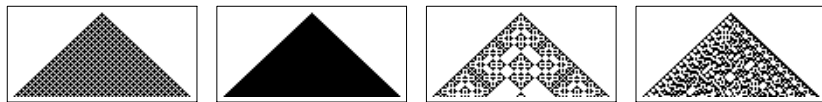
In practice I certainly do not expect that even parts of our universe where nothing much seems to be going on will actually have causal networks as simple as at the top of the next page. And in fact, as I mentioned at the end of the previous section, what I expect instead is that there will always tend to be all sorts of complicated and seemingly random behavior at small scales—though at larger scales this will typically get washed out to yield the kind of consistent average properties that we ordinarily associate with space.



A very simple substitution system whose causal network has slices that can be thought of as corresponding to a highly regular idealization of one-dimensional ordinary space. The rule effectively just sorts elements so that black ones come first, and yields the same causal network regardless of what updating scheme is used.

One of the defining features of space as we normally experience it is a certain locality that leads most things that happen at some particular position to be able at first to affect only things very near them.

Such locality is built into the basic structure of systems like cellular automata. For in such systems the underlying rules allow the color of a particular cell to affect only its immediate neighbors at each step. And this has the consequence that effects in such systems can spread only at a limited rate, as manifest for example in a maximum slope for the edges of patterns like those in the pictures below.



Examples of patterns produced by cellular automata, illustrating the fact discussed in Chapter 6 that the edge of each pattern has a maximum slope equal to one cell per step, corresponding to an absolute upper limit on the rate of information transmission—similar to the speed of light in physics.

In physics there also seems to be a maximum speed at which the effects of any event can spread: the speed of light, equal to about 300

million meters per second. And it is common in spacetime physics to draw “light cones” of the kind shown at the right to indicate the region that will be reached by a light signal emitted from a particular position in space at a particular time. So what is the analog of this in a causal network?

The answer is straightforward, for the very definition of a causal network shows that to see how the effects of a particular event spread one just has to follow the successive connections from it in the causal network.

But in the abstract there is no reason that these connections should lead to points that can in any way be viewed as nearby in space. Among the various kinds of underlying systems that I have studied in this book many have no particular locality in their basic rules. But the particular kinds of systems I have discussed for both strings and networks in the past few sections do have a certain locality, in that each individual replacement they make involves only a few nearby elements.

One might choose to consider systems like these just because it seems easier to specify their rules. But their locality also seems important in giving rise to anything that one can reasonably recognize as space.

For without it there will tend to be no particular way to match up corresponding parts in successive slices through the causal networks that are produced. And as a result there will not be the consistency between successive slices necessary to have a stable notion of space.

In the case of substitution systems for strings, locality of underlying replacement rules immediately implies overall locality of effects in the system. For the different elements in the system are always just laid out in a one-dimensional string, with the result that local replacement rules can only ever propagate effects to nearby elements in the string—much like in a one-dimensional cellular automaton.

If one is dealing with an underlying system based on networks, however, then the situation can be somewhat more complicated. For as we discussed several sections ago—and will discuss again in the final sections of this chapter—there will typically be only an approximate correspondence between the structure of the network and the structure of ordinary space. And so for example—as we will discuss later in connection with quantum phenomena—there may sometimes be a kind of thread that connects parts of the network that would not



Schematic illustration of a light cone in physics. Light emitted at a point in space will normally spread out with time into a cone, whose cross-section is shown schematically here.

normally be considered nearby in three-dimensional space. And so when clusters of nodes that are nearby with respect to connections on the network get updated, they can potentially propagate effects to what might be considered distant points in space.

Nevertheless, if a network is going to correspond to space as it seems to exist in our universe, such phenomena must not be too important—and in the end there must to a good approximation be the kind of straightforward locality that exists for example in the simple causal network of page 518.

In the next section I will discuss how actual physical entities like particles propagate in systems represented by causal networks. But ultimately the whole point of causal networks is that their connections represent all possible ways that effects propagate. Yet these connections are also what end up defining our notions of space and time in a system. And particularly in a causal network as regular as the one on page 518 one can then immediately view each connection in the causal network as corresponding to an effect propagating a certain distance in space during a certain interval in time.

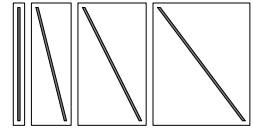
So what about a more complicated causal network? One might imagine that its connections could perhaps represent varying distances in space and varying intervals in time. But there is no independent way to work out distance in space or interval in time beyond looking at the connections in the causal network. So the only thing that ultimately makes sense is to measure space and time taking each connection in the causal network to correspond to an identical elementary distance in space and elementary interval in time.

One may guess that this elementary distance is around 10^{-35} meters, and that the elementary time interval is around 10^{-43} seconds. But whatever these values are, a crucial point is that their ratio must be a fixed speed, and we can identify this with the speed of light. So this means that in a sense every connection in a causal network can be viewed as representing the propagation of an effect at the speed of light.

And with this realization we are now close to being able to see how the kinds of systems I have discussed must almost inevitably succeed in reproducing the fundamental features of relativity theory.

But first we must consider the concept of motion.

To say that one is not moving means that one imagines one is in a sense sampling the same region of space throughout time. But if one is moving—say at a fixed speed—then this means that one imagines that the region of space one is sampling systematically shifts with time, as illustrated schematically in the simple pictures on the right.



Graphical representation in space and time of motion at fixed speeds.

But as we have seen in discussing causal networks, it is in general quite arbitrary how one chooses to match up space at different times. And in fact one can just view different states of motion as corresponding to different such choices: in each case one matches up space so as to treat the point one is at as being the same throughout time.

Motion at a fixed speed is then the simplest case—and the one emphasized in the so-called special theory of relativity. And at least in the context of a highly regular causal network like the one in the picture on page 518 there is a simple interpretation to this: it just corresponds to looking at slices at different angles through the causal network.

Successive parallel slices through the causal network in general correspond to successive states of the underlying system at successive moments in time. But there is nothing that determines in any absolute way the overall angle of these slices in pictures like those on page 518. And the point is that in fact one can interpret slices at different angles as corresponding to motion at different fixed speeds.

If the angle is so great that there are connections going up as well as down between slices, then there will be a problem. But otherwise it will always be the case that regardless of angle, successive slices must correspond to possible evolution histories for the underlying system.

One might have thought that states obtained from slices at different angles would inevitably be consistent only with different sets of underlying rules. But in fact this is not the case, and instead the exact same rules can reproduce slices at all angles. And this is a consequence of the fact that the substitution system on page 518 has the property of causal invariance—so that it gives the same causal network independent of the scheme used to apply its underlying rules.

It is slightly more complicated to represent uniform motion in causal networks that are not as regular as the one on page 518. But

whenever there is sufficient uniformity to give a stable structure to space one can still think of something like parallel slices at different angles as representing motion at different fixed speeds.

And the crucial point is that whenever the underlying system is causal invariant the exact same underlying rules will account for what one sees in slices at different angles. And what this means is that in effect the same rules will apply regardless of how fast one is going.

And the remarkable point is then that this is also what seems to happen in physics. For everyday experience—together with all sorts of detailed experiments—strongly support the idea that so long as there are no effects from acceleration or external forces, physical systems work exactly the same regardless of how fast they are moving.

At the outset it might not have seemed conceivable that any system which at some level just applies a fixed program to various underlying elements could successfully capture the phenomenon of motion. For certainly a system like a typical cellular automaton does not—since for example its effective rules for evolution at different angles will usually be quite different. But there are two crucial ideas that make motion work in the kinds of systems I am discussing here. First, that causal networks can represent everything that can be observed. And second, that with causal invariance different slices through a causal network can be produced by the same underlying rules.

Historically, the idea that physical processes should always be independent of overall motion goes back at least three hundred years. And from this idea one expects for example that light should always travel at its usual speed with respect to whatever emitted it. But what if one happens to be moving with respect to this emitter? Will the light then appear to be travelling at a different speed? In the case of sound it would. But what was discovered around the end of the 1800s is that in the case of light it does not. And it was essentially to explain this surprising fact that the special theory of relativity was developed.

In the past, however, there seemed to be no obvious underlying mechanism that could account for the validity of this basic theory. But now it turns out that the kinds of discrete causal network models that I have described almost inevitably end up being able to do this.

And essentially the reason for this is that—as I discussed above—each individual connection in any causal network must almost by definition represent propagation of effects at the speed of light. The overall structure of space that emerges may be complicated, and there may be objects that end up moving at all sorts of speeds. But at least locally the individual connections basically define the speed of light as a fixed maximum rate of propagation of any effect. And the point is that they do this regardless of how fast the source of an effect may be moving.

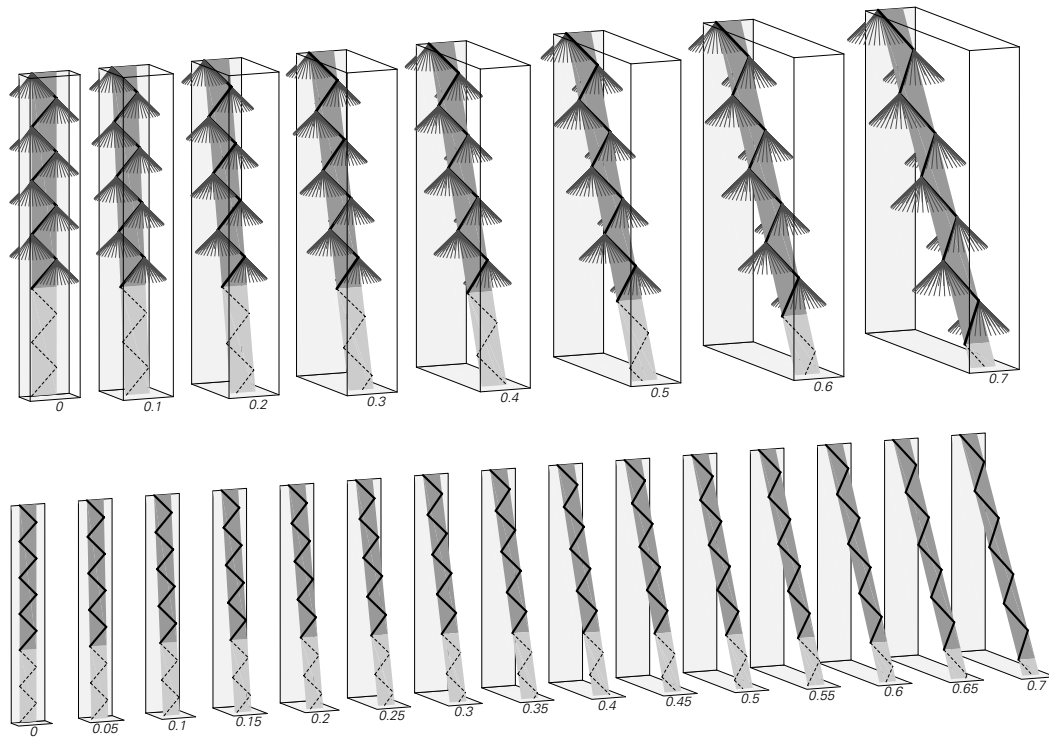
So from this one can use essentially standard arguments to derive all the various phenomena familiar from ordinary relativity theory. A typical example is time dilation, in which a fixed time interval for a system moving at some speed seems to correspond to a longer time interval for a system at rest. The picture on the next page shows schematically how this at first unexpected result arises.

The basic idea is to consider what happens when a system that can act as a simple clock moves at different speeds. At a traditional physics level one can think of the clock as having a photon of light bouncing backwards and forwards between mirrors a fixed distance apart. But more generally one can think of following criss-crossing connections that exist in some fixed fragment of a causal network.

In the picture on the next page time goes down the page. The internal mechanism of the clock is shown as a zig-zag black line—with each sweep of this line corresponding to the passage of one unit of time.

The black line is always assumed to be moving at the speed of light—so that it always lies on the surface of a light cone, as indicated in the top row of pictures. But then in successive pictures the whole clock is taken to move at increasing fractions of the speed of light.

The dark gray region in each picture represents a fixed amount of time for the clock—corresponding to a fixed number of sweeps of the black line. But as the pictures indicate, it is then essentially just a matter of geometry to see that this dark gray region will correspond to progressively larger amounts of time for a system at rest—in just the way predicted by the standard formula of relativistic time dilation.



A simple derivation of the classic phenomenon of relativistic time dilation. The pictures show the behavior of a very simple idealized clock going at different fractions of the speed of light. The clock can be thought of as consisting of a photon of light bouncing backwards and forwards between mirrors a fixed distance apart. (At a more general level in my approach it can also be thought of as a fragment of a causal network.) Time is shown going down the page, so that the photon in the clock traces out a zig-zag path. The fundamental assumption—that in my approach is just a consequence of basic properties of causal networks—is that the photon always goes at the speed of light, so that its path always lies on the surface of light cones like the ones in the top row of pictures. A fixed interval of time for the clock—as indicated by the length of the darker gray regions—corresponds to a progressively longer interval of time at rest. The amount of this time dilation is given by the classic relativistic formula $1/\sqrt{1-v^2/c^2}$, where v/c is the ratio of the speed of the clock to the speed of light. Such time dilation is routinely observed in particle accelerators—and has to be corrected for in GPS satellites. It leads to the so-called twin paradox in which less time will pass for a member of a twin going at high speed in a spacecraft than one staying at rest. The fact that time dilation is a general phenomenon not restricted to something like the simple clock shown relies in my approach on general properties of causal networks. Once the basic assumptions are established, the derivation of time dilation given here is no different in principle from the original one given in 1905, though I believe it is in many ways considerably clearer. Note that it is necessary to consider motion in two dimensions—so that the clock as a whole can be moving perpendicular to the path of the photon inside it. If these were parallel, one would inevitably get not just pure time dilation, but a mixture of it and length contraction.

Elementary Particles

There are some aspects of the universe—notably the structure of space and time—that present-day physics tends to assume are continuous. But over the past century it has at least become universally accepted that all matter is made up of identifiable discrete particles.

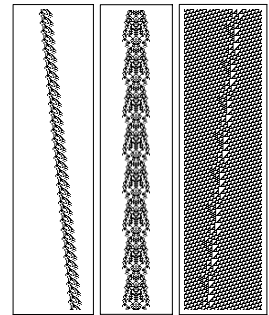
Experiments have found a fairly small number of fundamentally different kinds of particles, with electrons, photons, muons and the six basic types of quarks being a few examples. And it is one of the striking observed regularities of the universe that all particles of a given kind—say electrons—seem to be absolutely identical in their properties.

But what actually are particles? As far as present-day experiments can tell, electrons, for example, have zero size and no substructure. But particularly if space is discrete, it seems almost inevitable that electrons and other particles must be made up of more fundamental elements.

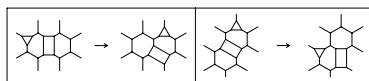
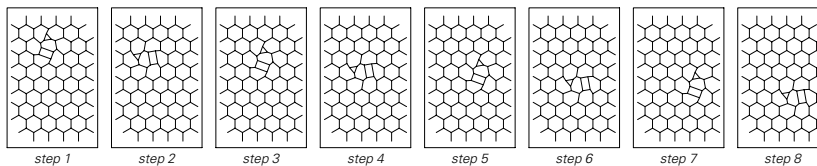
So how might this work? An immediate possibility that I suspect is actually not too far from the mark is that such particles are analogs of the localized structures that we saw earlier in this book in systems like the class 4 cellular automata shown on the right. And if this is so, then it means that at the lowest level, the rules for the universe need make no reference to particular particles. Instead, all the particles we see would just emerge as structures formed from more basic elements.

In networks it can be somewhat difficult to visualize localized structures. But the picture below nevertheless shows a simple example of how a localized structure can move across a regular planar network.

Both the examples on this page show structures that exist on very regular backgrounds. But to get any kind of realistic model for actual



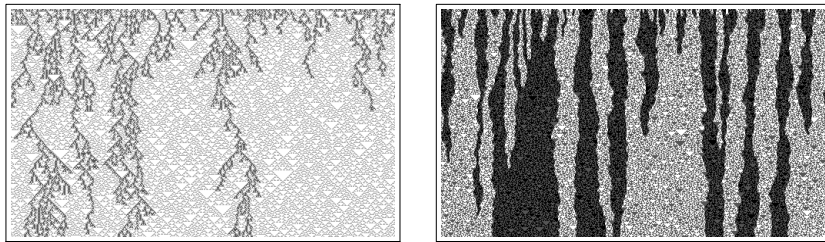
Typical examples of particle-like localized structures in class 4 cellular automata.



A particle-like localized structure in a network.

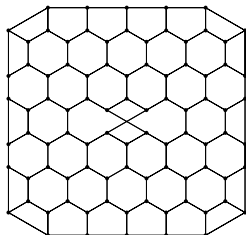
particles in physics one must consider structures on much more complicated and random backgrounds. For any network that has a serious chance of representing actual space—even a supposedly empty part—will no doubt show all sorts of seemingly random activity. So any localized structure that might represent a particle will somehow have to persist even on this kind of random background.

Yet at first one might think that such randomness would inevitably disrupt any kind of definite persistent structure. But the pictures below show two simple examples where it does not. In the first case, there are localized cracks that persist. And in the second case, there are two different types of regions, separated by boundaries that act like localized structures with definite properties, and persist until they annihilate.



Examples of one-dimensional cellular automata that support various forms of persistent structures even on largely random backgrounds. These are 3-color totalistic rules with codes 294 and 1893.

So what about networks? It turns out that here again it is possible to get definite structures that persist even in the presence of randomness. And to see an example of this consider setting up rules like those on page 509 that preserve the planarity of networks.

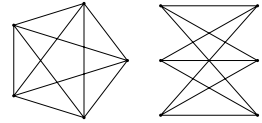


A network with a single irreducible crossing of lines.

Starting off with a network that is planar—so that it can be drawn flat on a page without any lines crossing—such rules can certainly give all sorts of complex and apparently random behavior. But the way the rules are set up, all the networks they produce must still be planar.

And if one starts off with a network like the one on the left that can only be drawn with lines crossing, then what will happen is that the non-planarity of the network will be preserved. But to what extent does this non-planarity correspond to a definite structure in the network?

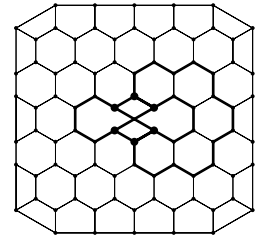
There are typically many different ways to draw a non-planar network, each with lines crossing in different places. But there is a fundamental result in graph theory that shows that if a network is not planar, then it must always be possible to identify in it a specific part that can be reduced to one of the two forms shown on the right—or just the second form for a network with three connections at each node.



The K_5 and $K_{3,3}$ forms that lead to non-planarity in networks.

So this implies that one can in fact meaningfully associate a definite structure with non-planarity. And while at some level the structure can be spread out in the network, the point is that it must always in effect have a localized core with the form shown on the right.

In general one can imagine having several pieces of non-planarity in a network—perhaps each pictured like a carrying handle. But if the underlying rules for the network preserve planarity then each of these pieces of non-planarity must on their own be persistent—and can in a sense only disappear through processes like annihilating with each other.



How $K_{3,3}$ is embedded in the network from the facing page.

So might these be like actual particles in physics?

In the realistic case of network rules for the universe, planarity as such is presumably not preserved. But observations in physics suggest that there are several quantities like electric charge that are conserved. And ultimately the values of these quantities must reflect properties of underlying networks that are preserved by network evolution rules.

And if these rules satisfy the constraint of causal invariance that I discussed in previous sections, then I suspect that this means that they will inevitably exhibit various additional features—perhaps notably including for example what is usually known as local gauge invariance.

But what is most relevant here is that it seems likely that—much as for non-planarity—nonzero values of quantities conserved by network evolution rules can be thought of as being associated with some sort of local structures or tangles of connections in the network. And I suspect that it is essentially such structures that define the cores of the various types of elementary particles that are seen in physics.

Before the results of this book it might have seemed completely implausible that anything like this could be correct. For independent of any specific arguments about networks and their evolution, traditional intuition would tend to make one think that the elaborate properties of

particles must inevitably be the result of an elaborate underlying setup. But what we have now seen over and over again in this book is that in fact it is perfectly possible to get phenomena of great complexity even with a remarkably simple underlying setup. And I suspect that particles in physics—with all their various properties and interactions—are just yet another example of this very general phenomenon.

One immediate thing that might seem to suggest that elementary particles must somehow be based on simple discrete structures is the fact that their values of quantities like electric charge always seem to be in simple rational ratios. In traditional particle physics this is explained by saying that many if not all particles are somehow just manifestations of the same underlying abstract object, related by a simple fixed group of symmetry operations. But in terms of networks one can imagine a much more explicit explanation: that there are just a simple discrete set of possible structures for the cores of particles—each perhaps related in some quite mechanical way by the group of symmetry operations.

But in addition to quantities like electric charge, another important intrinsic property of all particles is mass. And unlike for example electric charge the observed masses of elementary particles never seem to be in simple ratios—so that for example the muon is about 206.7683 times the mass of the electron, while the tau lepton is about 16.819 times the mass of the muon. But despite such results, it is still conceivable that there could in the end be simple relations between truly fundamental particle masses—since it turns out that the masses that have actually been observed in effect also include varying amounts of interaction energy.

A defining feature of any particle is that it can somehow move in space while maintaining its identity. In traditional physics, such motion has a straightforward mathematical representation, and it has not usually seemed meaningful to ask what might underlie it. But in the approach that I take here, motion is no longer such an intrinsic concept, and the motion of a particle must be thought of as a process that is made up of a whole sequence of explicit lower-level steps.

So at first, it might seem surprising that one can even set up a particular type of particle to move at different speeds. But from the discussion in the previous section it follows that this is actually an

almost inevitable consequence of having underlying rules that show causal invariance. For assuming that around the particle there is some kind of uniformity in the causal network—and thus in the apparent structure of space—taking slices through the causal network at an appropriate angle will always make any particle appear to be at rest. And the point is that causal invariance then implies that the same underlying rules can be used to update the network in all such cases.

But what happens if one has two particles that are moving with different velocities? What will the events associated with the second particle look like if one takes slices through the causal network so that the first particle appears to be at rest? The answer is that the more the second particle moves between successive slices, the more updating events must be involved. For in effect any node that was associated with the particle on either one slice or the next must be updated—and the more the particle moves, the less these will overlap. And in addition, there will inevitably appear to be an asymmetry in the pattern of events relative to whatever direction the particle is moving.

There are many subtleties here, and indeed to explain the details of what is going on will no doubt require quite a few new and rather abstract concepts. But the general picture that I believe will emerge is that when particles move faster they will appear to have more nodes associated with them.

Most likely the intrinsic properties of a particle—like its electric charge—will be associated with some sort of core that corresponds to a definite network structure involving a roughly fixed number of nodes. But I suspect that the apparent motion of the particle will be associated with a kind of coat that somehow interpolates from the core to the uniform background of surrounding space. With different slices through the causal network, the apparent size of this coat can change. But I suspect that the size of the coat in a particular case will somehow be related to the apparent energy and momentum of a particle in that case.

An important fact in traditional physics is that interactions between particles seem to conserve total energy and momentum. And conceivably the reason for this is that such interactions somehow tend to preserve the total number of network nodes. Indeed, perhaps in most

situations—save those associated with the overall expansion of the universe—the basic rules for the network at least on average just rearrange nodes and never change their number.

In traditional physics energy and momentum are always assumed to have continuous values. But just as in the case of position there is no contradiction with sufficiently small underlying discrete elements.

As I will discuss in the last section of this chapter, quantum mechanics tends to make one think of particles with higher momenta as being somehow progressively less spread out in space. So how can this be consistent with the idea that higher momentum is associated with having more nodes? Part of the answer probably has to do with the fact that outside the piece of the network that corresponds to the particle, the network presumably matches up to yield uniform space in much the same way as without the particle. And within the piece of the network corresponding to the particle, the effective structure of space may be very different—with for example more long-range connections added to reduce the effective overall distance.

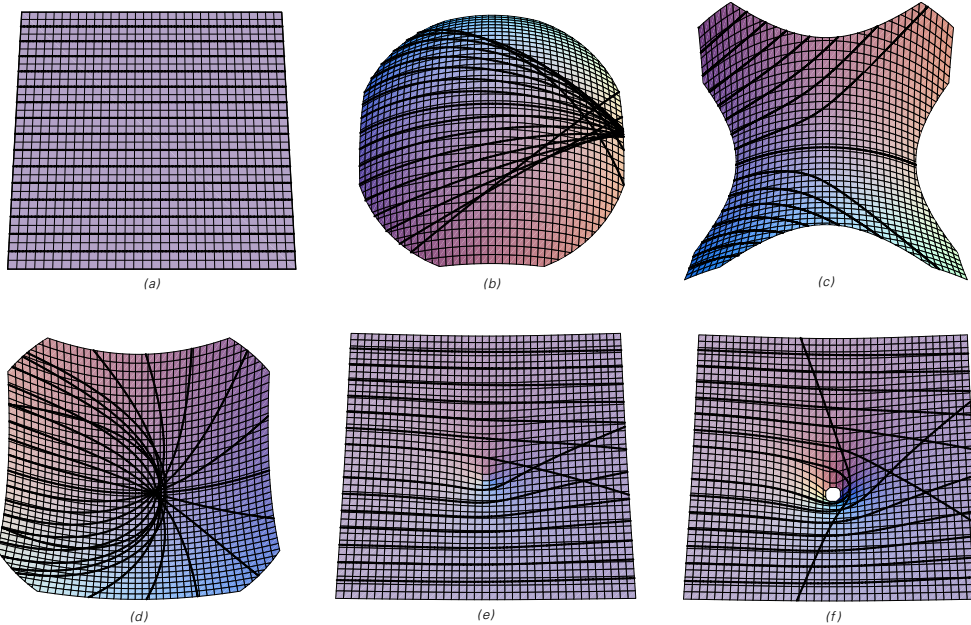
The Phenomenon of Gravity

At an opposite extreme from elementary particles one can ask how the universe behaves on the largest possible scales. And the most obvious effect on such scales is the phenomenon of gravity. So how then might this emerge from the kinds of models I have discussed here?

The standard theory of gravity for nearly a century has been general relativity—which is based on the idea of associating gravity with curvature in space, then specifying how this curvature relates to the energy and momentum of whatever matter is present.

Something like a magnetic field in general has different effects on objects made of different materials. But a key observation verified experimentally to considerable accuracy is that gravity has exactly the same effect on the motion of different objects, regardless of what those objects are made of. And it is this that allows one to think of gravity as a general feature of space—rather than for example as some type of force that acts specifically on different objects.

In the absence of any gravity or forces, our normal definition of space implies that when an object moves from one point to another, it always goes along a straight line, which corresponds to the shortest path. But when gravity is present, objects in general move on curved paths. Yet these paths can still be the shortest—or so-called geodesics—if one takes space to be curved. And indeed if space has appropriate curvature one can get all sorts of paths, as in the pictures below.



Examples of the effect of curvature in space on paths taken by objects. In each case all the paths shown start parallel, but do not remain so when there is curvature. The paths are geodesics which go the minimum distance on the surface to get to all the points they reach. (In general, the minimum may only be local.) Case (b) shows the top of a sphere, which is a surface of positive curvature. Case (c) shows the negatively curved surface $z = x^2 - y^2$, (d) a paraboloid $z = x^2 + y^2$, and (e, f) $z = 1/(r + \delta)$ —a rough analog of curvature in space produced by a sphere of mass.

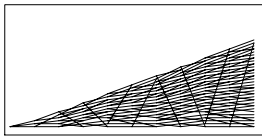
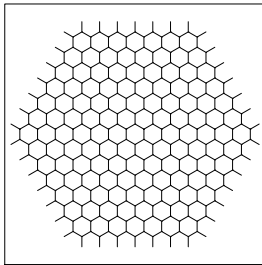
But in our actual universe what determines the curvature of space? The answer from general relativity is that the Einstein equations give conditions for the value of a particular kind of curvature in terms of the energy and momentum of matter that is present. And the point then is that the shortest paths in space with this curvature seem to be

consistent with those followed by objects moving under the influence of gravity associated with the given distribution of matter.

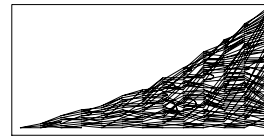
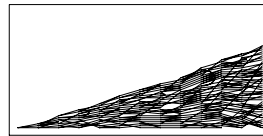
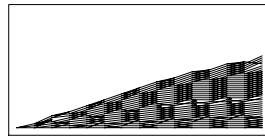
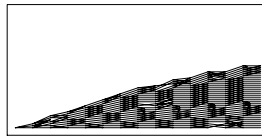
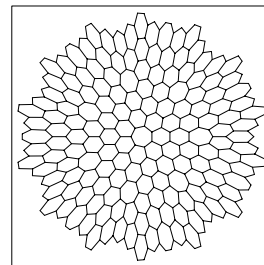
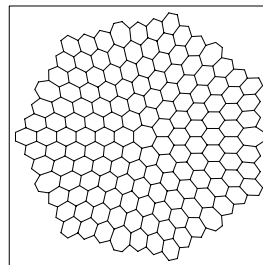
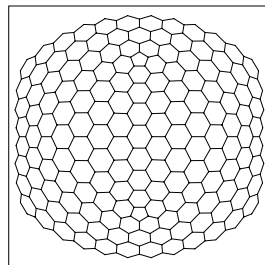
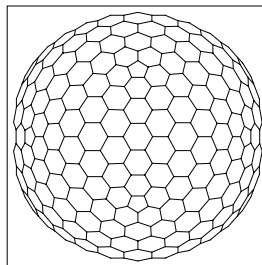
For a continuous surface—or in general a continuous space—the idea of curvature is a familiar one in traditional geometry. But if the universe is at an underlying level just a discrete network of nodes then how does curvature work? At some level the answer is that on large scales the discrete network must approximate continuous space.

But it turns out that one can actually also recognize curvature in the basic structure of a network. If one has a simple array of hexagons—as in the picture on the left—then this can readily be laid out flat on a two-dimensional plane. But what if one replaces some of these hexagons by pentagons? One still has a fundamentally two-dimensional surface. But if one tries to keep all edges the same length the surface will inevitably become curved—like a soccer ball or a geodesic dome.

So what this suggests is that in a network just changing the pattern of connections can in effect change the overall curvature. And indeed the pictures below show a succession of networks that in effect have curvatures with a range of negative and positive values.



A hexagonal array corresponding to flat two-dimensional space.



Networks with various limiting curvatures. If every region in the network is in effect a hexagon—as in the picture at the top of the page—then the network will behave as if it is flat. But if pentagons are introduced, as in the cases on the left, the network will increasingly behave as if it has positive curvature—like part of a sphere. And if heptagons are introduced, as in the cases on the right, the network will behave as if it has negative curvature. In the bottom row of pictures, the networks are laid out as on page 479, so that successive heights give the number of nodes at successive distances r from a particular node. In the limit of large r , this number is approximately $r^2(1 - k r^2 + \dots)$ where k turns out to be exactly proportional to the curvature.

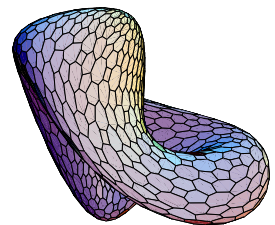
But how can we determine the curvature from the structure of each network? Earlier in this chapter we saw that if a network is going to correspond to ordinary space in some number of dimensions d , then this means that by going r connections from any given node one must reach about r^{d-1} nodes. But it turns out that when curvature is present it leads to a systematic correction to this.

In each of the pictures on the facing page the network shown can be thought of as corresponding to two-dimensional space. And this means that to a first approximation the number of nodes reached must increase linearly with r . But the bottom row of pictures show that there are corrections to this. And what happens is that when there is positive curvature—as in the pictures on the left—progressively fewer than r nodes end up being reached. But when there is negative curvature—as on the right—progressively more nodes end up being reached. And in general the leading correction to the number of nodes reached turns out to be proportional to the curvature multiplied by r^{d+1} .

So what happens in more than two dimensions? In general the result could be very complicated, and could for example involve all sorts of different forms of curvature and other characteristics of space. But in fact the leading correction to the number of nodes reached is always quite simple: it is just proportional to what is called the Ricci scalar curvature, multiplied by r^{d+1} . And already here this is some suggestion of general relativity—for the Ricci scalar curvature also turns out to be a central quantity in the Einstein equations.

But in trying to see a more detailed correspondence there are immediately a variety of complications. Perhaps the most obvious is that the traditional mathematical formulation of general relativity seems to rely on many detailed properties of continuous space. And while one expects that sufficiently large networks should in some sense act on average like continuous space, it is far from clear at first how the kinds of properties of relevance to general relativity will emerge.

If one starts, say, from an ordinary continuous surface, then it is straightforward to approximate it as in the picture on the right by a collection of flat faces. And one might think that the edges of these faces would define a network of the kind I have been discussing.



A surface approximated by flat faces whose edges form a trivalent network.

But in fact, such a network has vastly less information. For given just a set of connections between nodes, there is no obvious way even to know which of these connections should be associated with the same face—let alone to work out anything like angles between faces.

Yet despite this, it turns out that all the geometrical features that are ultimately of relevance to general relativity can actually be determined in large networks just from the connectivity of nodes.

One of these is the value of the so-called Ricci tensor, which in effect specifies how the Ricci scalar curvature is made up from different curvature components associated with different directions.

As indicated above, the scalar curvature associated with a network is directly related to how many nodes lie within successive distances r of a given node on the network—or in effect how many nodes lie within successive generalized spheres around that node. And it turns out that the projection of the Ricci tensor along a particular direction is then just related to the number of nodes that lie within a cylinder oriented in that direction. But even just defining a consistent direction in a network is not entirely straightforward. But one way to do it is simply to pick two points in the network, then to say that paths in the network are going in the same direction if they are segments of the same shortest path between those points. And with this definition, a region that approximates a cylinder can be formed just by setting up spheres with centers at every point on the path.

But there is now another issue to address: at least in its standard formulation general relativity is set up in terms of properties not of three-dimensional space but rather of four-dimensional spacetime. And this means that what is relevant are properties not so much of specific networks representing space, but rather of complete causal networks.

And one immediate feature of causal networks that differs from space networks is that their connections go only one way. But it turns out that this is exactly what one needs in order to set up the analog of a spacetime Ricci tensor. The idea is to start at a particular event in the causal network, then to form what is in effect a cone of events that can be reached from there. To define the spacetime Ricci tensor, one considers—as on page 516—a sequence of spacelike slices through this

cone and asks how the number of events that lie within the cone increases as one goes to successive slices. After t steps, the number of events reached will be proportional to t^d . But there is then a correction proportional to t^{d+2} , that has a coefficient that is a combination of the spacetime Ricci scalar and a projection of the spacetime Ricci tensor along what is in effect the time direction defined by the sequence of spacelike slices chosen.

So how does this relate to general relativity? It turns out that when there is no matter present the Einstein equations simply state that the spacetime Ricci tensor—and thus all of its projections—are exactly zero. There can still for example be higher-order curvature, but there can be no curvature at the level described by the Ricci tensor.

So what this means is that any causal network whose behavior obeys the Einstein equations must at the level of counting nodes in a cone have the same uniform structure as it would if it were going to correspond to ordinary flat space. As we saw a few sections ago, many underlying replacement rules end up producing networks that are for example too extensively connected to correspond to ordinary space in any finite number of dimensions. But I suspect that if one has replacement rules that are causal invariant and that in effect successfully maintain a fixed number of dimensions they will almost inevitably lead to behavior that follows something close to the Einstein equations.

Probably the situation is somewhat analogous to what we saw with fluid behavior in cellular automata in Chapter 8—that at least if there are underlying rules whose behavior is complicated enough to generate significant effective randomness, then almost whenever the rules lead to conservation of total particle number and momentum something close to the ordinary Navier-Stokes equation behavior emerges.

So what about matter?

As a first step, one can ask what effect the structure of space has on something like a particle—assuming that one can ignore the effect of the particle back on space. In traditional general relativity it is always assumed that a particle which is not interacting with anything else will move along a shortest path—or so-called geodesic—in space.

But what about an explicit particle of the kind we discussed in the previous section that exists as a structure in a network? Given two nodes in a network, one can always identify a shortest path from one to the other that goes along a sequence of individual connections in the network. But in a sense a structure that corresponds to a particle will normally not fit through this path. For usually the structure will involve many nodes, and thus typically require many connections going in more or less the same direction in order to be able to move across the network.

But if one assumes a certain uniformity in networks—and in particular in the causal network—then it still follows that particles of the kind that we discussed in the previous section will tend to move along geodesics. And whereas in traditional general relativity the idea of motion along geodesics is essentially an assumption, this can now in principle be derived explicitly from an underlying network model.

One might have thought that in the absence of matter there would be little to say about gravity—since after all the Einstein equations then say that there can be no curvature in space, at least of the kind described by the Ricci tensor. But it turns out that there can still be other kinds of curvature—described for example by the so-called Riemann tensor—and these can in fact lead to all sorts of phenomena. Examples include familiar ones like inverse-square gravitational fields around massive objects, as well as unfamiliar ones like gravitational waves.

But while the mathematical structure of general relativity is complicated enough that it is often difficult to see just where in spacetime effects come from, it is usually assumed that matter is somehow ultimately required to provide a source for gravity. And in the full Einstein equations the Ricci tensor need not be zero; instead it is specified at every point in space as being equal to a certain combination of energy and momentum density for matter at that point. So this means that to know what will happen even in phenomena primarily associated with gravity one typically has to know all sorts of properties of matter.

But why exactly does matter have to be introduced explicitly at all? It has been the assumption of traditional physics that even though gravity can be represented in terms of properties of space, other elements of our universe cannot. But in my approach everything just

emerges from the same underlying network—or in effect from the structure of space. And indeed even in traditional general relativity one can try avoiding introducing matter explicitly—for example by imagining that everything we call matter is actually made up of pure gravitational energy, or of something like gravitational waves.

But so far as one can tell, the details of this do not work out—so that at the level of general relativity there is no choice but to introduce matter explicitly. Yet I suspect that this is in effect just a sign of limitations in the Einstein equations and general relativity.

For while at a large scale these may provide a reasonable description of average behavior in a network, it is almost inevitable that closer to the scale of individual connections they will have to be modified. Yet presumably one can still use the Einstein equations on large scales if one introduces matter with appropriate properties as a way to represent small-scale effects in the network.

In the previous section I suggested that energy and momentum might in effect be associated with the presence of excess nodes in a network. And this now potentially seems to fit quite well with what we have seen in this section. For if the underlying rule for a network is going to maintain to a certain approximation the same average number of nodes as flat space, then it follows that wherever there are more nodes corresponding to energy and momentum, this must be balanced by something reducing the number of nodes. But such a reduction is exactly what is needed to correspond to positive curvature of the kind implied by the Einstein equations in the presence of ordinary matter.

Quantum Phenomena

From our everyday experience with objects that we can see and touch we develop a certain intuition about how things work. But nearly a century ago it became clear that when it comes to things like electrons some of this intuition is no longer correct. Yet there has developed an elaborate mathematical formalism in quantum theory that successfully reproduces much of what is observed. And while some aspects of this

formalism remain mysterious, it has increasingly come to be believed that any fundamental theory of physics must somehow be based on it.

Yet the kinds of programs I have discussed in this book are not in any obvious way set up to fit in with this formalism. But as we have seen a great many times in the course of the book, what emerges from a program can be very different from what is obvious in its underlying rules. And in fact it is my strong suspicion that the kinds of programs that I have discussed in the past few sections will actually in the end turn out to show many if not all the key features of quantum theory.

To see this, however, will not be easy. For the kinds of constructs that are emphasized in the standard formalism of quantum theory are very different from those immediately visible in the programs I have discussed. And ultimately the only reliable way to make contact will probably be to set up rather complete and realistic models of experiments—then gradually to see how limits and idealizations of these manage to match what is expected from the standard formalism. Yet from what we have seen in this chapter and earlier in this book there are already some encouraging signs that one can identify.

At first, though, things might not seem promising. For my model of particles such as electrons being persistent structures in a network might initially seem to imply that such particles are somehow definite objects just like ones familiar from everyday experience. But there are all sorts of phenomena in quantum theory that seem to indicate that electrons do not in fact behave like ordinary objects that have definite properties independent of us making observations of them.

So how can this be consistent? The basic answer is just that a network which represents our whole universe must also include us as observers. And this means that there is no way that we can look at the network from the outside and see the electron as a definite object. Instead, anything we deduce about the electron must come from processes that explicitly go on inside the network.

But this is not just an issue in studying things like electrons: it is actually a completely general feature of the models I have discussed. And in fact, as we saw earlier in this chapter, it is what allows them to support meaningful notions of even such basic concepts as time. At a

more formal level, it also implies that everything we can observe can be captured by a causal network. And as I will discuss a little below, I suspect that the idea of causal invariance for such a network will then be what turns out to account for some key features of quantum theory.

The basic picture of our universe that I have outlined in the past few sections is a network whose connections are continually updated according to some simple set of underlying rules. In the past one might have assumed that a system like this would be far too simple to correspond to our universe. But from the discoveries in this book we now know that even when the underlying rules for a system are simple, its overall behavior can still be immensely complex.

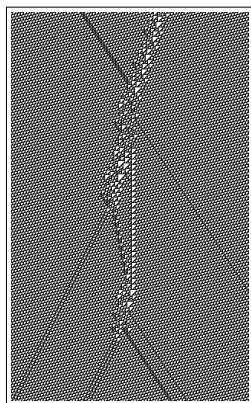
And at the lowest level what I expect is that even though the rules being applied are perfectly definite, the overall pattern of connections that will exist in the network corresponding to our universe will continually be rearranged in ways complicated enough to seem effectively random.

Yet on a slightly larger scale such randomness will then lead to a certain average uniformity. And it is then essentially this that I believe is responsible for maintaining something like ordinary space—with gradual variations giving rise to the phenomenon of gravity.

But superimposed on this effectively random background will then presumably also be some definite structures that persist through many updatings of the network. And it is these, I believe, that are what correspond to particles like electrons.

As I discussed in the last two sections, causal invariance of the underlying rules implies that such structures should be able to move at a range of uniform speeds through the background. Typically properties like charge will be associated with some specific pattern of connections at the core of the structure corresponding to a particle, while the energy and momentum of the particle will be associated with roughly the number of nodes in some outer region around the core.

So what about interactions? If the structures corresponding to different particles are isolated, then the underlying rules will make them persist. But if they somehow overlap, these same rules will usually make some different configuration of particles be produced.



A collision between localized structures in the rule 110 class 4 cellular automaton.

At some level the situation will no doubt be a little like in the evolution of a typical class 4 cellular automaton, as illustrated on the left. Given some initial set of persistent structures, these can interact to produce some intermediate pattern of behavior, which then eventually resolves into a final set of structures that again persist.

In the intermediate pattern of behavior one may also be able to identify some definite structures. Ones that do not last long can be very different from ones that would persist forever. But ones that last longer will tend to have properties progressively closer to genuinely persistent structures. And while persistent structures can be thought of as corresponding to real particles, intermediate structures are in many ways like the virtual particles of traditional particle physics.

So this means that a picture like the one on the left above can be viewed in a remarkably literal sense as being a spacetime diagram of particle interactions—a bit like a Feynman diagram from particle physics.

One immediate difference, however, is that in traditional particle physics one does not imagine a pattern of behavior as definite and determined as in the picture above. And indeed in my model for the universe it is already clear that there is more going on. For any process like the one in the picture above must occur on top of a background of apparently random small-scale rearrangements of the network. And in effect what this background does is to introduce a kind of random environment that can make many different detailed patterns of behavior occur with certain probabilities even with the same initial configuration of particles.

The idea that even a vacuum without particles will have a complicated and in some ways random form also exists in standard quantum field theory in traditional physics. The full mathematical structure of quantum field theory is far from completely worked out. But the basic notion is that for each possible type of particle there is some kind of continuous field that exists throughout space—with the presence of a particle corresponding to a simple type of structure in this field.

In general, the equations of quantum field theory seem to imply that there can be all sorts of complicated configurations in the field, even in the absence of actual particles. But as a first approximation, one can consider

just short-lived pairs of virtual particles and antiparticles. And in fact one can often do something similar for networks. For even in the planar networks discussed on page 527 a great many different arrangements of connections can be viewed as being formed from different configurations of nearby pairs of non-planar persistent structures.

Talking about a random background affecting processes in the universe immediately tends to suggest certain definite relations between probabilities for different processes. Thus for example, if there are two different ways that some process can occur, it suggests that the total probability for the whole process should be just the sum of the probabilities for the process to occur in the two different ways.

But the standard formalism of quantum theory says that this is not correct, and that in fact one has to look at so-called probability amplitudes, not ordinary probabilities. At a mathematical level, such amplitudes are analogous to ones for things like waves, and are in effect just numbers with directions. And what quantum theory says is that the probability for a whole process can be obtained by linearly combining the amplitudes for the different ways the process can occur, then looking at the square of the magnitude of the result—or the analog of intensity for something like a wave.

So how might this kind of mathematical procedure emerge from the types of models I have discussed? The answer seems complicated. For even though the procedure itself may sound straightforward, the constructs on which it operates are actually far from easy to define just on the basis of an underlying network—and I have seen no easy way to unravel the various limits and idealizations that have to be made.

Nevertheless, a potentially important point is that it is in some ways misleading to think of particles in a network as just interacting according to some definite rule, and being perturbed by what is in essence a random background. For this suggests that there is in effect a unique history to every particle interaction—determined by the initial conditions and the configuration that exists in the random background.

But the true picture is more complicated. For the sequence of updates to the underlying network can be made in any order—yet each order in effect gives a different detailed history for the network. But if

there is causal invariance, then ultimately all these different histories must in a sense be equivalent. And with this constraint, if one breaks some process into parts, there will typically be no simple way to describe how the effect of these parts combines together.

And for at least some purposes it may well make sense to think explicitly about different possible histories, combining something like amplitudes that one assigns to each of them. Yet quite how this might work will certainly depend on what feature of the network one tries to look at.

It has always been a major issue in quantum theory just how one tells what is happening with a particular particle like an electron. From our experience with everyday objects we might think that it should somehow be possible to do this without affecting the electron. But if the only things we have are particles, then to find out something about a given particle we inevitably have to have some other particle—say a photon of light—explicitly interact with it. And in this interaction the original particle will inevitably be affected in some way.

And in fact just one interaction will certainly not be enough. For we as humans cannot normally perceive individual particles. And indeed there usually have to be a huge number of particles doing more or less the same thing before we successfully register it.

Most often the way this is made to happen is by setting up some kind of detector that is initially in a state that is sufficiently unstable that just a single particle can initiate a whole cascade of consequences. And usually such a detector is arranged so that it evolves to one or another stable state that has sufficiently uniform properties that we can recognize it as corresponding to a definite outcome of a measurement.

At first, however, such evolution to an organized state might seem inconsistent with microscopic reversibility. But in fact—just as in so many other seemingly irreversible processes—all that is needed to preserve reversibility is that if one looks at sufficient details of the system there can be arbitrary and seemingly random behavior. And the point is just that in making conclusions about the result of a measurement we choose to ignore such details.

So even though the actual result that we take away from a measurement may be quite simple, many particles—and many events—

will always be involved in getting it. And in fact in traditional quantum theory no measurement can ultimately end up giving a definite result unless in effect an infinite number of particles are involved.

As I mentioned above, ordinary quantum processes can appear to follow different histories depending on what scheme is used to decide the order in which underlying rules are applied. But taking the idealized limit of a measurement in which an infinite number of particles are involved will probably in effect establish a single history.

And this implies that if one knew all of the underlying details of the network that makes up our universe, it should always be possible to work out the result of any measurement. I strongly believe that the initial conditions for the universe were quite simple. But like many of the processes we have seen in this book, the evolution of the universe no doubt intrinsically generates apparent randomness.

And the result is that most aspects of the network that represents the current state of our universe will seem essentially random. So this means that to know its form we would in essence have to sample every one of its details—which is certainly not possible if we have to use measurements that each involve a huge number of particles.

One might however imagine that as a first approximation one could take account of underlying apparent randomness just by saying that there are certain probabilities for particles to behave in particular ways. But one of the most often quoted results about foundations of quantum theory is that in practice there can be correlations observed between particles that seem impossible to account for in at least the most obvious kind of such a so-called hidden-variables theory.

For in particular, if one takes two particles that have come from a single source, then the result of a measurement on one of them is found in a sense to depend too much on what measurement gets done on the other—even if there is not enough time for information travelling at the speed of light to get from one to the other. And indeed this fact has often been taken to imply that quantum phenomena can ultimately never be the result of any definite underlying process of evolution.

But this conclusion depends greatly on traditional assumptions about the nature of space and of particles. And it turns out that for the kinds of models I have discussed here it in general no longer holds.

And the basic reason for this is that if the universe is a network then it can in a sense easily contain threads that continue to connect particles even when the particles get far apart in terms of ordinary space.

The picture that emerges is then of a background containing a very large number of connections that maintain an approximation to three-dimensional space, together with a few threads that in effect go outside of that space to make direct connections between particles.

If two particles get created together, it is reasonable to expect that the tangles that represent their cores will tend to have a few connections in common—and indeed this for example happens for lumps of non-planarity of the kind we discussed on page 527. But until there are interactions that change the structure of the cores, these common connections will then remain—and will continue to define a thread that goes directly from one particle to the other.

But there is immediately a slight subtlety here. For earlier in this chapter I discussed measuring distance on a network just by counting the minimum number of successive individual connections that one has to follow in order to get from one point to another. Yet if one uses this measure of distance then the distance between two particles will always tend to remain fixed as the number of connections in the thread.

But the point is that this measure of distance is in reality just a simple idealization of what is relevant in practice. For the only way we end up actually being able to measure physical distances is in effect by looking at the propagation of photons or other particles. Yet such particles always involve many nodes. And while they can get from one point to another through the large number of connections that define the background space, they cannot in a sense fit through a small number of connections in a thread. So this means that distance as we normally experience it is typically not affected by threads.

But it does not mean that threads can have no effect at all. And indeed what I suspect is that it is precisely the presence of threads that leads to the correlations that are seen in measurements on particles.

It so happens that the standard formalism of quantum theory provides a rather simple mathematical description of these correlations. And it is certainly far from obvious how this might emerge from detailed mechanisms associated with threads in a network. But the fact that this and other results seem simple in the standard formalism of quantum theory should not be taken to imply that they are in any sense particularly fundamental. And indeed my guess is that most of them will actually in the end turn out to depend on all sorts of limits and idealizations in quantum theory—and will emerge just as simple approximations to much more complex underlying behavior.

In its development since the early 1900s quantum theory has produced all sorts of elaborate results. And to try to derive them all from the kinds of models I have outlined here will certainly take an immense amount of work. But I consider it very encouraging that some of the most basic quantum phenomena seem to be connected to properties like causal invariance and the network structure of space that already arose in our discussion of quite different fundamental issues in physics.

And all of this supports my strong belief that in the end it will turn out that every detail of our universe does indeed follow rules that can be represented by a very simple program—and that everything we see will ultimately emerge just from running this program.

Fundamental Physics

The Notion of Reversibility

■ **Page 437 · Testing for reversibility.** To show that a cellular automaton is reversible it is sufficient to check that all configurations consisting of repetitions of different blocks have different successors. This can be done for blocks up to length n in a 1D cellular automaton with k colors using

```
ReversibleQ[rule_, k_, n_] := Catch[Do[
  If[Length[Union[Table[CAStep[rule, IntegerDigits[i, k, m]],
    {i, 0, k^m - 1}]]] ≠ k^m, Throw[False]], {m, n}]; True]
```

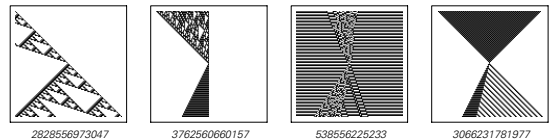
For $k=2, r=1$ it turns out that it suffices to test only up to $n=4$ (128 out of the 256 rules fail at $n=1$, 64 at $n=2$, 44 at $n=3$ and 14 at $n=4$); for $k=2, r=2$ it suffices to test up to $n=15$, and for $k=3, r=1$, up to $n=9$. But although these results suggest that in general it should suffice to test only up to $n=k^{2^r}$, all that has so far been rigorously proved is that $n=k^{2^r}(k^{2^r}-1)+2r+1$ (or $n=15$ for $k=2, r=1$) is sufficient.

For 2D cellular automata an analogous procedure can in principle be used, though there is no upper limit on the size of blocks that need to be tested, and in fact the question of whether a particular rule is reversible is directly equivalent to the tiling problem discussed on page 213 (compare page 942), and is thus formally undecidable.

■ **Numbers of reversible rules.** For $k=2, r=1$, there are 6 reversible rules, as shown on page 436. For $k=2, r=2$ there are 62 reversible rules, in 20 families inequivalent under symmetries, out of a total of 2^{32} or about 4 billion possible rules. For $k=3, r=1$ there are 1800 reversible rules, in 172 families. For $k=4, r=1$, some of the reversible rules can be constructed from the second-order cellular automata below. Note that for any k and r , no non-trivial totalistic rule can ever be reversible.

■ **Inverse rules.** Some reversible rules are self-inverse, so that applying the same rule twice yields the identity. Other rules come in distinct pairs. Most often a rule that involves r neighbors has an inverse that also involves at most r neighbors. But for both $k=2, r=2$ and $k=3, r=1$ there turn out to be reversible rules whose inverses involve larger

numbers of neighbors. For any given rule one can define the neighborhood size s to be the largest block of cells that is ever needed to determine the color of a single new cell. In general $s \leq 2r+1$, and for a simple identity or shift rule, $s=1$. For $k=2, r=1$, it then turns out that all the reversible rules and their inverses have $s=1$. For $k=2, r=2$, the reversible rules have values of s from 1 to 5, but their inverses have values \bar{s} from 1 to 6. There are only 8 rules (the inequivalent ones being 16740555 and 3327051468) where $\bar{s} > s$, and in each case $\bar{s}=6$ while $s=5$. For $k=3, r=1$, there are a total of 936 rules with this property: 576, 216 and 144 with $\bar{s}=4, 5$ and 6, and in all cases $s=3$. Examples with $\bar{s}=3, 4, 5$ and 6 are shown below. For arbitrary k and r , it is not clear what the maximum \bar{s} can be; the only bound rigorously established so far is $\bar{s} \leq r + 1/2 k^{2^{r+1}} (k^{2^r} - 1)$.



■ **Surjectivity and injectivity.** See page 959.

■ **Directional reversibility.** Even if successive time steps in the evolution of a cellular automaton do not correspond to an injective map, it is still possible to get an injective map by looking at successive lines at some angle in the spacetime evolution of the system. Examples where this works include the surjective rules 30 and 90.

■ **Page 437 · Second-order cellular automata.** Second-order elementary rules can be implemented using

```
CA2EvolveList[rule_List, {a_List, b_List}, t_Integer] :=
  Map[First, NestList[CA2Step[rule, #] &, {a, b}, t]]
CA2Step[rule_List, {a_, b_}] := {b, Mod[a + rule[[
  8 - (RotateLeft[b] + 2(b + 2 RotateRight[b]))], 2]}
```

where *rule* is obtained from the rule number using *IntegerDigits*[*n*, 2, 8].

The combination $Drop[list, -1] + 2Drop[list, 1]$ of the result from *CA2EvolveList* corresponds to evolution according to a first-order $k = 4, r = 1$ rule.

■ **History.** The concept of getting reversibility in a cellular automaton by having a second-order rule was apparently first suggested by Edward Fredkin around 1970 in the context of 2D systems—on the basis of an analogy with second-order differential equations in physics. Similar ideas had appeared in numerical analysis in the 1960s in connection with so-called symmetric or self-adjoint discrete approximations to differential equations.

■ **Page 438 • Properties.** The pattern from rule 67R with simple initial conditions grows irregularly, at an average rate of about 1 cell every 5 steps. The right-hand side of the pattern from rule 173R consists three triangles that repeat progressively larger at steps of the form $2(9^s - 1)$. Rule 90R has the property that of the diamond of cells at relative positions $\{-n, 0\}, \{0, -n\}, \{n, 0\}, \{0, n\}$ it is always true for any n that an even number are black.

■ **Page 439 • Properties.** The initial conditions used here have a single black cell on two successive initial steps. For rule 150R, however, there is no black cell on the first initial step. The pattern generated by rule 150R has fractal dimension $\log[2, 3 + \sqrt{17}] - 1$ or about 1.83. In rule 154R, each diagonal stripe is followed by at least one 0; otherwise, the positions of the stripes appear to be quite random, with a density around 0.44.

■ **Generalized additive rules.** Additive cellular automata of the kind discussed on page 952 can be generalized by allowing the new value of each cell to be obtained from combinations of cells on s previous steps. For rule 90 the combination c can be specified as $\{\{1, 0, 1\}\}$, while for rule 150R it can be specified as $\{\{0, 1, 0\}, \{1, 1, 1\}\}$. All generalized additive rules ultimately yield nested patterns. Starting with a list of the initial conditions for s steps, the configurations for the next s steps are given by

```
Append[Rest[list],
  Map[Mod[Apply[Plus, Flatten[c#]], 2] &, Transpose[
    Table[RotateLeft[list, {0, i}], {i, -r, r}], {3, 2, 1}]]]
```

where $r = (\text{Length}[\text{First}[c]] - 1)/2$.

Just as for ordinary additive rules on page 1091, an algebraic analysis for generalized additive rules can be given. The objects that appear are solutions to linear recurrences of order s , and in general involve s^{th} roots. For rule 150R, the configuration at step t as shown in the picture on page 439 is given by $(u^t - v^t)/\text{Sqrt}[4 + h^2]$, where $\{u, v\} = z/. \text{Solve}[z^2 == hz + 1]$ and $h = 1/x + 1 + x$. (See also page 1078.)

■ **Page 440 • Rule 37R.** Complicated structures are fairly easy to get with this rule. The initial condition $\{1, 0, 1\}$ with all cells 0 on the previous step yields a structure that repeats but only every 666 steps. The initial condition $\{\{0, 1, 1\}, \{1, 0, 0\}\}$ yields a pattern that grows sporadically for 3774 steps, then breaks into two repetitive structures. The typical background repeats every 3 steps.

■ **Classification of reversible rules.** In a reversible system it is possible with suitable initial conditions to get absolutely any arrangement of cells to appear at any step. Despite this, however, the overall spacetime pattern of cells is not arbitrary, but is instead determined by the underlying rules. If one starts with completely random initial conditions then class 2 and class 3 behavior are often seen. Class 1 behavior can never occur in a reversible system. Class 4 behavior can occur, as in rule 37R, but is typically obvious only if one starts say with a low density of black cells.

For arbitrary rules, difference patterns of the kind shown on page 250 can get both larger and smaller. In a reversible rule, such patterns can grow and shrink, but can never die out completely.

■ **Emergence of reversibility.** Once on an attractor, any system—even if it does not have reversible underlying rules—must in some sense show approximate reversibility. (Compare page 959.)

■ **Other reversible systems.** Reversible examples can be found of essentially all the types of systems discussed in this book. Reversible mobile automata can for instance be constructed using

```
Table[{IntegerDigits[i, 2, 3] -> If[First[#] == 0, {#, -1},
  {Reverse[#], 1}] &}[IntegerDigits[perm[[i]], 2, 3]], {i, 8}]
```

where $perm$ is an element of $Permutations[Range[8]]$. An example that exhibits complex behavior is:



Systems based on numbers are typically reversible whenever the mathematical operations they involve are invertible. Thus, for example, the system on page 121 based on successive multiplication by $3/2$ is reversible by using division by $3/2$. Page 905 gives another example of a reversible system based on numbers.

Multiway systems are reversible whenever both $a \rightarrow b$ and $b \rightarrow a$ are present as rules, so that the system corresponds mathematically to a semigroup. (See page 938.)

■ **Reversible computation.** Typical practical computers—and computer languages—are not even close to reversible: many inputs can lead to the same output, and there is no unique

way to undo the steps of a computation. But despite early confusion (see page 1020), it has been known since at least the 1970s that there is nothing in principle which prevents computation from being reversible. And indeed—just like with the cellular automata in this section—most of the systems in Chapter 11 that exhibit universal computation can readily be made reversible with only slight overhead.

Irreversibility and the Second Law of Thermodynamics

■ **Time reversal invariance.** The reversibility of the laws of physics implies that given the state of a physical system at a particular time, it is always possibly to work out uniquely both its future and its past. Time reversal invariance would further imply that the rules for going in each direction should be identical. To a very good approximation this appears to be true, but it turns out that in certain esoteric particle physics processes small deviations have been found. In particular, it was discovered in 1964 that the decay of the K^0 particle violated time reversal invariance at the level of about one part in a thousand. In current theories, this effect is not attributed any particularly fundamental origin, and is just assumed to be associated with the arbitrary setting of certain parameters. K^0 decay was for a long time the only example of time reversal violation that had explicitly been seen, although recently examples in B particle decays have probably also been seen. It also turns out that the only current viable theories of the apparent preponderance of matter over antimatter in the universe are based on the idea that a small amount of time reversal violation occurred in the decays of certain very massive particles in the very early universe.

The basic formalism used for particle physics assumes not only reversibility, but also so-called CPT invariance. This means that same rules should apply if one not only reverses the direction of time (T), but also simultaneously inverts all spatial coordinates (P) and conjugates all charges (C), replacing particles by antiparticles. In a certain mathematical sense, CPT invariance can be viewed as a generalization of relativistic invariance: with a speed faster than light, something close to an ordinary relativistic transformation is a CPT transformation.

Originally it was assumed that C, P and T would all separately be invariances, as they are in classical mechanics. But in 1957 it was discovered that in radioactive beta decay, C and P are in a sense each maximally violated: among other things, the correlation between spin and motion direction is exactly opposite for neutrinos and for antineutrinos that are emitted. Despite this, it was still assumed that CP and T

would be true invariances. But in 1964 these too were found to be violated. Starting with a pure beam of K^0 particles, it turns out that quantum mechanical mixing processes lead after about 10^{-8} seconds to a certain mixture of \bar{K}^0 particles—the antiparticles of the K^0 . And what effectively happens is that the amount of mixing differs by about 0.1% in the positive and negative time directions. (What is actually observed is a small probability for the long-lived component of a K^0 beam to decay into two rather than three pions. Some analysis is required to connect this with T violation.) Particle physics experiments so far support exact CPT invariance. Simple models of gravity potentially suggest CPT violation (as a consequence of deviations from pure special relativistic invariance), but such effects tend to disappear when the models are refined.

■ **History of thermodynamics.** Basic physical notions of heat and temperature were established in the 1600s, and scientists of the time appear to have thought correctly that heat is associated with the motion of microscopic constituents of matter. But in the 1700s it became widely believed that heat was instead a separate fluid-like substance. Experiments by James Joule and others in the 1840s put this in doubt, and finally in the 1850s it became accepted that heat is in fact a form of energy. The relation between heat and energy was important for the development of steam engines, and in 1824 Sadi Carnot had captured some of the ideas of thermodynamics in his discussion of the efficiency of an idealized engine. Around 1850 Rudolf Clausius and William Thomson (Kelvin) stated both the First Law—that total energy is conserved—and the Second Law of Thermodynamics. The Second Law was originally formulated in terms of the fact that heat does not spontaneously flow from a colder body to a hotter. Other formulations followed quickly, and Kelvin in particular understood some of the law's general implications. The idea that gases consist of molecules in motion had been discussed in some detail by Daniel Bernoulli in 1738, but had fallen out of favor, and was revived by Clausius in 1857. Following this, James Clerk Maxwell in 1860 derived from the mechanics of individual molecular collisions the expected distribution of molecular speeds in a gas. Over the next several years the kinetic theory of gases developed rapidly, and many macroscopic properties of gases in equilibrium were computed. In 1872 Ludwig Boltzmann constructed an equation that he thought could describe the detailed time development of a gas, whether in equilibrium or not. In the 1860s Clausius had introduced entropy as a ratio of heat to temperature, and had stated the Second Law in terms of the increase of this quantity. Boltzmann then showed that his

equation implied the so-called H Theorem, which states that a quantity equal to entropy in equilibrium must always increase with time. At first, it seemed that Boltzmann had successfully proved the Second Law. But then it was noticed that since molecular collisions were assumed reversible, his derivation could be run in reverse, and would then imply the opposite of the Second Law. Much later it was realized that Boltzmann's original equation implicitly assumed that molecules are uncorrelated before each collision, but not afterwards, thereby introducing a fundamental asymmetry in time. Early in the 1870s Maxwell and Kelvin appear to have already understood that the Second Law could not formally be derived from microscopic physics, but must somehow be a consequence of human inability to track large numbers of molecules. In responding to objections concerning reversibility Boltzmann realized around 1876 that in a gas there are many more states that seem random than seem orderly. This realization led him to argue that entropy must be proportional to the logarithm of the number of possible states of a system, and to formulate ideas about ergodicity. The statistical mechanics of systems of particles was put in a more general context by Willard Gibbs, beginning around 1900. Gibbs introduced the notion of an ensemble—a collection of many possible states of a system, each assigned a certain probability. He argued that if the time evolution of a single state were to visit all other states in the ensemble—the so-called ergodic hypothesis—then averaged over a sufficiently long time a single state would behave in a way that was typical of the ensemble. Gibbs also gave qualitative arguments that entropy would increase if it were measured in a “coarse-grained” way in which nearby states were not distinguished. In the early 1900s the development of thermodynamics was largely overshadowed by quantum theory and little fundamental work was done on it. Nevertheless, by the 1930s, the Second Law had somehow come to be generally regarded as a principle of physics whose foundations should be questioned only as a curiosity. Despite neglect in physics, however, ergodic theory became an active area of pure mathematics, and from the 1920s to the 1960s properties related to ergodicity were established for many kinds of simple systems. When electronic computers became available in the 1950s, Enrico Fermi and others began to investigate the ergodic properties of nonlinear systems of springs. But they ended up concentrating on recurrence phenomena related to solitons, and not looking at general questions related to the Second Law. Much the same happened in the 1960s, when the first simulations of hard sphere gases were led to concentrate on the specific phenomenon of long-time tails. And by the 1970s, computer experiments were mostly oriented towards ordinary

differential equations and strange attractors, rather than towards systems with large numbers of components, to which the Second Law might apply. Starting in the 1950s, it was recognized that entropy is simply the negative of the information quantity introduced in the 1940s by Claude Shannon. Following statements by John von Neumann, it was thought that any computational process must necessarily increase entropy, but by the early 1970s, notably with work by Charles Bennett, it became accepted that this is not so (see page 1018), laying some early groundwork for relating computational and thermodynamic ideas.

■ **Current thinking on the Second Law.** The vast majority of current physics textbooks imply that the Second Law is well established, though with surprising regularity they say that detailed arguments for it are beyond their scope. More specialized articles tend to admit that the origins of the Second Law remain mysterious. Most ultimately attribute its validity to unknown constraints on initial conditions or measurements, though some appeal to external perturbations, to cosmology or to unknown features of quantum mechanics.

An argument for the Second Law from around 1900, still reproduced in many textbooks, is that if a system is ergodic then it will visit all its possible states, and the vast majority of these will look random. But only very special kinds of systems are in fact ergodic, and even in such systems, the time necessary to visit a significant fraction of all possible states is astronomically long. Another argument for the Second Law, arising from work in the 1930s and 1940s, particularly on systems of hard spheres, is based on the notion of instability with respect to small changes in initial conditions. The argument suffers however from the same difficulties as the ones for chaos theory discussed in Chapter 6 and does not in the end explain in any real way the origins of randomness, or the observed validity of the Second Law.

With the Second Law accepted as a general principle, there is confusion about why systems in nature have not all dissipated into complete randomness. And often the rather absurd claim is made that all the order we see in the universe must just be a fluctuation—leaving little explanatory power for principles such as the Second Law.

■ **My explanation of the Second Law.** What I say in this book is not incompatible with much of what has been said about the Second Law before; it is simply that I make more definite some key points that have been left vague before. In particular, I use notions of computation to specify what kinds of initial conditions can reasonably be prepared, and what kinds of measurements can reasonably be made. In a sense

what I do is just to require that the operation of coarse graining correspond to a computation that is less sophisticated than the actual evolution of the system being studied. (See also Chapters 10 and 12.)

■ **Biological systems and Maxwell's demon.** Unlike most physical systems, biological systems typically seem capable of spontaneously organizing themselves. And as a result, even the original statements of the Second Law talked only about “inanimate systems”. In the mid-1860s James Clerk Maxwell then suggested that a demon operating at a microscopic level could reduce the randomness of a system such as a gas by intelligently controlling the motion of molecules. For many years there was considerable confusion about Maxwell's demon. There were arguments that the demon must use a flashlight that generates entropy. And there were extensive demonstrations that actual biological systems reduce their internal entropy only at the cost of increases in the entropy of their environment. But in fact the main point is that if the evolution of the whole system is to be reversible, then the demon must store enough information to reverse its own actions, and this limits how much the demon can do, preventing it, for example, from unscrambling a large system of gas molecules.

■ **Self-gravitating systems.** The observed existence of structures such as galaxies might lead one to think that any large number of objects subject to mutual gravitational attraction might not follow the Second Law and become randomized, but might instead always form orderly clumps. It is difficult to know, however, what an idealized self-gravitating system would do. For in practice, issues such as the limited size of a galaxy, its overall rotation, and the details of stellar collisions all seem to have major effects on the results obtained. (And it is presumably not feasible to do a small-scale experiment, say in Earth orbit.) There are known to be various instabilities that lead in the direction of clumping and core collapse, but how these weigh against effects such as the transfer of energy into tight binding of small groups of stars is not clear. Small galaxies such as globular clusters that contain less than a million stars seem to exhibit a certain uniformity which suggests a kind of equilibrium. Larger galaxies such as our own that contain perhaps 100 billion stars often have intricate spiral or other structure, whose origin may be associated with gravitational effects, or may be a consequence of detailed processes of star formation and explosion. (There is some evidence that older galaxies of a given size tend to develop more regularities in their structure.) Current theories of the early universe tend to assume that galaxies originally began to form as a result of density fluctuations of non-gravitational origin (and reflected

in the cosmic microwave background). But there is evidence that a widespread fractal structure develops—with a correlation function of the form $r^{-1.8}$ —in the distribution of stars in our galaxy, galaxies in clusters and clusters in superclusters, perhaps suggesting the existence of general overall laws for self-gravitating systems. (See also page 973.)

As mentioned on page 880, it so happens that my original interest in cellular automata around 1981 developed in part from trying to model the growth of structure in self-gravitating systems. At first I attempted to merge and generalize ideas from traditional areas of mathematical physics, such as kinetic theory, statistical mechanics and field theory. But then, particularly as I began to think about doing explicit computer simulations, I decided to take a different tack and instead to look for the most idealized possible models. And in doing this I quickly came up with cellular automata. But when I started to investigate cellular automata, I discovered some very remarkable phenomena, and I was soon led away from self-gravitating systems, and into the process of developing the much more general science in this book. Over the years, I have occasionally come back to the problem of self-gravitating systems, but I have never succeeded in developing what I consider to be a satisfactory approach to them.

■ **Cosmology and the Second Law.** In the standard big bang model it is assumed that all matter in the universe was initially in completely random thermal equilibrium. But such equilibrium implies uniformity, and from this it follows that the initial conditions for the gravitational forces in the universe must have been highly regular, resulting in simple overall expansion, rather than random expansion in some places and contraction in others. As I discuss on page 1026 I suspect that in fact the universe as a whole probably had what were ultimately very simple initial conditions, and it is just that the effective rules for the evolution of matter led to rapid randomization, whereas those for gravity did not.

■ **Alignment of time in the universe.** Evidence from astronomy clearly suggests that the direction of irreversible processes is the same throughout the universe. The reason for this is presumably that all parts of the universe are expanding—with the local consequence that radiation is more often emitted than absorbed, as evidenced by the fact that the night sky is dark. Olbers' paradox asks why one does not see a bright star in every direction in the night sky. The answer is that locally stars are clumped, and light from stars further away is progressively red-shifted to lower energy. Focusing a larger and larger distance away, the light one sees was emitted longer and longer ago. And eventually one sees light emitted when the universe was filled with hot opaque

gas—now red-shifted to become the 2.7K cosmic microwave background.

■ **Poincaré recurrence.** Systems of limited size that contain only discrete elements inevitably repeat their evolution after a sufficiently long time (see page 258). In 1890 Henri Poincaré established the somewhat less obvious fact that even continuous systems also always eventually get at least arbitrarily close to repeating themselves. This discovery led to some confusion in early interpretations of the Second Law, but the huge length of time involved in a Poincaré recurrence makes it completely irrelevant in practice.

■ **Page 446 • Billiards.** The discrete system I consider here is analogous to continuous so-called billiard systems consisting of circular balls in the plane. The simplest case involves one ball bouncing around in a region of a definite shape. In a rectangular region, the position is given by $\text{Mod}[at, \{w, h\}]$ and every point will be visited if the parameters have irrational ratios. In a region that contains fixed circular obstructions, the motion can become sensitively dependent on initial conditions. (This setup is similar to a so-called Lorentz gas.) For a system of balls in a region with cyclic boundaries, a complicated proof due to Yakov Sinai from the 1960s purports to show that every ball eventually visits every point in the region, and that certain simple statistical properties of trajectories are consistent with randomness. (See also page 971.)

■ **Page 449 • Entropy of particles in a box.** The number of possible states of a region of m cells containing q particles is $\text{Binomial}[m, q]$. In the large size limit, the logarithm of this can be approximated by $q \text{Log}[m/q]/m$.

■ **Page 457 • Periods in rule 37R.** With a system of size n , the maximum possible repetition period is 2^{2^n} . In actuality, however, the periods are considerably shorter. With all cells 0 on one step, and a block of nonzero cells on the next step, the periods are for example: $\{1\}$: 21; $\{1, 1\}$: $3n-8$; $\{1, 0, 1\}$: 666; $\{1, 1, 1\}$: $3n-8$; $\{1, 0, 0, 1\}$: irregular ($< 24n$; peaks at $6j+1$); $\{1, 0, 0, 1, 0, 1\}$: irregular ($\leq 2^n$; 857727 for $n=26$; 13705406 for $n=100$). With completely random initial conditions, there are great fluctuations, but a typical period is around $2^{n/3}$.

Conserved Quantities and Continuum Phenomena

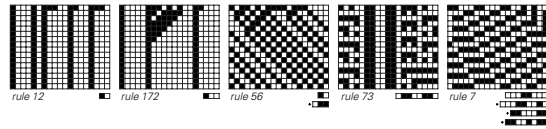
■ **Physics.** The quantities in physics that so far seem to be exactly conserved are: energy, momentum, angular momentum, electric charge, color charge, lepton number (as well as electron number, muon number and τ lepton number) and baryon number.

■ **Implementation.** Whether a k -color cellular automaton with range r conserves total cell value can be determined from

```
Catch[Do[
  (If[Apply[Plus, CAStep[rule, #] - #] != 0, Throw[False]] &)[
  IntegerDigits[i, k, m]], {m, w}, {i, 0, k^m - 1}]; True]
```

where w can be taken to be k^{2r} , and perhaps smaller. Among the 256 elementary cellular automata just 5 conserve total cell value. Among the 2^{32} $k=2$, $r=2$ rules 428 do, and of these 2 are symmetric, and 6 are reversible, and all these are just shift and identity rules.

■ **More general conserved quantities.** Some rules conserve not total numbers of cells with given colors, but rather total numbers of blocks of cells with given forms—or combinations of these. The pictures below show the simplest quantities of these kinds that end up being conserved by various elementary rules.



Among the 256 elementary rules, the total numbers that have conserved quantities involving at most blocks of lengths 1 through 10 are $\{5, 38, 66, 88, 102, 108, 108, 114, 118, 118\}$.

Rules that show complicated behavior usually do not seem to have conserved quantities, and this is true for example of rules 30, 90 and 110, at least up to blocks of length 10.

One can count the number of occurrences of each of the k^b possible blocks of length b in a given state using

```
BC[list.]:=
  With[{z = Map[FromDigits[#, k] &, Partition[list, b, 1, 1]]},
  Map[Count[z, #] &, Range[0, k^b - 1]]]
```

Conserved quantities of the kind discussed here are then of the form $q \cdot \text{BC}[a]$ where q is some fixed list. A way to find candidates for q is to compute

```
NullSpace[Table[With[{u = Table[Random[Integer,
  {0, k-1}], {m}], BC[CAStep[u] - BC[u]], {s}]]
```

for progressively larger m and s , and to see what lists continue to appear. For block size b , k^{b-1} lists will always appear as a result of trivial conserved quantities. (With $k=2$, for $b=1$, $\{1, 1\}$ represents conservation of the total number of cells, regardless of color, while for $b=2$, $\{1, 1, 1, 1\}$ represents the same thing, while $\{0, 1, -1, 0\}$ represents the fact that in going along in any state the number of black-to-white transitions must equal the number of white-to-black ones.) If more than k^{b-1} lists appear, however, then some must correspond to genuine non-trivial conserved quantities. To identify any such quantity with certainty, it turns out to be enough to look at the k^{b+2r-1} states where no block of length

$b + 2r - 1$ appears more than once (and perhaps even just some fairly small subset of these).

(See also page 981.)

■ **Other conserved quantities.** The conserved quantities discussed so far can all be thought of as taking values assigned to blocks of different kinds in a given state and then just adding them up as ordinary numbers. But one can also imagine using other operations to combine such values. Addition modulo n can be handled by inserting `Modulus → n` in `NullSpace` in the previous note. And doing this shows for example that rule 150 conserves the total number of black cells modulo 2. But in general not many additional conserved quantities are found in this way. One can also consider combining values of blocks by the multiplication operation in a group—and seeing whether the conjugacy class of the result is conserved.

■ **PDEs.** In the early 1960s it was discovered that certain nonlinear PDEs support an infinite number of distinct conserved quantities, associated with so-called integrability and the presence of solitons. Systematic methods now exist to find conserved quantities that are given by integrals of polynomials of any given degree in the dependent variables and their derivatives. Most randomly chosen PDEs appear, however, to have no such conserved quantities.

■ **Local conservation laws.** Whenever a system like a cellular automaton (or PDE) has a global conserved quantity there must always be a local conservation law which expresses the fact that every point in the system the total flux of the conserved quantity into a particular region must equal the rate of increase of the quantity inside it. (If the conserved quantity is thought of like charge, the flux is then current.) In any 1D $k = 2, r = 1$ cellular automaton, it follows from the basic structure of the rule that one can tell what the difference in values of a particular cell on two successive steps will be just by looking at the cell and its immediate neighbor on each side. But if the number of black cells is conserved, then one can compute this difference instead by defining a suitable flux, and subtracting its values on the left and right of the cell. What the flux should be depends on the rule. For rule 184, it can be taken to be 1 for each \blacksquare block, and to be 0 otherwise. For rule 170, it is 1 for both \square and \blacksquare . For rule 150, it is 1 for \square and \blacksquare , with all computations done modulo 2. In general, if the global conserved quantity involves blocks of size b , the flux can be computed by looking at blocks of size $b + 2r - 1$. What the values for these blocks should be can be found by solving a system of linear equations; that a solution must exist can be seen by looking at the de Bruijn network (see page 941), with nodes labelled by size $b + 2r - 1$ blocks,

and connections by value differences between size b blocks at the center of the possible size $b + 2r$ blocks. (Note that the same basic kind of setup works in any number of dimensions.)

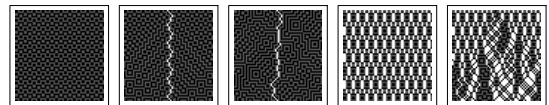
■ **Block cellular automata.** With a rule of the form $\{(1, 1) \rightarrow \{1, 1\}, \{1, 0\} \rightarrow \{1, 0\}, \{0, 1\} \rightarrow \{0, 0\}, \{0, 0\} \rightarrow \{0, 1\}\}$ the evolution of a block cellular automaton with blocks of size n can be implemented using

```
BCAEvolveList[{n_Integer, rule_}, init_, t_] :=
  FoldList[BCAStep[{n, rule}, #1, #2] &, init, Range[t]] /;
  Mod[Length[init], n] == 0
BCAStep[{n_, rule_}, a_, d_] := RotateRight[
  Flatten[Partition[RotateLeft[a, d], n] /. rule], d]
```

Starting with a single black cell, none of the $k = 2, n = 2$ block cellular automata generate anything beyond simple nested patterns. In general, there are $k^{n k^n}$ possible rules for block cellular automata with k colors and blocks of size n . Of these, $k^n!$ are reversible. For $k = 2$, the number of rules that conserve the total number of black cells can be computed from $q = \text{Binomial}[n, \text{Range}[0, n]]$ as $\text{Apply}[\text{Times}, q^q]$. The number of these rules that are also reversible is $\text{Apply}[\text{Times}, q!]$. In general, a block cellular automaton is reversible only if its rule simply permutes the k^n possible blocks.

Compressing each block into a single cell, and n steps into one, any block cellular automaton with k colors and block size n can be translated directly into an ordinary cellular automaton with k^n colors and range $r = n/2$.

■ **Page 461 · Block rules.** These pictures show the behavior of rule (c) starting from some special initial conditions.



The repetition period with a total of n cells can be 3^n steps. With random initial conditions, the period is typically up to about $3^{n/2}$. Starting with a block of q black cells, the period can get close to this. For $n = 20, q = 17$, for example, it is 31,300.

Note that even in rule (b) wraparound phenomena can lead to repetition periods that increase rapidly with n (e.g. 4820 for $n = 20, q = 15$), but presumably not exponentially.

In rule (d), the repetition periods can typically be larger than in rule (c): e.g. 803,780 for $n = 20, q = 13$.

■ **Page 464 · Limiting procedures.** Several different limiting procedures all appear to yield the same continuum behavior for the cellular automata shown here. In the pictures on this

page a large ensemble of different initial conditions is considered, and the density of each individual cell averaged over this ensemble is computed. In a more direct analogy to actual physical systems, one would consider instead a very large number of cells, then compute the density in a single state of the system by averaging over regions that contain many cells but are nevertheless small compared to the size of the whole system.

■ **PDE approximations.** Cellular automaton (d) in the main text can be viewed as minimal discrete approximations to the diffusion equation. The evolution of densities in the ensemble average is analogous to a traditional finite difference method with a real number at each site. The cellular automaton itself uses in effect a distributed representation of the density.

■ **Diffusion equation.** In an appropriate limit the density distribution for cellular automaton (d) appears to satisfy the usual diffusion equation $\partial_t f[x, t] = c \partial_{xx} f[x, t]$ discussed on page 163. The solution to this equation with an impulse initial condition is $\text{Exp}[-x^2/t]$, and with a block from $-a$ to a it is $(\text{Erf}[(a-x)/\sqrt{t}] + \text{Erf}[(a+x)/\sqrt{t}])/a$.

■ **Derivation of the diffusion equation.** With some appropriate assumptions, it is fairly straightforward to derive the usual diffusion equation from a cellular automaton. Let the density of black cells at position x and time t be $f[x, t]$, where this density can conveniently be computed by averaging over many instances of the system. If we assume that the density varies slowly with position and time, then we can make series expansions such as

$$f[x+dx, t] = f[x, t] + \partial_x f[x, t] dx + 1/2 \partial_{xx} f[x, t] dx^2 + \dots$$

where the coordinates are scaled so that adjacent cells are at positions $x-dx$, x , $x+dx$, etc. If we then assume perfect underlying randomness, the density at a particular position must be given in terms of the densities at neighboring positions on the previous step by

$$f[x, t+dt] = p_1 f[x-dx, t] + p_2 f[x, t] + p_3 f[x+dx, t]$$

Density conservation implies that $p_1 + p_2 + p_3 = 1$, while left-right symmetry implies $p_1 = p_3$. And from this it follows that

$$f[x, t+dt] = c(f[x-dx, t] + f[x+dx, t]) + (1-2c)f[x, t]$$

Performing a series expansion then yields

$$f[x, t+dt] + \partial_t f[x, t] dt = f[x, t] + c dx^2 \partial_{xx} f[x, t]$$

which in turn gives exactly the usual 1D diffusion equation $\partial_t f[x, t] = \xi \partial_{xx} f[x, t]$, where ξ is the diffusion coefficient for the system. I first gave this derivation in 1986, together with extensive generalizations.

■ **Page 464 · Non-standard diffusion.** To get ordinary diffusion behavior of the kind that occurs in gases—and is described by the diffusion equation—it is in effect necessary to have

perfect uncorrelated randomness, with no structure that persists too long. But for example in the rule (a) picture on page 463 there is in effect a block of solid that persists in the middle—so that no ordinary diffusion behavior is seen. In rule (c) there is considerable apparent randomness, but it turns out that there are also fluctuations that last too long to yield ordinary diffusion. And thus for example whenever there is a structure containing s identical cells (as on page 462), this typically takes about s^2 steps to decay away. The result is that on page 464 the limiting form of the average behavior does not end up being an ordinary Gaussian.

■ **Conservation of vector quantities.** Conservation of the total number of colored cells is analogous to conservation of a scalar quantity such as energy or particle number. One can also consider conservation of a vector quantity such as momentum which has not only a magnitude but also a direction. Direction makes little sense in 1D, but is meaningful in 2D. The 2D cellular automaton used as a model of an idealized gas on page 446 provides an example of a system that can be viewed as conserving a vector quantity. In the absence of fixed scatterers, the total fluxes of particles in the horizontal and the vertical directions are conserved. But in a sense there is too much conservation in this system, and there is no interaction between horizontal and vertical motions. This can be achieved by having more complicated underlying rules. One possibility is to use a hexagonal rather than square grid, thereby allowing six particle directions rather than four. On such a grid it is possible to randomize microscopic particle motions, but nevertheless conserve overall momenta. This is essentially the model used in my discussion of fluids on page 378.

Ultimate Models for the Universe

■ **History of ultimate models.** From the earliest days of Greek science until well into the 1900s, it seems to have often been believed that an ultimate model of the universe was not far away. In antiquity there were vague ideas about everything being made of elements like fire and water. In the 1700s, following the success of Newtonian mechanics, a common assumption seems to have been that everything (with the possible exception of light) must consist of tiny corpuscles with gravity-like forces between them. In the 1800s the notion of fields—and the ether—began to develop, and in the 1880s it was suggested that atoms might be knotted vortices in the ether (see page 1044). When the electron was discovered in 1897 it was briefly thought that it might be the fundamental constituent of everything. And later it was imagined that perhaps electromagnetic fields could underlie

everything. Then after the introduction of general relativity for the gravitational field in 1915, there were efforts, especially in the 1930s, to introduce extensions that would yield unified field theories of everything (see page 1028). By the 1950s, however, an increasing number of subatomic particles were being found, and most efforts at unification became considerably more modest. In the 1960s the quark model began to explain many of the particles that were seen. Then in the 1970s work in quantum field theory encouraged the use of gauge theories and by the late 1970s the so-called Standard Model had emerged, with the Weinberg-Salam $SU(2) \otimes U(1)$ gauge theory for weak interactions and electromagnetism, and the QCD $SU(3)$ gauge theory for strong interactions. The discoveries of the c quark, τ lepton and b quark were largely unexpected, but by the late 1970s there was widespread enthusiasm for the idea of a single “grand unified” gauge theory, based say on $SU(5)$, that would explain all forces except gravity. By the mid-1980s failure to observe expected proton decay cast doubts on simple versions of such models, and various possibilities based on supersymmetry and groups like $SO(10)$ were considered. Occasional attempts to construct quantum theories of gravity had been made since the 1930s, and in the late 1980s these began to be pursued more vigorously. In the mid-1980s the discovery that string theory could be given various mathematical features that were considered desirable made it emerge as the main hope for an ultimate “theory of everything”. But despite all sorts of elegant mathematical work, the theory remains rather distant from observed features of our universe. In some parts of particle physics, it is still sometimes claimed that an ultimate theory is not far away, but outside it generally seems to be assumed that physics is somehow instead an endless frontier—that will continue to yield a stream of surprising and increasingly complex discoveries forever—with no ultimate theory ever being found.

■ **Theological implications.** Some may view an ultimate model of the universe as “leaving no room for a god”, while others may view it as a direct reflection of the existence of a god. In any case, knowing a complete and ultimate model does make it impossible to have miracles or divine interventions that come from outside the laws of the universe—though working out what will happen on the basis of these laws may nevertheless be irreducibly difficult.

■ **Origins of physical models.** Considering the reputation of physics as an empirical science, it is remarkable how many significant theories were in fact first constructed on largely aesthetic grounds. Notable examples include Maxwell’s equations for electromagnetism (1880s), general relativity

(1915), the Dirac equation for relativistic electrons (1928), and QCD (early 1970s). This history makes it seem more plausible that one might be able to come up with an ultimate model of physics on largely aesthetic grounds, rather than mainly by working from detailed experimental observations.

■ **Simplicity in scientific models.** To curtail absurdly complicated early scientific models Occam’s razor principle that “entities should not be multiplied beyond necessity” was introduced in the 1300s. This principle has worked well in physics, where it has often proven to be the case, for example, that out of all possible terms in an equation the only ones that actually occur are the very simplest. But in a field like biology, the principle has usually been regarded as much less successful. For many complicated features are seen in biological organisms, and when there have been guesses of simple explanations for them, these have often turned out to be wrong. Much of what is seen is probably a reflection of complicated details of the history of biological evolution. But particularly after the discoveries in this book it seems likely that at least some of what appears complicated may actually be produced by very simple underlying programs—which perhaps occur because they were the first to be tried, or are the most efficient or robust. Outside of natural science, Occam’s principle can sometimes be useful—typically because simplicity is a good assumption in some aspect of human behavior or motivation. In looking at well-developed technological systems or human organizations simplicity is also quite often a reasonable assumption—since over the course of time parts that are complicated or difficult to understand will tend to have been optimized away.

■ **Numerology.** Ever since the Pythagoreans many attempts to find truly ultimate models of the universe have ended up centering on derivations of numbers that are somehow thought to be characteristic of the universe. In the past century, the emphasis has been on physical constants such as the fine structure constant $\alpha \approx 1/137.0359896$, and usually the idea is that such constants arise directly from counting objects of some specified type using traditional discrete mathematics. A notable effort along these lines was made by Arthur Eddington in the mid-1930s, and certainly over the past twenty or so years I have received a steady stream of mail presenting such notions with varying degrees of obscurity and mysticism. But while I believe that every feature of our universe does indeed come from an ultimate discrete model, I would be very surprised if the values of constants which happen to be easy for us to measure in the end turn out to be given by simple traditional mathematical formulas.

■ **Emergence of simple laws.** In statistical physics it is seen that universal and fairly simple overall laws often emerge

even in systems whose underlying molecular or other structure can be quite complicated. The basic origin of this phenomenon is the averaging effect of randomness discussed in Chapter 7 (technically, it is the survival only of leading operators at renormalization group fixed points). The same phenomenon is also seen in quantum field theory, where it is essentially a consequence of the averaging effect of quantum fluctuations, which have a direct mathematical analog to statistical physics.

■ **Apparent simplicity.** Given any rules it is always possible to develop a form of description in which these rules will be considered simple. But what is interesting to ask is whether the underlying rules of the universe will seem simple—or special, say in their elegance or symmetry—with respect to forms of description that we as humans currently use.

■ **Mechanistic models.** Until quite recently, it was generally assumed that if one were able to get at the microscopic constituents of the universe they would look essentially like small-scale analogs of objects familiar from everyday life. And so, for example, the various models of atoms from the end of the 1800s and beginning of the 1900s were all based on familiar mechanical systems. But with the rise of quantum mechanics it came to be believed throughout mainstream physics that any true fundamental model must be abstract and mathematical—and never ultimately amenable to any kind of direct mechanistic description. Occasionally there have been mechanistic descriptions used—as in the parton and bag models, and various continuum models of high-energy collisions—but they have typically been viewed only as convenient rough approximations. (Feynman diagrams may also seem superficially mechanistic, but are really just representations of quite abstract mathematical formulas.) And indeed since at least the 1960s mechanistic models have tended to carry the stigma of uninformed amateur science.

With the rise of computers there began to be occasional discussion—though largely outside of mainstream science—that the universe might have a mechanism related to computers. Since the 1950s science fiction has sometimes featured the idea that the universe or some part of it—such as the Earth—could be an intentionally created computer, or that our perception of the universe could be based on a computer simulation. Starting in the 1950s a few computer scientists considered the idea that the universe might have components like a computer. Konrad Zuse suggested that it could be a continuous cellular automaton; Edward Fredkin an ordinary cellular automaton (compare page 1027). And over the past few decades—normally in the context of amateur science—there have been a steady stream of systems like cellular automata constructed to have elements

reminiscent of observed particles or forces. From the point of view of mainstream physics, such models have usually seemed quite naive. And from what I say in the main text, no such literal mechanistic model can ever in the end realistically be expected to work. For if an ultimate model is going to be simple, then in a sense it cannot have room for all sorts of elements that are immediately recognizable in terms of everyday known physics. And instead I believe that what must happen relies on the phenomena discovered in this book—and involves the emergence of complex properties without any obvious underlying mechanistic set up. (Compare page 860.)

■ **The Anthropic Principle.** It is sometimes argued that the reason our universe has the characteristics it does is because otherwise an intelligence such as us could not have arisen to observe it. But to apply such an argument one must among other things assume that we can imagine all the ways in which intelligence could conceivably operate. Yet as we have seen in this book it is possible for highly complex behavior—ultimately not dissimilar to intelligence—to arise from simple programs in ways that we never came even close to imagining. And indeed, as we discuss in Chapter 12, it seems likely that above a fairly low threshold the vast majority of underlying rules can in fact in some way or another support arbitrarily complex computations—potentially allowing something one might call intelligence in a vast range of very different universes. (See page 822.)

■ **Physics versus mathematics.** Theoretical physics can be viewed as taking physical input in the form of models and then using mathematics to work out the consequences. If I am correct that there is a simple underlying program for the universe, then this means that theoretical physics must at some level have only a very small amount of true physical input—and the rest must in a sense all just be mathematics.

■ **Initial conditions.** To find the behavior of the universe one potentially needs to know not only its rule but also its initial conditions. Like the rule, I suspect that the initial conditions will turn out to be simple. And ultimately there should be traces of such simplicity in, say, the distribution of galaxies or the cosmic microwave background. But ideas like those on page 1055—as well as inflation—tend to suggest that we currently see only a tiny fraction of the whole universe, making it very difficult for example to recognize overall geometrical regularities. And it could also be that even though there might ultimately have been simple initial conditions, the current phase of our universe might be the result of some sequence of previous phases, and so effectively have much more complicated initial conditions. (Proposals discussed in quantum cosmology since the 1980s

that for example just involve requiring the universe to satisfy final but not initial boundary condition constraints do not fit well into my kinds of models.)

■ **Consequences of an ultimate model.** Even if one knows an ultimate model for the universe, there will inevitably be irreducible difficulty in working out all its consequences. Indeed, questions like “does there exist a way to transmit information faster than light?” may boil down to issues analogous to whether it is possible to construct a configuration that has a certain property in, say, the rule 110 cellular automaton. And while some such questions may be answered by fairly straightforward computational or mathematical means, there will be no upper bound on the amount of effort that it could take to answer any particular question.

■ **Meaning of the universe.** If the whole history of our universe can be obtained by following definite simple rules, then at some level this history has the same kind of character as a construct such as the digit sequence of π . And what this suggests is that it makes no more or less sense to talk about the meaning of phenomena in our universe as it does to talk about the meaning of phenomena in the digit sequence of π .

The Nature of Space

■ **History of discrete space.** The idea that matter might be made up of discrete particles existed in antiquity (see page 876), and occasionally the notion was discussed that space might also be discrete—and that this might for example be a way of avoiding issues like Zeno’s paradox. In 1644 René Descartes proposed that space might initially consist of an array of identical tiny discrete spheres, with motion then occurring through chains of these spheres going around in vortices—albeit with pieces being abraded off. But with the rise of calculus in the 1700s all serious fundamental models in physics began to assume continuous space. In discussing the notion of curved space, Bernhard Riemann remarked in 1854 that it would be easier to give a general mathematical definition of distance if space were discrete. But since physical theories seemed to require continuous space, the necessary new mathematics was developed and almost universally used—though for example in 1887 William Thomson (Kelvin) did consider a discrete foam-like model for the ether (compare page 988). Starting in 1930, difficulties with infinities in quantum field theory again led to a series of proposals that spacetime might be discrete. And indeed by the late 1930s this notion was fairly widely discussed as a possible inevitable feature of quantum mechanics. But there were problems with relativistic

invariance, and after ideas of renormalization developed in the 1940s, discrete space seemed unnecessary, and has been out of favor ever since. Some non-standard versions of quantum field theory involving discrete space did however continue to be investigated into the 1960s, and by then a few isolated other initiatives had arisen that involved discrete space. The idea that space might be defined by some sort of causal network of discrete elementary quantum events arose in various forms in work by Carl von Weizsäcker (ur-theory), John Wheeler (pregeometry), David Finkelstein (spacetime code), David Bohm (topochronology) and Roger Penrose (spin networks; see page 1055). General arguments for discrete space were also sometimes made—notably by Edward Fredkin, Marvin Minsky and to some extent Richard Feynman—on the basis of analogies to computers and in particular the idea that a given region of space should contain only a finite amount of information. In the 1980s approximation schemes such as lattice gauge theory and later Regge calculus (see page 1054) that take space to be discrete became popular, and it was occasionally suggested that versions of these could be exact models. There have been a variety of continuing initiatives that involve discrete space, with names like combinatorial physics—but most have used essentially mechanistic models (see page 1026), and none have achieved significant mainstream acceptance. Work on quantum gravity in the late 1980s and 1990s led to renewed interest in the microscopic features of spacetime (see page 1054). Models that involve discreteness have been proposed—most often based on spin networks—but there is usually still some form of continuous averaging present, leading for example to suggestions very different from mine that perhaps this could lead to the traditional continuum description through some analog of the wave-particle duality of elementary quantum mechanics. I myself became interested in the idea of completely discrete space in the mid-1970s, but I could not find a plausible framework for it until I started thinking about networks in the mid-1980s.

■ **Planck length.** Even in existing particle physics it is generally assumed that the traditional simple continuum description of space must break down at least below about the Planck length $\text{Sqrt}[\hbar G/c^3] \approx 2 \times 10^{-35}$ meters—since at this scale dimensional analysis suggests that quantum effects should be comparable in magnitude to gravitational ones.

■ **Page 472 · Symmetry.** A system like a cellular automaton that consists of a large number of identical cells must in effect be arranged like a crystal, and therefore must exhibit one of the limited number of possible crystal symmetries in any particular dimension, as discussed on page 929. And even a

generalized cellular automaton constructed say on a Penrose tiling still turns out to have a discrete spatial symmetry.

■ **Page 474 · Space and its contents.** A number of somewhat different ideas about space were discussed in antiquity. Around 375 BC Plato vaguely suggested that the universe might consist of large numbers of abstract polyhedra. A little later Aristotle proposed that space is set up so as to provide a definite place for everything—and in effect to force it there. But in geometry as developed by Euclid there was at least a mathematical notion of space as a kind of uniform background. And by sometime after 300 BC the Epicureans developed the idea of atoms of matter existing in a mostly featureless void of space. In the Middle Ages there was discussion about how the non-material character of God might fit in with ideas about space. In the early 1600s the concept of inertia developed by Galileo implied that space must have a certain fundamental uniformity. And with the formulation of mechanics by Isaac Newton in 1687 space became increasingly viewed as something purely abstract, quite different in character from material objects which exist in it. Philosophers had meanwhile discussed matter—as opposed to mind—being something characterized by having spatial extent. And for example in 1643 Thomas Hobbes suggested that the whole universe might be made of the same continuous stuff, with different densities of it corresponding to different materials, and geometry being just an abstract idealization of its properties. But in the late 1600s Gottfried Leibniz suggested instead that everything might consist of discrete monads, with space emerging from the pattern of relative distances between them. Yet with the success of Newtonian mechanics such ideas had by the late 1700s been largely forgotten—leading space almost always to be viewed just in simple abstract geometrical terms. The development of non-Euclidean geometry in the mid-1800s nevertheless suggested that even at the level of geometry space could in principle have a complicated structure. But in physics it was still assumed that space itself must have a standard fixed Euclidean form—and that everything in the universe must just exist in this space. By the late 1800s, however, it was widely believed that in addition to ordinary material objects, there must throughout space be a fluid-like ether with certain mechanical and electromagnetic properties. And in the 1860s it was even suggested that perhaps atoms might just correspond to knots in this ether (see page 1044). But this idea soon fell out of favor, and when relativity theory was introduced in 1905 it emphasized relations between material objects and in effect always treated space as just some kind of abstract background, with no real structure of its own. But in 1915 general relativity

introduced the idea that space could actually have a varying non-Euclidean geometry—and that this could represent gravity. Yet it was still assumed that matter was something different—that for example had to be represented separately by explicit terms in the Einstein equations. There were nevertheless immediate thoughts that perhaps at least electromagnetism could be like gravity and just arise from features of space. And in 1918 Hermann Weyl suggested that this could happen through local variations of scale or “gauge” in space, while in the 1920s Theodor Kaluza and Oskar Klein suggested that it could be associated with a fifth spacetime dimension of invisibly small extent. And from the 1920s to the 1950s Albert Einstein increasingly considered the possibility that there might be a unified field theory in which all matter would somehow be associated with the geometry of space. His main specific idea was to allow the metric of spacetime to be non-symmetric (see page 1052) and perhaps complex—with its additional components yielding electromagnetism. And he then tried to construct nonlinear field equations that would show no singularities, but would have solutions (perhaps analogous to the geons discussed on page 1054) that would exhibit various discrete features corresponding to particles—and perhaps quantum effects. But with the development of quantum field theory in the 1920s and 1930s most of physics again treated space as fixed and featureless—though now filled with various types of fields, whose excitations were set up to correspond to observed types of particles. Gravity has never fit very well into this framework. But it has always still been expected that in an ultimate quantum theory of gravity space will have to have a structure that is somehow like a quantum field. But when quantum gravity began to be investigated in earnest in the 1980s (see page 1054) most efforts concentrated on the already difficult problem of pure gravity—and did not consider how matter might enter. In the development of ordinary quantum field theories, supergravity theories studied in the 1980s did nominally support particles identified with gravitons, but were still formulated on a fixed background spacetime. And when string theory became popular in the 1980s the idea was again to have strings propagating in a background spacetime—though it turned out that for consistency this spacetime had to satisfy the Einstein equations. Consistency also typically required the basic spacetime to be 10-dimensional—with the reduction to observed 4D spacetime normally assumed to occur through restriction of the other dimensions to some kind of so-called Calabi-Yau manifold of small extent, associated excitations with various particles through an analog of the Kaluza-Klein mechanism. It has always been hoped that this kind of seemingly arbitrary setup would somehow automatically

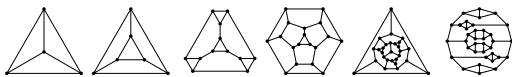
emerge from the underlying theory. And in the late 1990s there seemed to be some signs of this when dualities were discovered in various generalized string theories—notably for example between quantum particle excitations and gravitational black hole configurations. So while it remains impossible to work out all the consequences of string theories, it is conceivable that among the representations of such theories there might be ones in which matter can be viewed as just being associated with features of space.

Space as a Network

■ **Page 476 • Trivalent networks.** With n nodes and 3 connections at each node a network must always have an even number of nodes, and a total of $3n/2$ connections. Of all possible such networks, most large ones end up being connected. The number of distinct such networks for even n from 2 to 10 is {2, 5, 17, 71, 388}. If no self connections are allowed then these numbers become {1, 2, 6, 20, 91}, while if neither self nor multiple connections are allowed (yielding what are often referred to as cubic or 3-regular graphs), the numbers become {0, 1, 2, 5, 19, 85, 509, 4060, 41301, 510489}, or asymptotically $(6n)!/((3n)!(2n)!288^n e^2)$. (For symmetric graphs see page 1032.) If one requires the networks to be planar the numbers are {0, 1, 1, 3, 9, 32, 133, 681, 3893, 24809, 169206}. If one looks at subnetworks with dangling connections, the number of these up to size 10 is {2, 5, 7, 22, 43, 141, 373, 1270, 4053, 14671}, or {1, 1, 2, 6, 10, 29, 64, 194, 531, 1733} if no self or multiple connections are allowed (see also page 1039).

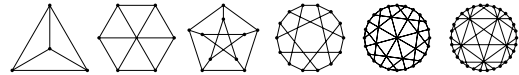
■ **Properties of networks.** Over the past century or so a variety of global properties of networks have been studied. Typical ones include:

- **Edge connectivity:** the minimum number of connections that must be removed to make the network disconnected.
- **Diameter:** the maximum distance between any two nodes in the network. The pictures below show the largest planar trivalent networks with diameters 1, 2 and 3, and the largest known ones with diameters 4, 5 and 6.



- **Circumference:** the length of the longest cycle in the network. Although difficult to determine in particular cases, many networks allow so-called Hamiltonian cycles that include every node. (Up to 8 nodes, all 8 trivalent networks have this property; up to 10 nodes 25 of 27 do.)

- **Girth:** the length of the shortest cycle in the network. The pictures below show the smallest trivalent networks with girths 3 through 8 (so-called cages). Girth can be relevant in seeing whether a particular cluster can ever occur in network.



- **Chromatic number:** the minimum of colors that can be assigned to nodes so that no adjacent nodes end up the same color. It follows from the Four-Color Theorem that the maximum for planar networks is 4. It turns out that for all trivalent networks the maximum is also 4, and is almost always 3.

■ **Regular polytopes.** In 3D, of the five regular polyhedra, only the tetrahedron, cube and dodecahedron have three edges meeting at each vertex, corresponding to a trivalent network. (Of the 13 additional Archimedean solids, 7 yield trivalent networks.) In 4D the six regular polytopes have 4, 4, 6, 8, 4 and 12 edges meeting at each vertex, and in higher dimensions the simplex ($d + 1$ vertices) and hypercube (2^d vertices) have d edges meeting at each vertex, while the cocube ($2d$ vertices) has $2(d - 1)$. (See also symmetric graphs on page 1032, and page 929.)

■ **Page 476 • Generalizations.** Almost any kind of generalized network can be emulated by a trivalent network just by introducing more nodes. As indicated in the main text, networks with more than three connections at each node can be emulated by combining nodes into groups, and looking only at the connections between groups. Networks with colored nodes can be emulated by representing each color of node by a fixed group of nodes. Going beyond ordinary networks, one can consider hypernetworks in which connections join not just pairs of nodes, but larger numbers of nodes. Such hypernetworks are specified by adjacency tensors rather than adjacency matrices. But it is possible to emulate any hypernetwork by having each generalized connection correspond to a group of connections in an ordinary trivalent network.

■ **Maintaining simple rules.** An important reason for considering models based solely on trivalent networks is that they allow simpler evolution rules to be maintained (see page 508). If nodes can have more than three connections, then they will often be able to evolve to have any number of connections—in which case one must give what is in effect an infinite set of rules to specify what to do for each number of connections.

■ **Page 477 · 3D network.** The 3D network (c) can be laid out in space using `Array[x[8[##]] &, {n, n, n}]` where

$$\begin{aligned} x_1[m : \{_, _ , _ \}] &:= \{x_1[m], x_1[m + 4], \\ &x_2[m + \{4, 2, 0\}], x_2[m + \{0, 6, 4\}]\} \\ x_1[m : \{_, _ , _ \}] &:= \text{Line}[\text{Map}[\# + m \&, \{\{1, 0, 0\}, \{1, 1, 1\}, \\ &\{0, 2, 1\}, \{1, 1, 1\}, \{3, 1, 3\}, \{3, 0, 4\}, \{3, 1, 3\}, \{4, 2, 3\}\}]] \\ x_2[\{i_, j_, k_\}] &:= \\ &x_1[\{-i - 4, -j - 2, k\}] /. \{a_, b_, c_\} \rightarrow \{-a, -b, c\} \end{aligned}$$

The resulting structure is a cubic array of blocks with each block containing 8 nodes. The shortest cycle that returns to a particular node turns out to involve 10 edges. The structure does not correspond to the way that chemical bonds are arranged in any common crystalline materials, probably because it would be likely to be mechanically unstable.

■ **Continuum limits.** For all everyday purposes a region in a network with enough nodes and an appropriate pattern of connections can act just like ordinary continuous space. But at a formal mathematical level this can happen rigorously only in an infinite limit. And in general, there is no reason to expect that all properties of the system (notably for example the existence of particles) will be preserved by taking such a limit. But in understanding the structure of space and comparing to ordinary continuous space it is convenient to imagine taking such a limit. Inevitably there are several scales involved, and one can only expect continuum behavior if one looks at scales intermediate between individual connections in the underlying network and the overall size of the whole network. Yet as I will discuss on pages 534 and 1050 even at such scales it is far from straightforward to see how all the various well-studied properties of ordinary continuous space (as embodied for example in the theory of manifolds) can emerge from discrete underlying networks.

■ **Page 478 · Definitions of distance.** Any measure of distance—whether in ordinary continuous space or elsewhere—takes a pair of points and yields a number. Several properties are normally assumed. First, that if the points are identical the distance is zero, and if they are different, it is a positive number. Second, that the distance between points A and B is the same as between B and A . And third, that the so-called triangle inequality holds, so that the distance AC is no greater than the sum of the distances AB and BC . With distance on a network defined as the length of shortest path between nodes one immediately gets all three of these properties. And even though all distances defined this way will be integers, they still make any network formally correspond in mathematical terms to a metric space (or strictly a path metric space). If the connections on the underlying network are one-way (as in causal networks) then one no longer necessarily gets the second property, and when

a continuum limit exists it can correspond to a (perhaps discontinuous) section through a fiber bundle rather than to a manifold. Note that as discussed on page 536 physical measures of distance will always end up being based not just on single paths in a network, but on the propagation of something like a particle, which typically in effect requires the presence of many paths. (See page 1048.)

■ **Page 478 · Definitions of dimension.** The most obvious way to define the dimension of a space is somehow to ask how many parameters—or coordinates—are needed to specify a point in it. But starting in the 1870s the discovery of constructs like space-filling curves (see page 1127) led to investigation of other definitions. And indeed there is some reason to believe that around 1884 Georg Cantor may have tried developing a definition based on essentially the idea that I use here of looking at growth rates of volumes of spheres (balls). But for standard continuous spaces this definition is hard to make robust—since unlike in discrete networks where one can define volume just by counting nodes, defining volume in a continuous space requires assigning a potentially arbitrary density function. And as a result, in the late 1800s and early 1900s other definitions of dimension were developed. What emerged as most popular is topological dimension, in which one fills space with overlapping balls, and asks what the minimum number that ever have to overlap at any point will be. Also considered was so-called Hausdorff dimension, which became popular in connection with fractals in the 1980s (see page 933), and which can have non-integer values. But for discrete networks the standard definitions for both topological and Hausdorff dimension give the trivial result 0. One can get more meaningful results by thinking about continuum limits, but the definition of dimension that I give in the main text seems much more straightforward. Even here, there are however some subtleties. For example, to find a definite volume growth rate one does still need to take some kind of limit—and one needs to avoid sampling too many or too few nodes in the network. And just as with fractal dimensions discussed on page 933 there are issues about whether a definite power law for the growth rate will emerge, and how one should average over results for different parts of the network. There are some alternative approaches to defining dimension in which some of these issues at least become less explicit. For example, one can imagine not just forming a ball on the network, but instead growing something like a cellular automaton, and seeing how big a pattern it produces after some number of steps. And similarly, one can for example look at the statistics of random walks on the network. A slightly different but still related approach is to study the

density of eigenvalues of the Laplace operator—which can also be thought of as measuring the number of solutions to equations giving linear constraints on numbers assigned to connected nodes. More sophisticated versions of this involve looking at invariants studied in topological field theory. And there are potentially also definitions based for example on considering geodesics and seeing how many linearly independent directions can be defined with them. (Note that given explicit coordinates, one can check whether one is in d or more dimensions by asking for all possible points

$Det[Table[(x[i]-x[j]).(x[i]-x[j]), \{i, d+3\}, \{j, d+3\}]] == 0$ and this should also work for sufficiently separated points on networks. Still another related approach is to consider coloring the edges of a network: if there are $d+1$ possible colors, all of which appear at every node, then it follows that d coordinates can consistently be assigned to each node.)

■ **Page 478 • Counting of nodes.** The number of nodes reached by going out to network distance r (with $r > 1$) from any node in the networks on page 477 is (a) $4r-4$, (b) $3r^2/2-3r/2+1$, and (c)

$$First[Select[4r^3/9+2r^2/3+\{2, 5/3, 5/3\}r-\{10/9, 1, -4/9\}, IntegerQ]]$$

In any trivalent network, the quantity $f[r]$ obtained by adding up the numbers of nodes reached by going distance r from each node must satisfy $f[0]=n$ and $f[1]=3n$, where n is the total number of nodes in the network. In addition, the limit of $f[r]$ for large r must be n^2 . The values of $f[r]$ for all other r will depend on the pattern of connections in the network.

■ **Page 479 • Cycle lengths.** The lengths of the shortest cycles (girths) of the networks on page 479 are (a) 3, (b) 5, (c) 4, (d) 4, (e) 3, (f) 5, (g) 6, (h) 10, (i) ∞ , (j) 3. Note that rules of the kind discussed on page 508 which involve replacing clusters of nodes can only apply when cycles in the cluster match those in the network.

■ **Page 479 • Volumes of spheres.** See page 1050.

■ **Page 480 • Implementation.** Networks are conveniently represented by assigning a number to each node, then having lists of rules which specify what nodes the connection from a particular node go to. The tetrahedron network from page 476 is for example given in this representation by

$$\{1 \rightarrow \{2, 3, 4\}, 2 \rightarrow \{1, 3, 4\}, 3 \rightarrow \{1, 2, 4\}, 4 \rightarrow \{1, 2, 3\}\}$$

The list of nodes reached by following up to n connections from node i are then given by

$$NodeLists[g_., i_., n_.] := NestList[Union[Flatten[# /. g]] &, {i}, n]$$

The network distance corresponding to the length of the shortest path between two nodes is given by

$$Distance[g_., \{i_., j_.\}] := Length[NestWhileList[Union[Flatten[# /. g]] &, {i}, !MemberQ[#, j] &]] - 1$$

■ **Finding layouts.** One way to lay out a network g so that network distances in it come as close as possible to ordinary distances in d -dimensional space, is just to search for values of the $x[i, k]$ which minimize a quantity such as

$$With[\{n = Length[g]\}, Apply[Plus, Flatten[(Table[Distance[g, \{i, j\}], \{i, n\}, \{j, n\}]^2 - Table[Sum[(x[i, k]-x[j, k])^2, \{k, d\}], \{i, n\}, \{j, n\}]^2)]]]$$

using for example $FindMinimum$ starting say with $x[1, _] \rightarrow 0$ and all the other $x[_., _] \rightarrow Random[.]$. Rarely is there a unique minimum that can be found, but the approach nevertheless seems to work fairly well whenever a good layout exists in a particular number of dimensions. One can imagine weighting different network distances differently, but usually I have found that equal weightings work best. If one ignores all constraints beyond network distance 1, then one is in effect just trying to build the network out of identical rigid rods. It turns out that this is almost always possible even in 2D (though not in 1D); the only exception is the tetrahedron network. And in fact very few trivalent structures are rigid, in the sense the angles between rods are uniquely determined. (In 3D, for example, this is true only for the tetrahedron.)

■ **Hamming distances.** In the so-called loop switching method of routing messages in communications systems one lays out a network on an m -dimensional Boolean hypercube so that the distance on the hypercube (equal to Hamming distance) agrees with distance in the network. It is known that to achieve this exactly, m must be at the least the number of either positive or negative eigenvalues of the distance matrix for the network, and can need to be as much as $n-1$, where n is the total number of nodes.

■ **Continuous mathematics.** Even though networks are discrete, it is conceivable that network-based models can also be formulated in terms of continuous mathematics, with a network-like structure emerging for example from the pattern of singularities or topology of continuous surfaces or functions.

The Relationship of Space and Time

■ **History.** The idea of representing time graphically like space has a long history—and was used for example by Nicholas Oresme in the mid-1300s. In the 1700s and 1800s the idea of position and time as just two coordinates was widespread in mathematical physics—and this then led to notions like “travelling in time” in H. G. Wells’s 1895 *The Time Machine*. The mathematical framework developed for relativity theory in the early 1900s (see page 1042) treated space and time very

symmetrically, leading popular accounts of the theory to emphasize a kind of fundamental equivalence between them and to try to make this seem inevitable through rather confusing thought experiments on such topics as idealized trains travelling near the speed of light.

In the context of traditional mathematical equations there has never been much reason to consider the possibility that space and time might be fundamentally different. For typically space and time are both just represented by abstract symbolic variables, and the formal process of solving equations as a function of position in space and as a function of time is essentially identical. But as soon as one tries to construct more explicit models of space and time one is immediately led to consider the possibility that they may be quite different.

■ **Page 482 · Discreteness in time.** In present-day physics, time, like space, is always assumed to be perfectly continuous. But experiments—the most direct of which are based on looking for quantization in the measured decay times of very short-lived particles—have only demonstrated continuity on scales longer than about 10^{-26} seconds, and there is nothing to say that on shorter scales time is not in fact discrete. (The possibility of a discrete quantum of time was briefly discussed in the 1920s when quantum mechanics was first being developed.)

■ **Page 483 · Network constraint systems.** Cases (a), (f) and (p) allow all networks that do not contain respectively cycles of length 1 (self-loops), cycles of length 3 or less, and cycles of length 5 or less. In cases where an infinite sequence of networks is allowed, there are typically particular subnetworks that can occur any number of times, making the sizes of allowed networks form arithmetic progressions. In cases (m), (n) and (o) respectively triangle, pentagon and square subnetworks can be repeated.

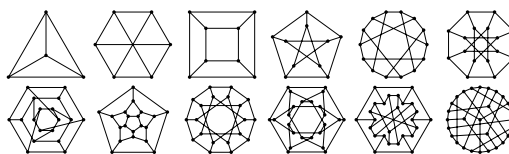
The main text excludes templates that have no dangling connections, and are thus themselves already complete networks. There are 5 such templates involving nodes out to distance one, but of these only 3 correspond to networks that satisfy the constraint that around each node the network has the same form as the template. Among templates involving nodes out to distance two there are 106 that have no dangling connections, and of these only 8 satisfy the constraints.

The main text considers only constraints based on a single template. One can also allow each node to have a neighborhood that corresponds to any of a set of templates. For templates involving nodes out to distance one, there are 13 minimal sets in the sense of page 941, of which only 6 contain just one template, 6 contain two and 1 contains three.

If one does allow dangling connections to be joined within a single template, the results are similar to those discussed so

far. There are 52 possible templates involving nodes out to distance two, of which 12 allow complete networks to be formed, none forced to be larger than 12 nodes. There are 46 minimal sets, with the largest containing 4 templates, but none forcing a network larger than 16 nodes.

■ **Symmetric graphs.** The constraints in a network constraint system require that the structure around each node agrees with a template that contains some number of nodes. A symmetric graph satisfies the same type of constraint, but with the template being the whole network. The pictures below show the smallest few symmetric graphs with 3 connections at each node (with up to 100 nodes there are still only 37 such graphs; compare page 1029).



■ **Cayley graphs.** As discussed on page 938, the structure of a group can be represented by a Cayley graph where nodes correspond to elements in the group, and connections specify results of multiplying by generators. The transitivity of group multiplication implies that Cayley graphs always have the property of being symmetric (see above). The number of connections at each node is fixed, and given by the number of distinct generators and inverses. In cases such as the tetrahedral group A_4 there are 3 connections at each node. The relations among the generators of a group can be thought of as constraints defining the Cayley graph. As mentioned on page 938, there are finite groups that have simple relations but at least very large Cayley graphs. For infinite groups, it is known (see page 938) that in most cases Cayley graphs are locally like trees, and so do not have finite dimension. It appears that only when the group is nilpotent (so that certain combinations of elements commute much as they do on a lattice) is there polynomial growth in the Cayley graph and thus finite dimension.

■ **Page 485 · Spacetime symmetric rules.** With $k=2$ and the neighborhoods shown here, only the additive rules 90R, 105R, 150R and 165R are space-time symmetric. For larger k and larger neighborhoods, there presumably begin to be non-additive rules with this property.

Time and Causal Networks

■ **Causal networks.** The idea of using networks to represent interdependencies of events seems to have developed with the systematization of manufacturing in the early 1900s—

notably in the work of Frank and Lillian Gilbreth—and has been popular since at least the 1940s. Early applications included switching circuits, logistics planning, decision analysis and general flowcharting. In the last few decades causal networks have been widely used in system specification methods such as Petri nets, as well as in schemes for medical and other diagnosis. Since at least the 1960s, causal networks have also been discussed as representations of connections between events in spacetime, particularly in quantum mechanics (see page 1027).

Causal networks like mine that are ultimately associated with some evolution or flow of activity always have certain properties. In particular, they can never contain loops, and thus correspond to directed acyclic graphs. And from this it follows for example that even the most circuitous path between two nodes must be of finite length.

Causal networks can also be viewed as Hasse diagrams of partially ordered sets, as discussed on page 1040.

■ **Implementation.** Given a list of successive positions of the active cell, as from `Map[Last, MAEvolveList[rule, init, t]]` (see page 887), the network can be generated using

```
MAToNet[list_] := Module[{u, j, k}, u[_] = ∞; Reverse[
  Table[j = list[[i]]; k = {u[j - 1], u[j], u[j + 1]}; u[j - 1] =
    u[j] = u[j + 1] = i; i → k, {i, Length[list], 1, -1}]]]
```

where nodes not yet found by explicit evolution are indicated by ∞.

■ **Page 488 · Mobile automata.** The special structure of mobile automata of the type used here leads to several special features in the causal networks derived from them. One of these is that every node always has exactly 3 incoming and 3 outgoing connections. Another feature is that there is always a path of doubled connections (associated with the active cell) that visits every node in some order. And in addition, the final network must always be planar—as it is whenever it is derived from the evolution of a local underlying 1D system.

■ **Computational compression.** In the model for time described here, it is noteworthy that in a sense an arbitrary amount of underlying computation can take place between successive moments in perceived time.

■ **Page 496 · 2D mobile automata.** As in 2D random walks, active cells in 2D mobile automata often do not return to positions they have visited before, with the result that no causal connections end up being created.

The Sequencing of Events in the Universe

■ **Implementation.** Sequential substitution systems in which only one replacement is ever done at each step can just be

implemented using `/.` as described on page 893. Substitution systems in which all replacements are done that are found to fit in a left-to-right scan can be implemented as follows

```
GSSEvolveList[rule_, s_, n_] :=
  NestList[GSSStep[rule, #] &, s, n]
GSSStep[rule_, s_] :=
  g[rule, s, f[StringPosition[s, Map[First, rule]]]]
f[{}] = {}; f[s_] := Fold[If[Last[Last[#1]] ≥ First[#2],
  #1, Append[#1, #2]] &, {First[s]}, Rest[s]]
g[rule_, s_, {}] := s; g[rule_, s_, pos_] := StringReplacePart[
  s, Map[StringTake[s, #] &, pos] /. rule, pos]
```

with rules given as `{"ABA" → "BAAB", "BBBB" → "AA"}`.

■ **Generating causal networks.** If every element generated in the evolution of a generalized substitution system is assigned a unique number, then events can be represented for example by `{4, 5} → {11, 12, 13}`—and from a list of such events a causal network can be built up using

```
With[{u = Map[First, list]}, MapIndexed[Function[
  {e, i}, First[i] → Map[If[# === {}, ∞, #][1, 1]] &][
  Position[u, #]] &, Last[e]], list]
```

■ **The sequential limit.** Even when the order of applying rules does not matter, using the scheme of a sequential substitution system will often give different results. If there is a tree of possible replacements (as in `"A" → "AA"`), then the sequential substitution system in a sense does depth-first recursion in the infinite tree, never returning from the single path it takes. Other schemes are closer to breadth-first recursion.

■ **Page 502 · Rule (b).** The maximum number of steps for which the rule can be applied occurs with initial conditions consisting of a white element followed by n black elements, and in this case the number of steps is $2^n + n$.

■ **String theory.** The sequences of symbols I call strings here have absolutely no direct connection to the continuous deformable 1D objects known as strings in string theory.

■ **String overlaps.** The total numbers of strings with length n and k colors that cannot overlap themselves are given by

$$a[0] = 1; a[n_] := k a[n - 1] - If[EvenQ[n], a[n/2], 0]$$

Up to reversal and interchange of A and B , the first few overlap-free strings with 2 colors are $A, AB, AAB, AAAB, AABB$.

The shortest pairs of strings of 2 elements with no self- or mutual overlaps are `{“A”, “B”}`, `{“AABB”, “AABAB”}`, `{“AABB”, “ABABB”}`; there are a total of 13 such pairs with strings up to length 5, and 85 with strings up to length 6.

The shortest non-overlapping triple of strings is `{“AAABB”, “ABABB”, “ABAABB”}` and its variants. There are a total of 36 such triples with no string having length more than 6.

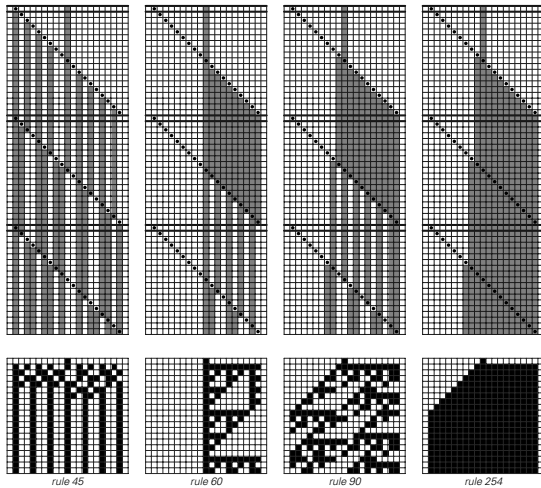
■ **Simulating mobile automata.** Given a mobile automaton like the one from page 73 with rules in the form used on page

887—and behavior of any complexity—the following will yield a causal-invariant substitution system that emulates it:

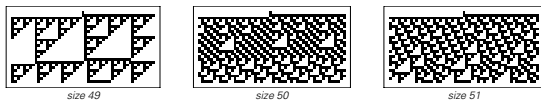
```
Map[StringJoin, Map[{"AAABB", "ABABB", "ABAABB"}][
  # + 1] &, Map[Insert[#[[1]], 2, 2] →
  Insert[#[[2, 1]], 2, 2 + #[[2, 2]]] &, rule], {2}], {2}]
```

■ **Sequential cellular automata.** Ordinary cellular automata are set up so that every cell is updated in parallel at each step, based on the colors of neighboring cells on the previous step. But in analogy with generalized substitution systems, one can also consider sequential cellular automata, in which cells are updated sequentially rather than in parallel. The behavior of such systems is usually very different from that of corresponding ordinary cellular automata, mainly because in sequential cellular automata the new color of a particular cell can depend on new rather than old colors of neighboring cells.

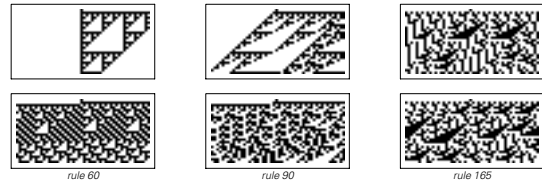
The pictures below show the behavior of several sequential cellular automata with $k = 2, r = 1$ elementary rules. In the top picture of each pair every individual update is indicated by a black dot. In the bottom picture each line represents one complete step of evolution, including one update of each cell. Note that in this representation, effects can propagate all the way across the system in a single step.



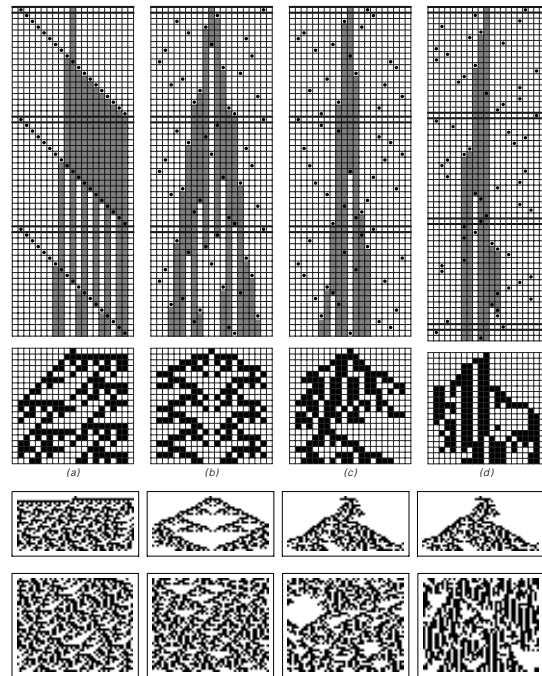
Size dependence. Because effects can propagate all the way across the system in a single step, the overall size, as well as boundary conditions, for the system can be significant after just a few steps, as illustrated in the pictures of rule 60 below.



Additive rules. Among elementary sequential cellular automata, those with additive rules turn out to yield some of the most complex behavior, as illustrated below. The top row shows evolution with the boundary forced to be white; the bottom row shows cyclic boundary conditions. Even though the basic rule is additive, there seems to be no simple traditional mathematical description of the results.



Updating orders. Somewhat different results are typically obtained if one allows different updating orders. For each complete update of a rule 90 sequential cellular automaton, the pictures below show results with (a) left-to-right scan, (b) random ordering of all cells, the same for each pass through the whole system, (c) random ordering of all cells, different for different passes, (d) completely random ordering, in which a particular cell can be updated twice before other cells have even been updated once.



History. Sequential cellular automata have a similar relationship to ordinary cellular automata as implicit updating schemes in finite difference methods have to explicit ones, or as infinite impulse response digital filters have to finite ones. There were several studies of sequential or asynchronous cellular automata done following my work on ordinary cellular automata in the early 1980s.

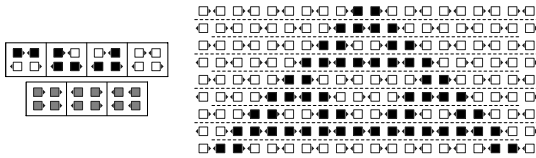
Implementation. The following will update triples of cells in the specified order by using the function f :

```
OrderedUpdate[f_, a_, order_] := Fold[ReplacePart[
  #1, f[Take[#1, {#2 - 1, #2 + 1}]], #2] &, a, order]
```

A random ordering of n cells corresponds to a random permutation of the form

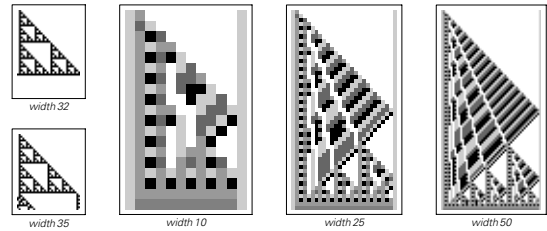
```
Fold[Insert[#1, #2, Random[Integer, Length[#1]] + 1] &,
  {}, Range[n]]
```

■ **Intrinsic synchronization in cellular automata.** Taking the rules for an ordinary cellular automaton and applying them sequentially will normally yield very different results. But it turns out that there are variants on cellular automata in which the rules can be applied in any order and the overall behavior obtained—or at least the causal network—is always the same. The picture below shows how this works for a simple block cellular automaton. The basic idea is that to each cell is added an arrow, and any pair of cells is updated only when their arrows point at each other. This in a sense forces cells to wait to be updated until the data they need is ready. Note that the rules can be thought of as replacements such as $A \langle B \rangle \rightarrow \langle AB \rangle$ for blocks of length 4 with 4 colors.



■ **“Firing squad” synchronization.** By choosing appropriate rules it is possible to achieve many forms of synchronization directly within cellular automata. One version posed as a problem by John Myhill in 1957 consists in setting up a rule in which all cells in a region go into a special state after exactly the same number of steps. The problem was first solved in the early 1960s; the solution using 6 colors and a minimal number of steps shown on the right below was found in 1988 by Jacques Mazoyer, who also determined that no similar 4-color solutions exist. Note that this solution in effect constructs a nested pattern of any width (it does this by optionally including or excluding one additional cell at each nesting level, using a mechanism related to the decimation systems of page 909). If one drops the requirement of cells

going into a special state, then even the 2-color elementary rule 60 shown on the left can be viewed as solving the problem—but only for widths that are powers of 2.



■ **Distributed computing.** Many of the basic issues about the progress of time in a universe consisting of many separate elements have analogs in the progress of computations that are distributed across many separate computing elements. In practice, such computations are most often done by requiring explicit synchronization of all elements at appropriate points, and implementing this using a mechanism that is outside of the computation. But more theoretical investigations of formal concurrent systems, temporal logics, dataflow systems, Petri nets and so on have led to ideas about distributed computing that are somewhat closer to the ones I discuss here for the universe. And, as it happens, in the mid-1980s I tried hard, though at the time without much success, to use updating rules for networks as the basis for a new kind of programming language intended for massively parallel computers.

Uniqueness and Branching in Time

■ **Page 506 · String transformations.** An example of a rule that allows one to go from any string of A 's and B 's to any other is $\{A \rightarrow AA, AA \rightarrow A, A \rightarrow B, B \rightarrow A\}$ (Compare page 1038.)

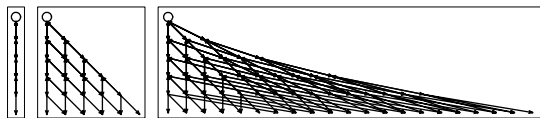
■ **Parallel universes.** The idea of parallel universes which somehow interact with each other has been much explored in science fiction. And one might think that if the history of each universe corresponds to one path in a multiway system then the convergence of paths might represent interactions between universes. But in fact, much as in the case of time travel, such connections do not represent additional observable effects; they simply imply consistency conditions, in this case between universes whose paths converge.

■ **Many-worlds models.** The notion of “many-figured time” has been discussed since the 1950s in the context of the many-worlds interpretation of quantum mechanics. There are some similarities to the multiway systems that I consider here. But an important difference is that while in the many-worlds

approach, branchings are associated with possible observation or measurement events, what I suggest here is that they could be an intrinsic feature of even the very lowest-level rules for the universe. (See also page 1063.)

■ **Spacetime networks from multiway systems.** The main text considers models in which the steps of evolution in a multiway system yield a succession of events in time. An alternative kind of model, somewhat analogous to the ones based on constraints on page 483, is to take the pattern of evolution of a multiway system to define directly a complete spacetime network. Instead of looking separately at strings produced at each step, one instead maintains just a single copy of each distinct string ever produced, and makes that correspond to a node in the network. Each node is then connected to the nodes associated with the strings reached by one application of the multiway rule, as on page 209.

It is fairly straightforward to generate in this way networks of any dimension. For example, starting with n A 's the rule $\{A \rightarrow AB, AB \rightarrow A\}$ yields a regular n -dimensional grid, as shown below.



If each node in a network is associated with a point in spacetime, then one slightly peculiar feature is that every such point would have an associated string—something like an encoded position coordinate. And it then becomes somewhat difficult to understand why different regions of spacetime seem to behave so similarly—and do not, for example, seem to depend on the details of their coordinates.

■ **Page 507 · Commuting operations.** If replacements on strings are viewed as mathematical operations, then when the replacements give the same result if applied in any order, the corresponding operations commute.

■ **Conditions for convergence.** One way to guarantee that there is convergence after one step is to require as in the previous section that blocks to be replaced cannot overlap with themselves or each other. And of the 196 possible rules involving two colors and blocks of length at most three, 112 have this property. But there are also an additional 20 rules which allow some overlap but which nevertheless yield convergence after one step. Examples are $AAA \rightarrow A$ and $AA \rightarrow ABA$. In these rules some of the elements essentially just supply context, but are not affected by the replacement. These elements can then overlap while not affecting the

result. Note that unless one excludes the context elements from events, paths in the multiway system will converge, but the causal networks on these paths will be locally slightly different.

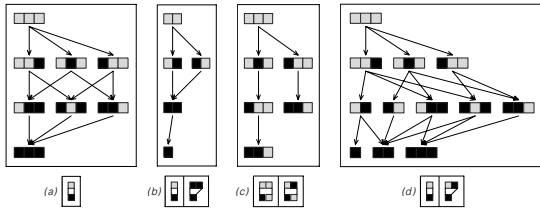
Much as in the previous section, even if paths do not converge for every possible string, it can still be true that paths converge for all strings that are actually generated from a particular initial string.

In general, one can consider convergence after any number of steps, requiring that any two strings which have a common ancestor must at some point also have a common successor. Note that a rule such as $\{A \rightarrow B, A \rightarrow C, B \rightarrow A, B \rightarrow D\}$ exhibits convergence for all paths that have diverged for only one step, but not for all those that have diverged for longer. In general it is formally undecidable whether a particular multiway system will eventually exhibit convergence of all paths.

■ **Confluence.** As mentioned on page 938, multiway systems have been studied in mathematical logic, typically under names such as rewrite systems, since the early 1900s. The property of path convergence discussed in the main text has been considered since the 1930s, usually under the name of confluence, or sometimes the Church-Rosser property. (Also considered is strong confluence—that paths can always converge in at most one step, and local confluence—that paths can converge after diverging for one step but not necessarily more. Early in its history confluence was most often studied for symbolic systems and lambda calculus rather than ordinary multiway systems.)

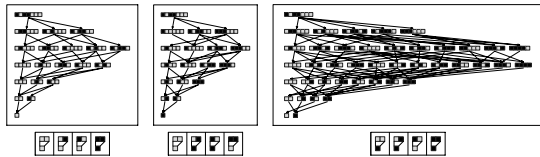
Confluence is important in defining a notion of equivalence for strings. One can say that two strings are equivalent if they can both be transformed to the same string by using the rules of the multiway system. And with such a definition, confluence is what is needed to obtain transitivity for equality, so that $p = q$ and $q = r$ implies $p = r$.

Most often confluence is studied in the context of terminating multiway systems—multiway systems in which eventually strings are produced to which no further replacements apply. If a terminating multiway system has the confluence property, then this implies that regardless of the path taken, a given string will always evolve to a unique string that can be thought of as giving a canonical or normal form for the original string. Examples (a) through (c) below have this property; (d) does not. In example (a), the canonical form is all elements black; in (b) it is a single black element, and in (c) all elements are black, except the last one, which is white if there were any initial white elements. Note that the first example on page 507 has a canonical form consisting of a sorted string.



The process of evaluation in mathematics or in a computer language such as *Mathematica* can be thought of as involving the application of a sequence of replacement rules. Only if these rules have the confluence property will the results always be unique, and independent of the order of rule application.

The evaluation of functions with attribute *Flat* in *Mathematica* provides an example of confluence. If *f* is *Flat*, then in evaluating $f[a, b, c]$ one can equally well start with $f[f[a, b], c]$ or $f[a, f[b, c]]$. Showing only the arguments to *f*, the pictures below illustrate how the flat functions *Xor* and *And* are confluent, while the non-flat function *Implies* is not.



■ **Completion.** If one has a multiway system that terminates but is not confluent then it turns out often to be possible to make it confluent by adding a finite set of new rules. Given a string p which gets transformed either to q or r by the original rules, one can always imagine adding a new rule $q \rightarrow r$ or $r \rightarrow q$ that makes the paths from p immediately converge. To do this explicitly for all possible p that can occur would however entail having infinitely many new rules. But as noted by Donald Knuth and Peter Bendix in 1970 it turns out often to be sufficient just iteratively to add new rules only for each so-called critical pair q, r that is obtained from strings p that represent minimal overlaps in the left-hand sides of the rules one has. To decide whether to add $q \rightarrow r$ or $r \rightarrow q$ in each case one can have some kind of ordering on strings. For the procedure to work this ordering must be such that the strings generated on successive steps in every possible evolution of the multiway system follow the ordering. A number of variations of the basic procedure—using different orderings and with different schemes for dropping redundant rules—have been proposed for systems arising in different kinds of applications. The original Knuth-Bendix procedure was for equations (of the form $a \leftrightarrow b$) had

the feature that it could terminate yet not give a confluent multiway system. But in the 1980s so-called unifying completion algorithms (see page 1158) were developed that—if they terminate—guarantee to give confluent systems. (The question of whether any procedure of this type will terminate in a particular case is nevertheless in general undecidable.)

The basic idea of so-called critical pair completion procedures has arisen several times—notably in the Gröbner basis approach of Bruno Buchberger from 1965 to finding canonical forms for systems of polynomials.

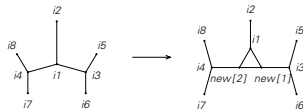
■ **Relationships between types of networks.** Each arrow on each path in a multiway system corresponds to a node in a causal network. Each element in each string in a multiway system corresponds to a connection in a causal network. Each complete string in a multiway system corresponds to a possible slice that goes through all connections across a causal network. Such a slice can be considered in traditional physics terms as a spacelike hypersurface (see page 1041).

Evolution of Networks

■ **Page 509 · Neighbor-independent rules.** Even though the same replacement is performed at each node at each step, the networks produced are not homogeneous. In the first case shown, the picture produced after t steps has $4 \times 3^{t-k-1}$ regions with 3×2^k edges. In the limit $t \rightarrow \infty$, the picture has the geometrical form of an Apollonian circle packing (see page 986). The number of nodes at distance up to r from a given node is at most $1 + \text{Sum}[c[i] + c[i - 1], \{i, n\}]$ where $c[_] := 2^{\text{DigitCount}[_], 2}$. In practice this number fluctuates greatly with r , making pictures like those on page 479 not exhibit smooth profiles. Averaged over all nodes, however, the number of nodes at distance up to r approximates $r^{\text{Log}[2, 3]}$, implying an effective dimension of $\text{Log}[2, 3]$. Note that there is no upper limit on the dimension that can be obtained with appropriate neighbor-independent rules.

■ **Implementation.** For many practical purposes the best representation for networks is the one given on page 1031. But in updating networks a particularly straightforward implementation of one scheme can be obtained if one uses instead a more explicit symbolic representation such as $u[1 \rightarrow v[2, 3, 4], 2 \rightarrow v[1, 3, 4], 3 \rightarrow v[1, 2, 4], 4 \rightarrow v[1, 2, 3]]$ This allows one to capture the basic character of networks by $\text{Attributes}[u] = \{\text{Flat}, \text{Orderless}\}; \text{Attributes}[v] = \text{Orderless}$ Updating rules can then be written in terms of ordinary *Mathematica* patterns. A slight complication is that the patterns have to include all nodes whose connections go to

nodes whose labels are changed by the update. The rule at the top of page 509 must therefore be written out as



and this corresponds to the *Mathematica* rule

```
u[i1_ -> v[i2_, i3_, i4_], i3_ -> v[i1_, i5_, i6_],
i4_ -> v[i1_, i7_, i8_]] -> u[i1 -> v[i2, new[1], new[2]],
new[1] -> v[i1, new[2], i3], new[2] -> v[i1, new[1], i4],
i3 -> v[new[1], i5, i6], i4 -> v[new[2], i7, i8]]
```

(Strictly there also need to be additional rules to cover where for example nodes 3 and 4 are actually the same.) With rules in this form the network update is simply

```
NetStep[rule_, net_] := Block[{new},
net /. rule /. new[n_] -> n + Apply[Max, Map[First, net]]]
```

Note that just as we discussed for strings on page 1033 the direct use of /. here corresponds to a particular scheme for applying the update rule.

■ **Identifying subnetworks.** The problem of finding where in a network a given subnetwork can occur turns out in general to be computationally difficult. For strings the analogous problem is straightforward, since in a string of length n one can ultimately just try each of the n possible starting points for the substring and see for which of them a match occurs. But for a network with n nodes, a similar procedure would require one to check n^k possible configurations in order to find out where a subnetwork of size k occurs. In practice, however, for fixed subnetworks, one can devise fairly efficient procedures. But the general problem of so-called subgraph isomorphism is formally NP-complete.

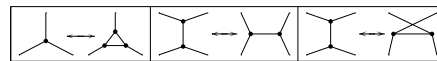
■ **Page 509 · Number of replacements.** The total number of distinct replacements that maintain planarity, involve clusters with up to five nodes and have from 3 to 7 dangling connections is {16, 8, 125, 24, 246}. Not maintaining planarity, the numbers are {14, 5, 13, 2, 2}. (See page 1039.)

■ **Cycles in networks.** See page 1031.

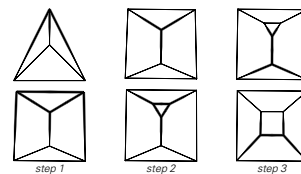
■ **Planar networks.** One feature of a planar network is that it is always possible to identify definite regions or faces bounded by connections in the network. And from Euler’s formula $f + n = e + 2$, it then follows that the average number of edges of each face is always $6(1 - 2/f)$, where f is the total number of faces. Note that with my definition of dimension for networks, the fact that a network is planar does not necessarily mean that it has been two-dimensional—and for example the networks on page 509 are not.

■ **Arbitrary transformations.** By applying the string transformation rules on page 1035 at appropriate locations, it

is possible to transform any string of A ’s and B ’s to any other. And the analog of this for networks is that by applying the rules shown below at appropriate locations it is possible to transform any network into any other. These rules correspond to the moves invented by James Alexander in 1923 in connection with transforming one knot into another. (Note that the first two rules suffice for all planar networks, and are sometimes called respectively T2 and T1.)



As an example, the pictures below show how a tetrahedron network can be transformed into a cube.



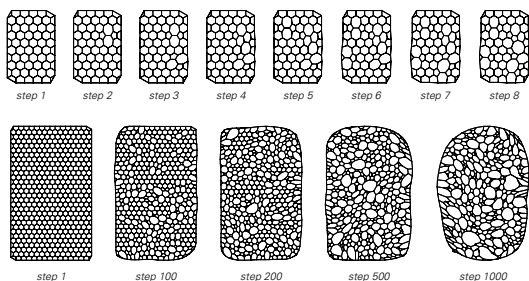
■ **Random networks.** One way to generate the connections for a “completely random” trivalent network with n nodes is just to apply a random permutation:

```
RandomNetwork[n_?EvenQ] := Partition[
Fold[Insert[#1, #2, Random[Integer, Length[#1]] + 1] &,
 {}, Floor[Range[1, n + 2/3, 1/3]]], 2]
```

Networks obtained in this way are usually connected, but will almost always contain self-loops and multiple edges. Properties of random networks are discussed on page 963. A convenient way to get somewhat random planar networks is from 2D Voronoi diagrams of the kind discussed on page 987.

■ **Random replacements.** As indicated in the note above, applying the second rule (T1, shown as (b) on page 511) at an appropriate sequence of positions can transform one planar network into any other with the same number of nodes. The pictures below show what happens if this rule is repeatedly applied at random positions in a network. Each time it is applied, the rule adds two edges to one face, and removes them from another. After many steps the pictures below show that faces with large numbers of edges appear. The average number of edges must always be 6 (see note above), but in a sufficiently large network the probability for a face to have n edges eventually approaches an equilibrium value of $8(n-2)(2n-3)!!(3/8)^n/n!$. (For large n this is approximately λ^n with $\lambda = 3/4$; if 1- and 2-edged regions are allowed then $\lambda = (3 + \sqrt{3})/6 \approx 0.79$.) There may be some easy way to derive such results, but so far it has only been done using fairly sophisticated techniques from quantum field theory developed in the late 1970s. The starting point is to look at a

ϕ^3 field theory with $SU(n)$ internal symmetry and to note that in the limit $n \rightarrow \infty$ what dominates are Feynman diagrams that have the structure of planar trivalent networks (see page 1040). And it then turns out that in zero spacetime dimensions the complete path integral for the theory can be evaluated exactly—yielding in effect a generating function for the number of possible networks. Parametric differentiation (to yield n -point correlation functions) then gives results for n -sided regions. Another result that has been derived is that the average total number $m[n]$ of edges of all faces around a given face with n edges is $7n + 3 + 9/(n + 1)$. Note that the networks obtained always have dimension 2 according to my definitions.



■ **Cellular structures.** There are many systems in nature that consist of assemblies of discrete regions—and the lines that define the interfaces between these regions form networks. In many cases the regions are fixed once established (compare page 988). But in other cases there is continuing evolution, as for example in soap and other foams and froths, grains in metals and perhaps some biological tissues. In 2D situations the lines between regions generically form a trivalent planar network. In a soap foam, the geometrical layout of this network is determined by surface tension forces—with connections meeting at 120° at each node, though being slightly curved and of different lengths. Pressure differences lead to diffusion of gas and on average to von Neumann’s Law that the area of an n -sided region changes linearly with time, at a rate proportional to $n - 6$. Typically the network topology of a foam continually rearranges itself through cascades of seemingly random T1 processes (rule (b) from page 511), with regions that reach zero size disappearing through T2 processes (reversed rule (a)). And as noted for example by Cyril Smith in the early 1950s there is a characteristic coarsening that occurs. Something similar is already visible in the pure T1 pictures in the note above. But results such as the so-called Aboav-Weaire law that $m[n]$ from the note above is in practice about $5n + c$ suggest that T2 processes are also important. (Processes like cell division

in 2D biological tissue in effect directly add connections to a network. But this can again be thought of as a combination of T1 and T2 processes, and in appropriate idealizations can lead to very similar results.)

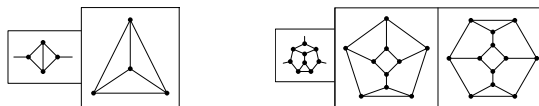
■ **Page 514 · Cluster numbers.** The following tables give the total numbers of distinct clusters—with number of nodes going across the page, and number of dangling connections going down. (See also page 1038.)

	1	2	3	4	5	6	7	8	9	10
0	0	0	0	1	0	2	0	5	0	19
1	0	0	0	0	1	0	4	0	19	0
2	0	0	0	1	0	5	0	23	0	132
3	1	0	1	0	3	0	15	0	91	0
4	0	1	0	2	0	9	0	54	0	390
5	0	0	1	0	4	0	22	0	166	0
6	0	0	0	2	0	9	0	63	0	551

	1	2	3	4	5	6	7	8	9	10
7	0	0	0	0	2	0	17	0	157	0
8	0	0	0	0	0	4	0	38	0	424
9	0	0	0	0	0	0	6	0	80	0
10	0	0	0	0	0	0	0	11	0	180
11	0	0	0	0	0	0	0	0	18	0
12	0	0	0	0	0	0	0	0	0	37

■ **Page 515 · Non-overlapping clusters.** The picture shows all distinct clusters with 3 dangling connections and 9 nodes that are not self-overlapping. The only smaller cluster with the same property is the trivial one with just a single node.

Most clusters that can overlap will be able to do so in an infinite number of possible networks. (One can see this by noting that they can overlap inside clusters with dangling connections, not just closed networks.) But there are some clusters that can overlap only in a few small networks. The pictures below show examples where this happens. The pictures in the main text still treat such clusters as non-overlapping.



If two clusters overlap, then this means that there is some network in which there are copies of these clusters that involve some of the same nodes. And it is possible to search for such a network by starting from a single node and then sequentially trying to take corresponding pieces from the two clusters.

■ **1- and 2-connection clusters.** Clusters with just one or two dangling connections can always in effect be thought of just as adding extra structure to single connections in a network. But this extra structure can be important in the application of other rules—and can for example emulate something like having multiple colors of connections.

■ **Connectedness.** It is not clear whether a network that represents the universe must remain globally connected, or whether pieces can break off. But any replacements that take connected clusters and yield connected clusters must always maintain the connectedness of any network.

■ **Reversibility.** By including both forward and backward versions of every transformation it is straightforward to set up reversible rules for network evolution. It is not clear, however, whether the basic rules for the universe are really reversible. It could well be that the apparent reversibility we see arises because the universe is effectively on an attractor, as discussed on page 1018. Note that if pieces of the universe can break off, but cannot reconnect, then there will inevitably be an irreversible loss of information.

■ **$1/n$ expansion.** If there are n possible colors for each connection in a network, then for large n it turns out that the vast majority of networks will be planar. This idea was used in the 1980s as a way of simplifying the Feynman diagrams to consider in QCD and other quantum field theories. (See page 1039.)

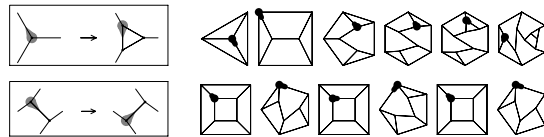
■ **Feynman diagrams.** In the standard approach to particle physics, possible interaction processes are represented by networks in which each node corresponds to an elementary interaction, and the nodes are joined by connections which correspond to the propagation of particles in spacetime. I can see no direct physical relationship between such diagrams and the networks I consider. However, at a mathematical level, the set of trivalent networks with n nodes formally corresponds to the set of n^{th} order Feynman diagrams in a ϕ^3 field theory. (Compare page 1039.)

■ **Chemical analogy.** The evolution of a network can be thought of as an idealized version of a chemical process in which molecules are networks of bonds. (See page 1193.)

■ **Symbolic representations.** Expressions in which common subexpressions are shared correspond to networks, as do collections of relations between objects representing nodes.

■ **Graph grammars.** The notion of generalizing substitutions for strings to the case of networks has been discussed in computer science since the 1960s—and a fair amount of formal work has been done on so-called graph grammars for specifying formal languages whose elements are networks. Even a good analog of regular languages has, however, not yet been found. But applications to constructing or verifying practical network-based system description schemes are quite often discussed. In mathematics rather little is usually done with anything but very trivial network substitutions. In mathematics, rather little is usually done with network substitutions, though the proof of the Four-Color Theorem in 1976 was for example based on showing that 300 or so possible replacement rules—if applied in an appropriate sequence—can transform any graph to have one of 1936 smaller subgraphs that require the same number of colors. (32 rules and 633 subgraphs are now known to be sufficient.)

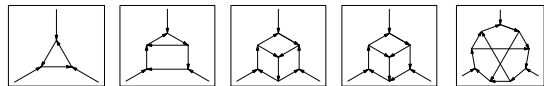
■ **Network mobile automata.** The analog of a mobile automaton can be defined for networks by setting up a single active node, then having rules which replace clusters of nodes around this active node, and move its position. The pictures below show two simple examples.



The total number of replacements that can be used in the rules of a network mobile automaton and which involve clusters with up to four nodes and have from 1 to 4 dangling connections is $\{14, 10, 2727, 781\}$. Despite looking at several hundred thousand cases I have not been able to find network mobile automata with especially complicated behavior.

Note that by having a cluster of nodes with a unique form it is possible to emulate a network mobile automaton using an ordinary network substitution system.

■ **Directed network systems.** If one adds directionality to the connections in a network it becomes particularly easy to set up rules for clusters of nodes that cannot overlap. For no two clusters whose dangling connections all point inwards can ever overlap, at least so long as neither of these clusters themselves contain subclusters whose dangling connections similarly all point inwards. The pictures below show a few examples of such clusters. Note that in a random network of n nodes, about $n/8$ such clusters typically occur.



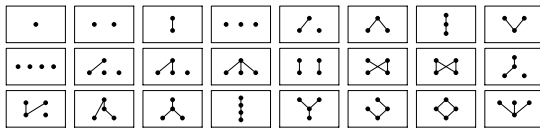
Space, Time and Relativity

■ **Page 516 • Posets.** The way I set things up, collections of events can be thought of as partially ordered sets (posets). If all events occurred in a definite sequence in time, this would define a total linear ordering for them. But with the setup I use, there is only a partial ordering of events, defined by causal connections. The causal networks I draw are so-called Hasse or order diagrams of the posets of events. If a connection goes directly from x to y in this network then x is said to cover y . And in general if there is a path from x to y then one writes $x > y$. The collection of all events that will lead to a given set of events (the union of their past light cones) is known as the filter of that set. Within a poset, there

can be sequences of elements that are totally ordered, and these are called chains. (The maximum length of any chain is sometimes called the dimension of a poset, but this is unrelated to the notions of dimension I consider.) There can also be sets of elements between which no ordering relations at all are defined, and these are called antichains.

Standard examples of posets include subsets of a set ordered by the subset relation, complex numbers ordered by magnitude, and integers ordered by divisibility. Posets first arose as general concepts in the late 1800s in connection with the development of mathematical logic, and to some extent abstract algebra. They became somewhat popular in the mid-1900s, both as formal generalizations in lattice theory, and as structures in various combinatorics applications. It was already noted in the 1920s that events in relativity theory formed posets.

The pictures below show the first few distinct possible Hasse diagrams for posets. For successive numbers of elements the total numbers of these are 1, 2, 5, 16, 63, 318, 2045, 16999, ...



■ **Page 517 · Spacelike slices.** The definition of spacelike slices used here is directly analogous to what is used in traditional relativity theory (typically under names like spacelike hypersurfaces and Cauchy surfaces). There will normally be many different possible choices of spacelike slices, but in all cases a particular such slice is set up to represent what can consistently be thought of as all of space at a given time. One definition of a spacelike slice is then a maximal set of points in which no pair are causally related (corresponding to a maximal antichain in a poset). Another definition (equivalent for any connected causal network) is that spacelike slices are what consistently divide a causal network into a past and a future. And an intermediate definition is that a spacelike slice contains points that are not themselves causally related, but which appear in either the past or the future of every other point. Given a spacelike slice in a causal network, it is always possible to construct another such slice by finding all those points whose immediate predecessors are all included either in the original slice or its predecessors.

■ **Page 518 · Speed of light.** In a vacuum the speed of light is 299,792,458 meters/second (and this is actually what is taken to define a meter). In materials light mostly travels

slower—basically because there are delays when it is absorbed and reemitted by atoms. In a first approximation, the slowdown factor is the refractive index. But particularly in materials which can amplify light a whole sequence of peculiar effects have been observed—and it is fairly subtle to account correctly for incoming and outgoing signals, and to show that at least no energy or information is transmitted faster than c . The standard mathematical framework of relativity theory implies that any massless particle must propagate at c in a vacuum—so that not only light but also gravitational waves presumably go at this speed (and the same is at least approximately true of neutrinos). The effective mass for massive particles increases by a factor $1/\text{Sqrt}[1 - v^2/c^2]$ at speed v , making it take progressively more energy to increase v . At a formal mathematical level it is possible to imagine tachyons which always travel faster than c . But the structure of modern physics would find it difficult to accommodate interactions between these and ordinary particles.

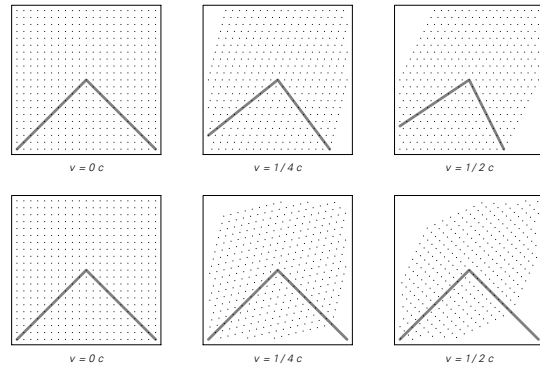
■ **Page 522 · History of relativity.** (See also page 1028.) The idea that mechanical processes should work the same regardless of how fast one is moving was expressed by Galileo in the early 1600s, particularly in connection with the motion of the Earth—and was incorporated in the laws of mechanics formulated by Isaac Newton in 1687. But when the wave theory of light finally became popular in the mid-1800s it seemed to imply that no similar principle could be true for light. For it was generally assumed that waves of light must correspond to explicit disturbances in a medium or ether that fills space. And it was thus expected that for example the apparent speed of light would depend on how fast one was moving with respect to this ether. And indeed in particular this was what the equations for electromagnetism developed by James Maxwell in the 1860s seemed to suggest. But in 1881 an experiment by Albert Michelson (repeated more accurately in 1887 as the Michelson-Morley experiment and now done to the 10^{-20} level) showed that in fact this was not correct. Already in 1882 George FitzGerald and Hendrik Lorentz noted that if there was a contraction in length by a factor $\text{Sqrt}[1 - v^2/c^2]$ in any object moving at speed v (with c being the speed of light) then this would explain the result. And in 1904 Lorentz pointed out that Maxwell's equations are formally invariant under a so-called Lorentz transformation of space and time coordinates (see note below). Then in 1905 Albert Einstein proposed his so-called special theory of relativity—which took as its basic postulates not only that the laws of mechanics and electrodynamics are independent of how fast one is moving, but that this is also true of the speed of light.

And while at first these postulates might seem incompatible, what Einstein showed was that they are not—at least if modifications are made to the basic laws of mechanics. In the few years that followed, various formulations of this result were given, with Hermann Minkowski in 1908 showing that it could be derived if one just assumes that space and time enter all physical laws together in a certain kind of 4D vector. In the late 1800s Ernst Mach had emphasized the idea of formulating science and particularly mechanics in terms only of concepts that can actually be measured by observers. And in this framework Einstein and others gave what seemed to be almost purely deductive arguments for relativity theory—with the result that it generally came to be assumed that there was no meaningful sense in which one could ever imagine deriving relativity from anything more fundamental. Yet as I discussed earlier in the chapter, if a complete theory of physics is to be as simple as possible, then most things like relativity theory must in effect be derived from more basic features of the theory—as I start to try to do in the main text of this section.

■ **Standard treatment.** In a standard treatment of relativity theory one way to begin is to consider setting up a square grid of points in space and time—and then to ask what kind of transformed grid corresponds to this same set of points if one is moving at some velocity v . At first one might assume that the answer would just be a grid that has been sheared by the simple transformation $\{t, x\} \rightarrow \{t, x - vt\}$, as in the first row of pictures below. And indeed for purposes of Newtonian mechanics this so-called Galilean transformation is exactly what is needed. But as the pictures below illustrate, it implies that light cones tip as v increases, so that the apparent speed of light changes, and for example Maxwell's equations must change their form. But the key point is that with an appropriate transformation that affects both space and time, the speed of light can be left the same. The necessary transformation is the so-called Lorentz transformation

$$\{t, x\} \rightarrow \{t - vx/c^2, x - vt\} / \text{Sqrt}[1 - v^2/c^2]$$

And from this the time dilation factor $1/\text{Sqrt}[1 - v^2/c^2]$ shown on page 524 follows, as well as the length contraction factor $\text{Sqrt}[1 - v^2/c^2]$. An important feature of the Lorentz transformation is that it preserves the quantity $c^2 t^2 - x^2$ —with the result that as v changes in the pictures below a given point in the grid traces out a hyperbola whose asymptotes lie on a light cone. Note that on a light cone $c^2 t^2 - x^2$ always vanishes. Note also that the intersection of the past and future light cones for two events separated by a distance x in space and t in time always has a volume proportional exactly to $c^2 t^2 - x^2$.



■ **Inferences from relativity.** The pictures on page 524 show that an idealized clock based on bouncing light between mirrors will exhibit relativistic time dilation. And from such derivations it is often assumed that the same result must hold for any possible clock system. But as a practical matter it does not. And indeed for example the clocks in GPS satellites are specifically set up so as to remove the effects of time dilation. And in the twin paradox one can certainly imagine that each twin could have an accelerometer whose readings they use to correct their clocks. Indeed, even when it comes to individual particles there are subtle effects associated with acceleration and radiation (see page 1062)—so that in the end not entirely clear that something like a biological system would actually in practice exhibit just standard time dilation.

One feature of relativity is that it implies that only relative motion is ultimately ever detectable. (This was also implied by Newtonian mechanics for purely mechanical systems.) And from this it is often concluded that there can be nothing like an ether that one can consider as defining an absolute state of rest in the universe. But in fact the cosmic microwave background in effect does exactly this. For in standard cosmological models it fills the universe, but is everywhere at rest relative to the global center of mass of the universe. And from the anisotropies we have observed in the microwave background it is thus possible to conclude that the Earth is moving at an absolute speed of about $c/10^3$ relative to the center of mass of the universe. In particle physics standard models also in effect introduce things that are assumed to be at rest relative to the center of mass of the universe. One example is the Higgs condensate discussed in connection with particle masses (see page 1047). Other possible examples include zero-point fluctuations in quantum fields.

Outside of science, relativity theory is sometimes given as evidence for various general ideas of cultural relativism (compare page 1131)—which have existed since well before

relativity theory in physics, and seem in the end to have no meaningful connection to it.

■ **Particle physics.** Relativity theory was originally formulated just for mechanics and electromagnetism. But its predictions like $E = mc^2$ were immediately applied for example to radioactivity, and soon it came to be assumed that the theory would work for any system at all—unless it involved gravity. So this has meant that in particle physics $c^2 t^2 - x^2 - y^2 - z^2$ is at some level the only quantity that ever appears. And to make mathematical work easier, what is very often done is to carry out the so-called Wick rotation $t \rightarrow it$ —so relativistic invariance is just independence on 4D orientation. (See page 1061.) But except in rather simple cases there is practically no evidence that results obtained after Wick rotation have anything to do with physical reality—and certainly the transformation removes some very basic phenomena such as particle propagation. One feature of it, however, is that it maps the equation for quantum mechanical time evolution into the equation for probabilities in statistical mechanics, with imaginary time corresponding to inverse temperature. And while it is conceivable that this mapping may have some deep significance, none has so far ever been identified.

■ **Time travel.** The idea that space and time are similar suggests that it might be possible to move backwards and forwards in time just like it is possible to move backwards and forwards in space. And indeed in the partial differential equations that define general relativity, it is formally possible for the motion of particles to achieve this, at least when there is sufficient negative energy density from matter or a cosmological constant. But even in this case there is no real progression in which one travels backwards in time. Instead, the possibility of motion that leads to earlier times simply implies a requirement of consistency between behavior at earlier and later times.

Elementary Particles

■ **Note for physicists.** My goal in the remainder of this chapter is not to present a specific ultimate model for physics, but rather to discuss at a fairly general level some features that I believe such a model will have, given the overall discoveries of this book, and the specific results I have described in this chapter. I am certainly aware that many physicists will want to know more details. But particularly in making contact with existing physics it is almost inevitable that all sorts of technical formalism will be needed—and to maintain balance in this book I have not included this here. (Given my own personal background in theoretical physics it will come as no

surprise that I have often used such formalism in the process of working out what I describe in these sections.)

■ **Page 525 • Types of particles.** Current particle physics identifies three basic types of known elementary particles: leptons, quarks and gauge bosons. The known leptons are the electron (e), muon (μ) and tau lepton (τ), and their corresponding neutrinos (ν_e, ν_μ, ν_τ). Quarks exist inside hadrons like the proton and pion, but never seem to occur as ordinary free particles. Six types are known: u, d, c (charm), s (strange), t (top), b . Gauge bosons are associated with forces. Those currently known are the photon (γ) for electromagnetism (QED), W and Z for so-called weak interactions, and the gluon (g) for QCD interactions between quarks. Gravitons associated with gravitational forces presumably also exist. In ordinary matter, the only particles that contribute in direct ways to everyday physical, chemical and even nuclear properties are electrons, photons and effectively u and d quarks, and gluons. (These, together presumably with some type of neutrino, are the only types of particles that never seem to decay.) The first reasonably direct observations of the various types of particles were as follows (some were predicted in advance): e (1897), γ (~1905), u, d (1914/~1970), μ (1937), s (1946), ν_e (1956), ν_μ (1962), c (1974), τ, ν_τ (1975), b (1977), g (~1979), W (1983), Z (1983), t (1995).

Most particles exist in several variations. Apart from the photon (and graviton), all have distinct antiparticles. Each quark has 3 possible color configurations; the gluon has 8. Most particles also have multiple spin states. Quarks and leptons have spin 1/2, yielding 2 spin states (neutrinos could have only 1 if they were massless). Gauge bosons normally have spin 1 (the graviton would have spin 2) yielding 3 spin states for massive ones. Real massless ones such as the photon always have just 2. (See page 1046.)

In the Standard Model the idea of spontaneous symmetry breaking (see page 1047) allows particles with different masses to be viewed as manifestations of single particles, and this is effectively done for W, Z, γ , as well as for each of the 3 so-called families of quarks and leptons: $u, d; c, s; t, b$ and $e, \nu_e; \mu, \nu_\mu; \tau, \nu_\tau$. Grand unified models typically do this for all known gauge bosons (except gravitons) and for corresponding families of quarks and leptons—and inevitably imply the existence of various additional particles more massive than those known, but with properties that are somehow intermediate. Some models also unify different families, and supersymmetric models unify quarks and leptons with gauge bosons.

■ **History.** The idea that matter—and light—might be made up of discrete particles was already discussed in antiquity

(see page 876). But it was only in the mid-1800s that there started to be real evidence for the existence of some kind of discrete atoms of matter. Yet at the time, the idea of fields was popular, and it was believed that the universe must be filled with a continuous fluid-like ether responsible at least for light and other electromagnetic phenomena. So for example following ideas of William Rankine from 1849 William Thomson (Kelvin) in 1867 suggested that perhaps atoms might be like knotted stable vortex rings in the ether—with different knots corresponding to different chemical elements. But though it initiated the mathematical classification of knots, and now has certain conceptual similarities to what I discuss in this book, the details of this model did not work out—and it had been largely abandoned even before the electron was discovered in 1897. Ernest Rutherford's work in the 1910s on scattering from atoms introduced the idea of an atomic nucleus, and after the discovery of the neutron in 1932 it became clear that the main constituents of nuclei were protons and neutrons. The positron and the muon were discovered in cosmic rays in the 1930s, followed in the 1940s by a handful of other particles. By the 1960s particle accelerators were finding large numbers of new particles every year. And the hypothesis was then suggested that all these particles might actually be composed of just three more fundamental particles that became known as quarks. An alternative so-called democratic or bootstrap hypothesis was also suggested: that somehow any particle could just be viewed as a composite of all others with the same overall properties—with everything being determined by consistency in the web of interactions between particles, and no particles in a sense being more fundamental than others. But by the early 1970s experiments on so-called deep inelastic scattering had given increasingly direct evidence for point-like constituents inside particles like protons—and by the mid-1970s these were routinely identified with quarks.

As soon as the electron was discovered there were questions about its possible size. For if its charge was distributed over a sphere of radius r , this was expected to lead to electrostatic repulsion energy proportional to $1/r$. And although it was suggested around 1900 that effects associated with this might account for the mass of the electron, this ran into problems with relativity theory, and it also remained mysterious just what might hold the electron together. (A late suggestion made in 1953 by Hendrik Casimir was that it could be forces associated with zero-point fluctuations in quantum fields—but at least with the simplest setup these turned out to have wrong sign.)

The development of quantum theory in the 1920s showed that discrete particles will inevitably exhibit continuous

wave-like features in their spatial distribution of probability amplitudes. But traditional quantum mechanics and quantum field theory are both normally formulated with the assumption that the basic particles they describe have zero intrinsic spatial size. Sometimes nonzero size is taken into account by inserting additional interaction parameters—as done in the 1950s with magnetic moments and form factors of protons and neutrons. But for example in quantum electrodynamics the definite assumption is made that electrons are intrinsically of zero size. Quantum fluctuations make any particle in an interacting field theory effectively be surrounded by virtual particles. Yet not unlike in classical electrodynamics having zero intrinsic size for the electron still immediately suggests that an electron should have infinite self-energy. In the 1930s ideas about avoiding this centered around modifying basic laws of electrodynamics or the structure of spacetime (see page 1027). But the development of renormalization in the 1940s showed that these infinities could in effect just be factored out. And by the 1960s a long series of successes in the predictions of QED had led to the almost universal belief that its assumption of point-like electrons must be correct. It was occasionally suggested that the muon might be some kind of composite object. But experiments seemed to indicate that it was in every way identical to the electron, except in mass. And although no reasonable explanation for its existence was found, it came to be generally assumed by the 1970s that it was just another point-like particle. And indeed—apart from few rare suggestions to the contrary—the same is now assumed throughout mainstream practical particle physics for all of the basic particles that appear in the Standard Model. (Actual experiments based on high-energy scattering and precision magnetic moment measurements have shown only that electrons and muons must have sizes smaller than about $\hbar c/(10 \text{ TeV}) \approx 10^{-20} \text{ m}$ —or about 10^{-5} times the size of a proton. One can make arguments that composite particles this small should have masses much larger than are observed—but it is easy to find theories that avoid these.)

In the 1980s superstring theory introduced the idea that particles might actually be tiny 1D strings—with different types of particles corresponding essentially just to strings in different modes of vibration. Since the 1960s it has been noted in many simplified quantum field theories that there can be a kind of duality in which a soliton or other extended field configuration in one representation becomes what acts like an elementary particle in another representation. And in the late 1990s there were indications that such phenomena could occur in generalized string theories—leading to suggestions of at least an abstract correspondence between

for example particles like electrons and gravitational configurations like black holes.

■ **Page 526 · Topological defects.** An idealized vortex in a 2D fluid involves velocity vectors that in effect wind around a point—and can never be unwound by making a series of small local perturbations. The result is a certain kind of stability that can be viewed as being of topological origin. One can classify forms of stability like this in terms of the mathematics of homotopy. Most common are point and line defects in vector fields, but more complicated defects can occur, notably in liquid crystals, models of condensates in the early universe, and certain nonlinear field theories. Analogs of homotopy can presumably be devised to represent certain forms of stability in systems like the networks I consider.

■ **Page 527 · Kuratowski's theorem.** Any network can be laid out in 3D space. (This is related to the Whitney embedding theorem that any d -dimensional manifold can be embedded in $(2d+1)$ -dimensional space.) When one says that a network is planar what one means is that it can be laid out in ordinary 2D space without any lines crossing. Kuratowski's theorem that planarity is associated with the absence of specific subgraphs in a network is an important result in graph theory established in the late 1920s. A subgraph is formally defined to be what one gets by selecting just some subset of connections in a network—and with this definition Kuratowski's theorem must allow extensions of K_5 and $K_{3,3}$ where extra nodes have been inserted in the middle of connections. (K_5 and $K_{3,3}$ are examples of so-called complete graphs, obtained by taking sets of specified numbers of nodes and connecting them in all possible ways.) Another approach is to consider reducing whole networks to so-called minors by deleting connections or merging connected nodes, and in this case Wagner's theorem shows that any non-planar network must be exactly reducible to either K_5 or $K_{3,3}$.

One can generalize the question of planarity to asking whether networks can be laid out on 2D surfaces with various topological structures—and in fact the genus of a graph can be defined to be the number of handles that must be added to a plane to embed the graph without crossings. But even on a torus it turns out that there is no finite set of (extended) subgraphs whose absence guarantees that a network can successfully be laid out. Nevertheless, if one considers minors a finite list does suffice—though for example on a torus it is known that at least 800 (and perhaps vastly more) are needed. (There is in fact a general theorem established since the 1980s that absolutely any list of networks—say for example ones that cannot be laid on a given surface—must actually in effect always all be reducible

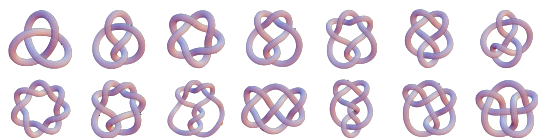
to some finite list of minors.) Note that finding the genus for a particular trivalent network is in general NP-complete.

■ **Page 527 · Gauge invariance.** It is often convenient to define quantities for which only differences or derivatives matter. In classical physics an example is electric potential, which can be shifted by any constant amount without affecting voltage differences or the electric field given by its gradient. In the mid-1800s the idea emerged of a vector potential whose curl gives the magnetic field, and it was soon recognized—notably by James Clerk Maxwell—that any function whose curl vanishes (and that can therefore normally be written as a gradient) could be added to the vector potential without affecting the magnetic field. By the end of the 1800s the general conditions on electromagnetic potentials for invariance of fields were known, though were not thought particularly significant. In 1918 Hermann Weyl tried to reproduce electromagnetism by adding the notion of an arbitrary scale or gauge to the metric of general relativity (see page 1028)—and noted the “gauge invariance” of his theory under simultaneous transformation of electromagnetic potentials and multiplication of the metric by a position-dependent factor. Following the introduction of the Schrödinger equation in quantum mechanics in 1926 it was almost immediately noticed that the equations for a charged particle in an electromagnetic field were invariant under gauge transformations in which the wave function was multiplied by a position-dependent phase factor. The idea then arose that perhaps some kind of gauge invariance could also be used as the basis for formulating theories of forces other than electromagnetism. And after a few earlier attempts, Yang-Mills theories were introduced in 1954 by extending the notion of a phase factor to an element of an arbitrary non-Abelian group. In the 1970s the Standard Model then emerged, based entirely on such theories. In mathematical terms, gauge theories can be viewed as describing fiber bundles in which connections between values of group elements in fibers at neighboring spacetime points are specified by gauge potentials—and curvatures correspond to gauge fields. (General relativity is in effect a special case in which the group elements are themselves related to spacetime coordinates.)

■ **Page 527 · Identifying particles.** In something like a class 4 cellular automaton it is quite straightforward to start enumerating possible persistent structures—as we saw in Chapter 6. But in a network system it can be much more difficult. Ultimately what one wants to do is to find what possible types of forms for local regions are inequivalent under the application of the underlying rules. But in general it may be undecidable even whether two such forms are actually equivalent (compare the notes below and on page

1051)—since to tell this one might need to be able to apply the rules infinitely many times. In specific cases, however, generalizations of concepts like planarity and homotopy may provide useful guides. And a first step may be to look at small closed networks and try to determine which of these can be transformed into each other by a given set of rules.

■ **Knot theory.** Somewhat analogous to the problem in the note above is the problem of classifying knots. The pictures below show some of the simplest distinct knots. But given presentations of two knots, no finite procedure is known that determines in general whether the knots are equivalent (or constructs a sequence of Reidemeister moves that transform one into the other). Quite probably this is in general undecidable, though since the 1920s a few polynomial invariants have been discovered—with recent ones being related to ideas from quantum field theory—that have allowed some progress to be made. (Even the problem of determining whether a knot specified by line segments is trivial is known to be NP-complete.)



■ **Page 528 · Charge quantization.** It is an observed fact that the electric and other charges of all particles are simple rational multiples of each other. In the context of electromagnetism alone, there would be no particular reason to expect this (unless magnetic monopoles exist). But as soon as different particles are related by a non-Abelian symmetry group, then the discreteness of the representations of such a group immediately implies that all charges must be rational multiples of each other.

■ **Spin.** Even when they appear to be of zero size, particles exhibit intrinsic angular momentum known as spin. The total spin is always a fixed multiple of the basic unit \hbar : $1/2$ for quarks and leptons, 1 for photons and other ordinary gauge bosons, 2 for gravitons, and in theory 0 for Higgs particles. (Observed mesons have spins up to perhaps 5 and nuclei up to more than 50 .) Particles of higher spin in effect require more information to specify their orientation (or polarization or its analog). And in the context of network models it could be that spin is somehow related to something as simple as the number of places at which the core of a particle is attached to the rest of the network. Spin values can be thought of as specifying which irreducible representation of the group of symmetries of spacetime is needed to describe a particle after momentum has been factored out. For ordinary massive

particles in d -dimensional space the group is $\text{Spin}(d)$, while for massless particles it is $E(d-1)$ (the Euclidean group). (For tachyons, it would be fundamentally non-compact, forcing continuous spin values.) For small transformations, $\text{Spin}(d)$ is just the ordinary rotation group $\text{SO}(d)$, but globally it is its universal cover, or $\text{SU}(2)$ in 3D. And this can be thought of as what allows half-integer spins, which must be described by spinors rather than vectors or tensors. Such objects have the property that they are not left invariant by 360° rotations, but only by 720° ones—a feature potentially fairly easy to reproduce with networks, perhaps even without definite integer dimensions. In the standard formalism of quantum field theory it can be shown that (above 2D) half-integer spins must always be associated with fermions (which for example satisfy the exclusion principle), and integer spins with bosons. (This spin-statistics connection also seems to hold for various kinds of objects defined by extended field configurations.)

■ **Page 528 · Particle masses.** The measured masses of known elementary particles in units of GeV (roughly equal to the proton mass) are: photon: 0 ; electron: 0.000510998902 ; muon: 0.1056583569 ; τ lepton: 1.77705 ; W : 80.4 ; Z : 91.19 . Recent evidence suggests a mass of about 10^{-11} GeV for at least one type of neutrino. Quarks and gluons presumably never occur as free particles, but still act in many ways as if they have definite masses. For all of them their confinement contributes perhaps 0.3 GeV of effective mass. Then there is also a direct mass: gluons 0 ; u : ~ 0.005 ; d ~ 0.01 ; s : ~ 0.2 ; c : 1.3 ; b : 4.4 ; t : 176 GeV. Note that among sets of particles that have the same quantum numbers—like d , s , b or γ , Z —mixing occurs that makes states of definite mass—that would propagate unchanged as free particles—differ by a unitary transformation from states that are left unchanged by interactions. When one sets up a quantum field theory one can typically in effect insert various mass parameters for particles. Self-interactions normally introduce formally infinite corrections—but if a theory is renormalizable then this means that there are only a limited number of independent such corrections, with the result that relations between masses of different particles are preserved. In quantum field theory any particle is always surrounded by a kind of cloud of virtual particles interacting with it. And following the Uncertainty Principle phenomena involving larger momentum scales will then to probe progressively smaller parts of this cloud—yielding different effective masses. (The masses tend to go up or down logarithmically with momentum scale—following so-called renormalization group equations.)

The Standard Model starts off with certain symmetries that force the masses of all ordinary particles to be zero. But then one assumes that nonzero masses are generated by spontaneous symmetry breaking. One starts by taking each particle to be coupled to a so-called Higgs field. Then one introduces self-interactions in this field so as to make its stable state be one that has constant nonzero value throughout the universe. But this means that as particles propagate, their interactions with the background give them an effective mass. And by having Higgs couplings be proportional to observed particle masses, it becomes inevitable that these will be the masses of particles. One prediction of the usual version of this mechanism for mass is that a definite Higgs particle should exist—which in the minimal Standard Model experiments should observe fairly soon. At times there have been hopes of so-called dynamical symmetry breaking giving the same effective results as the Higgs mechanism, but without an explicit Higgs field—perhaps through something similar to various phenomena in condensed matter physics. String theory, like the Standard Model, tends to start with zero mass particles—and then hopes that an appropriate Higgs-like mechanism will generate nonzero ones.

■ **More particles.** To produce more massive particles requires higher-energy particle collisions, and today's accelerators only allow one to search up to masses of perhaps 200 GeV. (Sufficiently stable particles could have survived from the early universe, and a few cosmic ray interactions in principle give higher energies—but are normally too rare to be useful.) I am not sure whether in my approach one should expect an infinite series of progressively more massive particles. The example of nonplanarity might suggest not, but even in the class 4 cellular automata discussed in Chapter 6 it is not clear whether fundamentally different progressively larger structures will appear forever. In quantum field theory particles of any mass can always in principle exist for short times in virtual form. But normally their effects decrease like powers of their mass—making them hard to measure. In two kinds of cases, however, this does not happen: one is so-called anomalies, the other interactions with the Higgs field, in which couplings are proportional to mass. In the minimal Standard Model it turns out to be impossible to get quarks or leptons with masses much above about 200 GeV without destabilizing the vacuum (a fact pointed out by David Politzer and me in 1979). But with more complicated models one can avoid this constraint. In supersymmetric models—and string theory—there are typically also all sorts of other types of particles, assumed to have high masses since they have not been observed. There is evidence against any more

than the three known generations of quarks and leptons in that the decay process $Z^0 \rightarrow \nu \bar{\nu}$ has a rate that rather accurately agrees with what is expected from just three types of low-mass neutrinos.

■ **Page 530 - Expansion of the universe.** See page 1055.

The Phenomenon of Gravity

■ **History.** With the Earth believed to be the center of the universe, gravity did not seem to require much explanation: it was just a force bringing things to a natural place. But with the advent of Copernican astronomy in the 1500s something more was needed. In the early 1600s Galileo noted that the force of gravity seems to depend only on the mass of an object, and not on any of its other features. In 1687 Isaac Newton then suggested a universal inverse square law of gravity between objects. In the 1700s and 1800s all sorts of celestial mechanics was done on the basis of this—with occasional observational anomalies being resolved for example by the discovery of new planets. Starting in the mid-1800s there were attempts to formulate gravity in the same way as electromagnetism—and in 1900 it was for example suggested that gravitational effects might propagate at the speed of light. Following his introduction of relativity theory in 1905, Albert Einstein began to seek a theory of gravity that would fit in with it. Ordinary special relativity has the feature that it assumes that systems behave the same regardless of their overall velocity—but not regardless of their acceleration. In 1907 Einstein then suggested the equivalence principle that gravity always locally has the same effect as an acceleration. (This principle requires only slightly more than Galileo's idea of the equivalence of gravitational and inertial mass, which has now been verified to the 10^{-12} level.) But by 1912 Einstein realized that if the effective laws of physics were somehow to remain the same in systems with different accelerations (or in different gravitational fields) then this would require a change in their perceived geometry. And building on ideas of differential geometry and tensor calculus from the late 1800s Einstein then began to formulate the concept that gravity is associated with curvature of space. In the late 1800s Ernst Mach had argued that phenomena like acceleration and rotation could ultimately be defined only relative to matter in the universe. And partly on this basis Einstein used the idea that curvature in space must be like a field produced by matter—leading eventually to his formulation in 1915 of the standard Einstein equations for general relativity. An immediate prediction of these was a deviation from the inverse square law, explaining an observed precession in the orbit of Mercury. After a dramatic verification in 1919 of predicted bending of light by

the Sun, general relativity began to be widely accepted. In the 1920s expansion of the universe was discovered, and this was seen to be consistent with general relativity. In the 1940s study of the evolution of stars then led to discussion of what became known as black holes. But for the most part general relativity was still viewed as being highly elegant though of little practical relevance. In the 1960s, however, more work began to be done on it. The discovery of the cosmic microwave background in 1965 led to increasing interest in cosmology. Precision tests—particularly with spacecraft—were designed. In calculations it was sometimes difficult to tell what was a genuine effect, and what was just a feature of the particular coordinates used. But a variety of increasingly abstract mathematical methods were developed, leading notably to general theorems about inevitability of singularities. Detailed calculations tended to require complicated symbolic tensor manipulation (with some associated problems being NP-complete), but with the development of computer algebra this gradually became more feasible—and by the mid-1970s approximate numerical methods were also being used. Various alternative formulations of general relativity were proposed, based for example on tetrads, spinors and twistors (and more recently on connection, loop and non-commutative geometry methods)—but none led to any great simplification. Meanwhile, there continued to be ever more accurate experimental tests of general relativity in the solar system—and at least in the weak gravitational fields available there (with metrics differing from the identity by at most one part in 10^6), all have worked out to around the 10^{-3} level. Starting in the 1960s, more and more ambitious gravitational wave detectors have been built—although none as yet have actually observed anything. Measurements done on a binary pulsar system are nevertheless consistent at a 10^{-3} level with the emission of gravitational radiation in a fairly strong gravitational field at the rate implied by general relativity. And since the 1980s there has been increasing conviction that at least indirect effects of black holes associated with very strong gravitational fields are being observed.

Over the years, some variants of general relativity have been proposed. At least when formulated in terms of tensors, none have quite the simplicity of the original theory—but some lead to rather different predictions, such as an absence of singularities like black holes. Ever since quantum theory began in the early 1900s there has been discussion of quantum gravity—and almost every major method developed for handling other quantum phenomena has been tried on gravity. Starting in the 1980s a variety of methods more specific to quantum gravity were also pursued, but none have yet had convincing success. (See page 1054.)

■ **Differential geometry.** Standard descriptions of properties like curvature—as used for example in general relativity—are normally based on differential geometry. In its usual formulation this assumes that space is continuous, and can always effectively be treated as some kind of deformed version of ordinary Euclidean space—thus forming what is known as a manifold. The result of this is that points in space can always be specified by lists of coordinates—although historically one of the objectives of differential geometry has been to find ways to define properties like curvature so that they do not depend on the choice of such coordinates. The geometrical properties of a space are in general specified by its so-called metric—and this metric allows one to compute quantities based on lengths and angles from coordinates. The metric can be written as a matrix g , defined so that the analog for infinitesimal vectors u and v of $u \cdot v$ in ordinary Euclidean space is $u \cdot g \cdot v$. (This is essentially equivalent to saying that infinitesimal arc length is related to infinitesimal coordinate distances by $ds^2 = g_{i,j} dx_i dx_j$.) In d dimensions the metric g for a so-called Riemannian space can in general be any $d \times d$ positive-definite symmetric matrix—and can vary with position. But for ordinary flat Euclidean space it is always just *IdentityMatrix*[d] (at least with Cartesian coordinates). Within say a surface whose points $\{x_1, x_2, \dots\}$ are obtained by evaluating an expression e as a function of parameters p (so that for example $e = \{x, y, f[x, y]\}$, $p = \{x, y\}$ for a *Plot3D* surface) the metric turns out to be given by

(Transpose[#].#&)[Outer[D,e,p]]

In ordinary Euclidean space a defining feature of geometry is that the shortest path between two points is a straight line. But in an arbitrary space things can be more complicated, and in general such a path will be a geodesic (see note below) which can have a more complicated form. If the coordinates along a path are given by an expression s (such as $\{t, 1+t, t^2\}$) that depends on a parameter t , and the metric at position p is $g[p]$, then the length of a path turns out to be

Integrate[Sqrt[$\partial_t s \cdot g[s] \cdot \partial_t s$], {t, t₁, t₂}]

and geodesics then correspond to paths that extremize this quantity. In ordinary Euclidean space, such paths are straight lines, so that the length of a path between points with lists of coordinates a and b is just the ordinary Euclidean distance *Sqrt*[($a-b$).($a-b$)]. But in general, even though geodesics are not straight lines their lengths can still be used to define a so-called geodesic distance—which turns out to have all the various properties of a distance discussed on page 1030.

If one draws a circle of radius r on a page, then the smaller r is, the more curved the circle will be—and one can define the

circle to have a constant curvature equal to $1/r$. If one draws a more general curve on a page, one can define its curvature at every point by seeing what size of circle fits it best at that point—or equivalently what the coefficients are in a quadratic approximation. (Compare page 418.) With a 2D surface in ordinary 3D space, one can imagine fitting quadrics (generalized ellipsoids). But these are now specified by two radii, yielding two principal curvatures. And in general these curvatures depend on the way the surface is laid out in 3D space. But a crucial point noted by Carl-Friedrich Gauss in the 1820s is that the product of such curvatures—the so-called Gaussian curvature—is always independent of how the surface is laid out, and can thus be viewed as intrinsic to the surface itself, and for example determined purely from the metric for the 2D space corresponding to the surface.

In a 2D space, intrinsic curvature is completely specified just by Gaussian curvature. In higher-dimensional spaces, there are more components, but in general they are all part of the so-called Riemann tensor—a rank-4 tensor introduced by Bernhard Riemann in 1854. (In *Mathematica*, the explicit form of such a tensor can be represented as a nested list for which *TensorRank[list] = 4*.) Several descriptions of the Riemann tensor can be given. One is based on looking at infinitesimal vectors u , v and w and asking how much w differs when transported two ways around the edges of a parallelogram, from x to $x+u+v$ via $x+u$ and via $x+v$. In ordinary flat space there is no difference, but in general the difference is a vector that is defined to be *Riemann.u.v.w.* (The *Riemann* that appears here is formally R_{ijk}^l .) Another description of the Riemann tensor is based on geodesics. In flat Euclidean space any two geodesics that start parallel always remain so. But a defining feature of general non-Euclidean spaces is that this is not in general so. And it turns out that the Riemann tensor is what determines the rate at which geodesics deviate from being parallel. Still another description of the Riemann tensor is as the coefficient of the quadratic terms in an expansion of the metric about a particular point, using so-called normal coordinates set up to make linear terms vanish. In general the Riemann tensor can always be computed from the metric, though it is somewhat complicated. If p is a list of coordinate parameters that appear in a d -dimensional metric g , then

$$Riemann = Table[\partial_{p[[i]]} \Gamma[[i, k]] - \partial_{p[[j]]} \Gamma[[j, k]] + \Gamma[[i, k]] . \Gamma[[j, k]] - \Gamma[[j, k]] . \Gamma[[i, k]], \{i, d\}, \{j, d\}, \{k, d\}]$$

where the so-called Christoffel symbol Γ_{ij}^k is

$$\Gamma = With[{gi = Inverse[g]}, Table[Sum[gi[[l, k]] (\partial_{p[[i]]} g[[i, l]] + \partial_{p[[j]]} g[[j, l]] - \partial_{p[[i]]} g[[i, j]]), \{l, d\}], \{i, d\}, \{j, d\}, \{k, d\}]/2$$

There are d^4 elements in the nested lists for *Riemann*, but symmetries and the so-called Bianchi identity reduce the

number of independent components to $1/12 d^2 (d^2 - 1)$ —or 20 for $d=4$. One can then compute the Ricci tensor ($R_{ik} = R_{ijk}^j$) using

$$RicciTensor = Map[Tr, Transpose[Riemann, \{1, 3, 2, 4\}], \{2\}]$$

and this has $1/2 d (d + 1)$ independent components in $d > 2$ dimensions. (The parts of the Riemann tensor not captured by the Ricci tensor correspond to the so-called Weyl tensor; for $d=2$ the Ricci tensor has only one independent component, equal to the negative of the Gaussian curvature.) Finally, the Ricci scalar curvature is given by

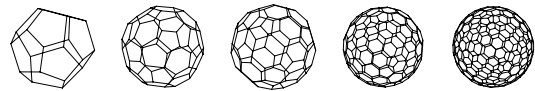
$$RicciScalar = Tr[RicciTensor . Inverse[g]]$$

■ **Page 531 • Geodesics.** On a sphere all geodesics are arcs of great circles. On a surface of constant negative curvature (like (c)) geodesics diverge exponentially, as noted in early work on chaos theory (see page 971). The path of a geodesic can in general be found by requiring that the analog of acceleration vanishes for it. In the case of a surface defined by $z = f[x, y]$ this is equivalent to solving

$$x''[t] = -(f^{(1,0)}[x[t], y[t]] (y'[t]^2 f^{(0,2)}[x[t], y[t]] + 2 x'[t] y'[t] f^{(1,1)}[x[t], y[t]] + x'[t]^2 f^{(2,0)}[x[t], y[t]])) / (1 + f^{(0,1)}[x[t], y[t]]^2 + f^{(1,0)}[x[t], y[t]]^2)$$

together with the corresponding equation for y'' , as already noted by Leonhard Euler in 1728 in connection with his development of the calculus of variations.

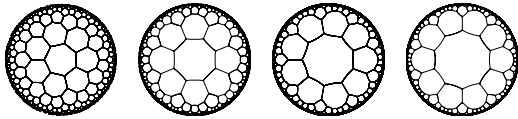
■ **Page 532 • Spherical networks.** One can construct networks of constant positive curvature by approximating the surface of a sphere—starting with a dodecahedron and adding hexagons. (Euler's theorem implies that at any stage there must always be exactly 12 pentagonal faces.) The following are examples with 20, 60, 80, 180 and 320 nodes:



The object with 60 nodes is a truncated icosahedron—the shape of a standard soccer ball, as well the shape of the fullerene molecule C_{60} . (Note that in C_{60} one of the connections at each node is always a double chemical bond, since carbon has valence 4.) Geodesic domes are typically duals of such networks—with three edges on each face.

■ **Hyperbolic networks.** Any surface that always has positive curvature must eventually close up to form something like a sphere. But a surface that has negative curvature (and no holes) must in some sense be infinite—more like cases (c) and (d) on page 412. Yet even in such a case one can always define coordinates that nominally allow the surface to be drawn in a finite way—and the Poincaré disk model used in the pictures below is the standard way of doing this. In ordinary flat space, regular polygons with more than 6

sides can never form a tessellation. But in a space with negative curvature this is possible for polygons with arbitrarily many sides—and the networks that result have been much studied as Cayley graphs of Fuchsian groups. One feature of these networks is that the number of nodes reached in them by following r connections always grows like 2^r . But if one intersperses hexagons in the networks (as in the main text) then one finds that for small r the number of nodes just grows like r^2 —as one would expect for something like a 2D surface. But if one tries to look at growth rates on scales that are not small compared to characteristic lengths associated with curvature then one again sees exponential growth—just as in the case of a uniform tessellation without hexagons.



■ **Page 533 • Sphere volumes.** In ordinary flat Euclidean space the area of a 2D circle is πr^2 , and the volume of a 3D sphere $4\pi r^3/3$. In general, the volume of a sphere in d -dimensional Euclidean space is $s[d]r^d$ where $s[d] = \pi^{d/2}/(d/2)!$ (the surface area is $d s[d]r^{d-1}$). (The function $s[d]$ has a maximum around $d = 5.26$, then decreases rapidly with d .)

If instead of flat space one considers a space defined by the surface of a 3D sphere—say with radius a —one can ask about areas of circles in this space. Such circles are no longer flat, but instead are like caps on the sphere—with a circle of radius r containing all points that are geodesic (great circle) distance less than r from its center. Such a circle has area

$$2\pi a^2 (1 - \text{Cos}[r/a]) = \pi r^2 (1 - r^2/(12a^2) + r^4/(360a^4) - \dots)$$

In the d -dimensional space corresponding to the surface of a $(d+1)$ -dimensional sphere of radius a , the volume of a d -dimensional sphere of radius r is similarly given by

$$d s[d] a^d \text{Integrate}[\text{Sin}[\theta]^{d-1}, \{\theta, 0, r/a\}] = s[d] r^d (1 - d(d-1)r^2/((6(d+2))a^2) + (d(5d^2 - 12d + 7))r^4/((360(d+4))a^4) + \dots)$$

where

$$\text{Integrate}[\text{Sin}[x]^{d-1}, x] = -\text{Cos}[x] \text{Hypergeometric2F1}[1/2, (2-d)/2, 3/2, \text{Cos}[x]^2]$$

In an arbitrary d -dimensional space the volume of a sphere can depend on position, but in general it is given by

$$s[d] r^d (1 - \text{RicciScalar} r^2/(6(d+2)) + \dots)$$

where the Ricci scalar curvature is evaluated at the position of the sphere. (The space corresponding to a $(d+1)$ -dimensional sphere has $\text{RicciScalar} = d(d-1)/a^2$.) The $d = 2$ version of this formula was derived in 1848; the general case in 1917 and 1939. Various derivations can be given. One can

start from the fact that the volume density in any space is given in terms of the metric by $\text{Sqrt}[\text{Det}[g]]$. But in normal coordinates the first non-trivial term in the expansion of the metric is proportional to the Riemann tensor, yet the symmetry of a spherical volume makes it inevitable that the Ricci scalar is the only combination of components that can appear at lowest order. To next order the result is

$$s[d] r^d (1 - \text{RicciScalar} r^2/(6(d+2)) + (5 \text{RicciScalar}^2 - 3 \text{RiemannNorm} + 8 \text{RicciNorm} - 18 \text{Laplacian}[\text{RicciScalar}]) r^4/(360(d+2)(d+4)) + \dots)$$

where the new quantities involved are

$$\text{RicciNorm} = \text{Norm}[\text{RicciTensor}, \{g, g\}]$$

$$\text{RiemannNorm} = \text{Norm}[\text{Riemann}, \{g, g, g, \text{Inverse}[g]\}]$$

$$\text{Norm}[t, gl] := \text{Tr}[\text{Flatten}[t \text{Dual}[t, gl]]]$$

$$\text{Dual}[t, gl] := \text{Fold}[\text{Transpose}[\#1, \text{Inverse}[\#2], \text{RotateLeft}[\text{Range}[\text{TensorRank}[t]]]] \&, t, \text{Reverse}[gl]]$$

$$\text{Laplacian}[f_] := \text{Inner}[D, \text{Sqrt}[\text{Det}[g]] (\text{Inverse}[g] \cdot \text{Map}[\partial_\mu f \&, p]), p]/\text{Sqrt}[\text{Det}[g]]$$

In general the series in r may not converge, but it is known that at least in most cases only flat space can give a result that shows no correction to the basic r^d form. It is also known that if the Ricci tensor is non-negative, then the volume never grows faster than r^d .

■ **Cylinder volumes.** In any d -dimensional space, the volume of a cylinder of length x and radius r whose direction is defined by a unit vector v turns out to be given by

$$s[d-1] r^{d-1} x (1 - (d-1)(\text{RicciScalar} - \text{RicciTensor} \cdot v \cdot v) r^2/(d+1) + \dots)$$

Note that what determines the volume of the cylinder is curvature orthogonal to its direction—and this is what leads to the combination of Ricci scalar and tensor that appears.

■ **Page 533 • Discrete spaces.** Most work with surfaces done on computers—whether for computer graphics, computer-aided design, solving boundary value problems or otherwise—makes use of discrete approximations. Typically surfaces are represented by collections of patches—with a simple mesh of triangles often being used. The triangles are however normally specified not so much by their network of connections as by the explicit coordinates of their vertices. And while there are various triangulation methods that for example avoid triangles with small angles, no standard method yields networks analogous to the ones I consider in which all triangle edges are effectively the same length.

In pure mathematics a basic idea in topology has been to look for finite or discrete ways to capture essential features of continuous surfaces and spaces. And as an early part of this Henri Poincaré in the 1890s introduced the concept of approximating manifolds by cell complexes consisting of collections of generalized polyhedra. By the 1920s there was

then extensive work on so-called combinatorial topology, in which spaces are thought of as being decomposed into abstract complexes consisting say of triangles, tetrahedra and higher-dimensional simplices. But while explicit coordinates and lengths are not usually discussed, it is still imagined that one knows more information than in the networks I consider: not only how vertices are connected by edges, but also how edges are arranged around faces, faces around volumes, and so on. And while in 2D and 3D it is possible to set up such an approximation to any manifold in this way, it turns out that at least in 5D and above it is not. Before the 1960s it had been hoped that in accordance with the Hauptvermutung of combinatorial topology it would be possible to tell whether a continuous mapping and thus topological equivalence exists between manifolds just by seeing whether subdivisions of simplicial complexes for them could be identical. But in the 1960s it was discovered that at least in 5D and above this will not always work. And largely as a result of this, there has tended to be less interest in ideas like simplicial complexes.

And indeed a crucial point for my discussion in the main text is that in formulating general relativity one actually does not appear to need all the structure of a simplicial complex. In fact, the only features of manifolds that ultimately seem relevant are ones that in appropriate limits are determined just from the connectivity of networks. The details of the limits are mathematically somewhat intricate (compare page 1030), but the basic approach is straightforward. One can find the volume of a sphere (geodesic ball) in a network just by counting the number of nodes out to a given network distance from a certain node. And from the limiting growth rate of this one can immediately get the Ricci scalar curvature—just as in the continuous case discussed above. To get the Ricci tensor one also needs a direction. But one can get this from a geodesic—which is in effect the analog of a straight line in the network. Note that unlike in a continuous space there is however usually no obvious way to continue a geodesic in a network. And in general, some—but not all—of the standard constructions used in continuous spaces can also immediately be used in networks. So for example it is straightforward to construct a triangle in a network: one just starts from a particular node, follows geodesics to two others, then joins these with a geodesic. But to extend the triangle into a parallelogram is not so easy—since there is no immediate notion of parallelism in the network. And this means that neither the Riemann tensor, nor a so-called Schild ladder for parallel transport, can readily be constructed.

Since the 1980s there has been increasing interest in formulating notions of continuous geometry for objects like Cayley graphs of groups—which are fundamentally discrete but have infinite limits analogous to continuous systems. (Compare page 938.)

■ **Manifold undecidability.** Given a particular set of network substitution rules there is in general no finite way to decide whether any sequence of such rules exists that will transform particular networks into each other. (Compare undecidability in multiway systems on page 779.) And although one might not expect it on the basis of traditional mathematical intuition, there is an analog of this even for topological equivalence of ordinary continuous manifolds. For the fundamental groups that represent how basic loops can be combined must be equivalent for equivalent manifolds. Yet it turns out that in 4D and above the fundamental group can have essentially any set of generators and relations—so that the undecidability of the word problem for arbitrary groups (see page 1141) implies undecidability of equivalence of manifolds. (In 2D it is straightforward to decide equivalence, and in 3D it is known that only some fundamental groups can be obtained—roughly because not all networks can be embedded in 2D—and it is expected that it will ultimately be possible to decide equivalence.)

■ **Non-integer dimensions.** Unlike in traditional differential geometry (and general relativity) my formulation of space as a network potentially allows concepts like curvature to be defined even outside of integers numbers of dimensions.

■ **Page 534 · Lorentzian spaces.** In ordinary Euclidean space distance is given by $\text{Sqrt}[x^2 + y^2 + z^2]$. In setting up relativity theory it is convenient (see page 1042) to define an analog of distance (so-called proper time) in 4D spacetime by $\text{Sqrt}[c^2 t^2 - x^2 - y^2 - z^2]$. And in terms of differential geometry such Minkowski space can be specified by the metric *DiagonalMatrix*[[+1, -1, -1, -1]] (now taking $c = 1$). To set up general relativity one then considers not Riemannian manifolds but instead Lorentzian ones in which the metric is not positive definite, but instead has the signature of Minkowski space.

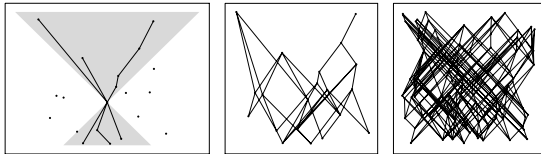
In such Lorentzian spaces, however, there is no useful immediate analog of a sphere. For given any point, even the light cone that corresponds to points at zero spacetime distance from it has an infinite volume. But with an appropriate definition one can still set up cones that have finite volume. To do this in general one starts by picking a vector e in a timelike direction, then normalizes it to be a unit vector so that $e \cdot g \cdot e = -1$. Then one defines a cone of height t whose apex is a given point to be those points whose displacement

vector v satisfies $0 > e \cdot g \cdot v > -t$ (and $0 > v \cdot g \cdot v$). And the volume of such a cone then turns out to be

$$\frac{s[d]t^{d+1}(1-t^2(d+1)(d \text{ RicciScalar} + 2(d+1)(\text{RicciTensor} \cdot e \cdot e)) / ((d+2)(d+3) + \dots) / (d+1))}{}$$

■ **Torsion.** In standard geometry, one assumes that the distance from one point to another is the same as the distance back, so that the metric tensor can be taken to be symmetric, and there is zero so-called torsion. But in for example a causal network, connections have definite directions, and there is in general no such symmetry. And if one looks at the volume of a cone this can then introduce a correction proportional to r . But as soon as there is enough uniformity to define a reasonable notion of static space, it seems that this effect must vanish. (Note that in pure mathematics there are several different uses of the word “torsion”. Here I use it to refer to the antisymmetric parts of the metric tensor.)

■ **Random causal networks.** If one assumes that there are events at random positions in continuous spacetime, then one can construct an effective causal network for them by setting up connections between each event and all events in its future light cone—then deleting connections that are redundant in the sense that they just provide shortcuts to events that could otherwise be reached by following multiple connections. The pictures below show examples of causal networks obtained in this way. The number of connections generally increases faster than linearly with the number of events. Most links end up being at angles that are close to the edge of the light cone.



■ **Page 534 • Einstein equations.** In the absence of matter, the standard statement of the Einstein equations is that all components of the Ricci tensor—and thus also the Ricci scalar—must be zero (or formally that $R_{ij} = 0$). But since the vanishing of all components of a tensor must be independent of the coordinates used, it follows that the vacuum Einstein equations are equivalent to the statement $\text{RicciTensor} \cdot e \cdot e = 0$ for all timelike unit vectors e —a statement that can readily be applied to networks of the kind I consider in the main text. (A related statement is that the 3D Ricci scalar curvature of all spacelike hypersurfaces must vanish wherever these have vanishing extrinsic curvature.)

Another way to state the Einstein equations—already discussed by David Hilbert in 1915—is as the constraint that the integral of $\text{RicciScalar} \sqrt{\text{Det}[g]}$ (the so-called Einstein-Hilbert action) be an extremum. (An idealized soap film or other minimal surface extremizes the integral of the intrinsic volume element $\sqrt{\text{Det}[g]}$, without a RicciScalar factor.) In the discrete Regge calculus that I mention on page 1054 this variational principle turns out to have a rather simple form.

The Einstein-Hilbert action—and the Einstein equations—can be viewed as having the simplest forms that do not ultimately depend on the choice of coordinates. Higher-order terms—say powers of the Ricci scalar curvature—could well arise from underlying network systems, but would not contribute noticeably except in very high gravitational fields.

Various physical interpretations can be given of the vanishing of the Ricci tensor implied by the ordinary vacuum Einstein equations. Closely related to my discussion of the absence of t^2 terms in volume growth for 4D spacetime cones is the statement that if one sets up a small 3D ball of comoving test particles then the volume it defines must have zero first and second derivatives with time.

Below 4D the vanishing of the Ricci tensor immediately implies the vanishing of all components of the Riemann tensor—so that the vacuum Einstein equations force space at least locally to have its ordinary flat form. (Even in 2D there can nevertheless still be non-trivial global topology—for example with flat space having its edges identified as on a torus. In the Euclidean case there were for a long time no non-trivial solutions to the Einstein equations known in any number of dimensions, but in the 1970s examples were found, including large families of Calabi-Yau manifolds.)

In the presence of matter, the typical formal statement of the full Einstein equations is $R_{\mu\nu} - R g_{\mu\nu} / 2 = 8 \pi G T_{\mu\nu} / c^4$, where $T_{\mu\nu}$ is the energy-momentum (stress-energy) tensor for matter and G is the gravitational constant. (An additional so-called cosmological term $\lambda g_{\mu\nu}$ is sometimes added on the right to adjust the effective overall energy density of the universe, and thus its expansion rate. Note that the equation can also be written $R_{\mu\nu} = 8 \pi G (T_{\mu\nu} - 1/2 T_{\mu}^{\mu} g_{\mu\nu}) / c^4$.) The μ, ν component of $T_{\mu\nu}$ gives the flux of the μ component of 4-momentum (whose components are energy and ordinary 3-momentum) in the ν direction. The fact that T_{00} is energy density implies that for static matter (where $E = m c^2$) the equation is in a sense a minimal extension of Poisson’s equation of Newtonian gravity theory. Note that conservation of energy and momentum implies that $T_{\mu\nu}$ must have zero divergence—a result guaranteed in the Einstein equations by the structure of the left-hand side.

In the variational approach to gravity mentioned above, the *RicciScalar* plays the role of a Lagrangian density for pure gravity—and in the presence of matter the Lagrangian density for matter must be added to it. At a physical level, the full Einstein equations can be interpreted as saying that the volume v of a small ball of comoving test particles satisfies

$$\partial_{tt} v[t]/v[t] = -1/2(\rho + 3p)$$

where ρ is the total energy density and p is the pressure averaged over all space directions.

To solve the full Einstein equations in any particular physical situation requires a knowledge of $T_{\mu\nu}$ —and thus of properties of matter such as the relation between pressure and energy density (equation of state). Quite a few global results about the formation of singularities and the absence of paths looping around in time can nevertheless be obtained just by assuming certain so-called energy conditions for $T_{\mu\nu}$. (A fairly stringent example is $0 \leq p \leq \rho/3$ —and whether this is actually true for non-trivial interacting quantum fields remains unclear.)

In their usual formulation, the Einstein equations are thought of as defining constraints on the structure of 4D spacetime. But at some level they can also be viewed as defining how 3D space evolves with time. And indeed the so-called initial value formulations constructed in the 1960s allow one to start with a 3D metric and various extrinsic curvatures defined for a 3D spacelike hypersurface, and then work out how these change on successive hypersurfaces. But at least in terms of tensors, the equations involved show nothing like the simplicity of the usual 4D Einstein equations. One can potentially view the causal networks that I discuss in the main text as providing another approach to setting up an initial value formulation of the Einstein equations.

■ **Page 536 • Pure gravity.** In the absence of matter, the Einstein equations always admit ordinary flat Minkowski space as a solution. But they also admit other solutions that in effect represent configurations of pure gravitational field. And in fact the 4D vacuum Einstein equations are already a sophisticated set of nonlinear partial differential equations that can support all sorts of complex behavior. Several tens of families of solutions to the equations have been found—some with obvious physical interpretations, others without.

Already in 1916 Karl Schwarzschild gave the solution for a spherically symmetric gravitational field. He imagined that this field itself existed in a vacuum—but that it was produced by a mass such as a star at its center. In its original form the metric becomes singular at radius $2Gm/c^2$ (or $3m$ km with m in solar masses). At first it was assumed that this would always be inside a star, where the vacuum Einstein equations

would not apply. But in the 1930s it was suggested that stars could collapse to concentrate their mass in a smaller radius. The singularity was then interpreted as an event horizon that separates the interior of a black hole from the ordinary space around it. In 1960 it was realized, however, that appropriate coordinates allowed smooth continuation across the event horizon—and that the only genuine singularity was infinite curvature at a single point at the center. Sometimes it was said that this must reflect the presence of a point mass, but soon it was typically just said to be a point at which the Einstein equations—for whatever reason—do not apply. Different choices of coordinates led to different apparent locations and forms of the singularity, and by the late 1970s the most common representation was just a smooth manifold with a topology reflecting the removal of a point—and without any specific reference to the presence of matter.

Appealing to ideas of Ernst Mach from the late 1800s it has often been assumed that to get curvature in space always eventually requires the presence of matter. But in fact even the vacuum Einstein equations for complete universes (with no points left out) have solutions that show curvature. If one assumes that space is both homogeneous and isotropic then it turns out that only ordinary flat Minkowski space is allowed. (When matter or a cosmological term is present one gets different solutions—that always expand or contract, and are much studied in cosmology.) If anisotropy is present, however, then there can be all sorts of solutions—classified for example as having different Bianchi symmetry types. And a variety of inhomogeneous solutions with no singularities are also known—an example being the 1962 Ozsváth-Schücking rotating vacuum. But in all cases the structure is too simple to capture much that seems relevant for our present universe.

One form of solution to the vacuum Einstein equations is a gravitational wave consisting of a small perturbation propagating through flat space. No solutions have yet been found that represent complete universes containing emitters and absorbers of such waves (or even for example just two massive bodies). But it is known that combinations of gravitational waves can be set up that will for example evolve to generate singularities. And I suspect that nonlinear interactions between such waves will also inevitably lead to the analog of turbulence for pure gravity. (Numerical simulations often show all sorts of complex behavior—but in the past this has normally been attributed just to the approximations used. Note that for example Bianchi type IX solutions for a complete universe show sensitive dependence on initial conditions—and no doubt this can also happen with nonlinear gravitational waves.)

As mentioned on page 1028, Albert Einstein considered the possibility that particles of matter might somehow just be localized structures in gravitational and electromagnetic fields. And in the mid-1950s John Wheeler studied explicit simple examples of such so-called geons. But in all cases they were found to be unstable—decaying into ordinary gravitational waves. The idea of having purely gravitational localized structures has also occasionally been considered—but so far no stable field configuration has been found. (And no purely repetitive solutions can exist.)

The equivalence principle (see page 1047) might suggest that anything with mass—or energy—should affect the curvature of space in the same way. But in the Einstein equations the energy-momentum tensor is not supposed to include contributions from the gravitational field. (There are alternative and seemingly inelegant theories of gravity that work differently—and notably do not yield black holes. The setup is also somewhat different in recent versions of string theory.) The very definition of energy for the gravitational field is not particularly straightforward in general relativity. But perhaps a definition could be found that would allow localized structures in the gravitational field to make effective contributions to the energy-momentum tensor that would mimic those from explicit particles of matter. Nevertheless, there are quite a few phenomena associated with particles that seem difficult to reproduce with pure gravity—at least say without extra dimensions. One example is parity violation; another is the presence of long-range forces other than gravity.

■ **Quantum gravity.** That there should be quantum effects in gravity was already noted in the 1910s, and when quantum field theory began to develop in the 1930s, there were immediately attempts to apply it to gravity. The first idea was to represent gravity as a field that exists in flat spacetime, and by analogy with photons in quantum electrodynamics to introduce gravitons (at one point identified with neutrinos). By the mid-1950s a path integral (see page 1061) based on the Einstein-Hilbert action had been constructed, and by the early 1960s Feynman diagram rules had been derived, and it had been verified that tree diagrams involving gravitons gave results that agreed with general relativity for small gravitational fields. But as soon as loop diagrams were considered, infinities began to appear. And unlike for quantum electrodynamics there did not seem to be only a finite number of these—that could be removed by renormalization. And in fact by 1973 gravity coupled to matter had been shown for certain not to be renormalizable—and the same was finally shown for pure gravity in 1986. There was an attempt in the 1970s and early 1980s to look

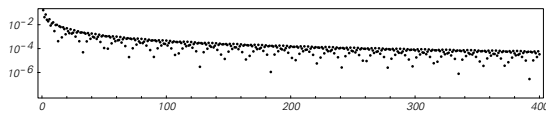
directly at the path integral—without doing an expansion in terms of Feynman diagrams. But despite the fact that at least in Euclidean spacetime a variety of seemingly relevant field configurations were identified, many mathematical difficulties were encountered. And in the late-1970s there began to be interest in the idea that supersymmetric field theories might make infinities associated with gravitons be cancelled by ones associated with other particles. But in the end this did not work out. And then in the mid-1980s one of the great attractions of string theory was that it seemed to support graviton excitations without the problem of infinities seen in point-particle field theories. But it had other problems, and to avoid these, supersymmetry had to be introduced, leading to the presence of many other particles that have so far not been observed. (See also page 1029.)

Starting in the 1950s a rather different approach to quantum gravity involved trying to find a representation of the structure of spacetime in which a quantum analog of the Einstein equations could be obtained by the formal procedure of canonical quantization (see page 1058). Yet despite a few signs of progress in the 1960s there was great difficulty in finding appropriately independent variables to use. In the late 1980s, however, it was suggested that variables could be used corresponding roughly to gravitational fluxes through loops in space. And in terms of these loop variables it was at least formally possible to write down a version of quantum gravity. Yet while this was found in the 1990s to have a correspondence with spin networks (see below), it has remained impossible to see just how it might yield ordinary general relativity as a limit.

Even if one assumes that spacetime is in a sense ultimately continuous one can imagine investigating quantum gravity by doing some kind of discrete approximation. And in 1961 Tullio Regge noted that for a simplicial complex (see page 1050) the Einstein-Hilbert action has a rather simple form in terms of angles between edges. Starting in the 1980s after the development of lattice gauge theories, simulations of random surfaces and higher-dimensional spaces set up in this way were done—often using so-called dynamic triangulation based on random sequences of generalized Alexander moves from page 1038. But there were difficulties with Lorentzian spaces, and when large-scale average behavior was studied, it seemed difficult to reproduce observed smooth spacetime. Analytical approaches (that happened to be like 0D string theory) were also found for 2D discrete spacetimes (compare page 1038)—but they were not successfully extended to higher dimensions.

Over the years, various attempts have been made to derive quantum gravity from fundamentally discrete models of

spacetime (compare page 1027). In recent times the most widely discussed have been spin networks—which despite their name ultimately seem to have fairly little to do with the systems I consider. Spin networks were introduced in 1964 by Roger Penrose as a way to set up an intrinsically quantum mechanical model of spacetime. A simple analog involves a 2D surface made out of triangles whose edges have integer lengths j_i . If one computes the product of $\text{Exp}[i(j_1 + j_2 - j_3)]$ for all triangles, then it turns out for example that this quantity is extremized exactly when the whole surface is flat. In 3D one imagines breaking space into tetrahedra whose edge lengths correspond to discrete quantum spin values. And in 1968 Tullio Regge and Giorgio Ponzano suggested—almost as an afterthought in technical work on $6j$ symbols—that the quantum probability amplitude for any form of space might perhaps be given by the product of $6j$ symbols for the spins on each tetrahedron. The $\text{SixJSymbol}[\{j_1, j_2, j_3\}, \{j_4, j_5, j_6\}]$ are slightly esoteric objects that correspond to recoupling coefficients for the 3D rotation group $\text{SO}(3)$, and that arose in 1940s studies of combinations of three angular momenta in atomic physics—and were often represented graphically as networks. For large j_i they are approximated by $\text{Cos}[\theta + \pi/4]/\text{Sqrt}[12\pi v]$, where v is the volume of the tetrahedron and θ is a deficit angle. And from this it turns out that limits of products of $6j$ symbols correspond essentially to $\text{Exp}[is]$, where s is the discrete form of the Einstein-Hilbert action—extremized by flat 3D space. (The picture below shows for example $\text{Abs}[\text{SixJSymbol}[\{j, j, j\}, \{j, j, j\}]]$. Note that for any j the $6j$ symbols can be given in terms of HypergeometricPFQ .)



In the early 1990s there was again interest in spin networks when the Turaev-Viro invariant for 3D spaces was discovered from a topological field theory involving triangulations weighted with $6j$ symbols of the quantum group $\text{SU}(2)_q$ —and it was seen that invariance under Alexander moves on the triangulation corresponded to the Biedenharn-Elliott identity for $6j$ symbols. In the mid-1990s it was then found that states in 3D loop quantum gravity (see above) could be represented in terms of spin networks—leading for example to quantization of all areas and volumes. In attempting extensions to 4D, spin foams have been introduced—and variously interpreted in terms of simplified Feynman diagrams, constructs in multidimensional category theory, and possible evolutions of spin networks. In all cases, however, spin networks and spin foams seem to be viewed

just as calculational constructs that must be evaluated and added together to get quantum amplitudes—quite different from my idea of associating an explicit evolution history for the universe with the evolution of a network.

■ **Cosmology.** On a large scale our universe appears to show a uniform expansion that makes progressively more distant galaxies recede from us at progressively higher speeds. In general relativity this is explained by saying that the initial conditions must have involved expansion—and that there is not enough in the way of matter or gravitational fields to produce the gravity to slow down this expansion too much. (Note that as soon as objects get gravitationally bound—like galaxies in clusters—there is no longer expansion between them.) The standard big bang model assumes that the universe starts with matter at what is in effect an arbitrarily high temperature. One issue much discussed in cosmology since the late 1970s is how the universe manages to be so uniform. Thermal equilibrium should eventually lead to uniformity—but different parts of the universe cannot come to equilibrium until there has at least been time for effects to propagate between them. Yet there seems for example to be overall uniformity in what we see if we look in opposite directions in the sky—even though extrapolating from the current rate of expansion there has not been enough time since the beginning of the universe for anything to propagate from one side to the other. But starting in the early 1980s it has been popular to think that early in its history the universe must have undergone a period of exponential expansion or so-called inflation. And what this would do is to take just a tiny region and make it large enough to correspond to everything we can now see in the universe. But the point is that a sufficiently tiny region will have had time to come to thermal equilibrium—and so will be approximately uniform, just as the cosmic microwave background is now observed to be. The actual process of inflation is usually assumed to reflect some form of phase transition associated with decreasing temperature of matter in the universe. Most often it is assumed that in the present universe a delicate balance must exist between energy density from a background Higgs field (see page 1047) and a cosmological term in the Einstein equations (see page 1052). But above a critical temperature thermal fluctuations should prevent the background from forming—leading to at least some period in which the universe is dominated by a cosmological term which yields exponential expansion. There tend to be various detailed problems with this scenario, but at least with a sufficiently complicated setup it seems possible to get results that are consistent with observations made so far.

In the context of the discoveries in this book, my expectation is that the universe started from a simple small network, then

progressively added more and more nodes as it evolved, until eventually on a large scale something corresponding to 4D spacetime emerged. And with this setup, the observed uniformity of the universe becomes much less surprising. Intrinsic randomness generation always tends to lead to a certain uniformity in networks. But the crucial point is that this will not take long to happen throughout any network if it is appropriately connected. Traditional models tend to assume that there are ultimately a fixed number of spacetime dimensions in the universe. And with this assumption it is inevitable that if the universe in a sense expands at the speed of light, then regions on opposite sides of it can essentially never share any common history. But in a network model the situation is different. The causal network always captures what happens. And in a case like page 518—with spacetime always effectively having a fixed finite dimension—points that are a distance t apart tend to have common ancestors only at least t steps back. But in a case like (a) on page 514—where spacetime has the structure of an exponentially growing tree—points a distance t apart typically have common ancestors just $\text{Log}[t]$ steps back. And in fact many kinds of causal networks—say associated with early randomly connected space networks—will inevitably yield common ancestors for distant parts of the universe. (Note that such phenomena presumably occur at around the Planck scale of 10^{19} GeV rather than at the 10^{15} GeV or lower scale normally discussed in connection with inflation. They can to some extent be captured in general relativity by imagining an effective spacetime dimension that is initially infinite, then gradually decreases to 4.)

Quantum Phenomena

■ **History.** In classical physics quantities like energy were always assumed to correspond to continuous variables. But in 1900 Max Planck noticed that fits to the measured spectrum of electromagnetic radiation produced by hot objects could be explained if there were discrete quanta of electromagnetic energy. And by 1910 work by Albert Einstein, notably on the photoelectric effect and on heat capacities of solids, had given evidence for discrete quanta of energy in both light and matter. In 1913 Niels Bohr then made the suggestion that the discrete spectrum of light emitted by hydrogen atoms could be explained as being produced by electrons making transitions between orbits with discrete quantized angular momenta. By 1920 ideas from celestial mechanics had been used to develop a formalism for quantized orbits which successfully explained various features of atoms and chemical elements. But it was not clear

how to extend the formalism say to a problem like propagation of light through a crystal. In 1925, however, Werner Heisenberg suggested a new and more general formalism that became known as matrix mechanics. The original idea was to imagine describing the state of an atom in terms of an array of amplitudes for virtual oscillators with each possible frequency. Particular conditions amounting to quantization were then imposed on matrices of transitions between these, and the idea was introduced that only certain kinds of amplitude combinations could ever be observed. In 1923 Louis de Broglie had suggested that just as light—which in optics was traditionally described in terms of waves—seemed in some respects to act like discrete particles, so conversely particles like electrons might in some respects act like waves. In 1926 Erwin Schrödinger then suggested a partial differential equation for the wave functions of particles like electrons. And when effectively restricted to a finite region, this equation allowed only certain modes, corresponding to discrete quantum states—whose properties turned out to be exactly the same as implied by matrix mechanics. In the late 1920s Paul Dirac developed a more abstract operator-based formalism. And by the end of the 1920s basic practical quantum mechanics was established in more or less the form it appears in textbooks today. In the period since, increasing computational capabilities have allowed coupled Schrödinger equations for progressively more particles to be solved (reasonably accurate solutions for hundreds of particles can now be found), allowing ever larger studies in atomic, molecular, nuclear and solid-state physics. A notable theoretical interest starting in the 1980s was so-called quantum chaos, in which it was found that modes (wave functions) in regions like stadiums that did not yield simple analytical solutions tended to show complicated and seemingly random forms.

Basic quantum mechanics is set up to describe how fixed numbers of particles behave—say in externally applied electromagnetic or other fields. But to describe things like fields one must allow particles to be created and destroyed. In the mid-1920s there was already discussion of how to set up a formalism for this, with an underlying idea again being to think in terms of virtual oscillators—but now one for each possible state of each possible one of any number of particles. At first this was just applied to a pure electromagnetic field of non-interacting photons, but by the end of the 1920s there was a version of quantum electrodynamics (QED) for interacting photons and electrons that is essentially the same as today. To find predictions from this theory a so-called perturbation expansion was made, with successive terms representing progressively more interactions, and each

having a higher power of the so-called coupling constant $\alpha \approx 1/137$. It was immediately noticed, however, that self-interactions of particles would give rise to infinities, much as in classical electromagnetism. At first attempts were made to avoid this by modifying the basic theory (see page 1044). But by the mid-1940s detailed calculations were being done in which infinite parts were just being dropped—and the results were being found to agree rather precisely with experiments. In the late 1940s this procedure was then essentially justified by the idea of renormalization: that since in all possible QED processes only three different infinities can ever appear, these can in effect systematically be factored out from all predictions of the theory. Then in 1949 Feynman diagrams were introduced (see note below) to represent terms in the QED perturbation expansion—and the rules for these rapidly became what defined QED in essentially all practical applications. Evaluating Feynman diagrams involved extensive algebra, and indeed stimulated the development of computer algebra (including my own interest in the field). But by the 1970s the dozen or so standard processes discussed in QED had been calculated to order α^2 —and by the mid-1980s the anomalous magnetic moment of the electron had been calculated to order α^4 , and nearly one part in a trillion (see note below).

But despite the success of perturbation theory in QED it did not at first seem applicable to other issues in particle physics. The weak interactions involved in radioactive beta decay seemed too weak for anything beyond lowest order to be relevant—and in any case not renormalizable. And the strong interactions responsible for holding nuclei together (and associated for example with exchange of pions and other mesons) seemed too strong for it to make sense to do an expansion with larger numbers of individual interactions treated as less important. So this led in the 1960s to attempts to base theories just on setting up simple mathematical constraints on the overall so-called S matrix defining the mapping from incoming to outgoing quantum states. But by the end of the 1960s theoretical progress seemed blocked by basic questions about functions of several complex variables, and predictions that were made did not seem to work well.

By the early 1970s, however, there was increasing interest in so-called gauge or Yang-Mills theories formed in essence by generalizing QED to operate not just with a scalar charge, but with charges viewed as elements of non-Abelian groups. In 1972 it was shown that spontaneously broken gauge theories of the kind needed to describe weak interactions were renormalizable—allowing meaningful use of perturbation theory and Feynman diagrams. And then in 1973 it was discovered that QCD—the gauge theory for quarks and

gluons with SU(3) color charges—was asymptotically free (it was known to be renormalizable), so that for processes probing sufficiently small distances, its effective coupling was small enough for perturbation theory. By the early 1980s first-order calculations of most basic QCD processes had been done—and by the 1990s second-order corrections were also known. Schemes for adding up all Feynman diagrams with certain very simple repetitive or other structures were developed. But despite a few results about large-distance analogs of renormalizability, the question of what QCD might imply for processes at larger distances could not really be addressed by such methods.

In 1941 Richard Feynman pointed out that amplitudes in quantum theory could be worked out by using path integrals that sum with appropriate weights contributions from all possible histories of a system. (The Schrödinger equation is like a diffusion equation in imaginary time, so the path integral for it can be thought of as like an enumeration of random walks. The idea of describing random walks with path integrals was discussed from the early 1900s.) At first the path integral was viewed mostly as a curiosity, but by the late 1970s it was emerging as the standard way to define a quantum field theory. Attempts were made to see if the path integral for QCD (and later for quantum gravity) could be approximated with a few exact solutions (such as instantons) to classical field equations. By the early 1980s there was then extensive work on lattice gauge theories in which the path integral (in Euclidean space) was approximated by randomly sampling discretized field configurations. But—I suspect for reasons that I discuss in the note below—such methods were never extremely successful. And the result is that beyond perturbation theory there is still no real example of a definitive success from standard relativistic quantum field theory. (In addition, even efforts in the context of so-called axiomatic field theory to set up mathematically rigorous formulations have run into many difficulties—with the only examples satisfying all proposed axioms typically in the end being field theories without any real interactions. In condensed matter physics there are nevertheless cases like the Kondo model where exact solutions have been found, and where the effective energy function for electrons happens to be roughly the same as in a relativistic theory.)

As mentioned on page 1044, ordinary quantum field theory in effect deals only with point particles. And indeed a recurring issue in it has been difficulty with constraints and redundant degrees of freedom—such as those associated with extended objects. (A typical goal is to find variables in which one can carry out what is known as canonical quantization: essentially applying the same straightforward

transformation of equations that happens to work in ordinary elementary quantum mechanics.) One feature of string theory and its generalizations is that they define presumably consistent quantum field theories for excitations of extended objects—though an analog of quantum field theory in which whole strings can be created and destroyed has not yet been developed.

When the formalism of quantum mechanics was developed in the mid-1920s there were immediately questions about its interpretation. But it was quickly suggested that given a wave function ψ from the Schrödinger equation $\text{Abs}[\psi]^2$ should represent probability—and essentially all practical applications have been based on this ever since. From a conceptual point of view it has however often seemed peculiar that a supposedly fundamental theory should talk only about probabilities. Following the introduction of the uncertainty principle and related formalism in the 1920s one idea that arose was that—in rough analogy to relativity theory—it might just be that there are only certain quantities that are observable in definite ways. But this was not enough, and by the 1930s it was being suggested that the validity of quantum mechanics might be a sign that whole new general frameworks for philosophy or logic were needed—a notion supported by the apparent need to bring consciousness into discussions about measurement in quantum mechanics (see page 1063). The peculiar character of quantum mechanics was again emphasized by the idealized experiment of Albert Einstein, Boris Podolsky and Nathan Rosen in 1935. But among most physicists the apparent lack of an ordinary mechanistic way to think about quantum mechanics ended up just being seen as another piece of evidence for the fundamental role of mathematical formalism in physics.

One way for probabilities to appear even in deterministic systems is for there to be hidden variables whose values are unknown. But following mathematical work in the early 1930s it was usually assumed that this could not be what was going on in quantum mechanics. In 1952 David Bohm did however manage to construct a somewhat elaborate model based on hidden variables that gave the same results as ordinary quantum mechanics—though involved infinitely fast propagation of information. In the early 1960s John Bell then showed that in any hidden variables theory of a certain general type there are specific inequalities that combinations of probabilities must satisfy (see page 1064). And by the early 1980s experiments had shown that such inequalities were indeed violated in practice—so that there were in fact correlations of the kind suggested by quantum mechanics. At first these just seemed like isolated esoteric effects, but by the mid-1990s they were being codified in the field of quantum

information theory, and led to constructions with names like quantum cryptography and quantum teleportation.

Particularly when viewed in terms of path integrals the standard formalism of quantum theory tends to suggest that quantum systems somehow do more computation in their evolution than classical ones. And after occasional discussion as early as the 1950s, this led by the late 1980s to extensive investigation of systems that could be viewed as quantum analogs of idealized computers. In the mid-1990s efficient procedures for integer factoring and a few other problems were suggested for such systems, and by the late 1990s small experiments on these were beginning to be done in various types of physical systems. But it is becoming increasingly unclear just how the idealizations in the underlying model really work, and to what extent quantum mechanics is actually in the end even required—as opposed, say, just to classical wave phenomena. (See page 1147.)

Partly as a result of discussions about measurement there began to be questions in the 1980s about whether ordinary quantum mechanics can describe systems containing very large numbers of particles. Experiments in the 1980s and 1990s on such phenomena as macroscopic superposition and Bose-Einstein condensation nevertheless showed that standard quantum effects still occur with trillions of atoms. But inevitably the kinds of general phenomena that I discuss in this book will also occur—leading to all sorts of behavior that at least cannot readily be foreseen just from the basic rules of quantum mechanics.

■ **Quantum effects.** Over the years, many suggested effects have been thought to be characteristic of quantum systems:

- Basic quantization (1913): mechanical properties of particles in effectively bounded systems are discrete;
- Wave-particle duality (1923): objects like electrons and photons can be described as either waves or particles;
- Spin (1925): particles can have intrinsic angular momentum even if they are of zero size;
- Non-commuting measurements (1926): one can get different results doing measurements in different orders;
- Complex amplitudes (1926): processes are described by complex probability amplitudes;
- Probabilism (1926): outcomes are random, though probabilities for them can be computed;
- Amplitude superposition (1926): there is a linear superposition principle for probability amplitudes;
- State superposition (1926): quantum systems can occur in superpositions of measurable states;

- Exclusion principle (1926): amplitudes cancel for fermions like electrons to go in the same state;
- Interference (1927): probability amplitudes for particles can interfere, potentially destructively;
- Uncertainty principle (1927): quantities like position and momenta have related measurement uncertainties;
- Hilbert space (1927): states of systems are represented by vectors of amplitudes rather than individual variables;
- Field quantization (1927): only discrete numbers of any particular kind of particle can in effect ever exist;
- Quantum tunnelling (1928): particles have amplitudes to go where no classical motion would take them;
- Virtual particles (1932): particles can occur for short times without their usual energy-momentum relation;
- Spinors (1930s): fermions show rotational invariance under $SU(2)$ rather than $SO(3)$;
- Entanglement (1935): separated parts of a system often inevitably behave in irreducibly correlated ways;
- Quantum logic (1936): relations between events do not follow ordinary laws of logic;
- Path integrals (1941): probabilities for behavior are obtained by summing contributions from many paths;
- Imaginary time (1947): statistical mechanics is like quantum mechanics in imaginary time;
- Vacuum fluctuations (1948): there are continual random field fluctuations even in the vacuum;
- Aharonov-Bohm effect (1959): magnetic fields can affect particles even in regions where they have zero strength;
- Bell's inequalities (1964): correlations between events can be larger than in any ordinary probabilistic system;
- Anomalies (1969): virtual particles can have effects that violate the original symmetries of a system;
- Delayed choice experiments (1978): whether particle or wave features are seen can be determined after an event;
- Quantum computing (1980s): there is the potential for fundamental parallelism in computations.

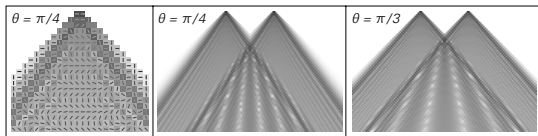
All of these effects are implied by the standard mathematical formalism of quantum theory. But it has never been entirely clear which of them are in a sense true defining features of quantum phenomena, and which are somehow just details. It does not help that most of the effects—at least individually—can be reproduced by mechanisms that seem to have little to do with the usual structure of quantum theory. So for example there will tend to be quantization whenever the

underlying elements of a system are discrete. Similarly, features like the uncertainty principle and path integrals tend to be seen whenever things like waves are involved. And probabilistic effects can arise from any of the mechanisms for randomness discussed in Chapter 7. Complex amplitudes can be thought of just as vector quantities. And it is straightforward to set up rules that will for example reproduce the detailed evolution of amplitudes according say to the Schrödinger equation (see note below). It is somewhat more difficult to set up a system in which such amplitudes will somehow directly determine probabilities. And indeed in recent times consequences of this—such as violations of Bell's inequalities—are what have probably most often been quoted as the most unique features of quantum systems. It is however notable that the vast majority of traditional applications of quantum theory do not seem to have anything to do with such effects. And in fact I do not consider it at all clear just what is really essential about them, and what is in the end just a consequence of the extreme limits that seem to need to be taken to get explicit versions of them.

■ **Reproducing quantum phenomena.** Given molecular dynamics it is much easier to see how to reproduce fluid mechanics than rigid-body mechanics—since to get rigid bodies with only a few degrees of freedom requires taking all sorts of limits of correlations between underlying molecules. And I strongly suspect that given a discrete underlying model of the type I discuss here it will similarly be much easier to reproduce quantum field theory than ordinary quantum mechanics. And indeed even with traditional formalism, it is usually difficult to see how quantum mechanics can be obtained as a limit of quantum field theory. (Classical limits are slightly easier: they tend to be associated with stationary features or caustics that occur at large quantum numbers—or coherent states that represent eigenstates of raising or particle creation operators. Note that the exclusion principle makes classical limits for fermions difficult—but crucial for the stability of bulk matter.)

■ **Discrete quantum mechanics.** While there are many issues in finding a complete underlying discrete model for quantum phenomena, it is quite straightforward to set up continuous cellular automata whose limiting behavior reproduces the evolution of probability amplitudes in standard quantum mechanics. One starts by assigning a continuous complex number value to each cell. Then given the list of such values the crucial constraint imposed by the standard formalism of quantum mechanics is unitarity: that the quantity $\text{Tr}[Abs[list]^2]$ representing total probability should be conserved. This is in a sense analogous to conservation of total density in diffusion processes. From

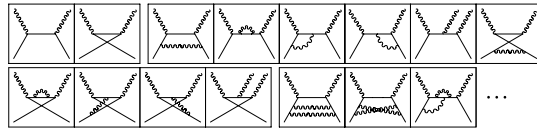
the discussion of page 1024 one can reproduce the 1D diffusion equation with a continuous block cellular automaton in which the new value of each block is given by $\{(1-\xi, \xi), \{\xi, 1-\xi\}\} \cdot \{a_1, a_2\}$. So in the case of quantum mechanics one can consider having each new block be given by $\{\{\text{Cos}[\theta], i\text{Sin}[\theta]\}, \{i\text{Sin}[\theta], \text{Cos}[\theta]\}\} \cdot \{a_1, a_2\}$. The pictures below show examples of behavior obtained with this rule. (Gray levels represent magnitude for each cell, and arrows phase.) And it turns out that in suitable limits one generally gets essentially the behavior expected from either the Dirac or Klein-Gordon equations for relativistic particles, or the Schrödinger equation for non-relativistic particles. (Versions of this were noticed by Richard Feynman in the 1940s in connection with his development of path integrals, and were pointed out again several times in the 1980s and 1990s.)



One might hope to be able to get an ordinary cellular automaton with a limited set of possible values by choosing a suitable θ . But in fact in non-trivial cases most of the cells generated at each step end up having distinct values. One can generalize the setup to more dimensions or to allow $n \times n$ matrices that are elements of $SU(n)$. Such matrices can be viewed in the context of ordinary quantum formalism as S matrices for elementary evolution events—and can in general represent interactions. (Note that all rules based on matrices are additive, reflecting the usual assumption of linearity at the level of amplitudes in quantum mechanics. Non-additive unitary rules can also be found. The analog of an external potential can be introduced by progressively changing values of certain cells at each step. Despite their basic setup the systems discussed here are not direct analogs of standard quantum spin systems, since these normally have local Hamiltonians and non-local evolution functions, while the systems here have local evolution functions but seem always to require non-local Hamiltonians.)

■ **Page 540 • Feynman diagrams.** The pictures below show a typical set of Feynman diagrams used to do calculations in QED—in this case for so-called Compton scattering of a photon by an electron. The straight lines in the diagrams represent electrons; the wavy ones photons. At some level each diagram can be thought of as representing a process in which an electron and photon come in from the left, interact in some way, then go out to the right. The incoming and

outgoing lines correspond to real particles that propagate to infinity. The lines inside each diagram correspond to virtual particles that in effect propagate only a limited distance, and have a distribution of energy-momentum and polarization properties that can differ from real particles. (Exchanges of virtual photons can be thought of as producing familiar electromagnetic forces; exchanges of virtual electrons as yielding an analog of covalent forces in chemistry.)



To work out the total probability for a process from Feynman diagrams, what one does is to find the expression corresponding to each diagram, then one adds these up, and squares the result. The first two blocks of pictures above show all the diagrams for Compton scattering that involve 2 or 3 photons—and contribute through order α^3 . Since for QED $\alpha \approx 1/137$, one might expect that this would give quite an accurate result—and indeed experiments suggest that it does. But the number of diagrams grows rapidly with order, and in fact the k^{th} order term can be about $(-1)^k \alpha^k (k/2)!$, yielding a series that formally diverges. In simpler examples where exact results are known, however, the first few terms typically still seem to give numerically accurate results for small α . (The high-order terms often seem to be associated with asymptotic series for things like $\text{Exp}[-1/\alpha]$.)

The most extensive calculation made so far in QED is for the magnetic moment of the electron. Ignoring parts that depend on particle masses the result (derived in successive orders from 1, 1, 7, 72, 891 diagrams) is

$$2(1 + \alpha/(2\pi) + (3\text{Zeta}[3])/4 - 1/2\pi^2 \text{Log}[2] + \pi^2/12 + 197/144)(\alpha/\pi)^2 + (83/72\pi^2 \text{Zeta}[3] - 215\text{Zeta}[5]/24 - 239\pi^4/2160 + 139\text{Zeta}[3]/18 + 25/18(24\text{PolyLog}[4, 1/2] + \text{Log}[2]^4 - \pi^2 \text{Log}[2]^2) - 298/9\pi^2 \text{Log}[2] + 17101\pi^2/810 + 28259/5184)(\alpha/\pi)^3 - 1.4(\alpha/\pi)^4 + \dots$$

or roughly

$$2. + 0.32\alpha - 0.067\alpha^2 + 0.076\alpha^3 - 0.029\alpha^4 + \dots$$

The comparative simplicity of the symbolic forms here (which might get still simpler in terms of suitable generalized polylogarithm functions) may be a hint that methods much more efficient than explicit Feynman diagram evaluation could be used. But it seems likely that there would be limits to this, and that in the end QED will exhibit the kind of computational irreducibility that I discuss in Chapter 12.

Feynman diagrams in QCD work at the formal level very much like those in QED—except that there are usually many more of them, and their numerical results tend to be larger,

with expansion parameters often effectively being $\alpha\pi$ rather than α/π . For processes with large characteristic momentum transfers in which the effective α in QCD is small, remarkably accurate results are obtained with first and perhaps second-order Feynman diagrams. But as soon as the effective α becomes larger, Feynman diagrams as such rapidly seem to stop being useful.

■ **Quantum field theory.** In standard approaches to quantum field theory one tends to think of particles as some kind of small perturbations in a field. Normally for calculations these perturbations are on their own taken to be plane waves of definite frequency, and indeed in many ways they are direct analogs of waves in classical field theories like those of electromagnetism or fluid mechanics. To investigate collisions between particles, one thus looks at what happens with multiple waves. In a system described by linear equations, there is always a simple superposition principle, and waves just pass through each other unchanged. But what in effect leads to non-trivial interactions between particles is the presence of nonlinearities. If these are small enough then it makes sense to do a perturbation expansion in which one approximates field configurations in terms of a succession of arrangements of ordinary waves—as in Feynman diagrams. But just as one cannot expect to capture fully turbulent fluid flow in terms of a few simple waves, so in general as soon as there is substantial nonlinearity it will no longer be sufficient just to do perturbation expansions. And indeed for example in QCD there are presumably many cases in which it is necessary to look at something closer to actual complete field configurations—and correlations in them.

The way the path integral for a quantum field theory works, each possible configuration of the field is in effect taken to make a contribution $\text{Exp}[is/\hbar]$, where s is the so-called action for the field configuration (given by the integral of the Lagrangian density—essentially a modified energy density), and \hbar is a basic scale factor for quantum effects (Planck's constant divided by 2π). In most places in the space of all possible field configurations, the value of s will vary quite quickly between nearby configurations. And assuming this variation is somehow random, the contributions of these nearby configurations will tend to cancel out. But inevitably there will be some places in the space where s is stationary (has zero variational derivative) with respect to changes in fields. And in some approximation the field configurations in these places can be expected to dominate the path integral. But it turns out that these field configurations are exactly the ones that satisfy the partial differential equations for the classical version of the field theory. (This is analogous to what happens for example in classical diffraction theory,

where there is an analog of the path integral—with \hbar replaced by inverse frequency—whose stationary points correspond through the so-called eikonal approximation to rays in geometrical optics.) In cases like QED and QCD the most obvious solutions to the classical equations are ones in which all fields are zero. And indeed standard perturbation theory is based on starting from these and then looking at the expansion of $\text{Exp}[is/\hbar]$ in powers of the coupling constant. But while this works for QED, it is only adequate for QCD in situations where the effective coupling is small. And indeed in other situations it seems likely that there will be all sorts of other solutions to the classical equations that become important. But apart from a few special cases with high symmetry, remarkably little is known about solutions to the classical equations even for pure gluon fields. No doubt the analog of turbulence can occur, and certainly there is sensitive dependence on initial conditions (even non-Abelian plane waves involve iterated maps that show this). Presumably much like in fluids there are various coherent structures such as color flux tubes and glueballs. But I doubt that states involving organized arrangements of these are common. And in general when there is strong coupling the path integral will potentially be dominated by large numbers of configurations not close to classical solutions.

In studying quantum field theories it has been common to consider effectively replacing time coordinates t by it to go from ordinary Minkowski space to Euclidean space (see page 1043). But while there is no problem in doing this at a formal mathematical level—and indeed the expressions one gets from Feynman diagrams can always be analytically continued in this way—what general correspondence there is for actual physical processes is far from clear. Formally continuing to Euclidean space makes path integrals easier to define with traditional mathematics, and gives them weights of the form $\text{Exp}[-\beta s]$ —analogous to constant temperature systems in statistical mechanics. Discretizing yields lattice gauge theories with energy functions involving for example $\text{Cos}[\theta_i - \theta_j]$ for color directions at adjacent sites. And Monte Carlo studies of such theories suggest all sorts of complex behavior, often similar in outline from what appears to occur in the corresponding classical field theories. (It seems conceivable that asymptotic freedom could lead to an analog of damping at small scales roughly like viscosity in turbulent fluids.)

One of the apparent implications of QCD is the confinement of quarks and gluons inside color-neutral hadrons. And at some level this is presumably a reflection of the fact that QCD forces get stronger rather than weaker with increasing distance. The beginnings of this are visible in perturbation

theory in the increase of the effective coupling with distance associated with asymptotic freedom. (In QED effective couplings decrease slightly with distance because fields get screened by virtual electron-positron pairs. The same happens with virtual quarks in QCD, but a larger effect is virtual gluon pairs whose color magnetic moments line up with a color field and serve to increase it.) At larger distances something like color flux tubes that act like elastic strings may form. But no detailed way to get confinement with purely classical gluon fields is known. In the quantum case, a sign of confinement would be exponential decrease with spacetime area of the average phase of color flux through so-called Wilson loops—and this is achieved if there is in a sense maximal randomness in field configurations. (Note that it is not inconceivable that the formal problem of whether quarks and gluons can ever escape to infinity starting from some given class of field configurations may in general be undecidable.)

■ **Vacuum fluctuations.** As an analog of the uncertainty principle, one of the implications of the basic formalism of quantum theory is that an ordinary quantum field can in a sense never maintain precisely zero value, but must always show certain fluctuations—even in what one considers the vacuum. And in terms of Feynman diagrams the way this happens is by virtual particle-antiparticle pairs of all types and all energy-momenta continually forming and annihilating at all points in the vacuum. Insofar as such vacuum fluctuations are always exactly the same, however, they presumably cannot be detected. (In the formalism of quantum field theory, they are usually removed by so-called normal ordering. But without this every mode of any quantum system will show a zero-point energy $\hbar\omega/2$ —positive in sign for bosons and negative for fermions, cancelling for perfect supersymmetry. Quite what gravitational effects such zero-point energy might have has never been clear.) If one somehow changes the space in which a vacuum exists, there can be directly observable effects of vacuum fluctuations. An example is the 1948 Casimir effect—in which the absence of low-energy (long wavelength) virtual particle pairs in the space between two metal plates (but not in the infinite space outside) leads to a small but measurable force of attraction between them. The different detailed patterns of modes of different fields in different spaces can lead to very different effective vacuum energies—often negative. And at least with the idealization of impermeable classical conducting boundaries one predicts (based on work of mine from 1981) the peculiar effect that closed cycles can be set up that systematically extract energy from vacuum fluctuations in a photon field.

If one has moving boundaries it turns out that vacuum fluctuations can in effect be viewed as producing real particles. And as known since the 1960s, the same is true for expanding universes. What happens in essence is that the modes of fields in different background spacetime structures differ to the point where zero-point excitations seem like actual particle excitations to a detector or observer calibrated to fields in ordinary fixed flat infinite spacetime. And in fact just uniform acceleration turns out to make detectors register real particles in a vacuum—in this case with a thermal spectrum at a temperature proportional to the acceleration. (Uniform rotation also leads to real particles, but apparently with a different spectrum.) As expected from the equivalence principle, a uniform gravitational field should produce the same effect. (Uniform electric fields lead in a formally similar way to production of charged particles.) And as pointed out by Stephen Hawking in 1974, black holes should also generate thermal radiation (at a temperature $\hbar c^3/(8\pi G k M)$). A common interpretation is that the radiated particles are somehow ones left behind when the other particle in a virtual pair goes inside the event horizon. (A similar explanation can be given for uniform acceleration—for which there is also an event horizon.) There has been much discussion of the idea that Hawking radiation somehow shows pure quantum states spontaneously turning into mixed ones, more or less as in quantum measurements. But presumably this is just a reflection of the idealization involved in looking at quantum fields in a fixed background classical spacetime. And indeed work in string theory in the mid-1990s may suggest ways in which quantum gravity configurations of black hole surfaces could maintain the information needed for the whole system to act as a pure state.

■ **Page 542 • Quantum measurement.** The basic mathematical formalism used in standard quantum theory to describe pure quantum processes deals just with vectors of probability amplitudes. Yet our everyday experience of the physical world is that we observe definite things to happen. And the way this is normally captured is by saying that when an observation is made the vector of amplitudes is somehow replaced by its projection s into a subspace corresponding to the outcome seen—with the probability of getting the outcome being taken to be determined by $s \cdot \text{Conjugate}[s]$.

At the level of pure quantum processes, the standard rules of quantum theory say that amplitudes should be added as complex numbers—with the result that they can for example potentially cancel each other, and generally lead to wave-like interference phenomena. But after an observation is made, it is in effect assumed that a system can be described by ordinary real-number probabilities—so that for example no

interference is possible. (At a formal level, results of pure quantum processes are termed pure quantum states, and are characterized by vectors of probability amplitudes; results of all possible observations are termed mixed states, and are in effect represented as mixtures of pure states.)

Ever since the 1930s there have been questions about just what should count as an observation. To explain everyday experience, conscious perception presumably always must. But it was not clear whether the operation of inanimate measuring devices of various kinds also should. And a major apparent problem was that if everything—including the measuring device—is supposed to be treated as part of the same quantum system, then all of it must follow the rules for pure quantum processes, which do not explicitly include any reduction of the kind supposed to occur in observations.

One approach to getting around this suggested in the late 1950s is the many-worlds interpretation (see page 1035): that there is in a sense a universal pure quantum process that involves all possible outcomes for every conceivable observation, and that represents the tree of all possible threads of history—but that in a particular thread, involving a particular sequence of tree branches, and representing a particular thread of experience for us, there is in effect a reduction in the pure quantum process at each branch point. Similar schemes have been popular in quantum cosmology since the early 1990s in connection with studying wave functions for the complete universe.

A quite different—and I think much more fruitful—approach is to consider analyzing actual potential measurement processes in the context of ordinary quantum mechanics. For even if one takes these processes to be pure quantum ones, what I believe is that in almost all cases appropriate idealized limits of them will reproduce what are in effect the usual rules for observations in quantum theory. A key point is that for one to consider something a reasonable measurement it must in a sense yield a definitive result. And in the context of standard quantum theory this means that somehow all the probability amplitudes associated with the measuring device must in effect be concentrated in specific outcomes—with no significant interference between different outcomes.

If one has just a few quantum particles—governed say by an appropriate Schrödinger equation—then presumably there can be no such concentration. But with a sufficiently large number of particles—and appropriate interactions—one expects that there can be. At first this might seem impossible. For the basic rules for pure quantum processes are entirely reversible (unitary). So one might think that if the evolution of a system leads to concentration of amplitudes, then it

should equally well lead to the reverse. But the crucial point is that while this may in principle be possible, it may essentially never happen in practice—just like classical reversible systems essentially never show behavior that goes against the Second Law of thermodynamics. As suggested by the main text, the details in the quantum measurement case are slightly more complicated—since to represent multiple outcomes measuring devices typically have to have the analogs of multiple equilibrium states. But the basic phenomena are ultimately very similar—and both are in effect based on the presence of microscopic randomness. (In a quantum system the randomness serves to give collections of complex numbers whose average is essentially always zero.)

This so-called decoherence approach was discussed in the 1930s, and finally began to become popular in the 1980s. But to make it work there needs to be some source of appropriate randomness. And almost without exception what has been assumed is that this must come through the first mechanism discussed in Chapter 7: that there is somehow randomness present in the environment that always gets into the system one is looking at. Various different specific mechanisms for this have been suggested, including ones based on ambient low-frequency photons, background quantum vacuum fluctuations and background spacetime metric fluctuations. (A somewhat related proposal involves quantum gravity effects in which irreversibility is assumed to be generated through analogs of the black hole processes mentioned in the previous note.) And indeed in recent practical experiments where simple pure quantum states have carefully been set up, they seem to be destroyed by randomness from the environment on timescales of at most perhaps microseconds. But this does not mean that in more complicated systems more characteristic of real measuring devices there may not be other sources of randomness that end up dominating.

One might imagine that a possibility would be the second mechanism for randomness from Chapter 7, based on ideas of chaos theory. For certainly in the standard formalism, quantum probability amplitudes are taken to be continuous quantities in which an arbitrary number of digits can be specified. But at least for a single particle, the Schrödinger equation is in all ways linear, and so it cannot support any kind of real sensitivity to initial conditions, or even to parameters. But when many particles are involved the situation can presumably be different, as it definitely can be in quantum field theory (see page 1061).

I suspect, however, that in fact the most important source of randomness in most cases will instead be the phenomenon of intrinsic randomness generation that I first discovered in systems like the rule 30 cellular automaton. Just like in so

many other areas, the emphasis on traditional mathematical methods has meant that for the most part fundamental studies have been made only on quantum systems that in the end turn out to have fairly simple behavior. Yet even within the standard formalism of quantum theory there are actually no doubt many closed systems that intrinsically manage to produce complex and seemingly random behavior even with very simple parameters and initial conditions. And in fact some clear signs of this were already present in studies of so-called quantum chaos in the 1980s—although most of the specific cases actually considered involved time-independent constraint satisfaction, not explicit time evolution. Curiously, what the Principle of Computational Equivalence suggests is that when quantum systems intrinsically produce apparent randomness they will in the end typically be capable of doing computations just as sophisticated as any other system—and in particular just as sophisticated as would be involved in conscious perception.

As a practical matter, mechanisms like intrinsic randomness generation presumably allow systems involving macroscopic numbers of particles to yield behavior in which interference becomes astronomically unlikely. But to reproduce the kind of exact reduction of probability amplitudes that is implied by the standard formalism of quantum theory inevitably requires taking the limit of an infinite system. Yet the Principle of Computational Equivalence suggests that the results of such a limit will typically be non-computable. (Using quantum field theory to represent infinite numbers of particles presumably cannot help; after appropriate analysis of the fairly sophisticated continuous mathematics involved, exactly the same computational issues should arise.)

It is often assumed that quantum systems should somehow easily be able to generate perfect randomness. But any sequence of bits one extracts must be deduced from a corresponding sequence of measurements. And certainly in practice—as mentioned on pages 303 and 970—correlations in the internal states of measuring devices between successive measurements will tend to lead to deviations from randomness. Whatever generates randomness and brings measuring devices back to equilibrium will eventually damp out such correlations. But insofar as measuring devices must formally involve infinite numbers of particles this process will formally require infinitely many steps. So this means that in effect an infinite computation is actually being done to generate each new bit. But with this amount of computation there are many ways to generate random bits. And in fact an infinite computation could even in principle produce algorithmic randomness (see page 1067) of the kind that is implicitly suggested by the

traditional continuous mathematical formalism of quantum theory. So what this suggests is that there may in the end be no clear way to tell whether randomness is coming from an underlying quantum process that is being measured, or from the actual process of measurement. And indeed when it comes to more realistic finite measuring devices I would not be surprised if most of the supposed quantum randomness they measure is actually more properly attributed to intrinsic randomness generation associated with their internal mechanisms.

■ **Page 543 · Bell's inequalities.** In classical physics one can set up light waves that are linearly polarized with any given orientation. And if these hit polarizing (“anti-glare”) filters whose orientation is off by an angle θ , then the waves transmitted will have intensity $\text{Cos}[\theta]^2$. In quantum theory the quantization of particle spin implies that any photon hitting a polarizing filter will always either just go through or be absorbed—so that in effect its spin measured relative to the orientation of the polarizer is either +1 or -1. A variety of atomic and other processes give pairs of photons that are forced to have total spin 0. And in what is essentially the Einstein-Podolsky-Rosen setup mentioned on page 1058 one can ask what happens if such photons are made to hit polarizers whose orientations differ by angle θ . In ordinary quantum theory, a straightforward calculation implies that the expected value of the product of the two measured spin values will be $-\text{Cos}[\theta]$. But now imagine instead that when each photon is produced it is assigned some “hidden variable” ϕ that in effect explicitly specifies the angle of its polarization. Then assume that a polarizer oriented at 0° will measure the spin of such a photon to have value $f[\phi]$ for some fixed function f . Now the expected value of the product of the two measured spin values is found just by averaging over ϕ as

$$\text{Integrate}[f[\phi]f[\theta - \phi], \{\phi, 0, 2\pi\}]/(2\pi)$$

A version of Bell's inequalities is then that this integral can decrease with θ no faster than $\theta/(2\pi) - 1$ —as achieved when $f = \text{Sign}$. (In 3D ϕ must be extended to a sphere, but the same final result holds.) Yet as mentioned on page 1058, actual experiments show that in fact the decrease with θ is more rapid—and is instead consistent with the quantum theory result $-\text{Cos}[\theta]$. So what this means is that there is in a sense more correlation between measurements made on separated photons than can apparently be explained by the individual photons carrying any kind of explicit hidden property. (In the standard formalism of quantum theory this is normally explained by saying that the two photons can only meaningfully be considered as part of a single “entangled” state. Note that because of the probabilistic nature of the

correlations it turns out to be impossible to use them to do anything that would normally be considered communicating information faster than the speed of light.)

A basic assumption in deriving Bell's inequalities is that the choice of polarizer angle for measuring one photon is not affected by the choice of angle for the other. And indeed experiments have been done which try to enforce this by choosing the angles for the polarizers only just before the photons reach them—and too close in time for a light signal to get from one to the other. Such experiments again show violations of Bell's inequalities. But inevitably the actually devices that work out choices of polarizer angles must be in causal contact as part of setting up the experiment. And although it seems contrived, it is thus at least conceivable that with a realistic model for their time evolution such devices could end up operating in just such a way as to yield observed violations of Bell's inequalities.

Another way to get violations of Bell's inequalities is to allow explicit instantaneous propagation of information. But traditional models involving for example a background quantum potential again seem quite contrived, and difficult to generalize to relativistic cases. The approach I discuss in the main text is quite different, in effect using the idea that in a network model of space there can be direct connections between particles that do not in a sense ever have to go through ordinary intermediate points in space.

When set up for pairs of particles, Bell's inequalities tend just to provide numerical constraints on probabilities. But for

triples of particles, it was noticed in the late 1980s that they can give constraints that force probabilities to be 0 or 1, implying that with the assumptions made, certain configurations of measurement results are simply impossible.

In quantum field theory the whole concept of measurement is much less developed than in quantum mechanics—not least because in field theory it is much more difficult to factor out subsystems, and so to avoid having to give explicit descriptions of measuring devices. But at least in axiomatic quantum field theory it is typically assumed that one can somehow measure expectation values of any suitably smeared product of field operators. (It is possible that these could be reconstructed from combinations of idealized scattering experiments). And to get a kind of analog of Bell's inequalities one can look at correlations defined by such expectation values for field operators at spacelike-separated points (too close in time for light signals to get from one to another). And it then turns out that even in the vacuum state the vacuum fluctuations that are present show nonzero such correlations—an analog of ordinary quantum mechanical entanglement. (In a non-interacting approximation these correlations turn out to be as large as is mathematically possible, but fall off exponentially outside the light cone, with exponents determined by the smallest particle mass or the measurement resolution.) In a sense, however, the presence of such correlations is just a reflection of the idealized way in which the vacuum state is set up—with each field mode determined all at once for the whole system.