

3 Introduction to computer science

In natural science, Nature has given us a world and we're just to discover its laws. In computers, we can stuff laws into it and create a world.

– Alan Kay

Our field is still in its embryonic stage. It's great that we haven't been around for 2000 years. We are still at a stage where very, very important results occur in front of our eyes.

– Michael Rabin, on computer science

Algorithms are the key concept of computer science. An algorithm is a precise recipe for performing some task, such as the elementary algorithm for adding two numbers which we all learn as children. This chapter outlines the modern theory of algorithms developed by computer science. Our fundamental model for algorithms will be the *Turing machine*. This is an idealized computer, rather like a modern personal computer, but with a simpler set of basic instructions, and an idealized unbounded memory. The apparent simplicity of Turing machines is misleading; they are very powerful devices. We will see that they can be used to execute any algorithm whatsoever, even one running on an apparently much more powerful computer.

The fundamental question we are trying to address in the study of algorithms is: what resources are required to perform a given computational task? This question splits up naturally into two parts. First, we'd like to understand what computational tasks are possible, preferably by giving explicit algorithms for solving specific problems. For example, we have many excellent examples of algorithms that can quickly sort a list of numbers into ascending order. The second aspect of this question is to demonstrate *limitations* on what computational tasks may be accomplished. For example, lower bounds can be given for the number of operations that must be performed by any algorithm which sorts a list of numbers into ascending order. Ideally, these two tasks – the finding of algorithms for solving computational problems, and proving limitations on our ability to solve computational problems – would dovetail perfectly. In practice, a significant gap often exists between the best techniques known for solving a computational problem, and the most stringent limitations known on the solution. The purpose of this chapter is to give a broad overview of the tools which have been developed to aid in the analysis of computational problems, and in the construction and analysis of algorithms to solve such problems.

Why should a person interested in *quantum* computation and *quantum* information spend time investigating *classical* computer science? There are three good reasons for this effort. First, classical computer science provides a vast body of concepts and techniques which may be reused to great effect in quantum computation and quantum information. Many of the triumphs of quantum computation and quantum information have come by combining existing ideas from computer science with novel ideas from quantum

mechanics. For example, some of the fast algorithms for quantum computers are based upon the Fourier transform, a powerful tool utilized by many classical algorithms. Once it was realized that quantum computers could perform a type of Fourier transform much more quickly than classical computers this enabled the development of many important quantum algorithms.

Second, computer scientists have expended great effort understanding what resources are required to perform a given computational task on a classical computer. These results can be used as the basis for a comparison with quantum computation and quantum information. For example, much attention has been focused on the problem of finding the prime factors of a given number. On a classical computer this problem is believed to have no ‘efficient’ solution, where ‘efficient’ has a meaning we’ll explain later in the chapter. What is interesting is that an efficient solution to this problem *is* known for quantum computers. The lesson is that, for this task of finding prime factors, there appears to be a *gap* between what is possible on a classical computer and what is possible on a quantum computer. This is both intrinsically interesting, and also interesting in the broader sense that it suggests such a gap may exist for a wider class of computational problems than merely the finding of prime factors. By studying this specific problem further, it may be possible to discern features of the problem which make it more tractable on a quantum computer than on a classical computer, and then act on these insights to find interesting quantum algorithms for the solution of other problems.

Third, and most important, there is learning to *think* like a computer scientist. Computer scientists think in a rather different style than does a physicist or other natural scientist. Anybody wanting a deep understanding of quantum computation and quantum information must learn to think like a computer scientist at least some of the time; they must instinctively know what problems, what techniques, and most importantly what problems are of greatest interest to a computer scientist.

The structure of this chapter is as follows. In Section 3.1 we introduce two models for computation: the Turing machine model, and the circuit model. The Turing machine model will be used as our fundamental model for computation. In practice, however, we mostly make use of the circuit model of computation, and it is this model which is most useful in the study of quantum computation. With our models for computation in hand, the remainder of the chapter discusses resource requirements for computation. Section 3.2 begins by overviewing the computational tasks we are interested in as well as discussing some associated resource questions. It continues with a broad look at the key concepts of *computational complexity*, a field which examines the time and space requirements necessary to solve particular computational problems, and provides a broad classification of problems based upon their difficulty of solution. Finally, the section concludes with an examination of the energy resources required to perform computations. Surprisingly, it turns out that the energy required to perform a computation can be made vanishingly small, provided one can make the computation reversible. We explain how to construct reversible computers, and explain some of the reasons they are important both for computer science and for quantum computation and quantum information. Section 3.3 concludes the chapter with a broad look at the entire field of computer science, focusing on issues of particular relevance to quantum computation and quantum information.

3.1 Models for computation

...algorithms are concepts that have existence apart from any programming language.

– Donald Knuth

What does it mean to have an *algorithm* for performing some task? As children we all learn a procedure which enables us to add together any two numbers, no matter how large those numbers are. This is an example of an algorithm. Finding a mathematically precise formulation of the concept of an algorithm is the goal of this section.

Historically, the notion of an algorithm goes back centuries; undergraduates learn Euclid's two thousand year old algorithm for finding the greatest common divisor of two positive integers. However, it wasn't until the 1930s that the fundamental notions of the modern theory of algorithms, and thus of computation, were introduced, by Alonzo Church, Alan Turing, and other pioneers of the computer era. This work arose in response to a profound challenge laid down by the great mathematician David Hilbert in the early part of the twentieth century. Hilbert asked whether or not there existed some algorithm which could be used, in principle, to solve all the problems of mathematics. Hilbert expected that the answer to this question, sometimes known as the *entscheidungsproblem*, would be yes.

Amazingly, the answer to Hilbert's challenge turned out to be no: there is no algorithm to solve all mathematical problems. To prove this, Church and Turing had to solve the deep problem of capturing in a mathematical definition what we mean when we use the intuitive concept of an algorithm. In so doing, they laid the foundations for the modern theory of algorithms, and consequently for the modern theory of computer science.

In this chapter, we use two ostensibly different approaches to the theory of computation. The first approach is that proposed by Turing. Turing defined a class of machines, now known as *Turing machines*, in order to capture the notion of an algorithm to perform a computational task. In Section 3.1.1, we describe Turing machines, and then discuss some of the simpler variants of the Turing machine model. The second approach is via the *circuit model* of computation, an approach that is especially useful as preparation for our later study of quantum computers. The circuit model is described in Section 3.1.2. Although these models of computation appear different on the surface, it turns out that they are equivalent. Why introduce more than one model of computation, you may ask? We do so because different models of computation may yield different insights into the solution of specific problems. Two (or more) ways of thinking about a concept are better than one.

3.1.1 Turing machines

The basic elements of a Turing machine are illustrated in Figure 3.1. A Turing machine contains four main elements: (a) a *program*, rather like an ordinary computer; (b) a *finite state control*, which acts like a stripped-down microprocessor, co-ordinating the other operations of the machine; (c) a *tape*, which acts like a computer memory; and (d) a *read-write tape-head*, which points to the position on the tape which is currently readable or writable. We now describe each of these four elements in more detail.

The *finite state control* for a Turing machine consists of a finite set of *internal states*,

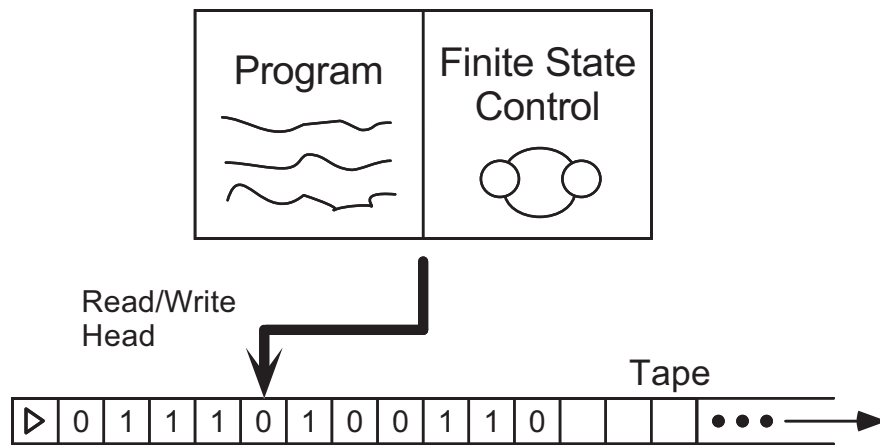


Figure 3.1. Main elements of a Turing machine. In the text, blanks on the tape are denoted by a ‘b’. Note the ▷ marking the left hand end of the tape.

q_1, \dots, q_m . The number m is allowed to be varied; it turns out that for m sufficiently large this does not affect the power of the machine in any essential way, so without loss of generality we may suppose that m is some fixed constant. The best way to think of the finite state control is as a sort of microprocessor, co-ordinating the Turing machine’s operation. It provides temporary storage off-tape, and a central place where all processing for the machine may be done. In addition to the states q_1, \dots, q_m , there are also two special internal states, labelled q_s and q_h . We call these the *starting state* and the *halting state*, respectively. The idea is that at the beginning of the computation, the Turing machine is in the starting state q_s . The execution of the computation causes the Turing machine’s internal state to change. If the computation ever finishes, the Turing machine ends up in the state q_h to indicate that the machine has completed its operation.

The Turing machine *tape* is a one-dimensional object, which stretches off to infinity in one direction. The tape consists of an infinite sequence of *tape squares*. We number the tape squares $0, 1, 2, 3, \dots$. The tape squares each contain one symbol drawn from some *alphabet*, Γ , which contains a finite number of distinct symbols. For now, it will be convenient to assume that the alphabet contains four symbols, which we denote by $0, 1, b$ (the ‘blank’ symbol), and \triangleright , to mark the left hand edge of the tape. Initially, the tape contains a \triangleright at the left hand end, a finite number of 0s and 1s, and the rest of the tape contains blanks. The *read-write tape-head* identifies a single square on the Turing machine tape as the square that is currently being accessed by the machine.

Summarizing, the machine starts its operation with the finite state control in the state q_s , and with the read-write head at the leftmost tape square, the square numbered 0. The computation then proceeds in a step by step manner according to the *program*, to be defined below. If the current state is q_h , then the computation has halted, and the *output* of the computation is the current (non-blank) contents of the tape.

A *program* for a Turing machine is a finite ordered list of *program lines* of the form $\langle q, x, q', x', s \rangle$. The first item in the program line, q , is a state from the set of internal states of the machine. The second item, x , is taken from the alphabet of symbols which may appear on the tape, Γ . The way the program works is that on each machine cycle, the Turing machine looks through the list of program lines in order, searching for a line $\langle q, x, \cdot, \cdot, \cdot \rangle$, such that the current internal state of the machine is q , and the symbol

being read on the tape is x . If it doesn't find such a program line, the internal state of the machine is changed to q_h , and the machine halts operation. If such a line is found, then that program line is *executed*. Execution of a program line involves the following steps: the internal state of the machine is changed to q' ; the symbol x on the tape is overwritten by the symbol x' , and the tape-head moves left, right, or stands still, depending on whether s is -1 , $+1$, or 0 , respectively. The only exception to this rule is if the tape-head is at the leftmost tape square, and $s = -1$, in which case the tape-head stays put.

Now that we know what a Turing machine is, let's see how it may be used to compute a simple function. Consider the following example of a Turing machine. The machine starts with a binary number, x , on the tape, followed by blanks. The machine has three internal states, q_1 , q_2 , and q_3 , in addition to the starting state q_s and halting state q_h . The program contains the following program lines (the numbers on the left hand side are for convenience in referring to the program lines in later discussion, and do not form part of the program):

- 1 : $\langle q_s, \triangleright, q_1, \triangleright, +1 \rangle$
- 2 : $\langle q_1, 0, q_1, b, +1 \rangle$
- 3 : $\langle q_1, 1, q_1, b, +1 \rangle$
- 4 : $\langle q_1, b, q_2, b, -1 \rangle$
- 5 : $\langle q_2, b, q_2, b, -1 \rangle$
- 6 : $\langle q_2, \triangleright, q_3, \triangleright, +1 \rangle$
- 7 : $\langle q_3, b, q_h, 1, 0 \rangle$.

What function does this program compute? Initially the machine is in the state q_s and at the left-most tape position so line 1, $\langle q_s, \triangleright, q_1, \triangleright, +1 \rangle$, is executed, which causes the tape-head to move right without changing what is written on the tape, but changing the internal state of the machine to q_1 . The next three lines of the program ensure that while the machine is in the state q_1 the tape-head will continue moving right while it reads either 0s (line 2) or 1s (line 3) on the tape, over-writing the tape contents with blanks as it goes and remaining in the state q_1 , until it reaches a tape square that is already blank, at which point the tape-head is moved one position to the left, and the internal state is changed to q_2 (line 4). Line 5 then ensures that the tape-head keeps moving left while blanks are being read by the tape-head, without changing the contents of the tape. This keeps up until the tape-head returns to its starting point, at which point it reads a \triangleright on the tape, changes the internal state to q_3 , and moves one step to the right (line 6). Line 7 completes the program, simply printing the number 1 onto the tape, and then halting.

The preceding analysis shows that this program computes the constant function $f(x) = 1$. That is, regardless of what number is input on the tape the number 1 is output. More generally, a Turing machine can be thought of as computing functions from the non-negative integers to the non-negative integers; the initial state of the tape is used to represent the input to the function, and the final state of the tape is used to represent the output of the function.

It seems as though we have gone to a very great deal of trouble to compute this simple function using our Turing machines. Is it possible to build up more complicated functions using Turing machines? For example, could we construct a machine such that when two numbers, x and y , are input on the tape with a blank to demarcate them, it will

output the sum $x + y$ on the tape? More generally, what class of functions is it possible to compute using a Turing machine?

It turns out that the Turing machine model of computation can be used to compute an enormous variety of functions. For example, it can be used to do all the basic arithmetical operations, to search through text represented as strings of bits on the tape, and many other interesting operations. Surprisingly, it turns out that a Turing machine can be used to simulate all the operations performed on a modern computer! Indeed, according to a thesis put forward independently by Church and by Turing, the Turing machine model of computation *completely captures* the notion of computing a function using an algorithm. This is known as the *Church–Turing thesis*:

The class of functions computable by a Turing machine corresponds exactly to the class of functions which we would naturally regard as being computable by an algorithm.

The Church–Turing thesis asserts an equivalence between a rigorous mathematical concept – function computable by a Turing machine – and the intuitive concept of what it means for a function to be computable by an algorithm. The thesis derives its importance from the fact that it makes the study of real-world algorithms, prior to 1936 a rather vague concept, amenable to rigorous mathematical study. To understand the significance of this point it may be helpful to consider the definition of a *continuous function* from real analysis. Every child can tell you what it means for a line to be continuous on a piece of paper, but it is far from obvious how to capture that intuition in a rigorous definition. Mathematicians in the nineteenth century spent a great deal of time arguing about the merits of various definitions of continuity before the modern definition of continuity came to be accepted. When making fundamental definitions like that of continuity or of computability it is important that good definitions be chosen, ensuring that one’s intuitive notions closely match the precise mathematical definition. From this point of view the Church–Turing thesis is simply the assertion that the Turing machine model of computation provides a good foundation for computer science, capturing the intuitive notion of an algorithm in a rigorous definition.

A priori it is not obvious that every function which we would intuitively regard as computable by an algorithm can be computed using a Turing machine. Church, Turing and many other people have spent a great deal of time gathering evidence for the Church–Turing thesis, and in sixty years no evidence to the contrary has been found. Nevertheless, it is possible that in the future we will discover in Nature a process which computes a function not computable on a Turing machine. It would be wonderful if that ever happened, because we could then harness that process to help us perform new computations which could not be performed before. Of course, we would also need to overhaul the definition of computability, and with it, computer science.

Exercise 3.1: (Non-computable processes in Nature) How might we recognize that a process in Nature computes a function not computable by a Turing machine?

Exercise 3.2: (Turing numbers) Show that single-tape Turing machines can each be given a number from the list $1, 2, 3, \dots$ in such a way that the number uniquely specifies the corresponding machine. We call this number the *Turing number* of the corresponding Turing machine. (*Hint*: Every positive integer has

a unique prime factorization $p_1^{a_1} p_2^{a_2} \dots p_k^{a_k}$, where p_i are distinct prime numbers, and a_1, \dots, a_k are non-negative integers.)

In later chapters, we will see that quantum computers also obey the Church–Turing thesis. That is, quantum computers can compute the same class of functions as is computable by a Turing machine. The difference between quantum computers and Turing machines turns out to lie in the *efficiency* with which the computation of the function may be performed – there are functions which can be computed much more efficiently on a quantum computer than is believed to be possible with a classical computing device such as a Turing machine.

Demonstrating in complete detail that the Turing machine model of computation can be used to build up all the usual concepts used in computer programming languages is beyond the scope of this book (see ‘History and further reading’ at the end of the chapter for more information). When specifying algorithms, instead of explicitly specifying the Turing machine used to compute the algorithm, we shall usually use a much higher level *pseudocode*, trusting in the Church–Turing thesis that this pseudocode can be translated into the Turing machine model of computation. We won’t give any sort of rigorous definition for pseudocode. Think of it as a slightly more formal version of English or, if you like, a sloppy version of a high-level programming language such as C++ or BASIC. Pseudocode provides a convenient way of expressing algorithms, without going into the extreme level of detail required by a Turing machine. An example use of pseudocode may be found in Box 3.2 on page 130; it is also used later in the book to describe quantum algorithms.

There are many variants on the basic Turing machine model. We might imagine Turing machines with different kinds of tapes. For example, one could consider two-way infinite tapes, or perhaps computation with tapes of more than one dimension. So far as is presently known, it is not possible to change any aspect of the Turing model in a way that is physically reasonable, and which manages to extend the class of functions computable by the model.

As an example consider a Turing machine equipped with multiple tapes. For simplicity we consider the two-tape case, as the generalization to more than two tapes is clear from this example. Like the basic Turing machine, a two-tape Turing machine has a finite number of internal states q_1, \dots, q_m , a start state q_s , and a halt state q_h . It has two tapes, each of which contain symbols from some finite alphabet of symbols, Γ . As before we find it convenient to assume that the alphabet contains four symbols, 0, 1, b and \triangleright , where \triangleright marks the left hand edge of each tape. The machine has two tape-heads, one for each tape. The main difference between the two-tape Turing machine and the basic Turing machine is in the program. Program lines are of the form $\langle q, x_1, x_2, q', x'_1, x'_2, s_1, s_2 \rangle$, meaning that if the internal state of the machine is q , tape one is reading x_1 at its current position, and tape two is reading x_2 at its current position, then the internal state of the machine should be changed to q' , x_1 overwritten with x'_1 , x_2 overwritten with x'_2 , and the tape-heads for tape one and tape two moved according to whether s_1 or s_2 are equal to +1, −1 or 0, respectively.

In what sense are the basic Turing machine and the two-tape Turing machine equivalent models of computation? They are equivalent in the sense that each computational model is able to *simulate* the other. Suppose we have a two-tape Turing machine which takes as input a bit string x on the first tape and blanks on the remainder of both tapes,

except the endpoint marker \triangleright . This machine computes a function $f(x)$, where $f(x)$ is defined to be the contents of the first tape after the Turing machine has halted. Rather remarkably, it turns out that given a two-tape Turing machine to compute f , there exists an equivalent single-tape Turing machine that is also able to compute f . We won't explain how to do this explicitly, but the basic idea is that the single-tape Turing machine *simulates* the two-tape Turing machine, using its single tape to store the contents of both tapes of the two-tape Turing machine. There is some computational overhead required to do this simulation, but the important point is that in principle it can always be done. In fact, there exists a Universal Turing machine (see Box 3.1) which can simulate any other Turing machine!

Another interesting variant of the Turing machine model is to introduce randomness into the model. For example, imagine that the Turing machine can execute a program line whose effect is the following: if the internal state is q and the tape-head reads x , then flip an unbiased coin. If the coin lands heads, change the internal state to q_{i_H} , and if it lands tails, change the internal state to q_{i_T} , where q_{i_H} and q_{i_T} are two internal states of the Turing machine. Such a program line can be represented as $\langle q, x, q_{i_H}, q_{i_T} \rangle$. However, even this variant doesn't change the essential power of the Turing machine model of computation. It is not difficult to see that we can simulate the effect of the above algorithm on a deterministic Turing machine by explicitly 'searching out' all the possible computational paths corresponding to different values of the coin tosses. Of course, this deterministic simulation may be far less efficient than the random model, but the key point for the present discussion is that the class of functions computable is not changed by introducing randomness into the underlying model.

Exercise 3.3: (Turing machine to reverse a bit string) Describe a Turing machine which takes a binary number x as input, and outputs the bits of x in reverse order. (*Hint:* In this exercise and the next it may help to use a multi-tape Turing machine and/or symbols other than $\triangleright, 0, 1$ and the blank.)

Exercise 3.4: (Turing machine to add modulo 2) Describe a Turing machine to add two binary numbers x and y modulo 2. The numbers are input on the Turing machine tape in binary, in the form x , followed by a single blank, followed by a y . If one number is not as long as the other then you may assume that it has been padded with leading 0s to make the two numbers the same length.

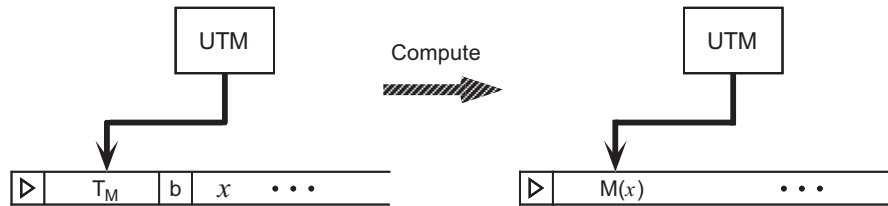
Let us return to Hilbert's entscheidungsproblem, the original inspiration for the founders of computer science. Is there an algorithm to decide all the problems of mathematics? The answer to this question was shown by Church and Turing to be no. In Box 3.2, we explain Turing's proof of this remarkable fact. This phenomenon of *undecidability* is now known to extend far beyond the examples which Church and Turing constructed. For example, it is known that the problem of deciding whether two topological spaces are topologically equivalent ('homeomorphic') is undecidable. There are simple problems related to the behavior of dynamical systems which are undecidable, as you will show in Problem 3.4. References for these and other examples are given in the end of chapter 'History and further reading'.

Besides its intrinsic interest, undecidability foreshadows a topic of great concern in computer science, and also to quantum computation and quantum information: the dis-

Box 3.1: The Universal Turing Machine

We've described Turing machines as containing three elements which may vary from machine to machine – the initial configuration of the tape, the internal states of the finite state control, and the program for the machine. A clever idea known as the *Universal Turing Machine* (UTM) allows us to fix the program and finite state control once and for all, leaving the initial contents of the tape as the only part of the machine which needs to be varied.

The Universal Turing Machine (see the figure below) has the following property. Let M be any Turing machine, and let T_M be the Turing number associated to machine M . Then on input of the binary representation for T_M followed by a blank, followed by any string of symbols x on the remainder of the tape, the Universal Turing Machine gives as output whatever machine M would have on input of x . Thus, the Universal Turing Machine is capable of simulating any other Turing machine!



The Universal Turing Machine is similar in spirit to a modern programmable computer, in which the action to be taken by the computer – the ‘program’ – is stored in memory, analogous to the bit string T_M stored at the beginning of the tape by the Universal Turing Machine. The data to be processed by the program is stored in a separate part of memory, analogous to the role of x in the Universal Turing Machine. Then some fixed hardware is used to run the program, producing the output. This fixed hardware is analogous to the internal states and the (fixed) program being executed by the Universal Turing Machine.

Describing the detailed construction of a Universal Turing Machine is beyond the scope of this book. (Though industrious readers may like to attempt the construction.) The key point is the existence of such a machine, showing that a single fixed machine can be used to run any algorithm whatsoever. The existence of a Universal Turing Machine also explains our earlier statement that the number of internal states in a Turing machine does not matter much, for provided that number m exceeds the number needed for a Universal Turing Machine, such a machine can be used to simulate a Turing machine with any number of internal states.

tion between problems which are easy to solve, and problems which are hard to solve. Undecidability provides the ultimate example of problems which are hard to solve – so hard that they are in fact impossible to solve.

Exercise 3.5: (Halting problem with no inputs) Show that given a Turing

machine M there is no algorithm to determine whether M halts when the input to the machine is a blank tape.

Exercise 3.6: (Probabilistic halting problem) Suppose we number the probabilistic Turing machines using a scheme similar to that found in Exercise 3.2 and define the probabilistic halting function $h_p(x)$ to be 1 if machine x halts on input of x with probability at least $1/2$ and 0 if machine x halts on input of x with probability less than $1/2$. Show that there is no probabilistic Turing machine which can output $h_p(x)$ with probability of correctness strictly greater than $1/2$ for all x .

Exercise 3.7: (Halting oracle) Suppose a *black box* is made available to us which takes a non-negative integer x as input, and then outputs the value of $h(x)$, where $h(\cdot)$ is the halting function defined in Box 3.2 on page 130. This type of black box is sometimes known as an *oracle* for the halting problem. Suppose we have a regular Turing machine which is augmented by the power to call the oracle. One way of accomplishing this is to use a two-tape Turing machine, and add an extra program instruction to the Turing machine which results in the oracle being called, and the value of $h(x)$ being printed on the second tape, where x is the current contents of the second tape. It is clear that this model for computation is more powerful than the conventional Turing machine model, since it can be used to compute the halting function. Is the halting problem for this model of computation undecidable? That is, can a Turing machine aided by an oracle for the halting problem decide whether a program for the Turing machine with oracle will halt on a particular input?

3.1.2 Circuits

Turing machines are rather idealized models of computing devices. Real computers are *finite* in size, whereas for Turing machines we assumed a computer of unbounded size. In this section we investigate an alternative model of computation, the *circuit model*, that is equivalent to the Turing machine in terms of computational power, but is more convenient and realistic for many applications. In particular the circuit model of computation is especially important as preparation for our investigation of quantum computers.

A circuit is made up of *wires* and *gates*, which carry information around, and perform simple computational tasks, respectively. For example, Figure 3.2 shows a simple circuit which takes as input a single bit, a . This bit is passed through a NOT gate, which flips the bit, taking 1 to 0 and 0 to 1. The wires before and after the NOT gate serve merely to carry the bit to and from the NOT gate; they can represent movement of the bit through space, or perhaps just through time.

More generally, a circuit may involve many input and output bits, many wires, and many logical gates. A *logic gate* is a function $f : \{0, 1\}^k \rightarrow \{0, 1\}^l$ from some fixed number k of *input bits* to some fixed number l of *output bits*. For example, the NOT gate is a gate with one input bit and one output bit which computes the function $f(a) = 1 \oplus a$, where a is a single bit, and \oplus is modulo 2 addition. It is also usual to make the convention that no loops are allowed in the circuit, to avoid possible instabilities, as illustrated in Figure 3.3. We say such a circuit is *acyclic*, and we adhere to the convention that circuits in the circuit model of computation be acyclic.

Box 3.2: The halting problem

In Exercise 3.2 you showed that each Turing machine can be uniquely associated with a number from the list $1, 2, 3, \dots$. To solve Hilbert's problem, Turing used this numbering to pose the *halting problem*: does the machine with Turing number x halt upon input of the number y ? This is a well posed and interesting mathematical problem. After all, it is a matter of some considerable interest to us whether our algorithms halt or not. Yet it turns out that there is no algorithm which is capable of solving the halting problem. To see this, Turing asked whether there is an algorithm to solve an even more specialized problem: does the machine with Turing number x halt upon input of the same number x ? Turing defined the *halting function*,

$$h(x) \equiv \begin{cases} 0 & \text{if machine number } x \text{ does not halt upon input of } x \\ 1 & \text{if machine number } x \text{ halts upon input of } x. \end{cases}$$

If there is an algorithm to solve the halting problem, then there surely is an algorithm to evaluate $h(x)$. We will try to reach a contradiction by supposing such an algorithm exists, denoted by $\text{HALT}(x)$. Consider an algorithm computing the function $\text{TURING}(x)$, with pseudocode

```

TURING(x)

y = HALT(x)
if y = 0 then
    halt
else
    loop forever
end if
  
```

Since HALT is a valid program, TURING must also be a valid program, with some Turing number, t . By definition of the halting function, $h(t) = 1$ if and only if TURING halts on input of t . But by inspection of the program for TURING , we see that TURING halts on input of t if and only if $h(t) = 0$. Thus $h(t) = 1$ if and only if $h(t) = 0$, a contradiction. Therefore, our initial assumption that there is an algorithm to evaluate $h(x)$ must have been wrong. We conclude that there is no algorithm allowing us to solve the halting problem.

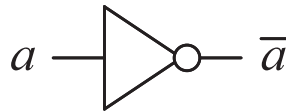


Figure 3.2. Elementary circuit performing a single NOT gate on a single input bit.

There are many other elementary logic gates which are useful for computation. A partial list includes the AND gate, the OR gate, the XOR gate, the NAND gate, and the NOR gate. Each of these gates takes two bits as input, and produces a single bit as output. The AND gate outputs 1 if and only if both of its inputs are 1. The OR gate outputs 1 if

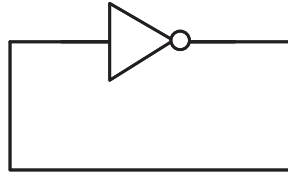


Figure 3.3. Circuits containing cycles can be unstable, and are not usually permitted in the circuit model of computation.

and only if at least one of its inputs is 1. The XOR gate outputs the sum, modulo 2, of its inputs. The NAND and NOR gates take the AND and OR, respectively, of their inputs, and then apply a NOT to whatever is output. The action of these gates is illustrated in Figure 3.4.

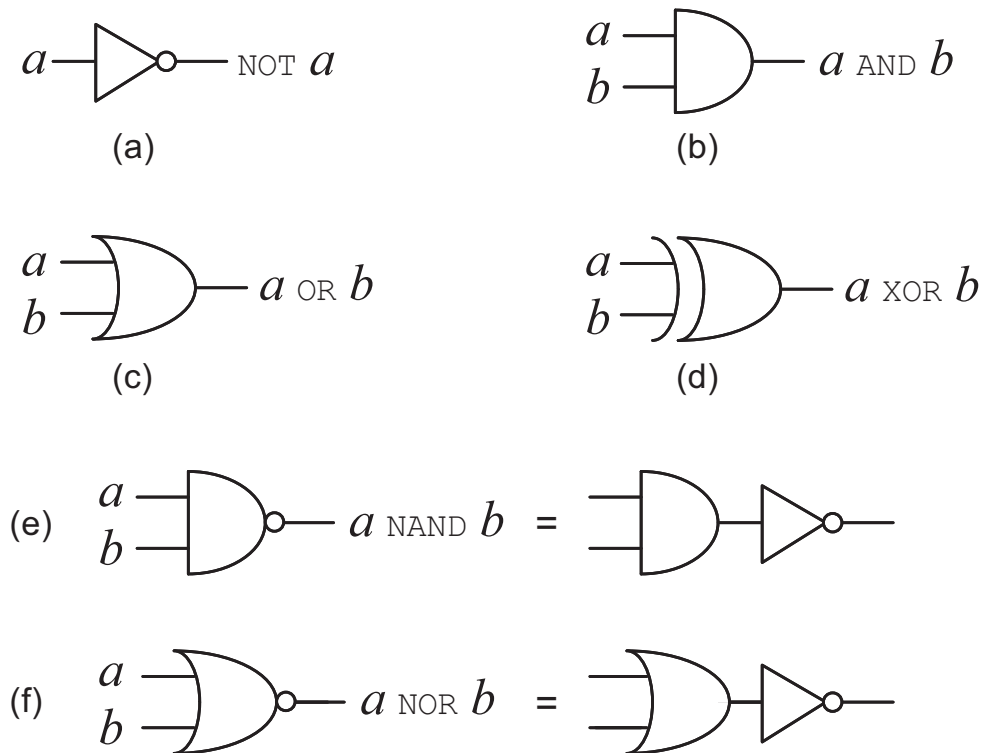


Figure 3.4. Elementary circuits performing the AND, OR, XOR, NAND, and NOR gates.

There are two important ‘gates’ missing from Figure 3.4, namely the FANOUT gate and the CROSSOVER gate. In circuits we often allow bits to ‘divide’, replacing a bit with two copies of itself, an operation referred to as FANOUT. We also allow bits to CROSSOVER, that is, the value of two bits are interchanged. A third operation missing from Figure 3.4, not really a logic gate at all, is to allow the preparation of extra *ancilla* or *work* bits, to allow extra working space during the computation.

These simple circuit elements can be put together to perform an enormous variety of computations. Below we’ll show that these elements can be used to compute any function whatsoever. In the meantime, let’s look at a simple example of a circuit which adds two n bit integers, using essentially the same algorithm taught to school-children around the

world. The basic element in this circuit is a smaller circuit known as a *half-adder*, shown in Figure 3.5. A half-adder takes two bits, x and y , as input, and outputs the sum of the bits $x \oplus y$ modulo 2, together with a carry bit set to 1 if x and y are both 1, or 0 otherwise.

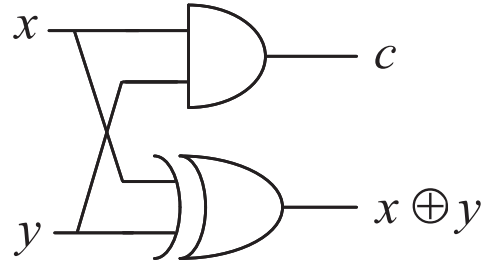


Figure 3.5. Half-adder circuit. The carry bit c is set to 1 when x and y are both 1, otherwise it is 0.

Two cascaded half-adders may be used to build a *full-adder*, as shown in Figure 3.6. A full-adder takes as input three bits, x , y , and c . The bits x and y should be thought of as data to be added, while c is a carry bit from an earlier computation. The circuit outputs two bits. One output bit is the modulo 2 sum, $x \oplus y \oplus c$ of all three input bits. The second output bit, c' , is a carry bit, which is set to 1 if two or more of the inputs is 1, and is 0 otherwise.

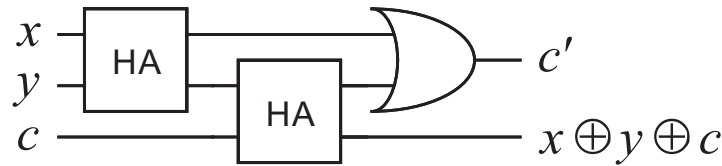


Figure 3.6. Full-adder circuit.

By cascading many of these full-adders together we obtain a circuit to add two n -bit integers, as illustrated in Figure 3.7 for the case $n = 3$.

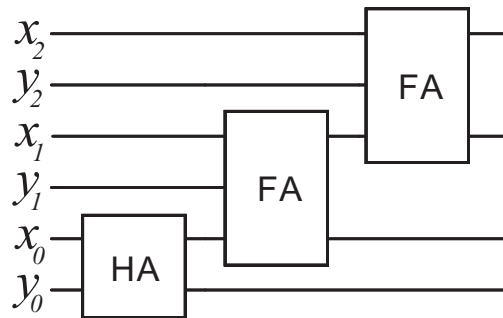


Figure 3.7. Addition circuit for two three-bit integers, $x = x_2x_1x_0$ and $y = y_2y_1y_0$, using the elementary algorithm taught to school-children.

We claimed earlier that just a few *fixed* gates can be used to compute *any* function $f : \{0, 1\}^n \rightarrow \{0, 1\}^m$ whatsoever. We will now prove this for the simplified case of a function $f : \{0, 1\}^n \rightarrow \{0, 1\}$ with n input bits and a single output bit. Such a function

is known as a *Boolean function*, and the corresponding circuit is a *Boolean circuit*. The general universality proof follows immediately from the special case of Boolean functions. The proof is by induction on n . For $n = 1$ there are four possible functions: the identity, which has a circuit consisting of a single wire; the bit flip, which is implemented using a single NOT gate; the function which replaces the input bit with a 0, which can be obtained by ANDing the input with a work bit initially in the 0 state; and the function which replaces the input with a 1, which can be obtained by ORing the input with a work bit initially in the 1 state.

To complete the induction, suppose that any function on n bits may be computed by a circuit, and let f be a function on $n + 1$ bits. Define n -bit functions f_0 and f_1 by $f_0(x_1, \dots, x_n) \equiv f(0, x_1, \dots, x_n)$ and $f_1(x_1, \dots, x_n) \equiv f(1, x_1, \dots, x_n)$. These are both n -bit functions, so by the inductive hypothesis there are circuits to compute these functions.

It is now an easy matter to design a circuit which computes f . The circuit computes both f_0 and f_1 on the last n bits of the input. Then, depending on whether the first bit of the input was a 0 or a 1 it outputs the appropriate answer. A circuit to do this is shown in Figure 3.8. This completes the induction.

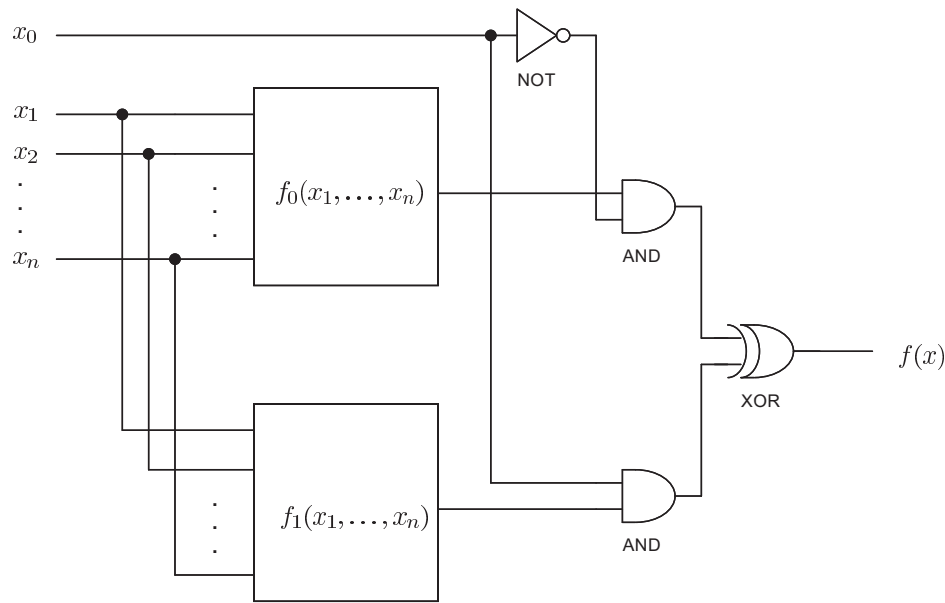


Figure 3.8. Circuit to compute an arbitrary function f on $n + 1$ bits, assuming by induction that there are circuits to compute the n -bit functions f_0 and f_1 .

Five elements may be identified in the universal circuit construction: (1) wires, which preserve the states of the bits; (2) ancilla bits prepared in standard states, used in the $n = 1$ case of the proof; (3) the FANOUT operation, which takes a single bit as input and outputs two copies of that bit; (4) the CROSSOVER operation, which interchanges the value of two bits; and (5) the AND, XOR, and NOT gates. In Chapter 4 we'll define the quantum circuit model of computation in a manner analogous to classical circuits. It is interesting to note that many of these five elements pose some interesting challenges when extending to the quantum case: it is not necessarily obvious that good quantum wires for the preservation of qubits can be constructed, even in principle, the FANOUT

operation cannot be performed in a straightforward manner in quantum mechanics, due to the no-cloning theorem (as explained in Section 1.3.5), and the AND and XOR gates are not invertible, and thus can't be implemented in a straightforward manner as unitary quantum gates. There is certainly plenty to think about in defining a quantum circuit model of computation!

Exercise 3.8: (Universality of NAND) Show that the NAND gate can be used to simulate the AND, XOR and NOT gates, provided wires, ancilla bits and FANOUT are available.

Let's return from our brief quantum digression, to the properties of classical circuits. We claimed earlier that the Turing machine model is equivalent to the circuit model of computation. In what sense do we mean the two models are equivalent? On the face of it, the two models appear quite different. The unbounded nature of a Turing machine makes them more useful for abstractly specifying what it is we mean by an algorithm, while circuits more closely capture what an actual physical computer does.

The two models are connected by introducing the notion of a *uniform circuit family*. A *circuit family* consists of a collection of circuits, $\{C_n\}$, indexed by a positive integer n . The circuit C_n has n input bits, and may have any finite number of extra work bits, and output bits. The output of the circuit C_n , upon input of a number x of at most n bits in length, is denoted by $C_n(x)$. We require that the circuits be *consistent*, that is, if $m < n$ and x is at most m bits in length, then $C_m(x) = C_n(x)$. The function computed by the circuit family $\{C_n\}$ is the function $C(\cdot)$ such that if x is n bits in length then $C(x) = C_n(x)$. For example, consider a circuit C_n that squares an n -bit number. This defines a family of circuits $\{C_n\}$ that computes the function, $C(x) = x^2$, where x is any positive integer.

It's not enough to consider unrestricted families of circuits, however. In practice, we need an algorithm to build the circuit. Indeed, if we don't place any restrictions on the circuit family then it becomes possible to compute all sorts of functions which we do not expect to be able to compute in a reasonable model of computation. For example, let $h_n(x)$ denote the halting function, restricted to values of x which are n bits in length. Thus h_n is a function from n bits to 1 bit, and we have proved there exists a circuit C_n to compute $h_n(\cdot)$. Therefore the circuit family $\{C_n\}$ computes the halting function! However, what prevents us from using this circuit family to solve the halting problem is that we haven't specified an algorithm which will allow us to build the circuit C_n for all values of n . Adding this requirement results in the notion of a uniform circuit family.

That is, a family of circuits $\{C_n\}$ is said to be a *uniform circuit family* if there is some algorithm running on a Turing machine which, upon input of n , generates a *description* of C_n . That is, the algorithm outputs a description of what gates are in the circuit C_n , how those gates are connected together to form a circuit, any ancilla bits needed by the circuit, FANOUT and CROSSOVER operations, and where the output from the circuit should be read out. For example, the family of circuits we described earlier for squaring n -bit numbers is certainly a uniform circuit family, since there is an algorithm which, given n , outputs a description of the circuit needed to square an n -bit number. You can think of this algorithm as the means by which an engineer is able to generate a description of (and thus build) the circuit for any n whatsoever. By contrast, a circuit family that is not uniform is said to be a *non-uniform* circuit family. There is no algorithm to construct

the circuit for arbitrary n , which prevents our engineer from building circuits to compute functions like the halting function.

Intuitively, a uniform circuit family is a family of circuits that can be generated by some reasonable algorithm. It can be shown that the class of functions computable by uniform circuit families is exactly the same as the class of functions which can be computed on a Turing machine. With this uniformity restriction, results in the Turing machine model of computation can usually be given a straightforward translation into the circuit model of computation, and vice versa. Later we give similar attention to issues of uniformity in the quantum circuit model of computation.

3.2 The analysis of computational problems

The analysis of computational problems depends upon the answer to three fundamental questions:

- (1) **What is a computational problem?** Multiplying two numbers together is a computational problem; so is programming a computer to exceed human abilities in the writing of poetry. In order to make progress developing a general theory for the analysis of computational problems we are going to isolate a special class of problems known as *decision problems*, and concentrate our analysis on those. Restricting ourselves in this way enables the development of a theory which is both elegant and rich in structure. Most important, it is a theory whose principles have application far beyond decision problems.
- (2) **How may we design algorithms to solve a given computational problem?** Once a problem has been specified, what algorithms can be used to solve the problem? Are there general techniques which can be used to solve wide classes of problems? How can we be sure an algorithm behaves as claimed?
- (3) **What are the minimal resources required to solve a given computational problem?** Running an algorithm requires the consumption of *resources*, such as time, space, and energy. In different situations it may be desirable to minimize consumption of one or more resource. Can we classify problems according to the resource requirements needed to solve them?

In the next few sections we investigate these three questions, especially questions 1 and 3. Although question 1, ‘what is a computational problem?’, is perhaps the most fundamental of the questions, we shall defer answering it until Section 3.2.3, pausing first to establish some background notions related to resource quantification in Section 3.2.1, and then reviewing the key ideas of *computational complexity* in Section 3.2.2.

Question 2, how to design good algorithms, is the subject of an enormous amount of ingenious work by many researchers. So much so that in this brief introduction we cannot even begin to describe the main ideas employed in the design of good algorithms. If you are interested in this beautiful subject, we refer you to the end of chapter ‘History and further reading’. Our closest direct contact with this subject will occur later in the book, when we study quantum algorithms. The techniques involved in the creation of quantum algorithms have typically involved a blend of deep existing ideas in algorithm design for classical computers, and the creation of new, wholly quantum mechanical techniques for algorithm design. For this reason, and because the spirit of quantum algorithm design

is so similar in many ways to classical algorithm design, we encourage you to become familiar with at least the basic ideas of algorithm design.

Question 3, what are the minimal resources required to solve a given computational problem, is the main focus of the next few sections. For example, suppose we are given two numbers, each n bits in length, which we wish to multiply. If the multiplication is performed on a single-tape Turing machine, how many computational steps must be executed by the Turing machine in order to complete the task? How much space is used on the Turing machine while completing the task?

These are examples of the type of resource questions we may ask. Generally speaking, computers make use of many different kinds of resources, however we will focus most of our attention on time, space, and energy. Traditionally in computer science, time and space have been the two major resource concerns in the study of algorithms, and we study these issues in Sections 3.2.2 through 3.2.4. Energy has been a less important consideration; however, the study of energy requirements motivates the subject of reversible classical computation, which in turn is a prerequisite for quantum computation, so we examine energy requirements for computation in some considerable detail in Section 3.2.5.

3.2.1 How to quantify computational resources

Different models of computation lead to different resource requirements for computation. Even something as simple as changing from a single-tape to a two-tape Turing machine may change the resources required to solve a given computational problem. For a computational task which is extremely well understood, like addition of integers, for example, such differences between computational models may be of interest. However, for a first pass at understanding a problem, we would like a way of quantifying resource requirements that is independent of relatively trivial changes in the computational model. One of the tools which has been developed to do this is the *asymptotic notation*, which can be used to summarize the *essential* behavior of a function. This asymptotic notation can be used, for example, to summarize the essence of how many time steps it takes a given algorithm to run, without worrying too much about the exact time count. In this section we describe this notation in detail, and apply it to a simple problem illustrating the quantification of computational resources – the analysis of algorithms for sorting a list of names into alphabetical order.

Suppose, for example, that we are interested in the number of gates necessary to add together two n -bit numbers. Exact counts of the number of gates required obscure the big picture: perhaps a specific algorithm requires $24n + 2\lceil \log n \rceil + 16$ gates to perform this task. However, in the limit of large problem size the only term which matters is the $24n$ term. Furthermore, we disregard constant factors as being of secondary importance to the analysis of the algorithm. The *essential* behavior of the algorithm is summed up by saying that the number of operations required scales like n , where n is the number of bits in the numbers being added. The asymptotic notation consists of three tools which make this notion precise.

The O ('big O ') notation is used to set *upper bounds* on the behavior of a function. Suppose $f(n)$ and $g(n)$ are two functions on the non-negative integers. We say ' $f(n)$ is in the class of functions $O(g(n))$ ', or just ' $f(n)$ is $O(g(n))$ ', if there are constants c and n_0 such that for all values of n greater than n_0 , $f(n) \leq cg(n)$. That is, for sufficiently large n , the function $g(n)$ is an upper bound on $f(n)$, up to an unimportant constant

factor. The big O notation is particularly useful for studying the worst-case behavior of *specific* algorithms, where we are often satisfied with an upper bound on the resources consumed by an algorithm.

When studying the behaviors of a *class* of algorithms – say the entire class of algorithms which can be used to multiply two numbers – it is interesting to set lower bounds on the resources required. For this the Ω ('big Omega') notation is used. A function $f(n)$ is said to be $\Omega(g(n))$ if there exist constants c and n_0 such that for all n greater than n_0 , $cg(n) \leq f(n)$. That is, for sufficiently large n , $g(n)$ is a lower bound on $f(n)$, up to an unimportant constant factor.

Finally, the Θ ('big Theta') notation is used to indicate that $f(n)$ behaves the same as $g(n)$ asymptotically, up to unimportant constant factors. That is, we say $f(n)$ is $\Theta(g(n))$ if it is both $O(g(n))$ and $\Omega(g(n))$.

Asymptotic notation: examples

Let's consider a few simple examples of the asymptotic notation. The function $2n$ is in the class $O(n^2)$, since $2n \leq 2n^2$ for all positive n . The function 2^n is $\Omega(n^3)$, since $n^3 \leq 2^n$ for sufficiently large n . Finally, the function $7n^2 + \sqrt{n} \log(n)$ is $\Theta(n^2)$, since $7n^2 \leq 7n^2 + \sqrt{n} \log(n) \leq 8n^2$ for all sufficiently large values of n . In the following few exercises you will work through some of the elementary properties of the asymptotic notation that make it a useful tool in the analysis of algorithms.

Exercise 3.9: Prove that $f(n)$ is $O(g(n))$ if and only if $g(n)$ is $\Omega(f(n))$. Deduce that $f(n)$ is $\Theta(g(n))$ if and only if $g(n)$ is $\Theta(f(n))$.

Exercise 3.10: Suppose $g(n)$ is a polynomial of degree k . Show that $g(n)$ is $O(n^l)$ for any $l \geq k$.

Exercise 3.11: Show that $\log n$ is $O(n^k)$ for any $k > 0$.

Exercise 3.12: ($n^{\log n}$ is super-polynomial) Show that n^k is $O(n^{\log n})$ for any k , but that $n^{\log n}$ is never $O(n^k)$.

Exercise 3.13: ($n^{\log n}$ is sub-exponential) Show that c^n is $\Omega(n^{\log n})$ for any $c > 1$, but that $n^{\log n}$ is never $\Omega(c^n)$.

Exercise 3.14: Suppose $e(n)$ is $O(f(n))$ and $g(n)$ is $O(h(n))$. Show that $e(n)g(n)$ is $O(f(n)h(n))$.

An example of the use of the asymptotic notation in quantifying resources is the following simple application to the problem of sorting an n element list of names into alphabetical order. Many sorting algorithms are based upon the 'compare-and-swap' operation: two elements of an n element list are compared, and swapped if they are in the wrong order. If this compare-and-swap operation is the only means by which we can access the list, how many such operations are required in order to ensure that the list has been correctly sorted?

A simple compare-and-swap algorithm for solving the sorting problem is as follows: (note that `compare-and-swap(j, k)` compares the list entries numbered j and k , and swaps them if they are out of order)

```

for j = 1 to n-1
  for k = j+1 to n
    compare-and-swap(j,k)
  end k
end j

```

It is clear that this algorithm correctly sorts a list of n names into alphabetical order. Note that the number of compare-and-swap operations executed by the algorithm is $(n-1) + (n-2) + \dots + 1 = n(n-1)/2$. Thus the number of compare-and-swap operations used by the algorithm is $\Theta(n^2)$. Can we do better than this? It turns out that we can. Algorithms such as ‘heapsort’ are known which run using $O(n \log n)$ compare-and-swap operations. Furthermore, in Exercise 3.15 you’ll work through a simple counting argument that shows any algorithm based upon the compare-and-swap operation requires $\Omega(n \log n)$ such operations. Thus, the sorting problem requires $\Theta(n \log n)$ compare-and-swap operations, in general.

Exercise 3.15: (Lower bound for compare-and-swap based sorts) Suppose an n element list is sorted by applying some sequence of compare-and-swap operations to the list. There are $n!$ possible initial orderings of the list. Show that after k of the compare-and-swap operations have been applied, at most 2^k of the possible initial orderings will have been sorted into the correct order. Conclude that $\Omega(n \log n)$ compare-and-swap operations are required to sort all possible initial orderings into the correct order.

3.2.2 Computational complexity

The idea that there won’t be an algorithm to solve it – this is something fundamental that won’t ever change – that idea appeals to me.

– Stephen Cook

Sometimes it is good that some things are impossible. I am happy there are many things that nobody can do to me.

– Leonid Levin

*It should not come as a surprise that our choice of polynomial algorithms as the mathematical concept that is supposed to capture the informal notion of ‘practically efficient computation’ is open to criticism from all sides. [...] Ultimately, our argument for our choice must be this: **Adopting polynomial worst-case performance as our criterion of efficiency results in an elegant and useful theory that says something meaningful about practical computation, and would be impossible without this simplification.***

– Christos Papadimitriou

What time and space resources are required to perform a computation? In many cases these are the most important questions we can ask about a computational problem. Problems like addition and multiplication of numbers are regarded as efficiently solvable because we have *fast* algorithms to perform addition and multiplication, which consume

little *space* when running. Many other problems have no known fast algorithm, and are effectively impossible to solve, not because we can't find an algorithm to solve the problem, but because all known algorithms consume such vast quantities of space or time as to render them practically useless.

Computational complexity is the study of the time and space resources required to solve computational problems. The task of computational complexity is to prove *lower bounds* on the resources required by the best possible algorithm for solving a problem, even if that algorithm is not explicitly known. In this and the next two sections, we give an overview of computational complexity, its major concepts, and some of the more important results of the field. Note that computational complexity is in a sense complementary to the field of algorithm design; ideally, the most efficient algorithms we could design would match perfectly with the lower bounds proved by computational complexity. Unfortunately, this is often not the case. As already noted, in this book we won't examine classical algorithm design in any depth.

One difficulty in formulating a theory of computational complexity is that different computational models may require different resources to solve the same problem. For instance, multiple-tape Turing machines can solve many problems substantially faster than single-tape Turing machines. This difficulty is resolved in a rather coarse way. Suppose a problem is specified by giving n bits as input. For instance, we might be interested in whether a particular n -bit number is prime or not. The chief distinction made in computational complexity is between problems which can be solved using resources which are bounded by a *polynomial* in n , or which require resources which grow faster than any polynomial in n . In the latter case we usually say that the resources required are *exponential* in the problem size, abusing the term exponential, since there are functions like $n^{\log n}$ which grow faster than any polynomial (and thus are 'exponential' according to this convention), yet which grow slower than any true exponential. A problem is regarded as *easy*, *tractable* or *feasible* if an algorithm for solving the problem using polynomial resources exists, and as *hard*, *intractable* or *infeasible* if the best possible algorithm requires exponential resources.

As a simple example, suppose we have two numbers with binary expansions $x_1 \dots x_{m_1}$ and $y_1 \dots y_{m_2}$, and we wish to determine the sum of the two numbers. The total size of the input is $n \equiv m_1 + m_2$. It's easy to see that the two numbers can be added using a number of elementary operations that scales as $\Theta(n)$; this algorithm uses a polynomial (indeed, linear) number of operations to perform its tasks. By contrast, it is believed (though it has never been proved!) that the problem of factoring an integer into its prime factors is intractable. That is, the belief is that there is no algorithm which can factor an arbitrary n -bit integer using $O(p(n))$ operations, where p is some fixed polynomial function of n . We will later give many other examples of problems which are believed to be intractable in this sense.

The polynomial versus exponential classification is rather coarse. In practice, an algorithm that solves a problem using $2^{n/1000}$ operations is probably more useful than one which runs in n^{1000} operations. Only for very large input sizes ($n \approx 10^8$) will the 'efficient' polynomial algorithm be preferable to the 'inefficient' exponential algorithm, and for many purposes it may be more practical to prefer the 'inefficient' algorithm.

Nevertheless, there are many reasons to base computational complexity primarily on the polynomial versus exponential classification. First, historically, with few exceptions, polynomial resource algorithms have been much faster than exponential algorithms. We

might speculate that the reason for this is lack of imagination: coming up with algorithms requiring n , n^2 or some other low degree polynomial number of operations is often much easier than finding a natural algorithm which requires n^{1000} operations, although examples like the latter do exist. Thus, the predisposition for the human mind to come up with relatively simple algorithms has meant that in practice polynomial algorithms usually do perform much more efficiently than their exponential cousins.

A second and more fundamental reason for emphasizing the polynomial versus exponential classification is derived from the *strong Church–Turing thesis*. As discussed in Section 1.1, it was observed in the 1960s and 1970s that probabilistic Turing machines appear to be the strongest ‘reasonable’ model of computation. More precisely, researchers consistently found that if it was possible to compute a function using k elementary operations in some model that was *not* the probabilistic Turing machine model of computation, then it was always possible to compute the same function in the probabilistic Turing machine model, using at most $p(k)$ elementary operations, where $p(\cdot)$ is some *polynomial* function. This statement is known as the *strong Church–Turing thesis*:

Strong Church–Turing thesis: *Any model of computation can be simulated on a probabilistic Turing machine with at most a polynomial increase in the number of elementary operations required.*

The strong Church–Turing thesis is great news for the theory of computational complexity, for it implies that attention may be restricted to the probabilistic Turing machine model of computation. After all, if a problem has no polynomial resource solution on a probabilistic Turing machine, then the strong Church–Turing thesis implies that it has no efficient solution on any computing device. Thus, the strong Church–Turing thesis implies that the entire theory of computational complexity will take on an elegant, model-independent form if the notion of efficiency is identified with polynomial resource algorithms, and this elegance has provided a strong impetus towards acceptance of the identification of ‘solvable with polynomial resources’ and ‘efficiently solvable’. Of course, one of the prime reasons for interest in quantum computers is that they cast into doubt the strong Church–Turing thesis, by enabling the efficient solution of a problem which is believed to be intractable on all classical computers, including probabilistic Turing machines! Nevertheless, it is useful to understand and appreciate the role the strong Church–Turing thesis has played in the search for a model-independent theory of computational complexity.

Finally, we note that, in practice, computer scientists are not only interested in the polynomial versus exponential classification of problems. This is merely the first and coarsest way of understanding how difficult a computational problem is. However, it is an exceptionally important distinction, and illustrates many broader points about the nature of resource questions in computer science. For most of this book, it will be our central concern in evaluating the efficiency of a given algorithm.

Having examined the merits of the polynomial versus exponential classification, we now have to confess that the theory of computational complexity has one remarkable outstanding failure: it seems very hard to prove that there are interesting classes of problems which require exponential resources to solve. It is quite easy to give non-constructive proofs that *most* problems require exponential resources (see Exercise 3.16, below), and furthermore many interesting problems are *conjectured* to require exponential resources for their solution, but rigorous proofs seem very hard to come by, at least with the present

state of knowledge. This failure of computational complexity has important implications for quantum computation, because it turns out that the computational power of quantum computers can be related to some major open problems in *classical* computational complexity theory. Until these problems are resolved, it cannot be stated with certainty how computationally powerful a quantum computer is, or even whether it is more powerful than a classical computer!

Exercise 3.16: (Hard-to-compute functions exist) Show that there exist Boolean functions on n inputs which require at least $2^n / \log n$ logic gates to compute.

3.2.3 Decision problems and the complexity classes P and NP

Many computational problems are most cleanly formulated as *decision problems* – problems with a yes or no answer. For example, is a given number m a prime number or not? This is the *primality* decision problem. The main ideas of computational complexity are most easily and most often formulated in terms of decision problems, for two reasons: the theory takes its simplest and most elegant form in this form, while still generalizing in a natural way to more complex scenarios; and historically computational complexity arose primarily from the study of decision problems.

Although most decision problems can easily be stated in simple, familiar language, discussion of the general properties of decision problems is greatly helped by the terminology of *formal languages*. In this terminology, a *language* L over the *alphabet* Σ is a subset of the set Σ^* of all (finite) strings of symbols from Σ . For example, if $\Sigma = \{0, 1\}$, then the set of binary representations of even numbers, $L = \{0, 10, 100, 110, \dots\}$ is a language over Σ .

Decision problems may be encoded in an obvious way as problems about languages. For instance, the primality decision problem can be encoded using the binary alphabet $\Sigma = \{0, 1\}$. Strings from Σ^* can be interpreted in a natural way as non-negative integers. For example, 0010 can be interpreted as the number 2. The language L is defined to consist of all binary strings such that the corresponding number is prime.

To solve the primality decision problem, what we would like is a Turing machine which, when started with a given number n on its input tape, eventually outputs some equivalent of ‘yes’ if n is prime, and outputs ‘no’ if n is not prime. To make this idea precise, it is convenient to modify our old Turing machine definition (of Section 3.1.1) slightly, replacing the halting state q_h with two states q_Y and q_N to represent the answers ‘yes’ and ‘no’ respectively. In all other ways the machine behaves in the same way, and it still halts when it enters the state q_Y or q_N . More generally, a language L is *decided* by a Turing machine if the machine is able to decide whether an input x on its tape is a member of the language of L or not, eventually halting in the state q_Y if $x \in L$, and eventually halting in the state q_N if $x \notin L$. We say that the machine has *accepted* or *rejected* x depending on which of these two cases comes about.

How quickly can we determine whether or not a number is prime? That is, what is the fastest Turing machine which decides the language representing the primality decision problem? We say that a problem is in **TIME**($f(n)$) if there exists a Turing machine which decides whether a candidate x is in the language in time $O(f(n))$, where n is the length of x . A problem is said to be solvable in *polynomial time* if it is in **TIME**(n^k) for some finite k . The collection of all languages which are in **TIME**(n^k), for some k , is denoted **P**. **P** is our first example of a *complexity class*. More generally, a complexity

class is defined to be a collection of languages. Much of computational complexity theory is concerned with the definition of various complexity classes, and understanding the relationship between different complexity classes.

Not surprisingly, there are problems which cannot be solved in polynomial time. Unfortunately, proving that any given problem can't be solved in polynomial time seems to be very difficult, although conjectures abound! A simple example of an interesting decision problem which is believed not to be in \mathbf{P} is the *factoring decision problem*:

FACToring: Given a composite integer m and $l < m$, does m have a non-trivial factor less than l ?

An interesting property of factoring is that if somebody claims that the answer is 'yes, m does have a non-trivial factor less than l ' then they can establish this by exhibiting such a factor, which can then be efficiently checked by other parties, simply by doing long-division. We call such a factor a *witness* to the fact that m has a factor less than l . This idea of an easily checkable witness is the key idea in the definition of the complexity class \mathbf{NP} , below. We have phrased factoring as a decision problem, but you can easily verify that the decision problem is equivalent to finding the factors of a number:

Exercise 3.17: Prove that a polynomial-time algorithm for finding the factors of a number m exists if and only if the factoring decision problem is in \mathbf{P} .

Factoring is an example of a problem in an important complexity class known as \mathbf{NP} . What distinguishes problems in \mathbf{NP} is that 'yes' instances of a problem can easily be verified with the aid of an appropriate witness. More rigorously, a language L is in \mathbf{NP} if there is a Turing machine M with the following properties:

- (1) If $x \in L$ then there exists a witness string w such that M halts in the state q_Y after a time polynomial in $|x|$ when the machine is started in the state x -blank- w .
- (2) If $x \notin L$ then for all strings w which attempt to play the role of a witness, the machine halts in state q_N after a time polynomial in $|x|$ when M is started in the state x -blank- w .

There is an interesting asymmetry in the definition of \mathbf{NP} . While we have to be able to quickly decide whether a possible witness to $x \in L$ is truly a witness, there is no such need to produce a witness to $x \notin L$. For instance, in the factoring problem, we have an easy way of proving that a given number has a factor less than m , but exhibiting a witness to prove that a number has no factors less than m is more daunting. This suggests defining \mathbf{coNP} , the class of languages which have witnesses to 'no' instances; obviously the languages in \mathbf{coNP} are just the complements of languages in \mathbf{NP} .

How are \mathbf{P} and \mathbf{NP} related? It is clear that \mathbf{P} is a subset of \mathbf{NP} . The most famous open problem in computer science is *whether or not there are problems in \mathbf{NP} which are not in \mathbf{P}* , often abbreviated as the $\mathbf{P} \neq \mathbf{NP}$ problem. Most computer scientists believe that $\mathbf{P} \neq \mathbf{NP}$, but despite decades of work nobody has been able to prove this, and the possibility remains that $\mathbf{P} = \mathbf{NP}$.

Exercise 3.18: Prove that if $\mathbf{coNP} \neq \mathbf{NP}$ then $\mathbf{P} \neq \mathbf{NP}$.

Upon first acquaintance it's tempting to conclude that the conjecture $\mathbf{P} \neq \mathbf{NP}$ ought to be pretty easy to resolve. To see why it's actually rather subtle it helps to see couple of

examples of problems that are in **P** and **NP**. We'll draw the examples from *graph theory*, a rich source of decision problems with surprisingly many practical applications. A *graph* is a finite collection of *vertices* $\{v_1, \dots, v_n\}$ connected by *edges*, which are pairs (v_i, v_j) of vertices. For now, we are only concerned with *undirected graphs*, in which the order of the vertices (in each edge pair) does not matter; similar ideas can be investigated for *directed graphs* in which the order of vertices does matter. A typical graph is illustrated in Figure 3.9.

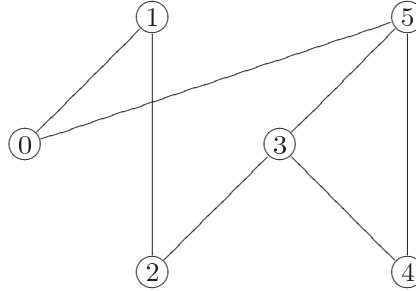


Figure 3.9. A graph.

A *cycle* in a graph is a sequence v_1, \dots, v_m of vertices such that each pair (v_j, v_{j+1}) is an edge, as is (v_1, v_m) . A *simple cycle* is a cycle in which none of the vertices is repeated, except for the first and last vertices. A *Hamiltonian cycle* is a simple cycle which visits every vertex in the graph. Examples of graphs with and without Hamiltonian cycles are shown in Figure 3.10.

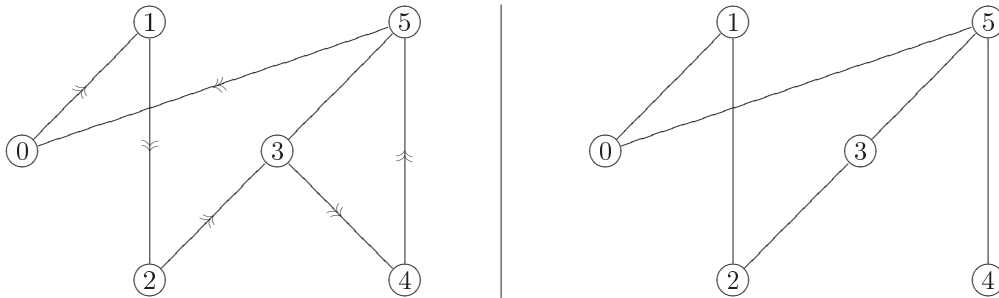


Figure 3.10. The graph on the left contains a Hamiltonian cycle, 0, 1, 2, 3, 4, 5, 0. The graph on the right contains no Hamiltonian cycle, as can be verified by inspection.

The *Hamiltonian cycle problem* (HC) is to determine whether a given graph contains a Hamiltonian cycle or not. HC is a decision problem in **NP**, since if a given graph has a Hamiltonian cycle, then that cycle can be used as an easily checkable witness. Moreover, HC has no known polynomial time algorithm. Indeed, HC is a problem in the class of so-called **NP-complete** problems, which can be thought of as the ‘hardest’ problems in **NP**, in the sense that solving HC in time t allows any other problem in **NP** to be solved in time $O(\text{poly}(t))$. This also means that if any **NP-complete** problem has a polynomial time solution then it will follow that **P** = **NP**.

There is a problem, the Euler cycle decision problem, which is superficially similar to HC, but which has astonishingly different properties. An *Euler cycle* is an ordering of the edges of a graph G so that every edge in the graph is visited exactly once. The Euler

cycle decision problem (EC) is to determine, given a graph G on n vertices, whether that graph contains an Euler cycle or not. EC is, in fact, exactly the same problem as HC, only the path visits edges, rather than vertices. Consider the following remarkable theorem, to be proven in Exercise 3.20:

Theorem 3.1: (Euler's theorem) A connected graph contains an Euler cycle if and only if every vertex has an even number of edges incident upon it.

Euler's theorem gives us a method for efficiently solving EC. First, check to see whether the graph is connected; this is easily done with $O(n^2)$ operations, as shown in Exercise 3.19. If the graph is not connected, then obviously no Euler cycle exists. If the graph is connected then for each vertex check whether there is an even number of edges incident upon the vertex. If a vertex is found for which this is not the case, then there is no Euler cycle, otherwise an Euler cycle exists. Since there are n vertices, and at most $n(n-1)/2$ edges, this algorithm requires $O(n^3)$ elementary operations. Thus EC is in P! Somehow, there is a structure present in the problem of visiting each edge that can be exploited to provide an efficient algorithm for EC, yet which does not seem to be reflected in the problem of visiting each vertex; it is not at all obvious why such a structure should be present in one case, but not in the other, if indeed it is absent for the HC problem.

Exercise 3.19: The REACHABILITY problem is to determine whether there is a path between two specified vertices in a graph. Show that REACHABILITY can be solved using $O(n)$ operations if the graph has n vertices. Use the solution to REACHABILITY to show that it is possible to decide whether a graph is connected in $O(n^2)$ operations.

Exercise 3.20: (Euler's theorem) Prove Euler's theorem. In particular, if each vertex has an even number of incident edges, give a constructive procedure for finding an Euler cycle.

The equivalence between the factoring decision problem and the factoring problem proper is a special instance of one of the most important ideas in computer science, an idea known as *reduction*. Intuitively, we know that some problems can be viewed as special instances of other problems. A less trivial example of reduction is the reduction of HC to the *traveling salesman* decision problem (TSP). The traveling salesman decision problem is as follows: we are given n cities $1, 2, \dots, n$ and a non-negative integer distance d_{ij} between each pair of cities. Given a distance d the problem is to determine if there is a tour of all the cities of distance less than d .

The reduction of HC to TSP goes as follows. Suppose we have a graph containing n vertices. We turn this into an instance of TSP by thinking of each vertex of the graph as a 'city' and defining the distance d_{ij} between cities i and j to be one if vertices i and j are connected, and the distance to be two if the vertices are unconnected. Then a tour of the cities of distance less than $n+1$ must be of distance n , and be a Hamiltonian cycle for the graph. Conversely, if a Hamiltonian cycle exists then a tour of the cities of distance less than $n+1$ must exist. In this way, given an algorithm for solving TSP, we can convert it into an algorithm for solving HC without much overhead. Two consequences can be inferred from this. First, if TSP is a tractable problem, then HC is also tractable. Second, if HC is hard then TSP must also be hard. This is an example of a general technique known

as *reduction*: we've reduced the problem HC to the problem TSP. This is a technique we will use repeatedly throughout this book.

A more general notion of reduction is illustrated in Figure 3.11. A language B is said to be *reducible* to another language A if there exists a Turing machine operating in polynomial time such that given as input x it outputs $R(x)$, and $x \in B$ if and only if $R(x) \in A$. Thus, if we have an algorithm for deciding A , then with a little extra overhead we can decide the language B . In this sense, the language B is essentially no more difficult to decide than the language A .

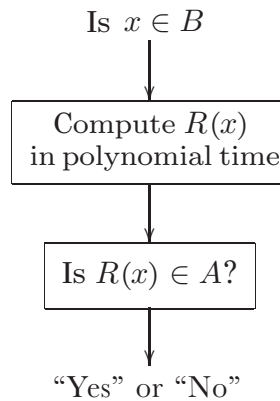


Figure 3.11. Reduction of B to A .

Exercise 3.21: (Transitive property of reduction) Show that if a language L_1 is reducible to the language L_2 and the language L_2 is reducible to L_3 then the language L_1 is reducible to the language L_3 .

Some complexity classes have problems which are *complete* with respect to that complexity class, meaning there is a language L in the complexity class which is the ‘most difficult’ to decide, in the sense that every other language in the complexity class can be reduced to L . Not all complexity classes have complete problems, but many of the complexity classes we are concerned with do have complete problems. A trivial example is provided by **P**. Let L be any language in **P** which is not empty or equal to the set of all words. That is, there exists a string x_1 such that $x_1 \notin L$ and a string x_2 such that $x_2 \in L$. Then any other language L' in **P** can be reduced to L using the following reduction: given an input x , use the polynomial time decision procedure to determine whether $x \in L'$ or not. If it is not, then set $R(x) = x_1$, otherwise set $R(x) = x_2$.

Exercise 3.22: Suppose L is complete for a complexity class, and L' is another language in the complexity class such that L reduces to L' . Show that L' is complete for the complexity class.

Less trivially, **NP** also contains complete problems. An important example of such a problem and the prototype for all other **NP**-complete problems is the *circuit satisfiability* problem or CSAT: given a Boolean circuit composed of AND, OR and NOT gates, is there an assignment of values to the inputs to the circuit that results in the circuit outputting 1, that is, is the circuit *satisfiable* for some input? The **NP**-completeness of CSAT is known as the *Cook–Levin theorem*, for which we now outline a proof.

Theorem 3.2: (Cook–Levin) CSAT is NP-complete.

Proof

The proof has two parts. The first part of the proof is to show that CSAT is in NP, and the second part is to show that any language in NP can be reduced to CSAT. Both parts of the proof are based on *simulation* techniques: the first part of the proof is essentially showing that a Turing machine can efficiently simulate a circuit, while the second part of the proof is essentially showing that a circuit can efficiently simulate a Turing machine. Both parts of the proof are quite straightforward; for the purposes of illustration we give the second part in some detail.

The first part of the proof is to show that CSAT is in NP. Given a circuit containing n circuit elements, and a potential witness w , it is obviously easy to check in polynomial time on a Turing machine whether or not w satisfies the circuit, which establishes that CSAT is in NP.

The second part of the proof is to show that any language $L \in \text{NP}$ can be reduced to CSAT. That is, we aim to show that there is a polynomial time computable reduction R such that $x \in L$ if and only if $R(x)$ is a satisfiable circuit. The idea of the reduction is to find a circuit which simulates the action of the machine M which is used to check instance-witness pairs, (x, w) , for the language L . The input variables for the circuit will represent the witness; the idea is that finding a witness which satisfies the circuit is equivalent to M accepting (x, w) for some specific witness w . Without loss of generality we may make the following assumptions about M to simplify the construction:

- (1) M 's tape alphabet is $\triangleright, 0, 1$ and the blank symbol.
- (2) M runs using time at most $t(n)$ and total space at most $s(n)$ where $t(n)$ and $s(n)$ are polynomials in n .
- (3) Machine M can actually be assumed to run using time *exactly* $t(n)$ for all inputs of size n . This is done by adding the lines $\langle q_Y, x, q_Y, x, 0 \rangle$, and $\langle q_N, x, q_N, x, 0 \rangle$ for each of $x = \triangleright, 0, 1$ and the blank, artificially halting the machine after exactly $t(n)$ steps.

The basic idea of the construction to simulate M is outlined in Figure 3.12. Each internal state of the Turing machine is represented by a single bit in the circuit. We name the corresponding bits $\tilde{q}_s, \tilde{q}_1, \dots, \tilde{q}_m, \tilde{q}_Y, \tilde{q}_N$. Initially, \tilde{q}_s is set to one, and all the other bits representing internal states are set to zero. Each square on the Turing machine tape is represented by three bits: two bits to represent the letter of the alphabet ($\triangleright, 0, 1$ or blank) currently residing on the tape, and a single ‘flag’ bit which is set to one if the read-write head is pointing to the square, and set to zero otherwise. We denote the bits representing the tape contents by $(u_1, v_1), \dots, (u_{s(n)}, v_{s(n)})$ and the corresponding flag bits by $f_1, \dots, f_{s(n)}$. Initially the u_j and v_j bits are set to represent the inputs x and w , as appropriate, while $f_1 = 1$ and all other $f_j = 0$. There is also a lone extra ‘global flag’ bit, F , whose function will be explained later. F is initially set to zero. We regard all the bits input to the circuit as fixed, except for those representing the witness w , which are the variable bits for the circuit.

The action of M is obtained by repeating $t(n)$ times a ‘simulation step’ which simulates the execution of a single program line for the Turing machine. Each simulation step may be broken up into a sequence of steps corresponding in turn to the respective program lines, with a final step which resets the global flag F to zero, as

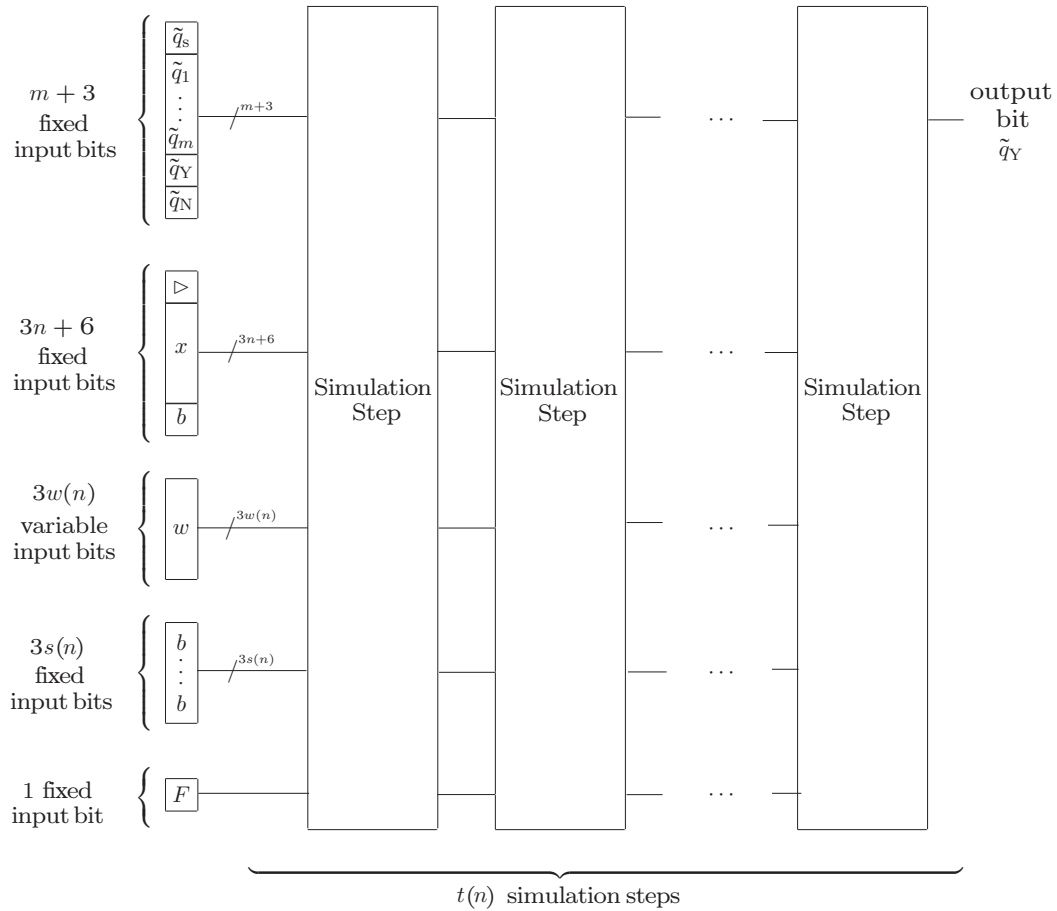


Figure 3.12. Outline of the procedure used to simulate a Turing machine using a circuit.

illustrated in Figure 3.13. To complete the simulation, we only need to simulate a program line of the form $\langle q_i, x, q_j, x', s \rangle$. For convenience, we assume $q_i \neq q_j$, but a similar construction works in the case when $q_i = q_j$. The procedure is as follows:

- (1) Check to see that $\tilde{q}_i = 1$, indicating that the current state of the machine is q_i .
- (2) For each tape square:
 - (a) Check to see that the global flag bit is set to zero, indicating that no action has yet been taken by the Turing machine.
 - (b) Check that the flag bit is set to one, indicating that the tape head is at this tape square.
 - (c) Check that the simulated tape contents at this point are x .
 - (d) If all conditions check out, then perform the following steps:
 1. Set $\tilde{q}_i = 0$ and $\tilde{q}_j = 1$.
 2. Update the simulated tape contents at this tape square to x' .
 3. Update the flag bit of this and adjacent 'squares' as appropriate, depending on whether $s = +1, 0, -1$, and whether we are at the left hand end of the tape.
 4. Set the global flag bit to one, indicating that this round of computation has been completed.

This is a fixed procedure which involves a constant number of bits, and by the universality result of Section 3.1.2 can be performed using a circuit containing a constant number of gates.

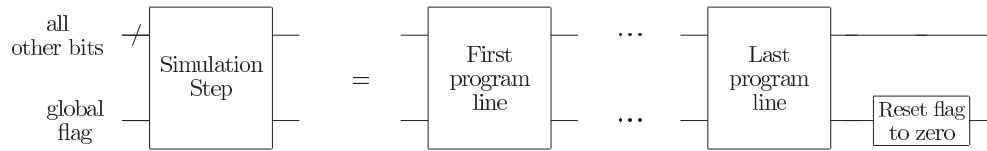


Figure 3.13. Outline of the simulation step used to simulate a Turing machine using a circuit.

The total number of gates in the entire circuit is easily seen to be $O(t(n)(s(n) + n))$, which is polynomial in size. At the end of the circuit, it is clear that $\tilde{q}_Y = 1$ if and only if the machine M accepts (x, w) . Thus, the circuit is satisfiable if and only if there exists w such that machine M accepts (x, w) , and we have found the desired reduction from L to CSAT. \square

CSAT gives us a foot in the door which enables us to easily prove that many other problems are **NP**-complete. Instead of directly proving that a problem is **NP**-complete, we can instead prove that it is in **NP** and that CSAT reduces to it, so by Exercise 3.22 the problem must be **NP**-complete. A small sample of **NP**-complete problems is discussed in Box 3.3. An example of another **NP**-complete problem is the *satisfiability* problem (SAT), which is phrased in terms of a Boolean formula. Recall that a Boolean formula φ is composed of the following elements: a set of Boolean variables, x_1, x_2, \dots ; Boolean connectives, that is, a Boolean function with one or two inputs and one output, such as \wedge (AND), \vee (OR), and \neg (NOT); and parentheses. The truth or falsity of a Boolean formula for a given set of Boolean variables is decided according to the usual laws of Boolean algebra. For example, the formula $\varphi = x_1 \vee \neg x_2$ has the satisfying assignment $x_1 = 0$ and $x_2 = 0$, while $x_1 = 0$ and $x_2 = 1$ is not a satisfying assignment. The satisfiability problem is to determine, given a Boolean formula φ , whether or not it is satisfiable by any set of possible inputs.

Exercise 3.23: Show that SAT is **NP**-complete by first showing that SAT is in **NP**, and then showing that CSAT reduces to SAT. (*Hint:* for the reduction it may help to represent each distinct wire in an instance of CSAT by different variables in a Boolean formula.)

An important restricted case of SAT is also **NP**-complete, the 3-satisfiability problem (3SAT), which is concerned with formulae in *3-conjunctive normal form*. A formula is said to be in *conjunctive normal form* if it is the AND of a collection of *clauses*, each of which is the OR of one or more *literals*, where a literal is an expression of the form x or $\neg x$. For example, the formula $(x_1 \vee \neg x_2) \wedge (x_2 \vee x_3 \vee \neg x_4)$ is in conjunctive normal form. A formula is in *3-conjunctive normal form* or *3-CNF* if each clause has exactly three literals. For example, the formula $(\neg x_1 \vee x_2 \vee \neg x_2) \wedge (\neg x_1 \vee x_3 \vee \neg x_4) \wedge (x_2 \vee x_3 \vee x_4)$ is in 3-conjunctive normal form. The 3-satisfiability problem is to determine whether a formula in 3-conjunctive normal form is satisfiable or not.

The proof that 3SAT is **NP**-complete is straightforward, but is a little too lengthy to justify inclusion in this overview. Even more than CSAT and SAT, 3SAT is in some sense

the NP-complete problem, and it is the basis for countless proofs that other problems are NP-complete. We conclude our discussion of NP-completeness with the surprising fact that 2SAT, the analogue of 3SAT in which every clause has two literals, can be solved in polynomial time:

Exercise 3.24: (2SAT has an efficient solution) Suppose φ is a Boolean formula in conjunctive normal form, in which each clause contains only two literals.

- (1) Construct a (directed) graph $G(\varphi)$ with directed edges in the following way: the vertices of G correspond to variables x_j and their negations $\neg x_j$ in φ . There is a (directed) edge (α, β) in G if and only if the clause $(\neg\alpha \vee \beta)$ or the clause $(\beta \vee \neg\alpha)$ is present in φ . Show that φ is not satisfiable if and only if there exists a variable x such that there are paths from x to $\neg x$ and from $\neg x$ to x in $G(\varphi)$.
- (2) Show that given a directed graph G containing n vertices it is possible to determine whether two vertices v_1 and v_2 are connected in polynomial time.
- (3) Find an efficient algorithm to solve 2SAT.

Box 3.3: A zoo of NP-complete problems

The importance of the class NP derives, in part, from the enormous number of computational problems that are known to be NP-complete. We can't possibly hope to survey this topic here (see 'History and further reading'), but the following examples, taken from many distinct areas of mathematics, give an idea of the delicious melange of problems known to be NP-complete.

- **CLIQUE (graph theory):** A clique in an undirected graph G is a subset of vertices, each pair of which is connected by an edge. The *size* of a clique is the number of vertices it contains. Given an integer m and a graph G , does G have a clique of size m ?
- **SUBSET-SUM (arithmetic):** Given a finite collection S of positive integers and a target t , is there any subset of S which sums to t ?
- **0-1 INTEGER PROGRAMMING (linear programming):** Given an integer $m \times n$ matrix A and an m -dimensional vector y with integer values, does there exist an n -dimensional vector x with entries in the set $\{0, 1\}$ such that $Ax \leq y$?
- **VERTEX COVER (graph theory):** A *vertex cover* for an undirected graph G is a set of vertices V' such that every edge in the graph has one or both vertices contained in V' . Given an integer m and a graph G , does G have a vertex cover V' containing m vertices?

Assuming that $P \neq NP$ it is possible to prove that there is a *non-empty* class of problems NPI (NP intermediate) which are neither solvable with polynomial resources, nor are NP-complete. Obviously, there are no problems known to be in NPI (otherwise we would know that $P \neq NP$) but there are several problems which are regarded as being likely candidates. Two of the strongest candidates are the factoring and graph isomorphism problems:

GRAPH ISOMORPHISM: Suppose G and G' are two undirected graphs over the vertices $V \equiv \{v_1, \dots, v_n\}$. Are G and G' isomorphic? That is, does there exist a one-to-one function $\varphi : V \rightarrow V$ such that the edge (v_i, v_j) is contained in G if and only if $(\varphi(v_i), \varphi(v_j))$ is contained in G' ?

Problems in **NPI** are interesting to researchers in quantum computation and quantum information for two reasons. First, it is desirable to find fast quantum algorithms to solve problems which are not in **P**. Second, many suspect that quantum computers will not be able to efficiently solve all problems in **NP**, ruling out **NP**-complete problems. Thus, it is natural to focus on the class **NPI**. Indeed, a fast quantum algorithm for factoring has been discovered (Chapter 5), and this has motivated the search for fast quantum algorithms for other problems suspected to be in **NPI**.

3.2.4 A plethora of complexity classes

We have investigated some of the elementary properties of some important complexity classes. A veritable pantheon of complexity classes exists, and there are many non-trivial relationships known or suspected between these classes. For quantum computation and quantum information, it is not necessary to understand all the different complexity classes that have been defined. However, it is useful to have some appreciation for the more important of the complexity classes, many of which have natural analogues in the study of quantum computation and quantum information. Furthermore, if we are to understand how powerful quantum computers are, then it behooves us to understand how the class of problems solvable on a quantum computer fits into the zoo of complexity classes which may be defined for classical computers.

There are essentially three properties that may be varied in the definition of a complexity class: the resource of interest (time, space, ...), the type of problem being considered (decision problem, optimization problem, ...), and the underlying computational model (deterministic Turing machine, probabilistic Turing machine, quantum computer, ...). Not surprisingly, this gives us an enormous range to define complexity classes. In this section, we briefly review a few of the more important complexity classes and some of their elementary properties. We begin with a complexity class defined by changing the resource of interest from time to *space*.

The most natural space-bounded complexity class is the class **PSPACE** of decision problems which may be solved on a Turing machine using a polynomial number of working bits, with no limitation on the amount of time that may be used by the machine (see Exercise 3.25). Obviously, **P** is included in **PSPACE**, since a Turing machine that halts after polynomial time can only traverse polynomially many squares, but it is also true that **NP** is a subset of **PSPACE**. To see this, suppose L is any language in **NP**. Suppose problems of size n have witnesses of size at most $p(n)$, where $p(n)$ is some polynomial in n . To determine whether or not the problem has a solution, we may sequentially test all $2^{p(n)}$ possible witnesses. Each test can be run in polynomial time, and therefore polynomial space. If we erase all the intermediate working between tests then we can test all the possibilities using polynomial space.

Unfortunately, at present it is not even known whether **PSPACE** contains problems which are not in **P**! This is a pretty remarkable situation – it seems fairly obvious that having unlimited time and polynomial spatial resources must be more powerful than having only a polynomial amount of time. However, despite considerable effort and in-

genuity, this has never been shown. We will see later that the class of problems solvable on a quantum computer in polynomial time is a subset of **PSPACE**, so proving that a problem efficiently solvable on a quantum computer is not efficiently solvable on a classical computer would establish that $\mathbf{P} \neq \mathbf{PSPACE}$, and thus solve a major outstanding problem of computer science. An optimistic way of looking at this result is that ideas from quantum computation might be useful in proving that $\mathbf{P} \neq \mathbf{PSPACE}$. Pessimistically, one might conclude that it will be a long time before anyone rigorously proves that quantum computers can be used to efficiently solve problems that are intractable on a classical computer. Even more pessimistically, it is possible that $\mathbf{P} = \mathbf{PSPACE}$, in which case quantum computers offer no advantage over classical computers! However, very few (if any) computational complexity theorists believe that $\mathbf{P} = \mathbf{PSPACE}$.

Exercise 3.25: ($\mathbf{PSPACE} \subseteq \mathbf{EXP}$) The complexity class **EXP** (for *exponential time*) contains all decision problems which may be decided by a Turing machine running in exponential time, that is time $O(2^{n^k})$, where k is any constant. Prove that $\mathbf{PSPACE} \subseteq \mathbf{EXP}$. (*Hint: If a Turing machine has l internal states, an m letter alphabet, and uses space $p(n)$, argue that the machine can exist in one of at most $lm^{p(n)}$ different states, and that if the Turing machine is to avoid infinite loops then it must halt before revisiting a state.*)

Exercise 3.26: ($\mathbf{L} \subseteq \mathbf{P}$) The complexity class **L** (for *logarithmic space*) contains all decision problems which may be decided by a Turing machine running in logarithmic space, that is, in space $O(\log(n))$. More precisely, the class **L** is defined using a two-tape Turing machine. The first tape contains the problem instance, of size n , and is a read-only tape, in the sense that only program lines which don't change the contents of the first tape are allowed. The second tape is a working tape which initially contains only blanks. The logarithmic space requirement is imposed on the second, working tape only. Show that $\mathbf{L} \subseteq \mathbf{P}$.

Does allowing more time or space give greater computational power? The answer to this question is yes in both cases. Roughly speaking, the *time hierarchy theorem* states that $\mathbf{TIME}(f(n))$ is a proper subset of $\mathbf{TIME}(f(n) \log^2(f(n)))$. Similarly, the *space hierarchy theorem* states that $\mathbf{SPACE}(f(n))$ is a proper subset of $\mathbf{SPACE}(f(n) \log(f(n)))$, where $\mathbf{SPACE}(f(n))$ is, of course, the complexity class consisting of all languages that can be decided with spatial resources $O(f(n))$. The hierarchy theorems have interesting implications with respect to the equality of complexity classes. We know that

$$\mathbf{L} \subseteq \mathbf{P} \subseteq \mathbf{NP} \subseteq \mathbf{PSPACE} \subseteq \mathbf{EXP}. \quad (3.1)$$

Unfortunately, although each of these inclusions is widely believed to be strict, none of them has ever been proved to be strict. However, the time hierarchy theorem implies that **P** is a strict subset of **EXP**, and the space hierarchy theorem implies that **L** is a strict subset of **PSPACE**! So we can conclude that at least one of the inclusions in (3.1) must be strict, although we do not know which one.

What should we do with a problem once we know that it is **NP**-complete, or that some other hardness criterion holds? It turns out that this is far from being the end of the story in problem analysis. One possible line of attack is to identify special cases of the problem which may be amenable to attack. For example, in Exercise 3.24 we saw that the 2SAT problem has an efficient solution, despite the **NP**-completeness of SAT.

Another approach is to change the type of problem which is being considered, a tactic which typically results in the definition of new complexity classes. For example, instead of finding exact solutions to an **NP**-complete problem, we can instead try to find good algorithms for finding *approximate* solutions to a problem. For example, the **VERTEX COVER** problem is an **NP**-complete problem, yet in Exercise 3.27 we show that it is possible to efficiently find an approximation to the minimal vertex cover which is correct to within a factor two! On the other hand, in Problem 3.6 we show that it is not possible to find approximations to solutions of TSP correct to within any factor, *unless* **P** = **NP**!

Exercise 3.27: (Approximation algorithm for **VERTEX COVER)** Let $G = (V, E)$ be an undirected graph. Prove that the following algorithm finds a vertex cover for G that is within a factor two of being a minimal vertex cover:

```

VC = ∅
E' = E
do until E' = ∅
    let (α, β) be any edge of E'
    VC = VC ∪ {α, β}
    remove from E' every edge incident on α or β
return VC.
```

Why is it possible to approximate the solution of one **NP**-complete problem, but not another? After all, isn't it possible to efficiently transform from one problem to another? This is certainly true, however it is not necessarily true that this transformation preserves the notion of a 'good approximation' to a solution. As a result, the computational complexity theory of approximation algorithms for problems in **NP** has a structure that goes beyond the structure of **NP** proper. An entire complexity theory of approximation algorithms exists, which unfortunately is beyond the scope of this book. The basic idea, however, is to define a notion of reduction that corresponds to being able to efficiently reduce one approximation problem to another, in such a way that the notion of good approximation is preserved. With such a notion, it is possible to define complexity classes such as **MAXSNP** by analogy to the class **NP**, as the set of problems for which it is possible to efficiently verify approximate solutions to the problem. Complete problems exist for **MAXSNP**, just as for **NP**, and it is an interesting open problem to determine how the class **MAXSNP** compares to the class of approximation problems which are efficiently solvable.

We conclude our discussion with a complexity class that results when the underlying model of computation itself is changed. Suppose a Turing machine is endowed with the ability to flip coins, using the results of the coin tosses to decide what actions to take during the computation. Such a Turing machine may only accept or reject inputs with a certain probability. The complexity class **BPP** (for *bounded-error probabilistic time*) contains all languages L with the property that there exists a probabilistic Turing machine M such that if $x \in L$ then M accepts x with probability at least $3/4$, and if $x \notin L$, then M rejects x with probability at least $3/4$. The following exercise shows that the choice of the constant $3/4$ is essentially arbitrary:

Exercise 3.28: (Arbitrariness of the constant in the definition of BPP) Suppose k is a fixed constant, $1/2 < k \leq 1$. Suppose L is a language such that there exists a Turing machine M with the property that whenever $x \in L$, M accepts x with probability at least k , and whenever $x \notin L$, M rejects x with probability at least k . Show that $L \in \mathbf{BPP}$.

Indeed, the *Chernoff bound*, discussed in Box 3.4, implies that with just a few repetitions of an algorithm deciding a language in **BPP** the probability of success can be amplified to the point where it is essentially equal to one, for all intents and purposes. For this reason, **BPP** even more than **P** is the class of decision problems which is usually regarded as being efficiently solvable on a classical computer, and it is the quantum analogue of **BPP**, known as **BQP**, that is most interesting in our study of quantum algorithms.

3.2.5 Energy and computation

Computational complexity studies the amount of time and space required to solve a computational problem. Another important computational resource is *energy*. In this section, we study the energy requirements for computation. Surprisingly, it turns out that computation, both classical and quantum, can in principle be done without expending any energy! Energy consumption in computation turns out to be deeply linked to the *reversibility* of the computation. Consider a gate like the NAND gate, which takes as input two bits, and produces a single bit as output. This gate is intrinsically *irreversible* because, given the output of the gate, the input is not uniquely determined. For example, if the output of the NAND gate is 1, then the input could have been any one of 00, 01, or 10. On the other hand, the NOT gate is an example of a *reversible* logic gate because, given the output of the NOT gate, it is possible to infer what the input must have been.

Another way of understanding irreversibility is to think of it in terms of information erasure. If a logic gate is irreversible, then some of the information input to the gate is lost irretrievably when the gate operates – that is, some of the information has been erased by the gate. Conversely, in a reversible computation, no information is ever erased, because the input can always be recovered from the output. Thus, saying that a computation is reversible is equivalent to saying that no information is erased during the computation.

What is the connection between energy consumption and irreversibility in computation? *Landauer's principle* provides the connection, stating that, in order to erase information, it is necessary to dissipate energy. More precisely, Landauer's principle may be stated as follows:

Landauer's principle (first form): Suppose a computer erases a single bit of information. The amount of energy dissipated into the environment is *at least* $k_B T \ln 2$, where k_B is a universal constant known as *Boltzmann's constant*, and T is the temperature of the environment of the computer.

According to the laws of thermodynamics, Landauer's principle can be given an alternative form stated not in terms of energy dissipation, but rather in terms of entropy:

Landauer's principle (second form): Suppose a computer erases a single bit of information. The entropy of the environment increases by *at least* $k_B \ln 2$, where k_B is Boltzmann's constant.

Justifying Landauer's principle is a problem of physics that lies beyond the scope of this

Box 3.4: BPP and the Chernoff bound

Suppose we have an algorithm for a decision problem which gives the correct answer with probability $1/2 + \epsilon$, and the wrong answer with probability $1/2 - \epsilon$. If we run the algorithm n times, then it seems reasonable to guess that the correct answer is whichever appeared most frequently. How reliably does this procedure work? The *Chernoff bound* is a simple result from elementary probability which answers this question.

Theorem 3.3: (The Chernoff bound) Suppose X_1, \dots, X_n are independent and identically distributed random variables, each taking the value 1 with probability $1/2 + \epsilon$, and the value 0 with probability $1/2 - \epsilon$. Then

$$p\left(\sum_{i=1}^n X_i \leq n/2\right) \leq e^{-2\epsilon^2 n}. \quad (3.2)$$

Proof

Consider any sequence (x_1, \dots, x_n) containing at most $n/2$ ones. The probability of such a sequence occurring is maximized when it contains $\lfloor n/2 \rfloor$ ones, so

$$p(X_1 = x_1, \dots, X_n = x_n) \leq \left(\frac{1}{2} - \epsilon\right)^{\frac{n}{2}} \left(\frac{1}{2} + \epsilon\right)^{\frac{n}{2}} \quad (3.3)$$

$$= \frac{(1 - 4\epsilon^2)^{\frac{n}{2}}}{2^n}. \quad (3.4)$$

There can be at most 2^n such sequences, so

$$p\left(\sum_i X_i \leq n/2\right) \leq 2^n \times \frac{(1 - 4\epsilon^2)^{\frac{n}{2}}}{2^n} = (1 - 4\epsilon^2)^{\frac{n}{2}}. \quad (3.5)$$

Finally, by calculus, $1 - x \leq \exp(-x)$, so

$$p\left(\sum_i X_i \leq n/2\right) \leq e^{-4\epsilon^2 n/2} = e^{-2\epsilon^2 n}. \quad (3.6)$$

□

What this tells us is that for fixed ϵ , the probability of making an error decreases *exponentially quickly* in the number of repetitions of the algorithm. In the case of **BPP** we have $\epsilon = 1/4$, so it takes only a few hundred repetitions of the algorithm to reduce the probability of error below 10^{-20} , at which point an error in one of the computer's components becomes much more likely than an error due to the probabilistic nature of the algorithm.

book – see the end of chapter ‘History and further reading’ if you wish to understand why Landauer’s principle holds. However, if we accept Landauer’s principle as given, then it raises a number of interesting questions. First of all, Landauer’s principle only provides a *lower bound* on the amount of energy that must be dissipated to erase information.

How close are existing computers to this lower bound? Not very, turns out to be the answer – computers circa the year 2000 dissipate roughly $500k_B T \ln 2$ in energy for each elementary logical operation.

Although existing computers are far from the limit set by Landauer's principle, it is still an interesting problem of principle to understand how much the energy consumption can be reduced. Aside from the intrinsic interest of the problem, a practical reason for the interest follows from Moore's law: if computer power keeps increasing then the amount of energy dissipated must also increase, unless the energy dissipated per operation drops at least as fast as the rate of increase in computing power.

If all computations could be done reversibly, then Landauer's principle would imply no lower bound on the amount of energy dissipated by the computer, since no bits at all are erased during a reversible computation. Of course, it is possible that some other physical principle might require that energy be dissipated during the computation; fortunately, this turns out not to be the case. But is it possible to perform universal computation without erasing any information? Physicists can cheat on this problem to see in advance that the answer to this question *must* be yes, because our present understanding of the laws of physics is that they are fundamentally reversible. That is, if we know the final state of a closed physical system, then the laws of physics allow us to work out the initial state of the system. If we believe that those laws are correct, then we must conclude that hidden in the irreversible logic gates like AND and OR, there must be some underlying reversible computation. But where is this hidden reversibility, and can we use it to construct manifestly reversible computers?

We will use two different techniques to give reversible circuit-based models capable of universal computation. The first model, a computer built entirely of billiard balls and mirrors, gives a beautiful concrete realization of the principles of reversible computation. The second model, based on a reversible logic gate known as the *Toffoli gate* (which we first encountered in Section 1.4.1), is a more abstract view of reversible computation that will later be of great use in our discussion of quantum computation. It is also possible to build reversible Turing machines that are universal for computation; however, we won't study these here, since the reversible circuit models turn out to be much more useful for quantum computation.

The basic idea of the billiard ball computer is illustrated in Figure 3.14. Billiard ball 'inputs' enter the computer from the left hand side, bouncing off mirrors and each other, before exiting as 'outputs' on the right hand side. The presence or absence of a billiard ball at a possible input site is used to indicate a logical 1 or a logical 0, respectively. The fascinating thing about this model is that it is manifestly reversible, insofar as its operation is based on the laws of classical mechanics. Furthermore, this model of computation turns out to be *universal* in the sense that it can be used to simulate an arbitrary computation in the standard circuit model of computation.

Of course, if a billiard ball computer were ever built it would be highly unstable. As any billiards player can attest, a billiard ball rolling frictionlessly over a smooth surface is easily knocked off course by small perturbations. The billiard ball model of computation depends on perfect operation, and the absence of external perturbations such as those caused by thermal noise. Periodic corrections can be performed, but information gained by doing this would have to be erased, requiring work to be performed. Expenditure of energy thus serves the purpose of reducing this susceptibility to noise, which is necessary for a practical, real-world computational machine. For the purposes of this introduction,

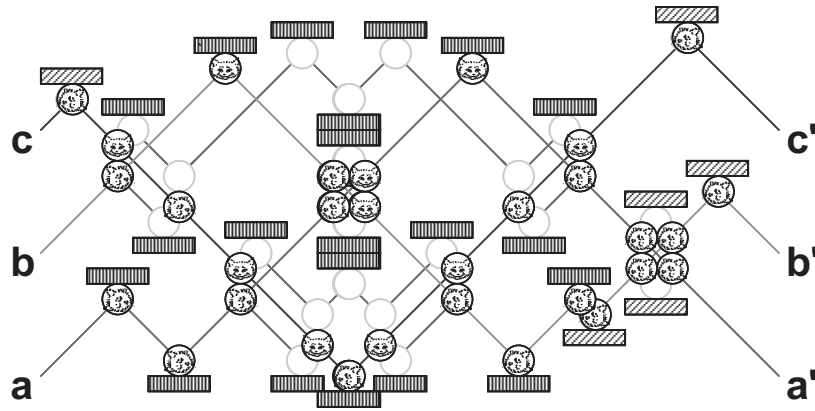


Figure 3.14. A simple billiard ball computer, with three input bits and three output bits, shown entering on the left and leaving on the right, respectively. The presence or absence of a billiard ball indicates a 1 or a 0, respectively. Empty circles illustrate potential paths due to collisions. This particular computer implements the Fredkin classical reversible logic gate, discussed in the text.

we will ignore the effects of noise on the billiard ball computer, and concentrate on understanding the essential elements of reversible computation.

The billiard ball computer provides an elegant means for implementing a reversible universal logic gate known as the *Fredkin gate*. Indeed, the properties of the Fredkin gate provide an informative overview of the general principles of reversible logic gates and circuits. The Fredkin gate has three input bits and three output bits, which we refer to as a, b, c and a', b', c' , respectively. The bit c is a *control bit*, whose value is not changed by the action of the Fredkin gate, that is, $c' = c$. The reason c is called the control bit is because it controls what happens to the other two bits, a and b . If c is set to 0 then a and b are left alone, $a' = a, b' = b$. If c is set to 1, a and b are swapped, $a' = b, b' = a$. The explicit truth table for the Fredkin gate is shown in Figure 3.15. It is easy to see that the Fredkin gate is reversible, because given the output a', b', c' , we can determine the inputs a, b, c . In fact, to recover the original inputs a, b and c we need only apply another Fredkin gate to a', b', c' :

Exercise 3.29: (Fredkin gate is self-inverse) Show that applying two consecutive Fredkin gates gives the same outputs as inputs.

Examining the paths of the billiard balls in Figure 3.14, it is not difficult to verify that this billiard ball computer implements the Fredkin gate:

Exercise 3.30: Verify that the billiard ball computer in Figure 3.14 computes the Fredkin gate.

In addition to reversibility, the Fredkin gate also has the interesting property that the number of 1s is *conserved* between the input and output. In terms of the billiard ball computer, this corresponds to the number of billiard balls going into the Fredkin gate being equal to the number coming out. Thus, it is sometimes referred to as being a *conservative* reversible logic gate. Such reversibility and conservative properties are interesting to a physicist because they can be motivated by fundamental physical principles.

Inputs			Outputs		
a	b	c	a'	b'	c'
0	0	0	0	0	0
0	0	1	0	0	1
0	1	0	0	1	0
0	1	1	1	0	1
1	0	0	1	0	0
1	0	1	0	1	1
1	1	0	1	1	0
1	1	1	1	1	1

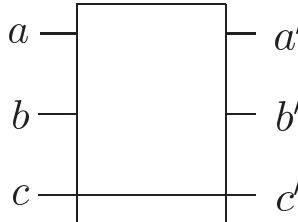


Figure 3.15. Fredkin gate truth table and circuit representation. The bits a and b are swapped if the control bit c is set, and otherwise are left alone.

ples. The laws of Nature are reversible, with the possible exception of the measurement postulate of quantum mechanics, discussed in Section 2.2.3 on page 84. The conservative property can be thought of as analogous to properties such as conservation of mass, or conservation of energy. Indeed, in the billiard ball model of computation the conservative property corresponds exactly to conservation of mass.

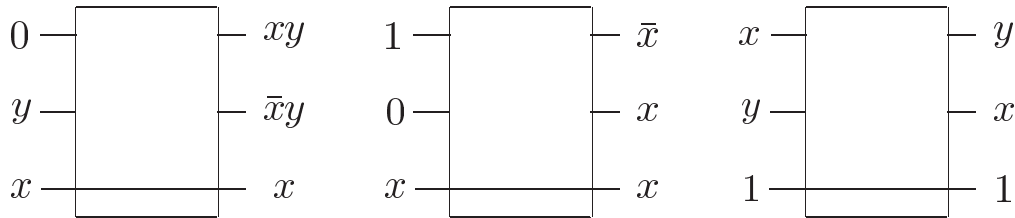


Figure 3.16. Fredkin gate configured to perform the elementary gates AND (left), NOT (middle), and a primitive routing function, the CROSSOVER (right). The middle gate also serves to perform the FANOUT operation, since it produces two copies of x at the output. Note that each of these configurations requires the use of extra ‘ancilla’ bits prepared in standard states – for example, the 0 input on the first line of the AND gate – and in general the output contains ‘garbage’ not needed for the remainder of the computation.

The Fredkin gate is not only reversible and conservative, it’s a universal logic gate as well! As illustrated in Figure 3.16, the Fredkin gate can be configured to simulate AND, NOT, CROSSOVER and FANOUT functions, and thus can be cascaded to simulate any classical circuit whatsoever.

To simulate irreversible gates such as AND using the Fredkin gate, we made use of two ideas. First, we allowed the input of ‘ancilla’ bits to the Fredkin gate, in specially prepared states, either 0 or 1. Second, the output of the Fredkin gate contained extraneous ‘garbage’ not needed for the remainder of the computation. These ancilla and garbage bits are not directly important to the computation. Their importance lies in the fact that they make the computation reversible. Indeed the irreversibility of gates like the AND and OR may be viewed as a consequence of the ancilla and garbage bits being ‘hidden’. Summarizing, given any classical circuit computing a function $f(x)$, we can build a reversible circuit made entirely of Fredkin gates, which on input of x , together with some ancilla bits

in a standard state a , computes $f(x)$, together with some extra ‘garbage’ output, $g(x)$. Therefore, we represent the action of the computation as $(x, a) \rightarrow (f(x), g(x))$.

We now know how to compute functions reversibly. Unfortunately, this computation produces unwanted garbage bits. With some modifications it turns out to be possible to perform the computation so that any garbage bits produced are in a *standard state*. This construction is crucial for quantum computation, because garbage bits whose value depends upon x will in general destroy the interference properties crucial to quantum computation. To understand how this works it is convenient to assume that the NOT gate is available in our repertoire of reversible gates, so we may as well assume that the ancilla bits a all start out as 0s, with NOT gates being added where necessary to turn the ancilla 0s into 1s. It will also be convenient to assume that the classical controlled-NOT gate is available, defined in a manner analogous to the quantum definition of Section 1.3.2, that is, the inputs (c, t) are taken to $(c, t \oplus c)$, where \oplus denotes addition modulo 2. Notice that $t = 0$ gives $(c, 0) \rightarrow (c, c)$, so the controlled-NOT can be thought of as a reversible copying gate or FANOUT, which leaves no garbage bits at the output.

With the additional NOT gates appended at the beginning of the circuit, the action of the computation may be written as $(x, 0) \rightarrow (f(x), g(x))$. We could also have added CNOT gates to the beginning of the circuit, in order to create a copy of x which is not changed during the subsequent computation. With this modification, the action of the circuit may be written

$$(x, 0, 0) \rightarrow (x, f(x), g(x)). \quad (3.7)$$

Equation (3.7) is a very useful way of writing the action of the reversible circuit, because it allows an idea known as *uncomputation* to be used to get rid of the garbage bits, for a small cost in the running time of the computation. The idea is the following. Suppose we start with a four register computer in the state $(x, 0, 0, y)$. The second register is used to store the result of the computation, and the third register is used to provide workspace for the computation, that is, the garbage bits $g(x)$. The use of the fourth register is described shortly, and we assume it starts in an arbitrary state y .

We begin as before, by applying a reversible circuit to compute f , resulting in the state $(x, f(x), g(x), y)$. Next, we use CNOTs to add the result $f(x)$ bitwise to the fourth register, leaving the machine in the state $(x, f(x), g(x), y \oplus f(x))$. However, all the steps used to compute $f(x)$ were reversible and did not affect the fourth register, so by applying the reverse of the circuit used to compute f we come to the state $(x, 0, 0, y \oplus f(x))$. Typically, we omit the ancilla 0s from the description of the function evaluation, and just write the action of the circuit as

$$(x, y) \rightarrow (x, y \oplus f(x)). \quad (3.8)$$

In general we refer to this modified circuit computing f as *the reversible circuit computing f* , even though in principle there are many other reversible circuits which could be used to compute f .

What resource overhead is involved in doing reversible computation? To analyze this question, we need to count the number of extra ancilla bits needed in a reversible circuit, and compare the gate counts with classical models. It ought to be clear that the number of gates in a reversible circuit is the same as in an irreversible circuit to within the constant factor which represents the number of Fredkin gates needed to simulate a single element of the irreversible circuit, and an additional factor of two for uncomputation, with an

overhead for the extra CNOT operations used in reversible computation which is linear in the number of bits involved in the circuit. Similarly, the number of ancilla bits required scales at most linearly with the number of gates in the irreversible circuit, since each element in the irreversible circuit can be simulated using a constant number of ancilla bits. As a result, natural complexity classes such as **P** and **NP** are the same no matter whether a reversible or irreversible model of computation is used. For more elaborate complexity classes like **PSPACE** the situation is not so immediately clear; see Problem 3.9 and ‘History and further reading’ for a discussion of some such subtleties.

Exercise 3.31: (Reversible half-adder) Construct a reversible circuit which, when two bits x and y are input, outputs $(x, y, c, x \oplus y)$, where c is the carry bit when x and y are added.

The Fredkin gate and its implementation using the billiard ball computer offers a beautiful paradigm for reversible computation. There is another reversible logic gate, the *Toffoli gate*, which is also universal for classical computation. While the Toffoli gate does not have quite the same elegant physical simplicity as the billiard ball implementation of the Fredkin gate, it will be more useful in the study of quantum computation. We have already met the Toffoli gate in Section 1.4.1, but for convenience we review its properties here.

The Toffoli gate has three input bits, a, b and c . a and b are known as the first and second *control bits*, while c is the *target bit*. The gate leaves both control bits unchanged, flips the target bit if both control bits are set, and otherwise leaves the target bit alone. The truth table and circuit representation for the Toffoli gate are shown in Figure 3.17.

Inputs			Outputs		
a	b	c	a'	b'	c'
0	0	0	0	0	0
0	0	1	0	0	1
0	1	0	0	1	0
0	1	1	0	1	1
1	0	0	1	0	0
1	0	1	1	0	1
1	1	0	1	1	1
1	1	1	1	1	0

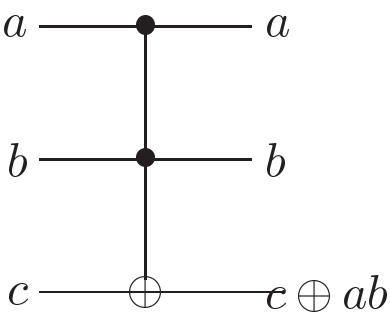


Figure 3.17. Truth table and circuit representation of the Toffoli gate.

How can the Toffoli gate be used to do universal computation? Suppose we wish to NAND the bits a and b . To do this using the Toffoli gate, we input a and b as control bits, and send in an ancilla bit set to 1 as the target bit, as shown in Figure 3.18. The NAND of a and b is output as the target bit. As expected from our study of the Fredkin gate, the Toffoli gate simulation of a NAND requires the use of a special ancilla input, and some of the outputs from the simulation are garbage bits.

The Toffoli gate can also be used to implement the FANOUT operation by inputting an ancilla 1 to the first control bit, and a to the second control bit, producing the output 1, a, a . This is illustrated in Figure 3.19. Recalling that NAND and FANOUT are together

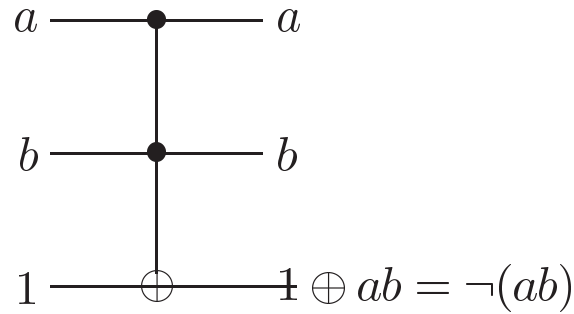


Figure 3.18. Implementing a NAND gate using a Toffoli gate. The top two bits represent the input to the NAND, while the third bit is prepared in the standard state 1, sometimes known as an *ancilla* state. The output from the NAND is on the third bit.

universal for computation, we see that an arbitrary circuit can be efficiently simulated using a reversible circuit consisting only of Toffoli gates and ancilla bits, and that useful additional techniques such as uncomputation may be achieved using the same methods as were employed with the Fredkin gate.

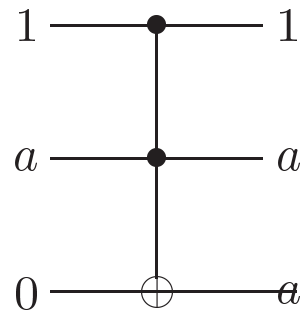


Figure 3.19. FANOUT with the Toffoli gate, with the second bit being the input to the FANOUT, and the other two bits standard ancilla states. The output from FANOUT appears on the second and third bits.

Our interest in reversible computation was motivated by our desire to understand the energy requirements for computation. It is clear that the noise-free billiard ball model of computation requires no energy for its operation; what about models based upon the Toffoli gate? This can only be determined by examining specific models for the computation of the Toffoli gate. In Chapter 7, we examine several such implementations, and it turns out that, indeed, the Toffoli gate can be implemented in a manner which does not require the expenditure of energy.

There is a significant caveat attached to the idea that computation can be done without the expenditure of energy. As we noted earlier, the billiard ball model of computation is highly sensitive to noise, and this is true of many other models of reversible computation. To nullify the effects of noise, some form of error-correction needs to be done. Such error-correction typically involves the performance of measurements on the system to determine whether the system is behaving as expected, or if an error has occurred. Because the computer's memory is finite, the bits used to store the measurement results utilized in error-correction must eventually be erased to make way for new measurement results. According to Landauer's principle, this erasure carries an associated energy cost

that must be accounted for when tallying the total energy cost of the computation. We analyze the energy cost associated with error-correction in more detail in Section 12.4.4.

What can we conclude from our study of reversible computation? There are three key ideas. First, reversibility stems from keeping track of every bit of information; irreversibility occurs only when information is lost or erased. Second, by doing computation reversibly, we obviate the need for energy expenditure during computation. All computations can be done, in principle, for zero cost in energy. Third, reversible computation can be done efficiently, without the production of garbage bits whose value depends upon the input to the computation. That is, if there is an irreversible circuit computing a function f , then there is an efficient simulation of this circuit by a reversible circuit with action $(x, y) \rightarrow (x, y \oplus f(x))$.

What are the implications of these results for physics, computer science, and for quantum computation and quantum information? From the point of view of a physicist or hardware engineer worried about heat dissipation, the good news is that, in principle, it is possible to make computation dissipation-free by making it reversible, although in practice energy dissipation is required for system stability and immunity from noise. At an even more fundamental level, the ideas leading to reversible computation also lead to the resolution of a century-old problem in the foundations of physics, the famous problem of *Maxwell's demon*. The story of this problem and its resolution is outlined in Box 3.5 on page 162. From the point of view of a computer scientist, reversible computation validates the use of irreversible elements in models of computation such as the Turing machine (since using them or not gives polynomially equivalent models). Moreover, since the physical world is fundamentally reversible, one can argue that complexity classes based upon reversible models of computation are more natural than complexity classes based upon irreversible models, a point revisited in Problem 3.9 and 'History and further reading'. From the point of view of quantum computation and quantum information, reversible computation is enormously important. To harness the full power of quantum computation, any classical subroutines in a quantum computation must be performed reversibly and without the production of garbage bits depending on the classical input.

Exercise 3.32: (From Fredkin to Toffoli and back again) What is the smallest number of Fredkin gates needed to simulate a Toffoli gate? What is the smallest number of Toffoli gates needed to simulate a Fredkin gate?

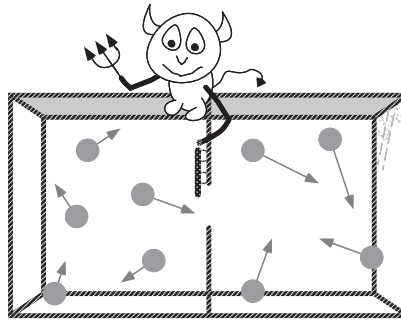
3.3 Perspectives on computer science

In a short introduction such as this chapter, it is not remotely possible to cover in detail all the great ideas of a field as rich as computer science. We hope to have conveyed to you something of what it means to *think* like a computer scientist, and provided a basic vocabulary and overview of some of the fundamental concepts important in the understanding of computation. To conclude this chapter, we briefly touch on some more general issues, in order to provide some perspective on how quantum computation and quantum information fits into the overall picture of computer science.

Our discussion has revolved around the Turing machine model of computation. How does the computational power of unconventional models of computation such as massively parallel computers, DNA computers and analog computers compare with the standard

Box 3.5: Maxwell's demon

The laws of thermodynamics govern the amount of work that can be performed by a physical system at thermodynamic equilibrium. One of these laws, the second law of thermodynamics, states that the *entropy* in a closed system can never decrease. In 1871, James Clerk Maxwell proposed the existence of a machine that apparently violated this law. He envisioned a miniature little 'demon', like that shown in the figure below, which could reduce the entropy of a gas cylinder initially at equilibrium by individually separating the fast and slow molecules into the two halves of the cylinder. This demon would sit at a little door at the middle partition. When a fast molecule approaches from the left side the demon opens a door between the partitions, allowing the molecule through, and then closes the door. By doing this many times the total entropy of the cylinder can be *decreased*, in apparent violation of the second law of thermodynamics.



The resolution to the Maxwell's demon paradox lies in the fact that the demon must perform *measurements* on the molecules moving between the partitions, in order to determine their velocities. The result of this measurement must be stored in the demon's memory. Because any memory is finite, the demon must eventually begin erasing information from its memory, in order to have space for new measurement results. By Landauer's principle, this act of erasing information increases the total entropy of the combined system – demon, gas cylinder, and their environments. In fact, a complete analysis shows that Landauer's principle implies that the entropy of the combined system is increased *at least as much* by this act of erasing information as the entropy of the combined system is decreased by the actions of the demon, thus ensuring that the second law of thermodynamics is obeyed.

Turing machine model of computation and, implicitly, with quantum computation? Let's begin with parallel computing architectures. The vast majority of computers in existence are serial computers, processing instructions one at a time in some central processing unit. By contrast, parallel computers can process more than one instruction at a time, leading to a substantial savings in time and money for some applications. Nevertheless, parallel processing does not offer any fundamental advantage over the standard Turing machine model when issues of efficiency are concerned, because a Turing machine can simulate a parallel computer with polynomially equivalent total physical resources – the total space and time used by the computation. What a parallel computer gains in time,

it loses in the total spatial resources required to perform the computation, resulting in a net of no essential change in the power of the computing model.

An interesting specific example of massively parallel computing is the technique of *DNA computing*. A strand of DNA, deoxyribonucleic acid, is a molecule composed of a sequence (a polymer) of four kinds of nucleotides distinguished by the bases they carry, denoted by the letter A (adenine), C (cytosine), G (guanine) and T (thymine). Two strands, under certain circumstances, can anneal to form a double strand, if the respective base pairs form complements of each other (A matches T and G matches C). The ends are also distinct and must match appropriately. Chemical techniques can be used to amplify the number of strands beginning or ending with specific sequences (polymerase chain reaction), separate the strands by length (gel electrophoresis), dissolve double strands into single strands (changing temperature and pH), read the sequence on a strand, cut strands at a specific position (restriction enzymes), and detect if a certain sequence of DNA is in a test tube. The procedure for using these mechanisms in a robust manner is rather involved, but the basic idea can be appreciated from an example.

The directed Hamiltonian path problem is a simple and equivalently hard variant of the Hamiltonian cycle problem of Section 3.2.2, in which the goal is to determine if a path exists or not between two specified vertices j_1 and j_N in a directed graph G of N vertices, entering each vertex exactly once, and following only allowed edge directions. This problem can be solved with a DNA computer using the following five steps, in which x_j are chosen to be unique sequences of bases (and \bar{x}_j their complements), DNA strands $x_j x_k$ encode edges, and strands $\bar{x}_j \bar{x}_j$ encode vertices. (1) Generate random paths through G , by combining a mixture of all possible vertex and edge DNA strands, and waiting for the strands to anneal. (2) Keep only the paths beginning with j_1 and ending with j_N , by amplifying only the double strands beginning with \bar{x}_{j_1} and ending with \bar{x}_{j_N} . (3) Select only paths of length N , by separating the strands according to their length. (4) Select only paths which enter each vertex at least once, by dissolving the DNA into single strands, and annealing with all possible vertex strands one at a time and filtering out only those strands which anneal. And (5) detect if any strands have survived the selection steps; if so, then a path exists, and otherwise, it does not. To ensure the answer is correct with sufficiently high probability, x_j may be chosen to contain many (≈ 30) bases, and a large number ($\approx 10^{14}$ or more are feasible) of strands are used in the reaction.

Heuristic methods are available to improve upon this basic idea. Of course, exhaustive search methods such as this only work as long as all possible paths can be generated efficiently, and thus the number of molecules used must grow exponentially as the size of the problem (the number of vertices in the example above). DNA molecules are relatively small and readily synthesized, and the huge number of DNA combinations one can fit into a test tube can stave off the exponential complexity cost increase for a while – up to a few dozen vertices – but eventually the exponential cost limits the applicability of this method. Thus, while DNA computing offers an attractive and physically realizable model of computation for the solution of certain problems, it is a classical computing technique and offers no essential improvement in principle over a Turing machine.

Analog computers offer a yet another paradigm for performing computation. A computer is analog when the physical representation of information it uses for computation is based on continuous degrees of freedom, instead of zeroes and ones. For example, a thermometer is an analog computer. Analog circuitry, using resistors, capacitors, and amplifiers, is also said to perform analog computation. Such machines have an infinite

resource to draw upon in the ideal limit, since continuous variables like position and voltage can store an unlimited amount of information. But this is only true in the absence of noise. The presence of a finite amount of noise reduces the number of *distinguishable states* of a continuous variable to a finite number – and thus restricts analog computers to the representation of a finite amount of information. In practice, noise reduces analog computers to being no more powerful than conventional digital computers, and through them Turing machines. One might suspect that quantum computers are just analog computers, because of the use of continuous parameters in describing qubit states; however, it turns out that the effects of noise on a quantum computer can effectively be *digitized*. As a result, their computational advantages remain even in the presence of a finite amount of noise, as we shall see in Chapter 10.

What of the effects of noise on digital computers? In the early days of computation, noise was a very real problem for computers. In some of the original computers a vacuum tube would malfunction every few minutes. Even today, noise is a problem for computational devices such as modems and hard drives. Considerable effort was devoted to the problem of understanding how to construct reliable computers from unreliable components. It was proven by von Neumann that this is possible with only a polynomial increase in the resources required for computation. Ironically, however, modern computers use none of those results, because the components of modern computers are fantastically reliable. Failure rates of 10^{-17} and even less are common in modern electronic components. For this reason, failures happen so rarely that the extra effort required to protect against them is not regarded as being worth making. On the other hand, we shall find that quantum computers are very delicate machines, and will likely require substantial application of error-correction techniques.

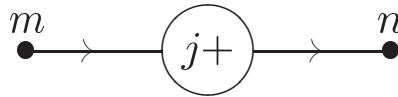
Different architectures may change the effects of noise. For example, if the effect of noise is ignored, then changing to a computer architecture in which many operations are performed in parallel may not change the number of operations which need to be done. However, a parallel system may be substantially more resistant to noise, because the effects of noise have less time to accumulate. Therefore, in a realistic analysis, the parallel version of an algorithm may have some substantial advantages over a serial implementation. Architecture design is a well developed field of study for classical computers. Hardly anything similar has been developed along the same lines for quantum computers, but the study of noise already suggests some desirable traits for future quantum computer architectures, such as a high level of parallelism.

A fourth model of computation is *distributed computation*, in which two or more spatially separated computational units are available to solve a computational problem. Obviously, such a model of computation is no more powerful than the Turing machine model in the sense that it can be efficiently simulated on a Turing machine. However, distributed computation gives rise to an intriguing new resource challenge: how best to utilize multiple computational units when the cost of *communication* between the units is high. This problem of distributed computation becomes especially interesting as computers are connected through high speed networks; although the total computational capacity of all the computers on a network might be extremely large, utilization of that potential is difficult. Most interesting problems do not divide easily into independent chunks that can be solved separately, and may frequently require global communication between different computational subsystems to exchange intermediate results or synchronize status. The field of *communication complexity* has been developed to address such issues, by

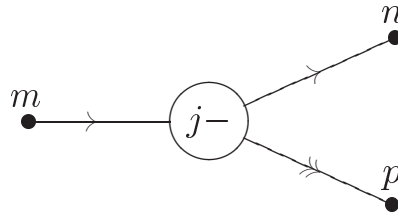
quantifying the cost of communication requirements in solving problems. When quantum resources are available and can be exchanged between distributed computers, the communication costs can sometimes be greatly reduced.

A recurring theme through these concluding thoughts and through the entire book is that despite the traditional independence of computer science from physical constraints, ultimately physical laws have tremendous impact not only upon how computers are realized, but also the class of problems they are capable of solving. The success of quantum computation and quantum information as a physically reasonable alternative model of computation questions closely held tenets of computer science, and thrusts notions of computer science into the forefront of physics. The task of the remainder of this book is to stir together ideas from these disparate fields, and to delight in what results!

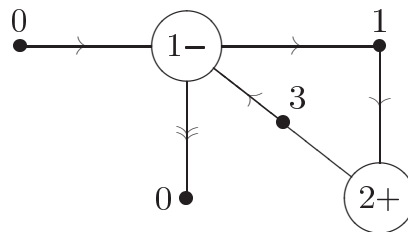
Problem 3.1: (Minsky machines) A *Minsky machine* consists of a finite set of registers, r_1, r_2, \dots, r_k , each capable of holding an arbitrary non-negative integer, and a *program*, made up of orders of one of two types. The first type has the form:



The interpretation is that at point m in the program register r_j is incremented by one, and execution proceeds to point n in the program. The second type of order has the form:



The interpretation is that at point m in the program, register r_j is decremented if it contains a positive integer, and execution proceeds to point n in the program. If register r_j is zero then execution simply proceeds to point p in the program. The *program* for the Minsky machine consists of a collection of such orders, of a form like:



The starting and all possible halting points for the program are conventionally labeled zero. This program takes the contents of register r_1 and adds them to register r_2 , while decrementing r_1 to zero.

- (1) Prove that all (Turing) computable functions can be computed on a Minsky machine, in the sense that given a computable function $f(\cdot)$ there is a Minsky machine program that when the registers start in the state $(n, 0, \dots, 0)$ gives as output $(f(n), 0, \dots, 0)$.
- (2) Sketch a proof that any function which can be computed on a Minsky machine, in the sense just defined, can also be computed on a Turing machine.

Problem 3.2: (Vector games) A *vector game* is specified by a finite list of vectors, all of the same dimension, and with integer co-ordinates. The game is to start with a vector x of non-negative integer co-ordinates and to add to x the first vector from the list which preserves the non-negativity of all the components, and to repeat this process until it is no longer possible. Prove that for any computable function $f(\cdot)$ there is a vector game which when started with the vector $(n, 0, \dots, 0)$ reaches $(f(n), 0, \dots, 0)$. (*Hint: Show that a vector game in $k + 2$ dimensions can simulate a Minsky machine containing k registers.*)

Problem 3.3: (Fractran) A *Fractran* program is defined by a list of positive rational numbers q_1, \dots, q_n . It acts on a positive integer m by replacing it by $q_i m$, where i is the least number such that $q_i m$ is an integer. If there is ever a time when there is no i such that $q_i m$ is an integer, then execution stops. Prove that for any computable function $f(\cdot)$ there is a Fractran program which when started with 2^n reaches $2^{f(n)}$ without going through any intermediate powers of 2. (*Hint: use the previous problem.*)

Problem 3.4: (Undecidability of dynamical systems) A Fractran program is essentially just a very simple dynamical system taking positive integers to positive integers. Prove that there is no algorithm to decide whether such a dynamical system ever reaches 1.

Problem 3.5: (Non-universality of two bit reversible logic) Suppose we are trying to build circuits using only one and two bit reversible logic gates, and ancilla bits. Prove that there are Boolean functions which cannot be computed in this fashion. Deduce that the Toffoli gate cannot be simulated using one and two bit reversible gates, even with the aid of ancilla bits.

Problem 3.6: (Hardness of approximation of TSP) Let $r \geq 1$ and suppose that there is an approximation algorithm for TSP which is guaranteed to find the shortest tour among n cities to within a factor r . Let $G = (V, E)$ be any graph on n vertices. Define an instance of TSP by identifying cities with vertices in V , and defining the distance between cities i and j to be 1 if (i, j) is an edge of G , and to be $\lceil r \rceil |V| + 1$ otherwise. Show that if the approximation algorithm is applied to this instance of TSP then it returns a Hamiltonian cycle for G if one exists, and otherwise returns a tour of length more than $\lceil r \rceil |V|$. From the NP-completeness of HC it follows that no such approximation algorithm can exist unless $\mathbf{P} = \mathbf{NP}$.

Problem 3.7: (Reversible Turing machines)

- (1) Explain how to construct a reversible Turing machine that can compute the same class of functions as is computable on an ordinary Turing machine. (*Hint: It may be helpful to use a multi-tape construction.*)

- (2) Give general space and time bounds for the operation of your reversible Turing machine, in terms of the time $t(x)$ and space $s(x)$ required on an ordinary single-tape Turing machine to compute a function $f(x)$.

Problem 3.8: (Find a hard-to-compute class of functions (Research)) Find a natural class of functions on n inputs which requires a super-polynomial number of Boolean gates to compute.

Problem 3.9: (Reversible PSPACE = PSPACE) It can be shown that the problem ‘quantified satisfiability’, or QSAT, is **PSPACE**-complete. That is, every other language in **PSPACE** can be reduced to QSAT in polynomial time. The language QSAT is defined to consist of all Boolean formulae φ in n variables x_1, \dots, x_n , and in conjunctive normal form, such that:

$$\exists x_1 \forall x_2 \exists x_3 \dots \forall x_n \varphi \text{ if } n \text{ is even;} \quad (3.9)$$

$$\exists x_1 \forall x_2 \exists x_3 \dots \exists x_n \varphi \text{ if } n \text{ is odd.} \quad (3.10)$$

Prove that a reversible Turing machine operating in polynomial space can be used to solve QSAT. Thus, the class of languages decidable by a computer operating reversibly in polynomial space is equal to **PSPACE**.

Problem 3.10: (Ancilla bits and efficiency of reversible computation) Let p_m be the m th prime number. Outline the construction of a reversible circuit which, upon input of m and n such that $n > m$, outputs the product $p_m p_n$, that is $(m, n) \rightarrow (p_m p_n, g(m, n))$, where $g(m, n)$ is the final state of the ancilla bits used by the circuit. Estimate the number of ancilla qubits your circuit requires. Prove that if a polynomial (in $\log n$) size reversible circuit can be found that uses $O(\log(\log n))$ ancilla bits then the problem of factoring a product of two prime numbers is in **P**.

History and further reading

Computer science is a huge subject with many interesting subfields. We cannot hope for any sort of completeness in this brief space, but instead take the opportunity to recommend a few titles of general interest, and some works on subjects of specific interest in relation to topics covered in this book, with the hope that they may prove stimulating.

Modern computer science dates to the wonderful 1936 paper of Turing^[Tur36]. The Church–Turing thesis was first stated by Church^[Chu36] in 1936, and was then given a more complete discussion from a different point of view by Turing. Several other researchers found their way to similar conclusions at about the same time. Many of these contributions and a discussion of the history may be found in a volume edited by Davis^[Dav65]. Provocative discussions of the Church–Turing thesis and undecidability may be found in Hofstadter^[Hof79] and Penrose^[Pen89].

There are many excellent books on algorithm design. We mention only three. First, there is the classic series by Knuth^[Knu97, Knu98a, Knu98b] which covers an enormous portion of computer science. Second, there is the marvelous book by Cormen, Leiserson, and Rivest^[CLR90]. This huge book contains a plethora of well-written material on many areas

of algorithm design. Finally, the book of Motwani and Raghavan^[MR95] is an excellent survey of the field of randomized algorithms.

The modern theory of computational complexity was especially influenced by the papers of Cook^[Coo71] and Karp^[Kar72]. Many similar ideas were arrived at independently in Russia by Levin^[Lev73], but unfortunately took time to propagate to the West. The classic book by Garey and Johnson^[GJ79] has also had an enormous influence on the field. More recently, Papadimitriou^[Pap94] has written a beautiful book that surveys many of the main ideas of computational complexity theory. Much of the material in this chapter is based upon Papadimitriou's book. In this chapter we considered only one type of reducibility between languages, polynomial time reducibility. There are many other notions of reductions between languages. An early survey of these notions was given by Ladner, Lynch and Selman^[LLS75]. The study of different notions of reducibility later blossomed into a subfield of research known as *structural complexity*, which has been reviewed by Balcázar, Diaz, and Gabarró^[BDG88a, BDG88b].

The connection between information, energy dissipation, and computation has a long history. The modern understanding is due to a 1961 paper by Landauer^[Lan61], in which Landauer's principle was first formulated. A paper by Szilard^[Szi29] and a 1949 lecture by von Neumann^[von66] (page 66) arrive at conclusions close to Landauer's principle, but do not fully grasp the essential point that it is the *erasure* of information that requires dissipation.

Reversible Turing machines were invented by Lecerf^[Lec63] and later, but independently, in an influential paper by Bennett^[Ben73]. Fredkin and Toffoli^[FT82] introduced reversible circuit models of computation. Two interesting historical documents are Barton's May, 1978 MIT 6.895 term paper^[Bar78], and Ressler's 1981 Master's thesis^[Res81], which contain designs for a reversible PDP-10! Today, reversible logic is potentially important in implementations of low-power CMOS circuitry^[YK95].

Maxwell's demon is a fascinating subject, with a long and intricate history. Maxwell proposed his demon in 1871^[Max71]. Szilard published a key paper in 1929^[Szi29] which anticipated many of the details of the final resolution of the problem of Maxwell's demon. In 1965 Feynman^[FLS65b] resolved a special case of Maxwell's demon. Bennett, building on Landauer's work^[Lan61], wrote two beautiful papers on the subject^[BBBW82, Ben87] which completed the resolution of the problem. An interesting book about the history of Maxwell's demon and its exorcism is the collection of papers by Leff and Rex^[LR90].

DNA computing was invented by Adleman, and the solution of the directed Hamiltonian path problem we describe is his^[Adl94]. Lipton has also shown how 3SAT and circuit satisfiability can be solved in this model^[Lip95]. A good general article is Adleman's *Scientific American* article^[Adl98]; for an insightful look into the universality of DNA operations, see Winfree^[Win98]. An interesting place to read about performing reliable computation in the presence of noise is the book by Winograd and Cowan^[WC67]. This topic will be addressed again in Chapter 10. A good textbook on computer architecture is by Hennessey, Goldberg, and Patterson.^[HGP96]

Problems 3.1 through 3.4 explore a line of thought originated by Minsky (in his beautiful book on computational machines^[Min67]) and developed by Conway^[Con72, Con86]. The Fractran programming language is certainly one of the most beautiful and elegant universal computational models known, as demonstrated by the following example, known

as PRIMEGAME^[Con86]. PRIMEGAME is defined by the list of rational numbers:

$$\frac{17}{91}, \frac{78}{85}, \frac{19}{51}, \frac{23}{38}, \frac{29}{33}, \frac{77}{29}, \frac{95}{23}, \frac{77}{19}, \frac{1}{17}, \frac{11}{13}, \frac{13}{11}, \frac{15}{2}, \frac{1}{7}, \frac{55}{1}. \quad (3.11)$$

Amazingly, when PRIMEGAME is started at 2, the other powers of 2 that appear, namely, $2^2, 2^3, 2^5, 2^7, 2^{11}, 2^{13}, \dots$, are precisely the prime powers of 2, with the powers stepping through the prime numbers, in order. Problem 3.9 is a special case of the more general subject of the spatial requirements for reversible computation. See the papers by Bennett^[Ben89], and by Li, Tromp and Vitanyi^[LV96, LTV98].

