

**FIGURE 4.10**  
The relationship among classes of languages

## 4.2

# UNDECIDABILITY

In this section, we prove one of the most philosophically important theorems of the theory of computation: There is a specific problem that is algorithmically unsolvable. Computers appear to be so powerful that you may believe that all problems will eventually yield to them. The theorem presented here demonstrates that computers are limited in a fundamental way.

What sorts of problems are unsolvable by computer? Are they esoteric, dwelling only in the minds of theoreticians? No! Even some ordinary problems that people want to solve turn out to be computationally unsolvable.

In one type of unsolvable problem, you are given a computer program and a precise specification of what that program is supposed to do (e.g., sort a list of numbers). You need to verify that the program performs as specified (i.e., that it is correct). Because both the program and the specification are mathematically precise objects, you hope to automate the process of verification by feeding these objects into a suitably programmed computer. However, you will be disappointed. The general problem of software verification is not solvable by computer.

In this section and in Chapter 5, you will encounter several computationally unsolvable problems. We aim to help you develop a feeling for the types of problems that are unsolvable and to learn techniques for proving unsolvability.

Now we turn to our first theorem that establishes the undecidability of a specific language: the problem of determining whether a Turing machine accepts a

given input string. We call it  $A_{\text{TM}}$  by analogy with  $A_{\text{DFA}}$  and  $A_{\text{CFG}}$ . But, whereas  $A_{\text{DFA}}$  and  $A_{\text{CFG}}$  were decidable,  $A_{\text{TM}}$  is not. Let

$$A_{\text{TM}} = \{\langle M, w \rangle \mid M \text{ is a TM and } M \text{ accepts } w\}.$$

### THEOREM 4.11

$A_{\text{TM}}$  is undecidable.

Before we get to the proof, let's first observe that  $A_{\text{TM}}$  is Turing-recognizable. Thus, this theorem shows that recognizers *are* more powerful than deciders. Requiring a TM to halt on all inputs restricts the kinds of languages that it can recognize. The following Turing machine  $U$  recognizes  $A_{\text{TM}}$ .

$U =$  "On input  $\langle M, w \rangle$ , where  $M$  is a TM and  $w$  is a string:

1. Simulate  $M$  on input  $w$ .
2. If  $M$  ever enters its accept state, *accept*; if  $M$  ever enters its reject state, *reject*."

Note that this machine loops on input  $\langle M, w \rangle$  if  $M$  loops on  $w$ , which is why this machine does not decide  $A_{\text{TM}}$ . If the algorithm had some way to determine that  $M$  was not halting on  $w$ , it could *reject* in this case. However, an algorithm has no way to make this determination, as we shall see.

The Turing machine  $U$  is interesting in its own right. It is an example of the *universal Turing machine* first proposed by Alan Turing in 1936. This machine is called universal because it is capable of simulating any other Turing machine from the description of that machine. The universal Turing machine played an important early role in the development of stored-program computers.

## THE DIAGONALIZATION METHOD

The proof of the undecidability of  $A_{\text{TM}}$  uses a technique called *diagonalization*, discovered by mathematician Georg Cantor in 1873. Cantor was concerned with the problem of measuring the sizes of infinite sets. If we have two infinite sets, how can we tell whether one is larger than the other or whether they are of the same size? For finite sets, of course, answering these questions is easy. We simply count the elements in a finite set, and the resulting number is its size. But if we try to count the elements of an infinite set, we will never finish! So we can't use the counting method to determine the relative sizes of infinite sets.

For example, take the set of even integers and the set of all strings over  $\{0,1\}$ . Both sets are infinite and thus larger than any finite set, but is one of the two larger than the other? How can we compare their relative size?

Cantor proposed a rather nice solution to this problem. He observed that two finite sets have the same size if the elements of one set can be paired with the elements of the other set. This method compares the sizes without resorting to counting. We can extend this idea to infinite sets. Here it is more precisely.

**DEFINITION 4.12**

Assume that we have sets  $A$  and  $B$  and a function  $f$  from  $A$  to  $B$ . Say that  $f$  is **one-to-one** if it never maps two different elements to the same place—that is, if  $f(a) \neq f(b)$  whenever  $a \neq b$ . Say that  $f$  is **onto** if it hits every element of  $B$ —that is, if for every  $b \in B$  there is an  $a \in A$  such that  $f(a) = b$ . Say that  $A$  and  $B$  are the **same size** if there is a one-to-one, onto function  $f: A \rightarrow B$ . A function that is both one-to-one and onto is called a **correspondence**. In a correspondence, every element of  $A$  maps to a unique element of  $B$  and each element of  $B$  has a unique element of  $A$  mapping to it. A correspondence is simply a way of pairing the elements of  $A$  with the elements of  $B$ .

Alternative common terminology for these types of functions is **injective** for one-to-one, **surjective** for onto, and **bijective** for one-to-one and onto.

**EXAMPLE 4.13**

Let  $\mathcal{N}$  be the set of natural numbers  $\{1, 2, 3, \dots\}$  and let  $\mathcal{E}$  be the set of even natural numbers  $\{2, 4, 6, \dots\}$ . Using Cantor's definition of size, we can see that  $\mathcal{N}$  and  $\mathcal{E}$  have the same size. The correspondence  $f$  mapping  $\mathcal{N}$  to  $\mathcal{E}$  is simply  $f(n) = 2n$ . We can visualize  $f$  more easily with the help of a table.

$n$	$f(n)$
1	2
2	4
3	6
$\vdots$	$\vdots$

Of course, this example seems bizarre. Intuitively,  $\mathcal{E}$  seems smaller than  $\mathcal{N}$  because  $\mathcal{E}$  is a proper subset of  $\mathcal{N}$ . But pairing each member of  $\mathcal{N}$  with its own member of  $\mathcal{E}$  is possible, so we declare these two sets to be the same size. ■

**DEFINITION 4.14**

A set  $A$  is **countable** if either it is finite or it has the same size as  $\mathcal{N}$ .

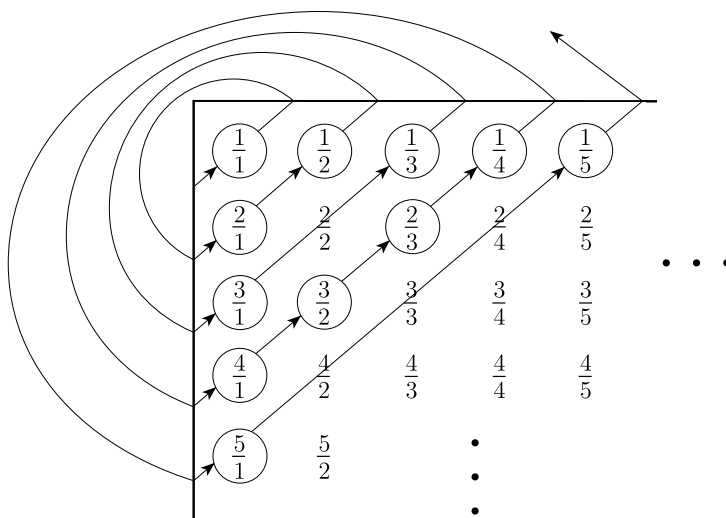
**EXAMPLE 4.15**

Now we turn to an even stranger example. If we let  $\mathcal{Q} = \{\frac{m}{n} \mid m, n \in \mathcal{N}\}$  be the set of positive rational numbers,  $\mathcal{Q}$  seems to be much larger than  $\mathcal{N}$ . Yet these two sets are the same size according to our definition. We give a correspondence with  $\mathcal{N}$  to show that  $\mathcal{Q}$  is countable. One easy way to do so is to list all the

elements of  $\mathcal{Q}$ . Then we pair the first element on the list with the number 1 from  $\mathcal{N}$ , the second element on the list with the number 2 from  $\mathcal{N}$ , and so on. We must ensure that every member of  $\mathcal{Q}$  appears only once on the list.

To get this list, we make an infinite matrix containing all the positive rational numbers, as shown in Figure 4.16. The  $i$ th row contains all numbers with numerator  $i$  and the  $j$ th column has all numbers with denominator  $j$ . So the number  $\frac{i}{j}$  occurs in the  $i$ th row and  $j$ th column.

Now we turn this matrix into a list. One (bad) way to attempt it would be to begin the list with all the elements in the first row. That isn't a good approach because the first row is infinite, so the list would never get to the second row. Instead we list the elements on the diagonals, which are superimposed on the diagram, starting from the corner. The first diagonal contains the single element  $\frac{1}{1}$ , and the second diagonal contains the two elements  $\frac{2}{1}$  and  $\frac{1}{2}$ . So the first three elements on the list are  $\frac{1}{1}$ ,  $\frac{2}{1}$ , and  $\frac{1}{2}$ . In the third diagonal, a complication arises. It contains  $\frac{3}{1}$ ,  $\frac{2}{2}$ , and  $\frac{1}{3}$ . If we simply added these to the list, we would repeat  $\frac{1}{1} = \frac{2}{2}$ . We avoid doing so by skipping an element when it would cause a repetition. So we add only the two new elements  $\frac{3}{1}$  and  $\frac{1}{3}$ . Continuing in this way, we obtain a list of all the elements of  $\mathcal{Q}$ .



**FIGURE 4.16**

A correspondence of  $\mathcal{N}$  and  $\mathcal{Q}$

After seeing the correspondence of  $\mathcal{N}$  and  $\mathcal{Q}$ , you might think that any two infinite sets can be shown to have the same size. After all, you need only demonstrate a correspondence, and this example shows that surprising correspondences do exist. However, for some infinite sets, no correspondence with  $\mathcal{N}$  exists. These sets are simply too big. Such sets are called *uncountable*.

The set of real numbers is an example of an uncountable set. A *real number* is one that has a decimal representation. The numbers  $\pi = 3.1415926\dots$  and

$\sqrt{2} = 1.4142135\dots$  are examples of real numbers. Let  $\mathcal{R}$  be the set of real numbers. Cantor proved that  $\mathcal{R}$  is uncountable. In doing so, he introduced the diagonalization method.

**THEOREM 4.17** .....

$\mathcal{R}$  is uncountable.

**PROOF** In order to show that  $\mathcal{R}$  is uncountable, we show that no correspondence exists between  $\mathcal{N}$  and  $\mathcal{R}$ . The proof is by contradiction. Suppose that a correspondence  $f$  existed between  $\mathcal{N}$  and  $\mathcal{R}$ . Our job is to show that  $f$  fails to work as it should. For it to be a correspondence,  $f$  must pair all the members of  $\mathcal{N}$  with all the members of  $\mathcal{R}$ . But we will find an  $x$  in  $\mathcal{R}$  that is not paired with anything in  $\mathcal{N}$ , which will be our contradiction.

The way we find this  $x$  is by actually constructing it. We choose each digit of  $x$  to make  $x$  different from one of the real numbers that is paired with an element of  $\mathcal{N}$ . In the end, we are sure that  $x$  is different from any real number that is paired.

We can illustrate this idea by giving an example. Suppose that the correspondence  $f$  exists. Let  $f(1) = 3.14159\dots$ ,  $f(2) = 55.55555\dots$ ,  $f(3) = \dots$ , and so on, just to make up some values for  $f$ . Then  $f$  pairs the number 1 with  $3.14159\dots$ , the number 2 with  $55.55555\dots$ , and so on. The following table shows a few values of a hypothetical correspondence  $f$  between  $\mathcal{N}$  and  $\mathcal{R}$ .

$n$	$f(n)$
1	3.14159...
2	55.55555...
3	0.12345...
4	0.50000...
$\vdots$	$\vdots$

We construct the desired  $x$  by giving its decimal representation. It is a number between 0 and 1, so all its significant digits are fractional digits following the decimal point. Our objective is to ensure that  $x \neq f(n)$  for any  $n$ . To ensure that  $x \neq f(1)$ , we let the first digit of  $x$  be anything different from the first fractional digit 1 of  $f(1) = 3.\underline{1}4159\dots$ . Arbitrarily, we let it be 4. To ensure that  $x \neq f(2)$ , we let the second digit of  $x$  be anything different from the second fractional digit 5 of  $f(2) = 55.5\underline{5}5555\dots$ . Arbitrarily, we let it be 6. The third fractional digit of  $f(3) = 0.12\underline{3}45\dots$  is 3, so we let  $x$  be anything different—say, 4. Continuing in this way down the diagonal of the table for  $f$ , we obtain all the digits of  $x$ , as shown in the following table. We know that  $x$  is not  $f(n)$  for any  $n$  because it differs from  $f(n)$  in the  $n$ th fractional digit. (A slight problem arises because certain numbers, such as  $0.1999\dots$  and  $0.2000\dots$ , are equal even though their decimal representations are different. We avoid this problem by never selecting the digits 0 or 9 when we construct  $x$ .)

$n$	$f(n)$	
1	3. <u>1</u> 4159...	$x = 0.4641 \dots$
2	55.5 <u>5</u> 555...	
3	0.12 <u>3</u> 45...	
4	0.500 <u>0</u> 0...	
$\vdots$	$\vdots$	

The preceding theorem has an important application to the theory of computation. It shows that some languages are not decidable or even Turing-recognizable, for the reason that there are uncountably many languages yet only countably many Turing machines. Because each Turing machine can recognize a single language and there are more languages than Turing machines, some languages are not recognized by any Turing machine. Such languages are not Turing-recognizable, as we state in the following corollary.

#### COROLLARY 4.18

Some languages are not Turing-recognizable.

**PROOF** To show that the set of all Turing machines is countable, we first observe that the set of all strings  $\Sigma^*$  is countable for any alphabet  $\Sigma$ . With only finitely many strings of each length, we may form a list of  $\Sigma^*$  by writing down all strings of length 0, length 1, length 2, and so on.

The set of all Turing machines is countable because each Turing machine  $M$  has an encoding into a string  $\langle M \rangle$ . If we simply omit those strings that are not legal encodings of Turing machines, we can obtain a list of all Turing machines.

To show that the set of all languages is uncountable, we first observe that the set of all infinite binary sequences is uncountable. An *infinite binary sequence* is an unending sequence of 0s and 1s. Let  $\mathcal{B}$  be the set of all infinite binary sequences. We can show that  $\mathcal{B}$  is uncountable by using a proof by diagonalization similar to the one we used in Theorem 4.17 to show that  $\mathcal{R}$  is uncountable.

Let  $\mathcal{L}$  be the set of all languages over alphabet  $\Sigma$ . We show that  $\mathcal{L}$  is uncountable by giving a correspondence with  $\mathcal{B}$ , thus showing that the two sets are the same size. Let  $\Sigma^* = \{s_1, s_2, s_3, \dots\}$ . Each language  $A \in \mathcal{L}$  has a unique sequence in  $\mathcal{B}$ . The  $i$ th bit of that sequence is a 1 if  $s_i \in A$  and is a 0 if  $s_i \notin A$ , which is called the *characteristic sequence* of  $A$ . For example, if  $A$  were the language of all strings starting with a 0 over the alphabet  $\{0,1\}$ , its characteristic sequence  $\chi_A$  would be

$$\begin{array}{lcl} \Sigma^* = \{ & \varepsilon, & 0, \quad 1, \quad 00, \quad 01, \quad 10, \quad 11, \quad 000, \quad 001, \quad \dots \} ; \\ A = \{ & & 0, \quad \quad 00, \quad 01, \quad \quad \quad 000, \quad 001, \quad \dots \} ; \\ \chi_A = & 0 & 1 \quad 0 \quad 1 \quad 1 \quad 0 \quad 0 \quad 1 \quad 1 \quad \dots \end{array}$$

The function  $f: \mathcal{L} \rightarrow \mathcal{B}$ , where  $f(A)$  equals the characteristic sequence of  $A$ , is one-to-one and onto, and hence is a correspondence. Therefore, as  $\mathcal{B}$  is uncountable,  $\mathcal{L}$  is uncountable as well.

Thus we have shown that the set of all languages cannot be put into a correspondence with the set of all Turing machines. We conclude that some languages are not recognized by any Turing machine.

---

## AN UNDECIDABLE LANGUAGE

Now we are ready to prove Theorem 4.11, the undecidability of the language

$$A_{\text{TM}} = \{\langle M, w \rangle \mid M \text{ is a TM and } M \text{ accepts } w\}.$$

**PROOF** We assume that  $A_{\text{TM}}$  is decidable and obtain a contradiction. Suppose that  $H$  is a decider for  $A_{\text{TM}}$ . On input  $\langle M, w \rangle$ , where  $M$  is a TM and  $w$  is a string,  $H$  halts and accepts if  $M$  accepts  $w$ . Furthermore,  $H$  halts and rejects if  $M$  fails to accept  $w$ . In other words, we assume that  $H$  is a TM, where

$$H(\langle M, w \rangle) = \begin{cases} \text{accept} & \text{if } M \text{ accepts } w \\ \text{reject} & \text{if } M \text{ does not accept } w. \end{cases}$$

Now we construct a new Turing machine  $D$  with  $H$  as a subroutine. This new TM calls  $H$  to determine what  $M$  does when the input to  $M$  is its own description  $\langle M \rangle$ . Once  $D$  has determined this information, it does the opposite. That is, it rejects if  $M$  accepts and accepts if  $M$  does not accept. The following is a description of  $D$ .

$D =$  “On input  $\langle M \rangle$ , where  $M$  is a TM:

1. Run  $H$  on input  $\langle M, \langle M \rangle \rangle$ .
2. Output the opposite of what  $H$  outputs. That is, if  $H$  accepts, *reject*; and if  $H$  rejects, *accept*.”

Don’t be confused by the notion of running a machine on its own description! That is similar to running a program with itself as input, something that does occasionally occur in practice. For example, a compiler is a program that translates other programs. A compiler for the language Python may itself be written in Python, so running that program on itself would make sense. In summary,

$$D(\langle M \rangle) = \begin{cases} \text{accept} & \text{if } M \text{ does not accept } \langle M \rangle \\ \text{reject} & \text{if } M \text{ accepts } \langle M \rangle. \end{cases}$$

What happens when we run  $D$  with its own description  $\langle D \rangle$  as input? In that case, we get

$$D(\langle D \rangle) = \begin{cases} \text{accept} & \text{if } D \text{ does not accept } \langle D \rangle \\ \text{reject} & \text{if } D \text{ accepts } \langle D \rangle. \end{cases}$$

No matter what  $D$  does, it is forced to do the opposite, which is obviously a contradiction. Thus, neither TM  $D$  nor TM  $H$  can exist.

---

Let's review the steps of this proof. Assume that a TM  $H$  decides  $A_{TM}$ . Use  $H$  to build a TM  $D$  that takes an input  $\langle M \rangle$ , where  $D$  accepts its input  $\langle M \rangle$  exactly when  $M$  does not accept its input  $\langle M \rangle$ . Finally, run  $D$  on itself. Thus, the machines take the following actions, with the last line being the contradiction.

- $H$  accepts  $\langle M, w \rangle$  exactly when  $M$  accepts  $w$ .
- $D$  rejects  $\langle M \rangle$  exactly when  $M$  accepts  $\langle M \rangle$ .
- $D$  rejects  $\langle D \rangle$  exactly when  $D$  accepts  $\langle D \rangle$ .

Where is the diagonalization in the proof of Theorem 4.11? It becomes apparent when you examine tables of behavior for TMs  $H$  and  $D$ . In these tables we list all TMs down the rows,  $M_1, M_2, \dots$ , and all their descriptions across the columns,  $\langle M_1 \rangle, \langle M_2 \rangle, \dots$ . The entries tell whether the machine in a given row accepts the input in a given column. The entry is *accept* if the machine accepts the input but is blank if it rejects or loops on that input. We made up the entries in the following figure to illustrate the idea.

	$\langle M_1 \rangle$	$\langle M_2 \rangle$	$\langle M_3 \rangle$	$\langle M_4 \rangle$	$\dots$
$M_1$	<i>accept</i>		<i>accept</i>		
$M_2$	<i>accept</i>	<i>accept</i>	<i>accept</i>	<i>accept</i>	
$M_3$					$\dots$
$M_4$	<i>accept</i>	<i>accept</i>			
$\vdots$			$\vdots$		

**FIGURE 4.19**

Entry  $i, j$  is *accept* if  $M_i$  accepts  $\langle M_j \rangle$

In the following figure, the entries are the results of running  $H$  on inputs corresponding to Figure 4.19. So if  $M_3$  does not accept input  $\langle M_2 \rangle$ , the entry for row  $M_3$  and column  $\langle M_2 \rangle$  is *reject* because  $H$  rejects input  $\langle M_3, \langle M_2 \rangle \rangle$ .

	$\langle M_1 \rangle$	$\langle M_2 \rangle$	$\langle M_3 \rangle$	$\langle M_4 \rangle$	$\dots$
$M_1$	<i>accept</i>	<i>reject</i>	<i>accept</i>	<i>reject</i>	
$M_2$	<i>accept</i>	<i>accept</i>	<i>accept</i>	<i>accept</i>	$\dots$
$M_3$	<i>reject</i>	<i>reject</i>	<i>reject</i>	<i>reject</i>	
$M_4$	<i>accept</i>	<i>accept</i>	<i>reject</i>	<i>reject</i>	
$\vdots$			$\vdots$		

**FIGURE 4.20**

Entry  $i, j$  is the value of  $H$  on input  $\langle M_i, \langle M_j \rangle \rangle$



In the following figure, we added  $D$  to Figure 4.20. By our assumption,  $H$  is a TM and so is  $D$ . Therefore, it must occur on the list  $M_1, M_2, \dots$  of all TMs. Note that  $D$  computes the opposite of the diagonal entries. The contradiction occurs at the point of the question mark where the entry must be the opposite of itself.

	$\langle M_1 \rangle$	$\langle M_2 \rangle$	$\langle M_3 \rangle$	$\langle M_4 \rangle$	$\dots$	$\langle D \rangle$	$\dots$
$M_1$	<u>accept</u>	reject	accept	reject		accept	
$M_2$	accept	<u>accept</u>	accept	accept	$\dots$	accept	$\dots$
$M_3$	reject	reject	<u>reject</u>	reject		reject	
$M_4$	accept	accept	reject	<u>reject</u>		accept	
$\vdots$			$\vdots$		$\ddots$		
$D$	reject	reject	accept	accept		<u>?</u>	
$\vdots$			$\vdots$				$\ddots$

**FIGURE 4.21**

If  $D$  is in the figure, a contradiction occurs at “?”

### A TURING-UNRECOGNIZABLE LANGUAGE

In the preceding section, we exhibited a language—namely,  $A_{\text{TM}}$ —that is undecidable. Now we exhibit a language that isn’t even Turing-recognizable. Note that  $A_{\text{TM}}$  will not suffice for this purpose because we showed that  $A_{\text{TM}}$  is Turing-recognizable (page 202). The following theorem shows that if both a language and its complement are Turing-recognizable, the language is decidable. Hence for any undecidable language, either it or its complement is not Turing-recognizable. Recall that the complement of a language is the language consisting of all strings that are not in the language. We say that a language is **co-Turing-recognizable** if it is the complement of a Turing-recognizable language.

### THEOREM 4.22

A language is decidable iff it is Turing-recognizable and co-Turing-recognizable.

In other words, a language is decidable exactly when both it and its complement are Turing-recognizable.

**PROOF** We have two directions to prove. First, if  $A$  is decidable, we can easily see that both  $A$  and its complement  $\bar{A}$  are Turing-recognizable. Any decidable language is Turing-recognizable, and the complement of a decidable language also is decidable.

## 7.4

## NP-COMPLETENESS

One important advance on the P versus NP question came in the early 1970s with the work of Stephen Cook and Leonid Levin. They discovered certain problems in NP whose individual complexity is related to that of the entire class. If a polynomial time algorithm exists for any of these problems, all problems in NP would be polynomial time solvable. These problems are called **NP-complete**. The phenomenon of NP-completeness is important for both theoretical and practical reasons.

On the theoretical side, a researcher trying to show that P is unequal to NP may focus on an NP-complete problem. If any problem in NP requires more than polynomial time, an NP-complete one does. Furthermore, a researcher attempting to prove that P equals NP only needs to find a polynomial time algorithm for an NP-complete problem to achieve this goal.

On the practical side, the phenomenon of NP-completeness may prevent wasting time searching for a nonexistent polynomial time algorithm to solve a particular problem. Even though we may not have the necessary mathematics to prove that the problem is unsolvable in polynomial time, we believe that P is unequal to NP. So proving that a problem is NP-complete is strong evidence of its nonpolynomiality.

The first NP-complete problem that we present is called the **satisfiability problem**. Recall that variables that can take on the values TRUE and FALSE are called **Boolean variables** (see Section 0.2). Usually, we represent TRUE by 1 and FALSE by 0. The **Boolean operations** AND, OR, and NOT, represented by the symbols  $\wedge$ ,  $\vee$ , and  $\neg$ , respectively, are described in the following list. We use the overbar as a shorthand for the  $\neg$  symbol, so  $\bar{x}$  means  $\neg x$ .

$$\begin{array}{lll} 0 \wedge 0 = 0 & 0 \vee 0 = 0 & \bar{0} = 1 \\ 0 \wedge 1 = 0 & 0 \vee 1 = 1 & \bar{1} = 0 \\ 1 \wedge 0 = 0 & 1 \vee 0 = 1 & \\ 1 \wedge 1 = 1 & 1 \vee 1 = 1 & \end{array}$$

A **Boolean formula** is an expression involving Boolean variables and operations. For example,

$$\phi = (\bar{x} \wedge y) \vee (x \wedge \bar{z})$$

is a Boolean formula. A Boolean formula is **satisfiable** if some assignment of 0s and 1s to the variables makes the formula evaluate to 1. The preceding formula is satisfiable because the assignment  $x = 0$ ,  $y = 1$ , and  $z = 0$  makes  $\phi$  evaluate to 1. We say the assignment *satisfies*  $\phi$ . The **satisfiability problem** is to test whether a Boolean formula is satisfiable. Let

$$SAT = \{\langle \phi \rangle \mid \phi \text{ is a satisfiable Boolean formula}\}.$$

Now we state a theorem that links the complexity of the SAT problem to the complexities of all problems in NP.

**THEOREM 7.27**

$SAT \in P$  iff  $P = NP$ .

Next, we develop the method that is central to the proof of this theorem.

**POLYNOMIAL TIME REDUCIBILITY**

In Chapter 5, we defined the concept of reducing one problem to another. When problem  $A$  reduces to problem  $B$ , a solution to  $B$  can be used to solve  $A$ . Now we define a version of reducibility that takes the efficiency of computation into account. When problem  $A$  is *efficiently* reducible to problem  $B$ , an efficient solution to  $B$  can be used to solve  $A$  efficiently.

**DEFINITION 7.28**

A function  $f: \Sigma^* \rightarrow \Sigma^*$  is a *polynomial time computable function* if some polynomial time Turing machine  $M$  exists that halts with just  $f(w)$  on its tape, when started on any input  $w$ .

**DEFINITION 7.29**

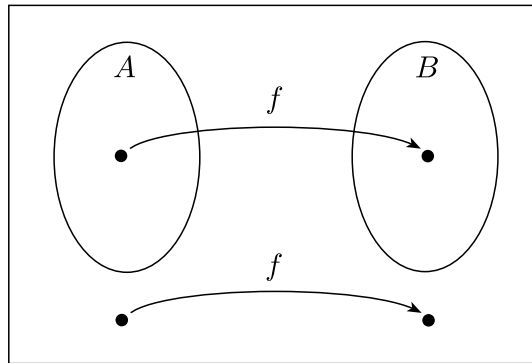
Language  $A$  is *polynomial time mapping reducible*,<sup>1</sup> or simply *polynomial time reducible*, to language  $B$ , written  $A \leq_P B$ , if a polynomial time computable function  $f: \Sigma^* \rightarrow \Sigma^*$  exists, where for every  $w$ ,

$$w \in A \iff f(w) \in B.$$

The function  $f$  is called the *polynomial time reduction* of  $A$  to  $B$ .

Polynomial time reducibility is the efficient analog to mapping reducibility as defined in Section 5.3. Other forms of efficient reducibility are available, but polynomial time reducibility is a simple form that is adequate for our purposes so we won't discuss the others here. Figure 7.30 illustrates polynomial time reducibility.

<sup>1</sup>It is called *polynomial time many-one reducibility* in some other textbooks.



**FIGURE 7.30**  
Polynomial time function  $f$  reducing  $A$  to  $B$

As with an ordinary mapping reduction, a polynomial time reduction of  $A$  to  $B$  provides a way to convert membership testing in  $A$  to membership testing in  $B$ —but now the conversion is done efficiently. To test whether  $w \in A$ , we use the reduction  $f$  to map  $w$  to  $f(w)$  and test whether  $f(w) \in B$ .

If one language is polynomial time reducible to a language already known to have a polynomial time solution, we obtain a polynomial time solution to the original language, as in the following theorem.

**THEOREM 7.31** .....

If  $A \leq_P B$  and  $B \in P$ , then  $A \in P$ .

**PROOF** Let  $M$  be the polynomial time algorithm deciding  $B$  and  $f$  be the polynomial time reduction from  $A$  to  $B$ . We describe a polynomial time algorithm  $N$  deciding  $A$  as follows.

$N =$  “On input  $w$ :

1. Compute  $f(w)$ .
2. Run  $M$  on input  $f(w)$  and output whatever  $M$  outputs.”

We have  $w \in A$  whenever  $f(w) \in B$  because  $f$  is a reduction from  $A$  to  $B$ . Thus,  $M$  accepts  $f(w)$  whenever  $w \in A$ . Moreover,  $N$  runs in polynomial time because each of its two stages runs in polynomial time. Note that stage 2 runs in polynomial time because the composition of two polynomials is a polynomial.

.....

Before demonstrating a polynomial time reduction, we introduce  $3SAT$ , a special case of the satisfiability problem whereby all formulas are in a special

form. A **literal** is a Boolean variable or a negated Boolean variable, as in  $x$  or  $\bar{x}$ . A **clause** is several literals connected with  $\vee$ s, as in  $(x_1 \vee \bar{x}_2 \vee \bar{x}_3 \vee x_4)$ . A Boolean formula is in **conjunctive normal form**, called a **cnf-formula**, if it comprises several clauses connected with  $\wedge$ s, as in

$$(x_1 \vee \bar{x}_2 \vee \bar{x}_3 \vee x_4) \wedge (x_3 \vee \bar{x}_5 \vee x_6) \wedge (x_3 \vee \bar{x}_6).$$

It is a **3cnf-formula** if all the clauses have three literals, as in

$$(x_1 \vee \bar{x}_2 \vee \bar{x}_3) \wedge (x_3 \vee \bar{x}_5 \vee x_6) \wedge (x_3 \vee \bar{x}_6 \vee x_4) \wedge (x_4 \vee x_5 \vee x_6).$$

Let  $3SAT = \{\langle \phi \rangle \mid \phi \text{ is a satisfiable 3cnf-formula}\}$ . If an assignment satisfies a cnf-formula, each clause must contain at least one literal that evaluates to 1.

The following theorem presents a polynomial time reduction from the  $3SAT$  problem to the  $CLIQUE$  problem.

### THEOREM 7.32

$3SAT$  is polynomial time reducible to  $CLIQUE$ .

**PROOF IDEA** The polynomial time reduction  $f$  that we demonstrate from  $3SAT$  to  $CLIQUE$  converts formulas to graphs. In the constructed graphs, cliques of a specified size correspond to satisfying assignments of the formula. Structures within the graph are designed to mimic the behavior of the variables and clauses.

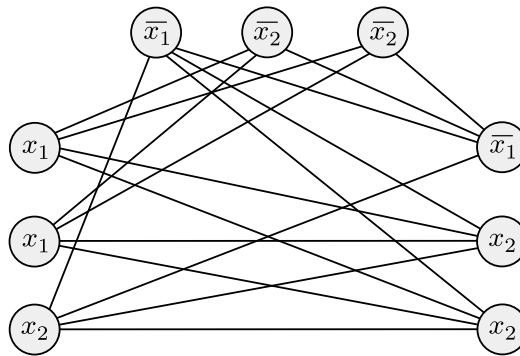
**PROOF** Let  $\phi$  be a formula with  $k$  clauses such as

$$\phi = (a_1 \vee b_1 \vee c_1) \wedge (a_2 \vee b_2 \vee c_2) \wedge \cdots \wedge (a_k \vee b_k \vee c_k).$$

The reduction  $f$  generates the string  $\langle G, k \rangle$ , where  $G$  is an undirected graph defined as follows.

The nodes in  $G$  are organized into  $k$  groups of three nodes each called the **triples**,  $t_1, \dots, t_k$ . Each triple corresponds to one of the clauses in  $\phi$ , and each node in a triple corresponds to a literal in the associated clause. Label each node of  $G$  with its corresponding literal in  $\phi$ .

The edges of  $G$  connect all but two types of pairs of nodes in  $G$ . No edge is present between nodes in the same triple, and no edge is present between two nodes with contradictory labels, as in  $x_2$  and  $\bar{x}_2$ . Figure 7.33 illustrates this construction when  $\phi = (x_1 \vee x_1 \vee x_2) \wedge (\bar{x}_1 \vee \bar{x}_2 \vee \bar{x}_2) \wedge (\bar{x}_1 \vee x_2 \vee x_2)$ .

**FIGURE 7.33**

The graph that the reduction produces from

$$\phi = (x_1 \vee x_1 \vee x_2) \wedge (\overline{x_1} \vee \overline{x_2} \vee \overline{x_2}) \wedge (\overline{x_1} \vee x_2 \vee x_2)$$

Now we demonstrate why this construction works. We show that  $\phi$  is satisfiable iff  $G$  has a  $k$ -clique.

Suppose that  $\phi$  has a satisfying assignment. In that satisfying assignment, at least one literal is true in every clause. In each triple of  $G$ , we select one node corresponding to a true literal in the satisfying assignment. If more than one literal is true in a particular clause, we choose one of the true literals arbitrarily. The nodes just selected form a  $k$ -clique. The number of nodes selected is  $k$  because we chose one for each of the  $k$  triples. Each pair of selected nodes is joined by an edge because no pair fits one of the exceptions described previously. They could not be from the same triple because we selected only one node per triple. They could not have contradictory labels because the associated literals were both true in the satisfying assignment. Therefore,  $G$  contains a  $k$ -clique.

Suppose that  $G$  has a  $k$ -clique. No two of the clique's nodes occur in the same triple because nodes in the same triple aren't connected by edges. Therefore, each of the  $k$  triples contains exactly one of the  $k$  clique nodes. We assign truth values to the variables of  $\phi$  so that each literal labeling a clique node is made true. Doing so is always possible because two nodes labeled in a contradictory way are not connected by an edge and hence both can't be in the clique. This assignment to the variables satisfies  $\phi$  because each triple contains a clique node and hence each clause contains a literal that is assigned TRUE. Therefore,  $\phi$  is satisfiable.

Theorems 7.31 and 7.32 tell us that if *CLIQUE* is solvable in polynomial time, so is *3SAT*. At first glance, this connection between these two problems appears quite remarkable because, superficially, they are rather different. But polynomial time reducibility allows us to link their complexities. Now we turn to a definition that will allow us similarly to link the complexities of an entire class of problems.

## DEFINITION OF NP-COMPLETENESS

**DEFINITION 7.34**

A language  $B$  is *NP-complete* if it satisfies two conditions:

1.  $B$  is in NP, and
2. every  $A$  in NP is polynomial time reducible to  $B$ .

**THEOREM 7.35**

If  $B$  is NP-complete and  $B \in P$ , then  $P = NP$ .

**PROOF** This theorem follows directly from the definition of polynomial time reducibility.

**THEOREM 7.36**

If  $B$  is NP-complete and  $B \leq_P C$  for  $C$  in NP, then  $C$  is NP-complete.

**PROOF** We already know that  $C$  is in NP, so we must show that every  $A$  in NP is polynomial time reducible to  $C$ . Because  $B$  is NP-complete, every language in NP is polynomial time reducible to  $B$ , and  $B$  in turn is polynomial time reducible to  $C$ . Polynomial time reductions compose; that is, if  $A$  is polynomial time reducible to  $B$  and  $B$  is polynomial time reducible to  $C$ , then  $A$  is polynomial time reducible to  $C$ . Hence every language in NP is polynomial time reducible to  $C$ .

**THE COOK–LEVIN THEOREM**

Once we have one NP-complete problem, we may obtain others by polynomial time reduction from it. However, establishing the first NP-complete problem is more difficult. Now we do so by proving that *SAT* is NP-complete.

**THEOREM 7.37**

*SAT* is NP-complete.<sup>2</sup>

This theorem implies Theorem 7.27.

<sup>2</sup>An alternative proof of this theorem appears in Section 9.3.

**PROOF IDEA** Showing that *SAT* is in NP is easy, and we do so shortly. The hard part of the proof is showing that any language in NP is polynomial time reducible to *SAT*.

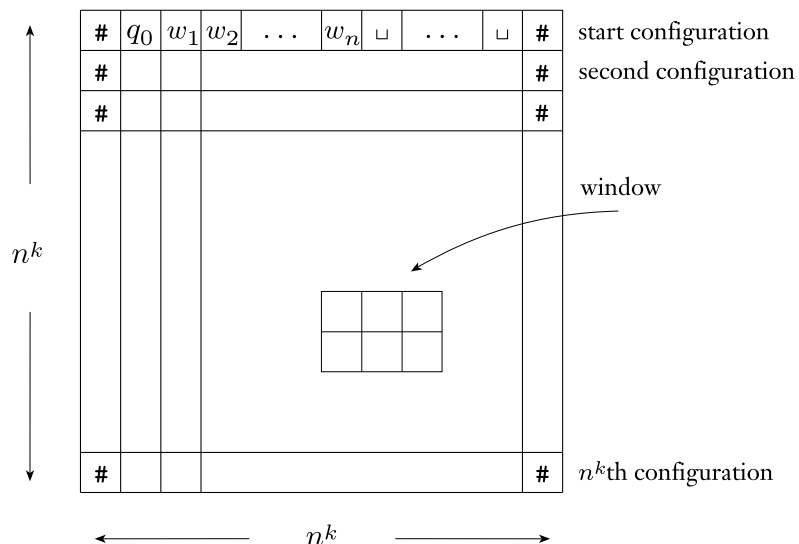
To do so, we construct a polynomial time reduction for each language  $A$  in NP to *SAT*. The reduction for  $A$  takes a string  $w$  and produces a Boolean formula  $\phi$  that simulates the NP machine for  $A$  on input  $w$ . If the machine accepts,  $\phi$  has a satisfying assignment that corresponds to the accepting computation. If the machine doesn't accept, no assignment satisfies  $\phi$ . Therefore,  $w$  is in  $A$  if and only if  $\phi$  is satisfiable.

Actually constructing the reduction to work in this way is a conceptually simple task, though we must cope with many details. A Boolean formula may contain the Boolean operations AND, OR, and NOT, and these operations form the basis for the circuitry used in electronic computers. Hence the fact that we can design a Boolean formula to simulate a Turing machine isn't surprising. The details are in the implementation of this idea.

**PROOF** First, we show that *SAT* is in NP. A nondeterministic polynomial time machine can guess an assignment to a given formula  $\phi$  and accept if the assignment satisfies  $\phi$ .

Next, we take any language  $A$  in NP and show that  $A$  is polynomial time reducible to *SAT*. Let  $N$  be a nondeterministic Turing machine that decides  $A$  in  $n^k$  time for some constant  $k$ . (For convenience, we actually assume that  $N$  runs in time  $n^k - 3$ ; but only those readers interested in details should worry about this minor point.) The following notion helps to describe the reduction.

A **tableau** for  $N$  on  $w$  is an  $n^k \times n^k$  table whose rows are the configurations of a branch of the computation of  $N$  on input  $w$ , as shown in the following figure.



**FIGURE 7.38**  
A tableau is an  $n^k \times n^k$  table of configurations



For convenience later, we assume that each configuration starts and ends with a # symbol. Therefore, the first and last columns of a tableau are all #s. The first row of the tableau is the starting configuration of  $N$  on  $w$ , and each row follows the previous one according to  $N$ 's transition function. A tableau is **accepting** if any row of the tableau is an accepting configuration.

Every accepting tableau for  $N$  on  $w$  corresponds to an accepting computation branch of  $N$  on  $w$ . Thus, the problem of determining whether  $N$  accepts  $w$  is equivalent to the problem of determining whether an accepting tableau for  $N$  on  $w$  exists.

Now we get to the description of the polynomial time reduction  $f$  from  $A$  to  $SAT$ . On input  $w$ , the reduction produces a formula  $\phi$ . We begin by describing the variables of  $\phi$ . Say that  $Q$  and  $\Gamma$  are the state set and tape alphabet of  $N$ , respectively. Let  $C = Q \cup \Gamma \cup \{\#\}$ . For each  $i$  and  $j$  between 1 and  $n^k$  and for each  $s$  in  $C$ , we have a variable,  $x_{i,j,s}$ .

Each of the  $(n^k)^2$  entries of a tableau is called a **cell**. The cell in row  $i$  and column  $j$  is called  $cell[i, j]$  and contains a symbol from  $C$ . We represent the contents of the cells with the variables of  $\phi$ . If  $x_{i,j,s}$  takes on the value 1, it means that  $cell[i, j]$  contains an  $s$ .

Now we design  $\phi$  so that a satisfying assignment to the variables does correspond to an accepting tableau for  $N$  on  $w$ . The formula  $\phi$  is the AND of four parts:  $\phi_{\text{cell}} \wedge \phi_{\text{start}} \wedge \phi_{\text{move}} \wedge \phi_{\text{accept}}$ . We describe each part in turn.

As we mentioned previously, turning variable  $x_{i,j,s}$  on corresponds to placing symbol  $s$  in  $cell[i, j]$ . The first thing we must guarantee in order to obtain a correspondence between an assignment and a tableau is that the assignment turns on exactly one variable for each cell. Formula  $\phi_{\text{cell}}$  ensures this requirement by expressing it in terms of Boolean operations:

$$\phi_{\text{cell}} = \bigwedge_{1 \leq i, j \leq n^k} \left[ \left( \bigvee_{s \in C} x_{i,j,s} \right) \wedge \left( \bigwedge_{\substack{s, t \in C \\ s \neq t}} (\overline{x_{i,j,s}} \vee \overline{x_{i,j,t}}) \right) \right].$$

The symbols  $\bigwedge$  and  $\bigvee$  stand for iterated AND and OR. For example, the expression in the preceding formula

$$\bigvee_{s \in C} x_{i,j,s}$$

is shorthand for

$$x_{i,j,s_1} \vee x_{i,j,s_2} \vee \cdots \vee x_{i,j,s_l}$$

where  $C = \{s_1, s_2, \dots, s_l\}$ . Hence  $\phi_{\text{cell}}$  is actually a large expression that contains a fragment for each cell in the tableau because  $i$  and  $j$  range from 1 to  $n^k$ . The first part of each fragment says that at least one variable is turned on in the corresponding cell. The second part of each fragment says that no more than one variable is turned on (literally, it says that in each pair of variables, at least one is turned off) in the corresponding cell. These fragments are connected by  $\wedge$  operations.

The first part of  $\phi_{\text{cell}}$  inside the brackets stipulates that at least one variable that is associated with each cell is on, whereas the second part stipulates that no more than one variable is on for each cell. Any assignment to the variables that satisfies  $\phi$  (and therefore  $\phi_{\text{cell}}$ ) must have exactly one variable on for every cell. Thus, any satisfying assignment specifies one symbol in each cell of the table. Parts  $\phi_{\text{start}}$ ,  $\phi_{\text{move}}$ , and  $\phi_{\text{accept}}$  ensure that these symbols actually correspond to an accepting tableau as follows.

Formula  $\phi_{\text{start}}$  ensures that the first row of the table is the starting configuration of  $N$  on  $w$  by explicitly stipulating that the corresponding variables are on:

$$\begin{aligned}\phi_{\text{start}} = & x_{1,1,\#} \wedge x_{1,2,q_0} \wedge \\ & x_{1,3,w_1} \wedge x_{1,4,w_2} \wedge \dots \wedge x_{1,n+2,w_n} \wedge \\ & x_{1,n+3,\sqcup} \wedge \dots \wedge x_{1,n^k-1,\sqcup} \wedge x_{1,n^k,\#}.\end{aligned}$$

Formula  $\phi_{\text{accept}}$  guarantees that an accepting configuration occurs in the tableau. It ensures that  $q_{\text{accept}}$ , the symbol for the accept state, appears in one of the cells of the tableau by stipulating that one of the corresponding variables is on:

$$\phi_{\text{accept}} = \bigvee_{1 \leq i, j \leq n^k} x_{i,j,q_{\text{accept}}}.$$

Finally, formula  $\phi_{\text{move}}$  guarantees that each row of the tableau corresponds to a configuration that legally follows the preceding row's configuration according to  $N$ 's rules. It does so by ensuring that each  $2 \times 3$  window of cells is legal. We say that a  $2 \times 3$  window is **legal** if that window does not violate the actions specified by  $N$ 's transition function. In other words, a window is legal if it might appear when one configuration correctly follows another.<sup>3</sup>

For example, say that  $a$ ,  $b$ , and  $c$  are members of the tape alphabet, and  $q_1$  and  $q_2$  are states of  $N$ . Assume that when in state  $q_1$  with the head reading an  $a$ ,  $N$  writes a  $b$ , stays in state  $q_1$ , and moves right; and that when in state  $q_1$  with the head reading a  $b$ ,  $N$  nondeterministically either

1. writes a  $c$ , enters  $q_2$ , and moves to the left, or
2. writes an  $a$ , enters  $q_2$ , and moves to the right.

Expressed formally,  $\delta(q_1, a) = \{(q_1, b, R)\}$  and  $\delta(q_1, b) = \{(q_2, c, L), (q_2, a, R)\}$ . Examples of legal windows for this machine are shown in Figure 7.39.

<sup>3</sup>We could give a precise definition of **legal window** here, in terms of the transition function. But doing so is quite tedious and would distract us from the main thrust of the proof argument. Anyone desiring more precision should refer to the related analysis in the proof of Theorem 5.15, the undecidability of the Post Correspondence Problem.

(a) 

a	$q_1$	b
$q_2$	a	c

(b) 

a	$q_1$	b
a	a	$q_2$

(c) 

a	a	$q_1$
a	a	b

(d) 

#	b	a
#	b	a

(e) 

a	b	a
a	b	$q_2$

(f) 

b	b	b
c	b	b

**FIGURE 7.39**  
Examples of legal windows

In Figure 7.39, windows (a) and (b) are legal because the transition function allows  $N$  to move in the indicated way. Window (c) is legal because, with  $q_1$  appearing on the right side of the top row, we don't know what symbol the head is over. That symbol could be an  $a$ , and  $q_1$  might change it to a  $b$  and move to the right. That possibility would give rise to this window, so it doesn't violate  $N$ 's rules. Window (d) is obviously legal because the top and bottom are identical, which would occur if the head weren't adjacent to the location of the window. Note that  $\#$  may appear on the left or right of both the top and bottom rows in a legal window. Window (e) is legal because state  $q_1$  reading a  $b$  might have been immediately to the right of the top row, and it would then have moved to the left in state  $q_2$  to appear on the right-hand end of the bottom row. Finally, window (f) is legal because state  $q_1$  might have been immediately to the left of the top row, and it might have changed the  $b$  to a  $c$  and moved to the left.

The windows shown in the following figure aren't legal for machine  $N$ .

(a)

a	b	a
a	a	a

(b)

a	$q_1$	b
$q_2$	a	a

(c)

b	$q_1$	b
$q_2$	b	$q_2$

**FIGURE 7.40**  
Examples of illegal windows

In window (a), the central symbol in the top row can't change because a state wasn't adjacent to it. Window (b) isn't legal because the transition function specifies that the  $b$  gets changed to a  $c$  but not to an  $a$ . Window (c) isn't legal because two states appear in the bottom row.

**CLAIM 7.41** .....

If the top row of the tableau is the start configuration and every window in the tableau is legal, each row of the tableau is a configuration that legally follows the preceding one.

We prove this claim by considering any two adjacent configurations in the tableau, called the upper configuration and the lower configuration. In the upper configuration, every cell that contains a tape symbol and isn't adjacent to a state symbol is the center top cell in a window whose top row contains no states. Therefore, that symbol must appear unchanged in the center bottom of the window. Hence it appears in the same position in the bottom configuration.

The window containing the state symbol in the center top cell guarantees that the corresponding three positions are updated consistently with the transition function. Therefore, if the upper configuration is a legal configuration, so is the lower configuration, and the lower one follows the upper one according to  $N$ 's rules. Note that this proof, though straightforward, depends crucially on our choice of a  $2 \times 3$  window size, as Problem 7.41 shows.

Now we return to the construction of  $\phi_{\text{move}}$ . It stipulates that all the windows in the tableau are legal. Each window contains six cells, which may be set in a fixed number of ways to yield a legal window. Formula  $\phi_{\text{move}}$  says that the settings of those six cells must be one of these ways, or

$$\phi_{\text{move}} = \bigwedge_{1 \leq i < n^k, 1 \leq j < n^k} (\text{the } (i, j)\text{-window is legal}).$$

The  $(i, j)$ -window has  $\text{cell}[i, j]$  as the upper central position. We replace the text “the  $(i, j)$ -window is legal” in this formula with the following formula. We write the contents of six cells of a window as  $a_1, \dots, a_6$ .

$$\bigvee_{\substack{a_1, \dots, a_6 \\ \text{is a legal window}}} (x_{i,j-1,a_1} \wedge x_{i,j,a_2} \wedge x_{i,j+1,a_3} \wedge x_{i+1,j-1,a_4} \wedge x_{i+1,j,a_5} \wedge x_{i+1,j+1,a_6})$$

Next, we analyze the complexity of the reduction to show that it operates in polynomial time. To do so, we examine the size of  $\phi$ . First, we estimate the number of variables it has. Recall that the tableau is an  $n^k \times n^k$  table, so it contains  $n^{2k}$  cells. Each cell has  $l$  variables associated with it, where  $l$  is the number of symbols in  $C$ . Because  $l$  depends only on the TM  $N$  and not on the length of the input  $n$ , the total number of variables is  $O(n^{2k})$ .

We estimate the size of each of the parts of  $\phi$ . Formula  $\phi_{\text{cell}}$  contains a fixed-size fragment of the formula for each cell of the tableau, so its size is  $O(n^{2k})$ . Formula  $\phi_{\text{start}}$  has a fragment for each cell in the top row, so its size is  $O(n^k)$ . Formulas  $\phi_{\text{move}}$  and  $\phi_{\text{accept}}$  each contain a fixed-size fragment of the formula for each cell of the tableau, so their size is  $O(n^{2k})$ . Thus,  $\phi$ 's total size is  $O(n^{2k})$ . That bound is sufficient for our purposes because it shows that the size of  $\phi$  is polynomial in  $n$ . If it were more than polynomial, the reduction wouldn't have any chance of generating it in polynomial time. (Actually, our estimates are low by a factor of  $O(\log n)$  because each variable has indices that can range up to  $n^k$  and so may require  $O(\log n)$  symbols to write into the formula, but this additional factor doesn't change the polynomiality of the result.)

To see that we can generate the formula in polynomial time, observe its highly repetitive nature. Each component of the formula is composed of many nearly

identical fragments, which differ only at the indices in a simple way. Therefore, we may easily construct a reduction that produces  $\phi$  in polynomial time from the input  $w$ .

---

Thus, we have concluded the proof of the Cook–Levin theorem, showing that *SAT* is NP-complete. Showing the NP-completeness of other languages generally doesn’t require such a lengthy proof. Instead, NP-completeness can be proved with a polynomial time reduction from a language that is already known to be NP-complete. We can use *SAT* for this purpose; but using *3SAT*, the special case of *SAT* that we defined on page 302, is usually easier. Recall that the formulas in *3SAT* are in conjunctive normal form (cnf) with three literals per clause. First, we must show that *3SAT* itself is NP-complete. We prove this assertion as a corollary to Theorem 7.37.

#### COROLLARY 7.42

---

*3SAT* is NP-complete.

**PROOF** Obviously *3SAT* is in NP, so we only need to prove that all languages in NP reduce to *3SAT* in polynomial time. One way to do so is by showing that *SAT* polynomial time reduces to *3SAT*. Instead, we modify the proof of Theorem 7.37 so that it directly produces a formula in conjunctive normal form with three literals per clause.

Theorem 7.37 produces a formula that is already almost in conjunctive normal form. Formula  $\phi_{\text{cell}}$  is a big AND of subformulas, each of which contains a big OR and a big AND of ORs. Thus,  $\phi_{\text{cell}}$  is an AND of clauses and so is already in cnf. Formula  $\phi_{\text{start}}$  is a big AND of variables. Taking each of these variables to be a clause of size 1, we see that  $\phi_{\text{start}}$  is in cnf. Formula  $\phi_{\text{accept}}$  is a big OR of variables and is thus a single clause. Formula  $\phi_{\text{move}}$  is the only one that isn’t already in cnf, but we may easily convert it into a formula that is in cnf as follows.

Recall that  $\phi_{\text{move}}$  is a big AND of subformulas, each of which is an OR of ANDs that describes all possible legal windows. The distributive laws, as described in Chapter 0, state that we can replace an OR of ANDs with an equivalent AND of ORs. Doing so may significantly increase the size of each subformula, but it can only increase the total size of  $\phi_{\text{move}}$  by a constant factor because the size of each subformula depends only on  $N$ . The result is a formula that is in conjunctive normal form.

Now that we have written the formula in cnf, we convert it to one with three literals per clause. In each clause that currently has one or two literals, we replicate one of the literals until the total number is three. In each clause that has more than three literals, we split it into several clauses and add additional variables to preserve the satisfiability or unsatisfiability of the original.

For example, we replace clause  $(a_1 \vee a_2 \vee a_3 \vee a_4)$ , wherein each  $a_i$  is a literal, with the two-clause expression  $(a_1 \vee a_2 \vee z) \wedge (\bar{z} \vee a_3 \vee a_4)$ , wherein  $z$  is a new