

## 1 Definition of a Turing machine

Turing machines are an abstract model of computation. They provide a precise, formal definition of what it means for a function to be computable. Many other definitions of computation have been proposed over the years — for example, one could try to formalize precisely what it means to run a program in Java on a computer with an infinite amount of memory — but it turns out that all known definitions of computation agree on what is computable and what is not. The Turing Machine definition seems to be the simplest, which is why we present it here. The key features of the Turing machine model of computation are:

1. A finite amount of internal state.
2. An infinite amount of external data storage.
3. A program specified by a finite number of instructions in a predefined language.
4. Self-reference: the programming language is expressive enough to write an interpreter for its own programs.

Models of computation with these key features tend to be equivalent to Turing machines, in the sense that the distinction between computable and uncomputable functions is the same in all such models.

A Turing machine can be thought of as a finite state machine sitting on an infinitely long tape containing symbols from some finite alphabet  $\Sigma$ . Based on the symbol it's currently reading, and its current state, the Turing machine writes a new symbol in that location (possibly the same as the previous one), moves left or right or stays in place, and enters a new state. It may also decide to halt and, optionally, to output “yes” or “no” upon halting. The machine's *transition function* is the “program” that specifies each of these actions (overwriting the current symbol, moving left or right, entering a new state, optionally halting and outputting an answer) given the current state and the symbol the machine is currently reading.

**Definition 1.** A Turing machine is specified by a finite alphabet  $\Sigma$ , a finite set of states  $K$  with a special element  $s$  (the starting state), and a *transition function*  $\delta : K \times \Sigma \rightarrow (K \cup \{\text{halt, yes, no}\}) \times \Sigma \times \{\leftarrow, \rightarrow, -\}$ . It is assumed that  $\Sigma$ ,  $K$ ,  $\{\text{halt, yes, no}\}$ , and  $\{\leftarrow, \rightarrow, -\}$  are disjoint sets, and that  $\Sigma$  contains two special elements  $\triangleright, \sqcup$  representing the start and end of the tape, respectively. We require that for every  $q \in K$ , if  $\delta(q, \triangleright) = (p, \sigma, d)$  then  $\sigma = \triangleright$  and  $d \neq \leftarrow$ . In other words, the machine never tries to overwrite the leftmost symbol on its tape nor to move to the left of it.<sup>1</sup>

---

<sup>1</sup>One could instead define Turing machines as having a doubly-infinite tape, with the ability to move arbitrarily far both left and right. Our choice of a singly-infinite tape makes certain definitions more convenient, but does not limit the computational power of Turing machines. A Turing machine as defined here can easily simulate one with a doubly-infinite tape by using the even-numbered positions on its tape to simulate the non-negative positions on the doubly-infinite tape, and using the odd-numbered positions on its tape to simulate the negative positions on the doubly-infinite tape. This simulation requires small modifications to the transition diagram of the Turing machine. Essentially, every move on the doubly-infinite tape gets simulated with two consecutive moves on the single-infinite tape, and this requires a couple of extra “bridging states” that are used for passing through the middle location while doing one of these two-step moves. The interested reader may fill in the details.

Note that a Turing machine is not prevented from overwriting the rightmost symbol on its tape or moving to the right of it. In fact, this capability is necessary in order for Turing machines to perform computations that require more space than is given in their original input string.

Having defined the *specification* of a Turing machine, we must now pin down a definition of *how they operate*. This has been informally described above, but it's time to make it formal. That begins with formally defining the configuration of the Turing machine at any time (the state of its tape, as well as the machine's own state and its position on the tape) and the rules for how its configuration changes over time.

**Definition 2.** The set  $\Sigma^*$  is the set of all finite sequences of elements of  $\Sigma$ . When an element of  $\Sigma^*$  is denoted by a letter such as  $x$ , then the elements of the sequence  $x$  are denoted by  $x_0, x_1, x_2, \dots, x_{n-1}$ , where  $n$  is the length of  $x$ . The length of  $x$  is denoted by  $|x|$ .

A *configuration* of a Turing machine is an ordered triple  $(x, q, k) \in \Sigma^* \times K \times \mathbb{N}$ , where  $x$  denotes the string on the tape,  $q$  denotes the machine's current state, and  $k$  denotes the position of the machine on the tape. The string  $x$  is required to begin with  $\triangleright$  and end with  $\sqcup$ . The position  $k$  is required to satisfy  $0 \leq k < |x|$ .

If  $M$  is a Turing machine and  $(x, q, k)$  is its configuration at any point in time, then its configuration  $(x', q', k')$  at the following point in time is determined as follows. Let  $(p, \sigma, d) = \delta(q, x_k)$ . The string  $x'$  is obtained from  $x$  by changing  $x_k$  to  $\sigma$ , and also appending  $\sqcup$  to the end of  $x$ , if  $k = |x| - 1$ . The new state  $q'$  is equal to  $p$ , and the new position  $k'$  is equal to  $k - 1, k + 1$ , or  $k$  according to whether  $d$  is  $\leftarrow, \rightarrow$ , or  $-$ , respectively. We express this relation between  $(x, q, k)$  and  $(x', q', k')$  by writing  $(x, q, k) \xrightarrow{M} (x', q', k')$ .

A *computation* of a Turing machine is a sequence of configurations  $(x_i, q_i, k_i)$ , where  $i$  runs from 0 to  $T$  (allowing for the case  $T = \infty$ ) that satisfies:

- The machine starts in a valid starting configuration, meaning that  $q_0 = s$  and  $k_0 = 0$ .
- Each pair of consecutive configurations represents a valid transition, i.e. for  $0 \leq i < T$ , it is the case that  $(x_i, q_i, k_i) \xrightarrow{M} (x_{i+1}, q_{i+1}, k_{i+1})$ .
- If  $T = \infty$ , we say that the computation *does not halt*.
- If  $T < \infty$ , we require that  $q_T \in \{\text{halt}, \text{yes}, \text{no}\}$  and we say that the computation *halts*.

If  $q_T = \text{yes}$  (respectively,  $q_T = \text{no}$ ) we say that the computation outputs “yes” (respectively, outputs “no”). If  $q_T = \text{halt}$  then the output of the computation is defined to be the string obtained by removing  $\triangleright$  and  $\sqcup$  from the start and end of  $x_T$ . In all three cases, the output of the computation is denoted by  $M(x)$ , where  $x$  is the input, i.e. the string  $x_0$  without the initial  $\triangleright$  and final  $\sqcup$ . If the computation does not halt, then its output is undefined and we write  $M(x) = \nearrow$ .

---

The upshot of this discussion is that one must standardize on either a singly-infinite or doubly-infinite tape, for the sake of making a precise definition, but the choice has no effect on the computational power of the model that is eventually defined. As we go through the other parts of the definition of Turing machines, we will encounter many other examples of this phenomenon: details that must be standardized for the sake of precision, but where the precise choice of definition has no bearing on the distinction between computable and uncomputable.

## 2 Examples of Turing machines

**Example 1.** As our first example, let's construct a Turing machine that takes a binary string and appends 0 to the left side of the string. The machine has four states:  $s, r_0, r_1, \ell$ . State  $s$  is the starting state, in state  $r_0$  and  $r_1$  it is moving right and preparing to write a 0 or 1, respectively, and in state  $\ell$  it is moving left.

The state  $s$  will be used only for getting started: thus, we only need to define how the Turing machine behaves when reading  $\triangleright$  in state  $s$ . The states  $r_0$  and  $r_1$  will be used, respectively, for writing 0 and writing 1 while remembering the overwritten symbol and moving to the right. Finally, state  $\ell$  is used for returning to the left side of the tape without changing its contents. This plain-English description of the Turing machine implies the following transition function. For brevity, we have omitted from the table the lines corresponding to pairs  $(q, \sigma)$  such that the Turing machine can't possibly be reading  $\sigma$  when it is in state  $q$ .

$q$	$\sigma$	$\delta(q, \sigma)$		
		state	symbol	direction
$s$	$\triangleright$	$r_0$	$\triangleright$	$\rightarrow$
$r_0$	0	$r_0$	0	$\rightarrow$
$r_0$	1	$r_1$	0	$\rightarrow$
$r_0$	$\sqcup$	$\ell$	0	$\leftarrow$
$r_1$	0	$r_0$	1	$\rightarrow$
$r_1$	1	$r_1$	1	$\rightarrow$
$r_1$	$\sqcup$	$\ell$	1	$\leftarrow$
$\ell$	0	$\ell$	0	$\leftarrow$
$\ell$	1	$\ell$	1	$\leftarrow$
$\ell$	$\triangleright$	halt	$\triangleright$	—

**Example 2.** Using similar ideas, we can design a Turing machine that takes a binary integer  $n$  (with the digits written in order, from the most significant digits on the left to the least significant digits on the right) and outputs the binary representation of its successor, i.e. the number  $n + 1$ .

It is easy to see that the following rule takes the binary representation of  $n$  and outputs the binary representation of  $n + 1$ . Find the rightmost occurrence of the digit 0 in the binary representation of  $n$ , change this digit to 1, and change every digit to the right of it from 1 to 0. The only exception is if the binary representation of  $n$  does not contain the digit 0; in that case, one should change every digit from 1 to 0 and prepend the digit 1.

$q$	$\sigma$	$\delta(q, \sigma)$		
		state	symbol	direction
$s$	$\triangleright$	$r$	$\triangleright$	$\rightarrow$
$r$	0	$r$	0	$\rightarrow$
$r$	1	$r$	1	$\rightarrow$
$r$	$\sqcup$	$\ell$	$\sqcup$	$\leftarrow$
$\ell$	0	$t$	1	$\leftarrow$
$\ell$	1	$\ell$	0	$\leftarrow$
$\ell$	$\triangleright$	prepend:: $s$	$\triangleright$	—
$t$	0	$t$	0	$\leftarrow$
$t$	1	$t$	1	$\leftarrow$
$t$	$\triangleright$	halt	$\triangleright$	—
prepend:: $s$	$\triangleright$	prepend:: $s$	$\triangleright$	$\rightarrow$
prepend:: $s$	0	prepend:: $r$	1	$\rightarrow$
prepend:: $r$	0	prepend:: $r$	0	$\rightarrow$
prepend:: $r$	$\sqcup$	prepend:: $\ell$	0	$\leftarrow$
prepend:: $\ell$	0	prepend:: $\ell$	0	$\leftarrow$
prepend:: $\ell$	1	prepend:: $\ell$	1	$\leftarrow$
prepend:: $\ell$	$\triangleright$	halt	$\triangleright$	—

Thus, the Turing machine for computing the binary successor function works as follows: it uses one state  $r$  to initially scan from left to right without modifying any of the digits, until it encounters the symbol  $\sqcup$ . At that point it changes into a new state  $\ell$  in which it moves to the left, changing any 1's that it encounters to 0's, until the first time that it encounters a symbol other than 1. (This may happen before encountering any 1's.) If that symbol is 0, it changes it to 1 and enters a new state  $t$  that moves leftward to  $\triangleright$  and then halts. On the other hand, if the symbol is  $\triangleright$ , then the original input consisted exclusively of 1's. In that case, it prepends

a 1 to the input using a subroutine very similar to Example 1. We'll refer to the states in that subroutine as `prepend::s`, `prepend::r`, `prepend::l`.

**Example 3.** We can compute the binary predecessor function in roughly the same way. We must take care to define the value of the predecessor function when its input is the number 0, since we haven't yet specified how negative numbers should be represented on a Turing machine's tape using the alphabet  $\{\triangleright, \sqcup, 0, 1\}$ . Rather than specify a convention for representing negative numbers, we will simply define the value of the predecessor function to be 0 when its input is 0. Also, for convenience, we will allow the output of the binary predecessor function to have any number of initial copies of the digit 0.

$q$	$\sigma$	$\delta(q, \sigma)$		
		state	symbol	direction
$s$	$\triangleright$	$z$	$\triangleright$	$\rightarrow$
$z$	0	$z$	0	$\rightarrow$
$z$	1	$r$	1	$\rightarrow$
$z$	$\sqcup$	$t$	$\sqcup$	$\leftarrow$
$r$	0	$r$	0	$\rightarrow$
$r$	1	$r$	1	$\rightarrow$
$r$	$\sqcup$	$\ell$	$\sqcup$	$\leftarrow$
$\ell$	0	$\ell$	1	$\leftarrow$
$\ell$	1	$t$	0	$\leftarrow$
$t$	0	$t$	0	$\leftarrow$
$t$	1	$t$	1	$\leftarrow$
$t$	$\triangleright$	halt	$\triangleright$	—

Our program to compute the binary predecessor of  $n$  thus begins with a test to see if  $n$  is equal to 0. The machine moves to the right, remaining in state  $z$  unless it sees the digit 1. If it reaches the end of its input in state  $z$ , then it simply rewinds to the beginning of the tape. Otherwise, it decrements the input in much the same way that the preceding example incremented it: in this case we use state  $\ell$  to change the trailing 0's to 1's until we encounter the rightmost 1, and then we enter state  $t$  to rewind to the beginning of the tape.

## 2.1 Pseudocode for Turing machines

Transforming even simple algorithms into Turing machine state transition diagrams can be an unbearably cumbersome process. It is even more unbearable to read a state transition diagram and try to deduce a plain-English description of what the algorithm accomplishes. It is desirable to express higher-level specifications of Turing machines using pseudocode.

Of course, the whole point of defining Turing machines is to get away from the *informal* notion of “algorithm” as exemplified by pseudocode. In this section we will discuss a straightforward, direct technique for transforming a certain type of pseudocode into Turing machine state diagrams. The pseudocode to which this transformation applies must satisfy the following restrictions:

1. The program has a finite, predefined number of variables.
2. Each variable can take only a finite, predefined set of possible values.
3. Recursive function calls are not allowed.

The word “predefined” above means, “specified by the program, not depending on the input.” Note that this is really quite restrictive. For example, our pseudocode is not allowed to use pointers to locations on the tape. (Such variables could not be limited to a predefined finite set of values.) Integer-valued variables are similarly off limits, unless there are predefined lower and upper bounds on the values of said variables.

Given pseudocode that meets the above restrictions, it is straightforward to transform it into an actual Turing machine. Suppose that the pseudocode has  $L$  lines and  $n$  variables  $v_1, \dots, v_n$ . For  $1 \leq i \leq n$  let  $V_i$  denote the predefined finite set of values for  $v_i$ . The state set of the

Turing machine is  $\{1, \dots, L\} \times V_1 \times \dots \times V_n$ . A state of the machine thus designates the line number that is currently being executed in the pseudocode, along with the current value of each variable. State transitions change the program counter in predictable ways (usually incrementing it, except when processing an if-then-else statement or reaching the end of a loop) and they may also change the values of  $v_1, \dots, v_n$  in ways specified by the pseudocode.

Note that the limitation to variables with a predefined finite set of values is forced on us by the requirement that the Turing machine should have a finite set of states. A variable with an unbounded number of possible values cannot be represented using a state set of some predefined, finite size. More subtly, a recursive algorithm must maintain a stack of program counters; if this stack can potentially grow to unbounded depth, then once again it is impossible to represent it using a finite state set.

Despite these syntactic limitations on pseudocode, we can still use Turing machines to execute algorithms that have recursion, integer-valued variables, pointers, arrays, and so on, because we can store any data we want on the Turing machine's tape, which has infinite capacity. We merely need to be careful to specify how this data is stored and retrieved. If necessary, this section of the notes could have been expanded to include standardized conventions for using the tape to store stacks of program counters, integer-valued variables, pointers to locations on the tape, or any other desired feature. We would then be permitted to use those more advanced features of pseudocode in the remainder of these notes. However, for our purposes it suffices to work with the very limited form of pseudocode that obeys the restrictions (1)-(3) listed above.

### 3 Universal Turing machines

The key property of Turing machines, and all other equivalent models of computation, is *universality*: there is a single Turing machine  $U$  that is capable of simulating any other Turing machine — even those with vastly more states than  $U$ . In other words, one can think of  $U$  as a “Turing machine interpreter,” written in the language of Turing machines. This capability for self-reference (the language of Turing machines is expressive enough to write an interpreter for itself) is the source of the surprising versatility of Turing machines and other models of computation. It is also the Pandora's Box that allows us to come up with undecidable problems, as we shall see in the following section.

#### 3.1 Describing Turing machines using strings

To define a universal Turing machine, we must first explain what it means to give a “description” of one Turing machine as the input to another one. For example, we must explain how a single Turing machine with bounded alphabet size can read the description of a Turing machine with a much larger alphabet.

To do so, we will make the following assumptions. For a Turing machine  $M$  with alphabet  $\Sigma$  and state set  $K$ , let

$$\ell = \lceil \log_2(|\Sigma| + |K| + 6) \rceil.$$

We will assume that each element of  $\Sigma \cup K \cup \{\text{halt}, \text{yes}, \text{no}\} \cup \{\leftarrow, \rightarrow, -\}$  is identified with a distinct binary string in  $\{0, 1\}^\ell$ . For example, we could always assume that  $\Sigma$  consists of the numbers  $0, 1, 2, \dots, |\Sigma| - 1$  represented as  $\ell$ -bit binary numbers (with initial 0's prepended, if necessary, to make the binary representation exactly  $\ell$  digits long), that  $K$  consists of the

numbers  $|\Sigma|, |\Sigma| + 1, \dots, |\Sigma| + |K| - 1$  (with  $|\Sigma|$  representing the starting state  $s$ ) and that  $\{\text{halt, yes, no}\} \cup \{\leftarrow, \rightarrow, -\}$  consists of the numbers  $|\Sigma| + |K|, \dots, |\Sigma| + |K| + 5$ .

Now, the description of Turing machine  $M$  is defined to be a finite string in the alphabet  $\{0, 1, '(', ')', ',', '\}\}$ , determined as follows.

$$M = m, n, \delta_1 \delta_2 \dots \delta_N$$

where  $m$  is the number  $|\Sigma|$  in binary,  $n$  is the number  $|K|$  in binary, and each of  $\delta_1, \dots, \delta_N$  is a string encoding one of the transition rules that make up the machine's transition function  $\delta$ . Each such rule is encoded using a string

$$\delta_i = (q, \sigma, p, \tau, d)$$

where each of the five parts  $q, \sigma, p, \tau, d$  is a string in  $\{0, 1\}^\ell$  encoding an element of  $\Sigma \cup K \cup \{\text{halt, yes, no}\} \cup \{\leftarrow, \rightarrow, -\}$  as described above. The presence of  $\delta_i$  in the description of  $M$  indicates that when  $M$  is in state  $q$  and reading symbol  $\sigma$ , then it transitions into state  $p$ , writes symbol  $\tau$ , and moves in direction  $d$ . In other words, the string  $\delta_i = (q, \sigma, p, \tau, d)$  in the description of  $M$  indicates that  $\delta(q, \sigma) = (p, \tau, d)$ .

### 3.2 Definition of a universal Turing machine

A universal Turing machine is a Turing machine  $U$  with alphabet  $\{0, 1, '(', ')', ',', '\}\}$ . It takes an input of the form  $\triangleright M; x \sqcup$ , where  $M$  is a valid description of a Turing machine and  $x$  is a string in the alphabet of  $M$ , encoded using  $\ell$ -bit blocks as described earlier. (If its input fails to match this specification, the universal Turing machine  $U$  is allowed to behave arbitrarily.)

The computation of  $U$ , given input  $\triangleright M; x \sqcup$ , has the same termination status — halting in state “halt”, “yes”, or “no”, or never halting — as the computation of  $M$  on input  $x$ . Furthermore, if  $M$  halts on input  $x$  (and, consequently,  $U$  halts on input  $\triangleright M; x \sqcup$ ) then the string on  $U$ 's tape at the time it halts is equal to the string on  $M$ 's tape at the time it halts, again translated into binary using  $\ell$ -bit blocks as specified above.

It is far from obvious that a universal Turing machine exists. In particular, such a machine must have a finite number of states, yet it must be able to simulate a computation performed by a Turing machine with a much greater number of states. In Section 3.4 we will describe how to construct a universal Turing machine. First, it is helpful to extend the definition of Turing machines to allow multiple tapes. After describing this extension we will indicate how a multi-tape Turing machine can be simulated by a single-tape machine (at the cost of a slower running time).

### 3.3 Multi-tape Turing machines

A multi-tape Turing machine with  $k$  tapes is defined in nearly the same way as a conventional single-tape Turing machine, except that it stores  $k$  strings simultaneously, maintains a position in each of the  $k$  strings, updates these  $k$  positions independently (i.e. it can move left in one string while moving right in another), and its state changes are based on the entire  $k$ -tuple of symbols that it reads at any point in time. More formally, a  $k$ -tape Turing machine has a finite alphabet  $\Sigma$  and state set  $K$  as before, but its transition function is

$$\delta : K \times \Sigma^k \rightarrow (K \cup \{\text{halt, yes, no}\}) \times (\Sigma \times \{\leftarrow, \rightarrow, -\})^k,$$

the difference being that the transition is determined by the entire  $k$ -tuple of symbols that it is reading, and that the instruction for the next action taken by the machine must include the symbol to be written, and the direction of motion, for *each* of the  $k$  strings.

The purpose of this section is to show that any multi-tape Turing machine  $M$  can be simulated with a single-tape Turing machine  $M^\circ$ . To store the contents of the  $k$  strings simultaneously, as well as the position of  $M$  within each of these strings, our single-tape machine  $M^\circ$  will have alphabet  $\Sigma^\circ = \Sigma^k \times \{0, 1\}^k$ . The  $k$ -tuple of symbols  $(\sigma_1, \dots, \sigma_k) \in \Sigma^k$  at the  $i^{\text{th}}$  location on the tape is interpreted to signify the symbols at the  $i^{\text{th}}$  location on each of the  $k$  tapes used by machine  $M$ . The  $k$ -tuple of binary values (“flags”)  $(p_1, \dots, p_k) \in \{0, 1\}^k$  at the  $i^{\text{th}}$  location on the tape represents whether the position of machine  $M$  on each of its  $k$  tapes is currently at location  $i$ . (1 = present, 0 = absent.) Our definition of Turing machine requires that the alphabet must contain two special symbols  $\triangleright, \sqcup$ . In the case of the alphabet  $\Sigma^\circ$ , the special symbol  $\triangleright$  is interpreted to be identical with  $(\triangleright)^k \times (0)^k$ , and  $\sqcup$  is interpreted to be identical with  $(\sqcup)^k \times (0)^k$ .

Machine  $M^\circ$  runs in a sequence of *simulation rounds*, each divided into  $k$  phases. The purpose of each round is to simulate one step of  $M$ , and the purpose of phase  $i$  in a round is to simulate what happens in the  $i^{\text{th}}$  string of  $M$  during that step. At the start of a phase,  $M^\circ$  is always at the left edge of its tape, on the symbol  $\triangleright$ . It makes a forward-and-back pass through its string, locating the flag that indicates the position of  $M$  on its  $i^{\text{th}}$  tape and updating that symbol (and possibly the adjacent one) to reflect the way that  $M$  updates its  $i^{\text{th}}$  string during the corresponding step of its execution.

The pseudocode presented in Algorithm 1 describes the simulation. See Section 2.1 above for a discussion of how to interpret this type of pseudocode as a precise specification of a Turing machine. Note that although the number  $k$  is not a specific finite number (such as 1024), it is still a “predefined, finite number” from the standpoint of interpreting this piece of pseudocode. The reason is that  $k$  depends only on the number of tapes of the machine  $M$  that we are simulating; it does not depend on the input to Algorithm 1.

---

**Algorithm 1** Simulation of multi-tape Turing machine  $M$  by single-tape Turing machine  $M^\circ$ .

---

```
1:  $q \leftarrow s$  // Initialize state.
2: repeat
3:   // Simulation round
4:   // First, find out what symbols  $M$  is seeing.
5:   repeat
6:     Move right without changing contents of string.
7:     for  $i = 1, \dots, k$  do
8:       if  $i^{\text{th}}$  flag at current location equals 1 then
9:          $\sigma_i \leftarrow i^{\text{th}}$  symbol at current location.
10:      end if
11:    end for
12:  until  $M^\circ$  reaches  $\sqcup$ 
13:  repeat
14:    Move left without changing contents of string.
15:  until  $M^\circ$  reaches  $\triangleright$ 
16:  // Now,  $\sigma_1, \dots, \sigma_k$  store the  $k$  symbols that  $M$  is seeing.
17:  // Evaluate state transition function for machine  $M$ .
18:   $q', (\sigma'_1, d_1), (\sigma'_2, d_2), \dots, (\sigma'_k, d_k) \leftarrow \delta(q, \sigma_1, \dots, \sigma_k)$ 
19:  for  $i = 1, \dots, k$  do
20:    // Phase  $i$ 
21:    repeat
22:      Move right without changing the contents of the string.
23:    until a location whose  $i^{\text{th}}$  flag is 1 is reached.
24:    Change the  $i^{\text{th}}$  symbol at this location to  $\sigma'_i$ .
25:    if  $d_i = \leftarrow$  then
26:      Change  $i^{\text{th}}$  flag to 0 at this location.
27:      Move one step left
28:      Change  $i^{\text{th}}$  flag to 1.
29:    else if  $d_i = \rightarrow$  then
30:      Change  $i^{\text{th}}$  flag to 0 at this location.
31:      Move one step right.
32:      Change  $i^{\text{th}}$  flag to 1.
33:    end if
34:    repeat
35:      Move left without changing the contents of the string.
36:    until the symbol  $\triangleright$  is reached.
37:  end for
38:   $q \leftarrow q'$ 
39: until  $q \in \{\text{halt}, \text{yes}, \text{no}\}$ 
40: // If the simulation reaches this line,  $q \in \{\text{halt}, \text{yes}, \text{no}\}$ .
41: if  $q = \text{yes}$  then
42:   Transition to “yes” state.
43: else if  $q = \text{no}$  then
44:   Transition to “no” state.
45: else
46:   Transition to “halt” state.
47: end if
```

---



### 3.4 Construction of a universal Turing machine

We now proceed to describe the construction of a universal Turing machine  $U$ . Taking advantage of Section 3.3, we can describe  $U$  as a 5-tape Turing machine; the existence of a single-tape universal Turing machine then follows from the general simulation presented in that section.

Our universal Turing machine has four tapes:

- the input tape: a read-only tape containing the input, which is never overwritten;
- the description tape: a tape containing the description of  $M$ , which is written once at initialization time and never overwritten afterwards;
- the working tape: a tape whose contents correspond to the contents of  $M$ 's tape, translated into  $\ell$ -bit blocks of binary symbols separated by commas, as the computation proceeds.
- the state tape: a tape describing the current state of  $M$ , encoded as an  $\ell$ -bit block of binary symbols.
- the special tape: a tape containing the binary encodings of the special states {halt,yes,no} and the “directional symbols”  $\{\leftarrow, \rightarrow, -\}$ .

The state tape solves the mystery of how a machine with a finite number of states can simulate a machine with many more states: it encodes these states using a tape that has the capacity to hold an unbounded number of symbols.

It would be too complicated to write down a diagram of the entire state transition function of a universal Turing machine, but we can describe it in plain English and pseudocode. The machine begins with an initialization phase in which it performs the following tasks:

- copy the description of  $M$  onto the description tape,
- copy the input string  $x$  onto the working tape, inserting commas between each  $\ell$ -bit block;
- copy the starting state of  $M$  onto the state tape;
- write the identifiers of the special states and directional symbols onto the special tape;
- move each cursor to the leftmost position on its respective tape.

After this initialization, the universal Turing machine executes its main loop. Each iteration of the main loop corresponds to one step in the computation executed by  $M$  on input  $x$ . At the start of any iteration of  $U$ 's main loop, the working tape and state tape contain the binary encodings of  $M$ 's tape contents and its state at the start of the corresponding step in  $M$ 's computation. Also, when  $U$  begins an iteration of its main loop, the cursor on each of its tapes except the working tape is at the leftmost position, and the cursor on the working tape is at the location corresponding to the position of  $M$ 's cursor at the start of the corresponding step in its computation. (In other words,  $U$ 's working tape cursor is pointing to the comma preceding the binary encoding of the symbol that  $M$ 's cursor is pointing to.)

The first step in an iteration of the main loop is to check whether it is time to halt. This is done by reading the contents of  $M$ 's state tape and comparing it to the binary encoding of the states “halt”, “yes”, and “no”, which are stored on the special tape. Assuming that  $M$  is not in one of the states “halt”, “yes”, “no”, it is time to simulate one step in the execution of  $M$ . This is done by working through the segment of the description tape that contains the strings  $\delta_1, \delta_2, \dots, \delta_N$  describing  $M$ 's transition function. The universal Turing machine moves through these strings in order from left to right. As it encounters each pair  $(q, \sigma)$ , it checks whether  $q$  is

identical to the  $\ell$ -bit string on its state tape and  $\sigma$  is identical to the  $\ell$ -bit string on its working tape. These comparisons are performed one bit at a time, and if either of the comparisons fails, then  $U$  rewinds its state tape cursor back to the  $\triangleright$  and it rewinds its working tape cursor back to the comma that marked its location at the start of this iteration of the main loop. It then moves its description tape cursor forward to the description of the next rule in  $M$ 's transition function. When it finally encounters a pair  $(q, \sigma)$  that matches its current state tape and working tape, then it moves forward to read the corresponding  $(p, \tau, d)$ , and it copies  $p$  onto the state tape,  $\tau$  onto the working tape, and then moves its working tape cursor in the direction specified by  $d$ . Finally, to end this iteration of the main loop, it rewinds the cursors on its description tape and state tape back to the leftmost position.

It is worth mentioning that the use of five tapes in the universal Turing machine is overkill. In particular, there is no need to save the input  $\triangleright M; x \sqcup$  on a separate read-only input tape. As we have seen, the input tape is never used after the end of the initialization stage. Thus, for example, we can skip the initialization step of copying the description of  $M$  from the input tape to the description tape; instead, after the initialization finishes, we can treat the input tape henceforward as if it were the description tape.

---

**Algorithm 2** Universal Turing machine, initialization.

---

```

1: // Copy the transition function of  $M$  from the input tape to the description tape.
2: Move right past the first and second commas on the input tape.
3: while not reading ';' on input tape do
4:   Read input tape symbol, copy to description tape, move right on both tapes.
5: end while
6: // Now, write the identifiers of the halting states and the "direction of motion symbols" on
   the special tape.
7: Move to start of input tape.
8: Using binary addition subroutine, write  $m + n$  in binary on working tape.
9: for  $i = 0, 1, 2, 3, 4, 5$  do
10:   On special tape, write binary representation of  $m + n + i$ , followed by ';'.
11:   // This doesn't require storing  $m + n$  in memory, because it's stored on the working tape.
12: end for
13: // Copy the input string  $x$  onto the working tape, inserting commas between each  $\ell$ -bit block.
14: Write  $m + n + 6$  in binary on the state tape. // In order to store the value of  $\ell$ .
15: Starting from left edge of input tape, move right until ';' is reached.
16: Move to left edge of working tape and state tape.
17: Move one step right on state tape.
18: repeat
19:   while state tape symbol is not  $\sqcup$  do
20:     Move right on input tape, working tape, and state tape.
21:     Copy symbol from input tape to working tape.
22:   end while
23:   Write ',' on working tape.
24:   Rewind to left edge of state tape, then move one step right.
25: until input tape symbol is  $\sqcup$ 
26: Copy  $m$  from input tape to state tape.
27: // Done with initialization!

```

---

---

**Algorithm 3** Universal Turing machine, main loop.

---

```
1: // Infinite loop. Each loop iteration simulates one step of  $M$ .
2: Move to the left edge of all tapes.
3: // First check if we need to halt.
4:  $match \leftarrow \text{TRUE}$ 
5: repeat
6:   Move right on special and state tapes.
7:   if symbols don't match then
8:      $match \leftarrow \text{FALSE}$ 
9:   end if
10: until reaching ',' on special tape
11: if  $match = \text{TRUE}$  then
12:   Enter "halt" state.
13: end if
14: Repeat the above steps two more times, with "yes", "no" in place of "halt".
15: Move back to left edge of description tape and state tape.
16: repeat
17:   Move right on description tape to next occurrence of '(' or  $\sqcup$ .
18:    $match \leftarrow \text{TRUE}$ 
19:   repeat
20:     Move right on description and state tapes.
21:     if symbols don't match then
22:        $match \leftarrow \text{FALSE}$ 
23:     end if
24:   until reaching ',' on description tape
25:   repeat
26:     Move right on description and working tapes.
27:     if symbols don't match then
28:        $match \leftarrow \text{FALSE}$ 
29:     end if
30:   until reaching ',' on description tape
31:   if  $match = \text{TRUE}$  then
32:     Move to left edge of state tape.
33:     Move left on working tape until ',' is reached.
34:   end if
35:   repeat
36:     Move right on description and state tapes.
37:     Copy symbol from description to state tape.
38:   until reaching ',' on description tape
39:   repeat
40:     Move right on description and working tapes.
41:     Copy symbol from description to working tape.
42:   until reaching ',' on description tape
43:   Move right on special tape past three occurrences of ','; stop at the fourth.
44:    $match \leftarrow \text{TRUE}$ 
45:   repeat
46:     Move right on description and special tapes.
47:     if symbols don't match then
48:        $match \leftarrow \text{FALSE}$ 
49:     end if
50:   until reaching ',' on special tape
51:   if  $match = \text{TRUE}$  then
52:     repeat
53:       Move left on working tape.
54:       until reaching ',' or  $\triangleright$ 
55:     end if
56:      $match \leftarrow \text{TRUE}$ 
57:     repeat
58:       Move right on description and special tapes.
59:       if symbols don't match then
60:          $match \leftarrow \text{FALSE}$ 
61:       end if
62:     until reaching ',' on special tape
63:     if  $match = \text{TRUE}$  then
64:       repeat
65:         Move right on working tape.
66:         until reaching ',' or  $\triangleright$ 
67:       end if
68:     until Reading  $\sqcup$  on description tape.
```

---

## 4 Undecidable problems

In this section we will see that there exist computational problems that are too difficult to be solved by any Turing machine. Since Turing machines are universal enough to represent any algorithm running on a deterministic computer, this means there are problems too difficult to be solved by any algorithm.

### 4.1 Definitions

To be precise about the notion of what we mean by “computational problems” and what it means for a Turing machine to “solve” a problem, we define problems in terms of *languages*, which correspond to decision problems with a yes/no answer. We define two notions of “solving” a problem specified by a language  $L$ . The first of these definitions (“deciding  $L$ ”) corresponds to what we usually mean when we speak of solving a computational problem, i.e. terminating and outputting a correct yes/no answer. The second definition (“accepting  $L$ ”) is a “one-sided” definition: if the answer is “yes”, the machine must halt and provide this answer after a finite amount of time; if the answer is “no”, the machine need not ever halt and provide this answer. Finally, we give some definitions that apply to computational problems where the goal is to output a string rather than just a simple yes/no answer.

**Definition 3.** Let  $\Sigma_0 = \Sigma \setminus \{\triangleright, \sqcup\}$ . A *language* is any set of strings  $L \subseteq \Sigma_0^*$ . Suppose  $M$  is a Turing machine and  $L$  is a language.

1.  $M$  *decides*  $L$  if every computation of  $M$  halts in the “yes” or “no” state, and  $L$  is the set of strings occurring in starting configurations that lead to the “yes” state.
2.  $L$  is *decidable* if there is a machine  $M$  that decides  $L$ .
3.  $M$  *accepts*  $L$  if  $L$  is the set of strings occurring in starting configurations that lead  $M$  to halt in the “yes” state.
4.  $L$  is *recursively enumerable* if there is a machine  $M$  that accepts  $L$ .
5.  $M$  *computes* a given function  $f : \Sigma_0^* \rightarrow \Sigma_0^*$  if every computation of  $M$  halts, and for all  $x \in \Sigma_0^*$ , the computation with starting configuration  $(\triangleright x \sqcup, s, 0)$  ends in configuration  $(\triangleright f(x) \sqcup \cdots \sqcup, \text{halt}, 0)$ , where  $\sqcup \cdots \sqcup$  denotes any sequence of one or more repetitions of the symbol  $\sqcup$ .
6.  $f$  is a *computable function* if there is a Turing machine  $M$  that computes  $f$ .

### 4.2 Undecidability via counting

One simple explanation for the existence of undecidable languages is via a counting argument: there are simply too many languages, and not enough Turing machines to decide them all! This can be formalized using the distinction between countable and uncountable sets.

**Definition 4.** An infinite set is *countable* if and only if there is a one-to-one correspondence between its elements and the natural numbers. Otherwise it is said to be *uncountable*.

**Lemma 1.** If  $\Sigma$  is a finite set then  $\Sigma^*$  is countable.

*Proof.* If  $|\Sigma| = 1$  then a string in  $\Sigma^*$  is uniquely determined by its length and this defines a one-to-one correspondence between  $\Sigma^*$  and the natural numbers. Otherwise, without loss of generality,  $\Sigma$  is equal to the set  $\{0, 1, \dots, b-1\}$  for some positive integer  $b > 1$ . Every natural number has an expansion in base  $b$  which is a finite-length string of elements of  $\Sigma$ . This gives a one-to-one correspondence between natural numbers and elements of  $\Sigma^*$  beginning with a non-zero element of  $\Sigma$ . To get a full one-to-one correspondence, we need one more trick: every positive integer can be uniquely represented in the form  $2^s \cdot (2t-1)$  where  $s$  and  $t$  are natural numbers and  $t > 0$ . We can map the positive integer  $2^s \cdot (2t+1)$  to the string consisting of  $s$  occurrences of 0 followed by the base- $b$  representation of  $t$ . This gives a one-to-one correspondence between positive natural numbers and nonempty strings in  $\Sigma^*$ . If we map 0 to the empty string, we have defined a full one-to-one correspondence.  $\square$

**Lemma 2.** *Let  $X$  be a countable set and let  $2^X$  denote the set of all subsets of  $X$ . The set  $2^X$  is uncountable.*

*Proof.* Consider any function  $f : X \rightarrow 2^X$ . We will construct an element  $T \in 2^X$  such that  $T$  is not equal to  $f(x)$  for any  $x \in X$ . This proves that there is no one-to-one correspondence between  $X$  and  $2^X$ ; hence  $2^X$  is uncountable.

Let  $x_0, x_1, x_2, \dots$  be a list of all the elements of  $X$ , indexed by the natural numbers. (Such a list exists because of our assumption that  $X$  is countable.) The construction of the set  $T$  is best explained via a diagram, in which the set  $f(x)$ , for every  $x \in X$  is represented by a row of 0's and 1's in an infinite table. This table has columns indexed by  $x_0, x_1, \dots$ , and the row labeled  $f(x_i)$  has a 0 in the column labeled  $x_j$  if and only if  $x_j$  belongs to the set  $f(x_i)$ .

	$x_0$	$x_1$	$x_2$	$x_3$	$\dots$
$f(x_0)$	0	1	0	0	$\dots$
$f(x_1)$	1	1	0	1	$\dots$
$f(x_2)$	0	0	1	0	$\dots$
$f(x_3)$	1	0	1	1	$\dots$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	

To construct the set  $T$ , we look at the main diagonal of this table, whose  $i$ -th entry specifies whether  $x_i \in f(x_i)$ , and we flip each bit. This produces a new sequence of 0's and 1's encoding a subset of  $X$ . In fact, the subset can be succinctly defined by

$$T = \{x \in X \mid x \notin f(x)\}. \quad (1)$$

We see that  $T$  cannot be equal to  $f(x)$  for any  $x \in X$ . Indeed, if  $T = f(x)$  and  $x \in T$  then this contradicts the fact that  $x \notin f(x)$  for all  $x \in T$ ; similarly, if  $T = f(x)$  and  $x \notin T$  then this contradicts the fact that  $x \in f(x)$  for all  $x \notin T$ .  $\square$

**Remark 1.** Actually, the same proof technique shows that there is never a one-to-one correspondence between  $X$  and  $2^X$  for any set  $X$ . We can always define the “diagonal set”  $T$  via (1) and argue that the assumption  $T = f(x)$  leads to a contradiction. The idea of visualizing the function  $f$  using a two-dimensional table becomes more strained when the number of rows and columns of the table is uncountable, but this doesn't interfere with the validity of the argument based directly on defining  $T$  via (1).

**Theorem 3.** *For every alphabet  $\Sigma$  there is a language  $L \subseteq \Sigma^*$  that is not recursively enumerable.*

*Proof.* The set of languages  $L \subseteq \Sigma^*$  is uncountable by Lemmas 1 and 2. The set of Turing machines with alphabet  $\Sigma$  is countable because each such Turing machine has a description which is a finite-length string of symbols in the alphabet  $\{0, 1, '(', ')', ';'\}$ . Therefore there are strictly more languages than there are Turing machines, so there are languages that are not accepted by any Turing machine.  $\square$

### 4.3 Undecidability via diagonalization

The proof of Theorem 3 is quite unsatisfying because it does not provide any example of an interesting language that is not recursively enumerable. In this section we will repeat the “diagonal argument” from the proof of Theorem 2, this time in the context of Turing machines and languages, to obtain a more interesting example of a set that is not recursively enumerable.

**Definition 5.** For a Turing machine  $M$ , we define  $L(M)$  to be the set of all strings accepted by  $M$ :

$$L(M) = \{x \mid M \text{ halts and outputs “yes” on input } x\}.$$

Suppose  $\Sigma$  is a finite alphabet, and suppose we have specified a mapping  $\phi$  from  $\{0, 1, '(', ')', ';'\}$  to strings of some fixed length in  $\Sigma^*$ , so that each Turing machine has a description in  $\Sigma^*$  obtained by taking its standard description, using  $\phi$  to map each symbol to  $\Sigma^*$ , and concatenating the resulting sequence of strings. For every  $x \in \Sigma^*$  we will now define a language  $L(x) \subseteq \Sigma_0^*$  as follows. If  $x$  is the description of a Turing machine  $M$  then  $L(x) = L(M)$ ; otherwise,  $L(x) = \emptyset$ .

We are now in a position to repeat the diagonal construction from the proof of Theorem 2. Consider an infinite two-dimensional table whose rows and columns are indexed by elements of  $\Sigma_0^*$ .

	$x_0$	$x_1$	$x_2$	$x_3$	$\dots$
$L(x_0)$	0	1	0	0	$\dots$
$L(x_1)$	1	1	0	1	$\dots$
$L(x_2)$	0	0	1	0	$\dots$
$L(x_3)$	1	0	1	1	$\dots$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	

**Definition 6.** The *diagonal language*  $D \subseteq \Sigma_0^*$  is defined by

$$D = \{x \in \Sigma_0^* \mid x \notin L(x)\}.$$

**Theorem 4.** *The diagonal language  $D$  is not recursively enumerable.*

*Proof.* The proof is exactly the same as the proof of Theorem 2. Assume, by way of contradiction, that  $D = L(x)$  for some  $x$ . Either  $x \in D$  or  $x \notin D$  and we obtain a contradiction in both cases. If  $x \in D$  then this violates the fact that  $x \notin L(x)$  for all  $x \in D$ . If  $x \notin D$  then this violates the fact that  $x \in L(x)$  for all  $x \notin D$ .  $\square$

**Corollary 5.** *There exists a language  $L$  that is recursively enumerable but its complement is not.*

*Proof.* Let  $L$  be the complement of the diagonal language  $D$ . A Turing machine  $M$  that accepts  $L$  can be described as follows: given input  $x$ , construct the string  $x; x$  and run a universal Turing machine on this input. It is clear from the definition of  $L$  that  $L = L(M)$ , so  $L$  is recursively enumerable. We have already seen that its complement,  $D$ , is not recursively enumerable.  $\square$

**Remark 2.** Unlike Theorem 3, there is no way to obtain Corollary 5 using a simple counting argument.

**Corollary 6.** *There exists a language  $L$  that is recursively enumerable but not decidable.*

*Proof.* From the definition of a decidable language, it is clear that the complement of a decidable language is decidable. For the recursively enumerable language  $L$  in Corollary 5, we know that the complement of  $L$  is not even recursively enumerable (hence, *a fortiori*, also not decidable) and this implies that  $L$  is not decidable.  $\square$

## 4.4 The halting problem

Theorem 4 gave an explicit example of a language that is not recursively enumerable — hence not decidable — but it is still a rather unnatural example, so perhaps one might hope that every *interesting* computational problem is decidable. Unfortunately, this is not the case!

**Definition 7.** The *halting problem* is the problem of deciding whether a given Turing machine halts when presented with a given input. In other words, it is the language  $H$  defined by

$$H = \{M; x \mid M \text{ is a valid Turing machine description, and } M \text{ halts on input } x\}.$$

**Theorem 7.** *The halting problem is not decidable.*

*Proof.* We give a proof by contradiction. Given a Turing machine  $M_H$  that decides  $H$ , we will construct a Turing machine  $M_D$  that accepts  $D$ , contradicting Theorem 4. The machine  $M_D$  operates as follows. Given input  $x$ , it constructs the string  $x; x$  and runs  $M_H$  on this input until the step when  $M_H$  is just about to output “yes” or “no”. At that point in the computation, instead of outputting “yes” or “no”,  $M_D$  does the following. If  $M_H$  is about to output “no” then  $M_D$  instead outputs “yes”. If  $M_H$  is about to output “yes” then  $M_D$  instead runs a universal Turing machine  $U$  on input  $x; x$ . Just before  $U$  is about to halt, if it is about to output “yes”, then  $M_D$  instead outputs “no”. Otherwise  $M_D$  outputs “yes”. (Note that it is not possible for  $U$  to run forever without halting, because of our assumption that  $M_H$  decides the halting problem and that it outputs “yes” on input  $x; x$ .)

By construction, we see that  $M_D$  always halts and outputs “yes” or “no”, and that its output is “no” if and only if  $M_H(x; x) = \text{yes}$  and  $U(x; x) = \text{yes}$ . In other words,  $M_D(x) = \text{no}$  if and only if  $x$  is the description of a Turing machine that halts and outputs “yes” on input  $x$ . Recalling our definition of  $L(x)$ , we see that another way of saying this is as follows:  $M_D(x) = \text{no}$  if and only if  $x \in L(x)$ . We now have the following chain of “if and only if” statements:

$$x \in L(M_D) \iff M_D(x) = \text{yes} \iff M_D(x) \neq \text{no} \iff x \notin L(x) \iff x \in D.$$

We have proved that  $L(M_D) = D$ , i.e.  $M_D$  is a Turing machine that accepts  $D$ , as claimed.  $\square$

**Remark 3.** It is easy to see that  $H$  is recursively enumerable. In fact, if  $U$  is any universal Turing machine then  $H = L(U)$ .

## 4.5 Rice's Theorem

Theorem 7 shows that, unfortunately, there exist interesting languages that are not decidable. In particular, the question of whether a given Turing machine halts on a given input is not decidable. Unfortunately, the situation is much worse than this! Our next theorem shows that essentially any non-trivial property of Turing machines is undecidable. (Actually, this is an overstatement; the theorem will show that any non-trivial property of *the languages accepted by Turing machines* is undecidable. Non-trivial properties of the Turing machine itself — e.g., does it run for more than 100 steps when presented with the input string  $\triangleright\sqcup$  — may be decidable.)

**Definition 8.** Language  $\mathcal{L} \subseteq \Sigma_0^*$  is a *Turing machine I/O property* if it is the case that for all pairs  $x, y \in \Sigma^*$  such that  $L(x) = L(y)$ , we have  $x \in \mathcal{L} \Leftrightarrow y \in \mathcal{L}$ . We say that  $\mathcal{L}$  is *non-trivial* if both of the sets  $\mathcal{L}$  and  $\overline{\mathcal{L}} = \Sigma_0^* \setminus \mathcal{L}$  are nonempty.

**Theorem 8** (Rice's Theorem). *If  $\mathcal{L}$  is a non-trivial Turing machine I/O property, then  $\mathcal{L}$  is undecidable.*

*Proof.* Consider any  $x_1$  such that  $L(x_1) = \emptyset$ . We can assume without loss of generality that  $x_1 \notin \mathcal{L}$ . The reason is that a language is decidable if and only if its complement is decidable. Thus, if  $x_1 \in \mathcal{L}$  then we can replace  $\mathcal{L}$  with  $\overline{\mathcal{L}}$  and continue with the rest of proof. In the end we will have proven that  $\overline{\mathcal{L}}$  is undecidable from which it follows that  $\mathcal{L}$  is also undecidable.

Since  $\mathcal{L}$  is non-trivial, there is also a string  $x_0 \in \mathcal{L}$ . This  $x_0$  must be a description of a Turing machine  $M_0$ . (If  $x_0$  were not a Turing machine description, then it would be the case that  $L(x_0) = \emptyset = L(x_1)$ , and hence that  $x_0 \notin \mathcal{L}$  since  $x_1 \notin \mathcal{L}$  and  $\mathcal{L}$  is a Turing machine I/O property.)

We are now ready to prove Rice's Theorem by contradiction. Given a Turing machine  $M_{\mathcal{L}}$  that decides  $\mathcal{L}$  we will construct a Turing machine  $M_H$  that decides the halting problem, in contradiction to Theorem 7.

The construction of  $M_H$  is as follows. On input  $M; x$ , it transforms the pair  $M; x$  into the description of another Turing machine  $M^*$ , and then it feeds this description into  $M_{\mathcal{L}}$ . The definition of  $M^*$  is a little tricky. On input  $y$ , machine  $M^*$  does the following. First it runs  $M$  on input  $x$ , without overwriting the string  $y$ . If  $M$  ever halts, then instead of halting  $M^*$  enters the second phase of its execution, which consists of running  $M_0$  on  $y$ . (Recall that  $M_0$  is a Turing machine whose description  $x_0$  belongs to  $\mathcal{L}$ .) If  $M$  never halts, then  $M^*$  also never halts.

This completes the construction of  $M_H$ . To recap, when  $M_H$  is given an input  $M; x$  it first transforms the pair  $M; x$  into the description of a related Turing machine  $M^*$ , then it runs  $M_{\mathcal{L}}$  on the input consisting of the description of  $M^*$  and it outputs the same answer that  $M_{\mathcal{L}}$  outputs. There are two things we still have to prove.

1. The function that takes the string  $M; x$  and outputs the description of  $M^*$  is a computable function, i.e. there is a Turing machine that can transform  $M; x$  into the description of  $M^*$ .
2. Assuming  $M_{\mathcal{L}}$  decides  $\mathcal{L}$ , then  $M_H$  decides the halting problem.

The first of these facts is elementary but tedious.  $M^*$  needs to have a bunch of extra states that append a special symbol (say,  $\sharp$ ) to the end of its input and then write out the string  $x$  after the special symbol,  $\sharp$ . It also has the same states as  $M$  with the same transition function, with two modifications: first, this modified version of  $M$  treats the symbol  $\sharp$  exactly as if it were  $\triangleright$ . (This ensures that  $M$  will remain on the right side of the tape and will not modify the copy of  $y$  that sits to the left of the  $\sharp$  symbol.) Finally, whenever  $M$  would halt, the modified version



of  $M$  enters a special set of states that move left to the  $\sharp$  symbol, overwrite this symbol with  $\sqcup$ , continue moving left until the symbol  $\triangleright$  is reached, and then run the machine  $M_0$ .

Now let's prove the second fact — that  $M_H$  decides the halting problem, assuming  $M_{\mathcal{L}}$  decides  $\mathcal{L}$ . Suppose we run  $M_H$  on input  $M; x$ . If  $M$  does not halt on  $x$ , then the machine  $M^*$  constructed by  $M_H$  never halts on any input  $y$ . (This is because the first thing  $M^*$  does on input  $y$  is to simulate the computation of  $M$  on input  $x$ .) Thus, if  $M$  does not halt on  $x$  then  $L(M^*) = \emptyset = L(x_1)$ . Recalling that  $\mathcal{L}$  is a Turing machine I/O property and that  $x_1 \notin \mathcal{L}$ , this means that the description of  $M^*$  also does not belong to  $\mathcal{L}$ , so  $M_{\mathcal{L}}$  outputs “no,” which means  $M_H$  outputs “no” on input  $M; x$  as desired. On the other hand, if  $M$  halts on input  $x$ , then the machine  $M^*$  constructed by  $M_H$  behaves as follows on any input  $y$ : it first spends a finite amount of time running  $M$  on input  $x$ , then ignores the answer and runs  $M_0$  on  $y$ . This means that  $M^*$  accepts input  $y$  if and only if  $y \in L(M_0)$ , i.e.  $L(M^*) = L(M_0) = L(x_0)$ . Once again using our assumption that  $\mathcal{L}$  is a Turing machine I/O property, and recalling that  $x_0 \in \mathcal{L}$ , this means that  $M_{\mathcal{L}}$  outputs “yes” on input  $M^*$ , which means that  $M_H$  outputs “yes” on input  $M; x$ , as it should.  $\square$

## 5 Nondeterminism

Up until now, the Turing machines we have been discussing have all been *deterministic*, meaning that there is only one valid computation starting from any given input. Nondeterministic Turing machines are defined in the same way as their deterministic counterparts, except that instead of a *transition function* associating one and only one  $(p, \tau, d)$  to each (state, symbol) pair  $(q, \sigma)$ , there is a *transition relation* (which we will again denote by  $\delta$ ) consisting of any number of ordered 5-tuples  $(q, \sigma, p, \tau, d)$ . If  $(q, \sigma, p, \tau, d)$  belongs to  $\delta$ , it means that when observing symbol  $\sigma$  in state  $q$ , one of the allowed behaviors of the nondeterministic machine is to enter state  $p$ , write  $\tau$ , and move in direction  $d$ . If a nondeterministic machine has more than one allowed behavior in a given configuration, it is allowed to choose one of them arbitrarily. Thus, the relation  $(x, q, k) \xrightarrow{M} (x', q', k')$  is interpreted to mean that configuration  $(x', q', k')$  is obtained from  $(x, q, k)$  by applying any one of the allowed rules  $(q, x_k, p, \tau, d)$  associated to (state, symbol) pair  $(q, x_k)$  in the transition relation  $\delta$ . Given this interpretation of the relation  $\xrightarrow{M}$ , we define a computation of a nondeterministic Turing machine is exactly as in Definition 2.

**Definition 9.** If  $M$  is a nondeterministic Turing machine and  $x \in \Sigma_0^*$  is a string, we say that  $M$  *accepts*  $x$  if there exists a computation of  $M$  starting with input  $\triangleright x \sqcup$  that ends by halting in the “yes” state. The set of all strings accepted by  $M$  is denoted by  $L(M)$ .

This interpretation of “ $M$  accepts  $x$ ” illustrates the big advantage of nondeterministic Turing machines over deterministic ones. In effect, a nondeterministic Turing machine is able to try out many possible computations in parallel and to accept its input if any one of these computations accepts it. No physical computer could ever aspire to this sort of unbounded parallelism, and thus nondeterministic Turing machines are, in some sense, a pure abstraction that does not correspond to any physically realizable computer. (This is in contrast to deterministic Turing machines, which are intended as a model of physically realizable computation, despite the fact that they make idealized assumptions — such as an infinitely long tape — that are intended only as approximations to the reality of computers having a finite but effectively unlimited amount of storage.) However, nondeterminism is a useful abstraction in computer science for a couple of reasons.

1. As we will see in Section 5.1, if there is a nondeterministic Turing machine that accepts language  $L$  then there is also a deterministic Turing machine that accepts  $L$ . Thus, if one ignores running time, nondeterministic Turing machines have no more power than deterministic ones.
2. The question of whether nondeterministic computation can be simulated deterministically with only a *polynomial* increase in running time is the P vs. NP question, perhaps the deepest open question in computer science.
3. Nondeterminism can function as a useful abstraction for computation with an untrusted external source of advice: the nondeterministic machine's transitions are guided by the advice from the external source, but it remains in control of the decision whether to accept its input or not.

## 5.1 Nondeterministic Turing machines and deterministic verifiers

One useful way of looking at nondeterministic computation is by relating it to the notion of a *verifier* for a language.

**Definition 10.** Let  $\Sigma$  be a language not containing the symbol ‘;’. A *verifier* for a language  $L \subseteq \Sigma^*$  is a deterministic Turing machine with alphabet  $\Sigma \cup \{; \}$ , such that

$$L = \{x \mid \exists y \text{ s.t. } V(x; y) = \text{yes}\}.$$

We sometimes refer to the string  $y$  in the verifier's input as the *evidence*.

We say that  $V$  is a *polynomial-time verifier* for  $L$  if  $V$  is a verifier for  $L$  and there exists a polynomial function  $p$  such that for all  $x \in L$  there exists a  $y$  such that  $V(x; y)$  outputs “yes” after at most  $p(|x|)$  steps. (Note that if any such  $y$  exists, then there is one such  $y$  satisfying  $|y| \leq p(|x|)$ , since the running time bound prevents  $V$  from reading any symbols of  $y$  beyond the first  $p(|x|)$  symbols.)

The following theorem details the close relationship between nondeterministic Turing machines and deterministic verifiers.

**Theorem 9.** *The following three properties of a language  $L$  are equivalent.*

1.  $L$  is recursively enumerable, i.e. there exists a deterministic Turing machine that accepts  $L$ .
2. There exists a nondeterministic Turing machine that accepts  $L$ .
3. There exists a verifier for  $L$ .

*Proof.* A deterministic Turing machine that accepts  $L$  is a special case of a nondeterministic Turing machine that accepts  $L$ , so it is trivial that (1) implies (2). To see that (2) implies (3), suppose that  $M$  is a nondeterministic Turing machine that accepts  $L$ . The transition relation of  $M$  is a finite set of 5-tuples and we can number the elements of this set as  $1, 2, \dots, K$  for some  $K$ . The verifier's evidence will be a string  $y$  encoding a sequence of elements of  $\{1, \dots, K\}$ , the sequence of transitions that  $M$  undergoes in a computation that leads to accepting  $x$ . The verifier  $V$  operates as follows. Given input  $x; y$ , it simulates a computation of  $M$  on input  $x$ . In every step of the simulation,  $V$  consults the evidence string  $y$  to obtain the next transition rule. If it is not a valid transition rule for the current configuration (for example because it applies to a state other than the current state of  $M$  in the simulation) then  $V$  instantly outputs “no”,

otherwise it performs the indicated transition and moves to the next simulation step. Having defined the verifier in this way, it is clear that there exists  $y$  such that  $V(x; y) = \text{yes}$  if and only if there exists a computation of  $M$  that accepts  $x$ ; in other words,  $V$  is a verifier for  $L(M)$ , as desired.

To see that (3) implies (1), suppose that  $V$  is a verifier for  $L$ . The following algorithm describes a deterministic Turing machine that accepts  $L$ .

- 1: Let  $x$  denote the input string.
- 2: **for**  $k = 1$  to  $\infty$  **do**
- 3:   **for all**  $y \in \Sigma_0^*$  such that  $|y| \leq k$  **do**
- 4:     Simulate  $V(x; y)$ , artificially terminating after  $k$  steps unless  $V$  halts earlier.
- 5:     **if**  $V(x; y) = \text{yes}$  **then** output “yes.”
- 6:   **end for**
- 7: **end for**

Clearly, if this algorithm outputs “yes” then there exists a  $y$  such that  $V(x; y) = \text{yes}$  and therefore  $x \in L$ . Conversely, if  $x \in L$  then there exists a  $y$  such that  $V(x; y) = \text{yes}$ . Letting  $t$  denote the number of steps that  $V$  executes when processing input  $x; y$ , we see that our algorithm will output “yes” during the outer-loop iteration in which  $k = \max\{|y|, t\}$ .  $\square$

Recall that we defined the complexity class NP, earlier in the semester, to be the set of all languages that have a deterministic polynomial-time verifier. The following definition and theorem provide an alternate definition of NP in terms of nondeterministic polynomial-time computation.

**Definition 11.** A *polynomial-time nondeterministic Turing machine* is a nondeterministic Turing machine  $M$  such that for every input  $x \in \Sigma_0^*$  and every computation of  $M$  starting with input  $x$ , the computation halts after at most  $p(|x|)$  steps, where  $p$  is a polynomial function called the *running time* of  $M$ .

**Theorem 10.** A language  $L \subseteq \Sigma_0^*$  has a deterministic polynomial-time verifier if and only if there is a polynomial-time nondeterministic Turing machine that accepts  $L$ .

*Proof.* If  $L$  has a deterministic polynomial-time verifier  $V$ , and  $p$  is a polynomial such that for all  $x \in L$  there exists  $y \in \Sigma_0^*$  such that  $V(x; y)$  outputs “yes” in  $p(|x|)$  or fewer steps, then we can design a nondeterministic Turing machine  $M$  that accepts  $L$ , as follows. On input  $x$ ,  $M$  begins by moving to the end of the string  $x$ , writing ‘;’, moving  $p(|x|)$  steps further to the right, and then going into a special state during which it moves to the left and (nondeterministically) writes arbitrary symbols in  $\Sigma_0 \cup \{\sqcup\}$  until it encounters the symbol ‘;’. Upon encountering ‘;’ it moves left to  $\triangleright$  and then simulates the deterministic verifier  $V$ , artificially terminating the simulation after  $p(|x|)$  steps unless it ends earlier.  $M$  outputs “yes” if and only if  $V$  outputs “yes” in this simulation.

Conversely, if  $M$  is a polynomial-time nondeterministic Turing machine that accepts  $L$ , then we can construct a polynomial-time verifier for  $L$  using exactly the same construction that was used in the proof of Theorem 9 to transform a nondeterministic Turing machine into a verifier. The reader is invited to check that the verifier given by that transformation is a polynomial-time verifier, as long as  $M$  is a polynomial-time nondeterministic Turing machine.  $\square$

## 5.2 The Cook-Levin Theorem

In this section we prove that 3SAT is NP-complete. Technically, it is easier to work with the language CNF-SAT consisting of all satisfiable Boolean formulas in conjunctive normal form. In

other words, CNF-SAT is defined in exactly the same way as 3SAT except that a clause is allowed to contain any number of literals.

It is clear that CNF-SAT and 3SAT belong to NP: a verifier merely needs to take a proposed truth assignment and go through each clause, checking that at least one of its literals is satisfied by the proposed assignment.

There is an easy reduction from CNF-SAT to 3SAT. Suppose we are given a SAT instance with variables  $x_1, \dots, x_n$  and clauses  $C_1, \dots, C_m$ . For each clause  $C_j$  that is a disjunction of  $k(j) > 3$  literals, we let  $\ell_1, \dots, \ell_{k(j)}$  be the literals in  $C_j$  and we transform  $C_j$  into a conjunction of  $k(j) - 2$  clauses of size 3, as follows. First we create auxiliary variables  $z_{j,1}, \dots, z_{j,k(j)-3}$ , then we represent  $C_j$  using the following conjunction:

$$(\ell_1 \vee \ell_2 \vee z_{j,1}) \wedge (\overline{z_{j,1}} \vee \ell_3 \vee z_{j,2}) \wedge (\overline{z_{j,2}} \vee \ell_4 \vee z_{j,3}) \wedge (\overline{z_{j,3}} \vee \ell_5 \vee z_{j,4}) \wedge \dots \wedge (\overline{z_{j,k(j)-3}} \vee \ell_{k(j)-1} \vee \ell_{k(j)}).$$

It is an exercise to see that a given truth assignment of variables  $x_1, \dots, x_n$  satisfies  $C_j$  if and only if there exists a truth assignment of  $z_{j,1}, \dots, z_{j,k(j)-3}$  that, when combined with the given assignment of  $x_1, \dots, x_n$ , satisfies all of the clauses in the conjunction given above. Similarly, we can use auxiliary variables to replace each clause having  $k < 3$  literals with a conjunction of clauses having exactly three literals. Specifically, we transform a clause  $C_j$  containing a single literal  $\ell_1$  into the conjunction

$$(\ell_1 \vee z_j \vee z'_j) \wedge (\ell_1 \vee \overline{z_j} \vee z'_j) \wedge (\ell_1 \vee z_j \vee \overline{z'_j}) \wedge (\ell_1 \vee \overline{z_j} \vee \overline{z'_j}),$$

and we transform a clause  $C_j = (\ell_1 \vee \ell_2)$  into the conjunction

$$(\ell_1 \vee \ell_2 \vee z_j) \wedge (\ell_1 \vee \ell_2 \vee \overline{z_j}).$$

Once again it is an exercise to see that a given truth assignment of variables  $x_1, \dots, x_n$  satisfies  $C_j$  if and only if there exists a truth assignment of the auxiliary variables that, when combined with the given assignment of  $x_1, \dots, x_n$ , satisfies all of the clauses in the conjunctions given above.

To complete the proof of the Cook-Levin Theorem, we now show that every language in NP can be reduced, in polynomial time, to CNF-SAT.

**Theorem 11.** *If  $L \in \Sigma_0^*$  is in NP then there exists a polynomial-time reduction from  $L$  to CNF-SAT.*

*Proof.* Suppose  $V$  is a polynomial-time verifier for  $L$  and  $p$  is a polynomial function such that for every  $x \in L$  there exists  $y$  such that  $V(x; y)$  outputs “yes” after at most  $p(|x|)$  steps. The computation of  $V$  can be expressed using a rectangular table with  $p(|x|)$  rows and columns, with each row representing a configuration of  $V$ . Each entry of the table is an ordered pair in  $(K \cup \{*\}) \times \Sigma$ , where  $K$  is the state set of  $V$ . The meaning of an entry  $(q, \sigma)$  in the  $i^{\text{th}}$  row and  $j^{\text{th}}$  column of the table is that during step  $i$  of the computation,  $V$  was in state  $q$ , visiting location  $j$ , and reading symbol  $\sigma$ . If the entry is  $(*, \sigma)$  it means that  $\sigma$  was stored at that location during step  $i$ , but the location of  $V$  was elsewhere.

To reduce  $L$  to CNF-SAT, we use a set of variables representing the contents of this table and a set of clauses asserting that the table represents a valid computation of  $V$  that ends in the “yes” state. This is surprisingly easy to do. For each table entry  $(i, j)$  where  $0 \leq i, j \leq p(|x|)$ , and each pair  $(q, \sigma) \in (K \cup \{ast\}) \times \Sigma$ , we have a Boolean variable  $u_{i,j}^{q,\sigma}$ . The clauses are defined as follows.

1. **The verifier starts in state  $s$  on the left side of its tape.** This is represented by a clause consisting of a single literal,  $u_{0,0}^{s,\triangleright}$ .

2. **The input to the verifier is of the form  $x; y$ .** For  $1 \leq j \leq |x| + 1$ , if  $\sigma$  is the  $j^{\text{th}}$  symbol of the string “ $x$ ,” there is a clause consisting of the single literal  $u_{0,j}^{*,\sigma}$ .
3. **There is exactly one entry in each position of the table.** For each pair  $(i, j)$  with  $0 \leq i, j \leq p(|x|)$ , there is one clause  $\bigvee_{q,\sigma} u_{i,j}^{q,\sigma}$  asserting that at least one entry occurs in position  $i, j$ , and a set of clauses  $\bar{u}_{i,j}^{q,\sigma} \vee \bar{u}_{i,j}^{q',\sigma'}$  asserting that for any distinct pairs  $(q, \sigma)$  and  $(q', \sigma')$ , at least one of them *does not* occur in position  $i, j$ .
4. **The table entries obey the transition rules of the verifier.** For each 5-tuple  $(q, \sigma, p, \tau, d)$  such that  $\delta(q, \sigma) = (p, \tau, d)$ , and each position  $i, j$  such that  $0 \leq i \leq p(|x|)$ ,  $0 \leq j < p(|x|)$ , there is a set of clauses ensuring that if the verifier is in state  $q$ , reading symbol  $\sigma$ , in location  $j$ , at time  $i$ , then its state and location at time  $i + 1$ , as well as the tape contents, are consistent with the specified transition rule. For example, if  $\delta(q, \sigma) = (p, \tau, \rightarrow)$  then the statement, “If, at time  $i$ , the verifier is in state  $q$ , reading symbol  $\sigma$ , at location  $j$ , then the symbol  $\tau$  should occur at location  $j$  at time  $i + 1$  and the verifier should be absent from that location,” is expressed by the Boolean formula  $u_{i,j}^{q,\sigma} \Rightarrow u_{i+1,j}^{*,\tau}$  which is equivalent to the disjunction  $\bar{u}_{i,j}^{q,\sigma} \vee u_{i+1,j}^{*,\tau}$ . The statement, “If, at time  $i$ , the verifier is in state  $q$ , reading symbol  $\sigma$ , at location  $j$ , and the symbol  $\sigma'$  occurs at location  $j + 1$ , then at time  $i + 1$  the verifier should be in state  $p$  reading symbol  $\sigma'$  at location  $j + 1$ ,” is expressed by the Boolean formula  $(u_{i,j}^{q,\sigma} \wedge u_{i,j+1}^{*,\sigma'}) \Rightarrow u_{i+1,j+1}^{p,\sigma'}$  which is equivalent to the disjunction  $\bar{u}_{i,j}^{q,\sigma} \vee \bar{u}_{i,j+1}^{*,\sigma'} \vee u_{i+1,j+1}^{p,\sigma'}$ . Note that we need one such clause for every possible neighboring symbol  $\sigma'$ .
5. **The verifier enters the “yes” state during the computation.** This is expressed by a single disjunction  $\bigvee_{\sigma,i,j} u_{i,j}^{\text{yes},\sigma}$ .

The reader may check that the total number of clauses is  $O(p(|x|)^2)$ , where the  $O(\cdot)$  masks constants that depend on the size of the alphabet  $\Sigma$  and of the verifier’s state set.  $\square$