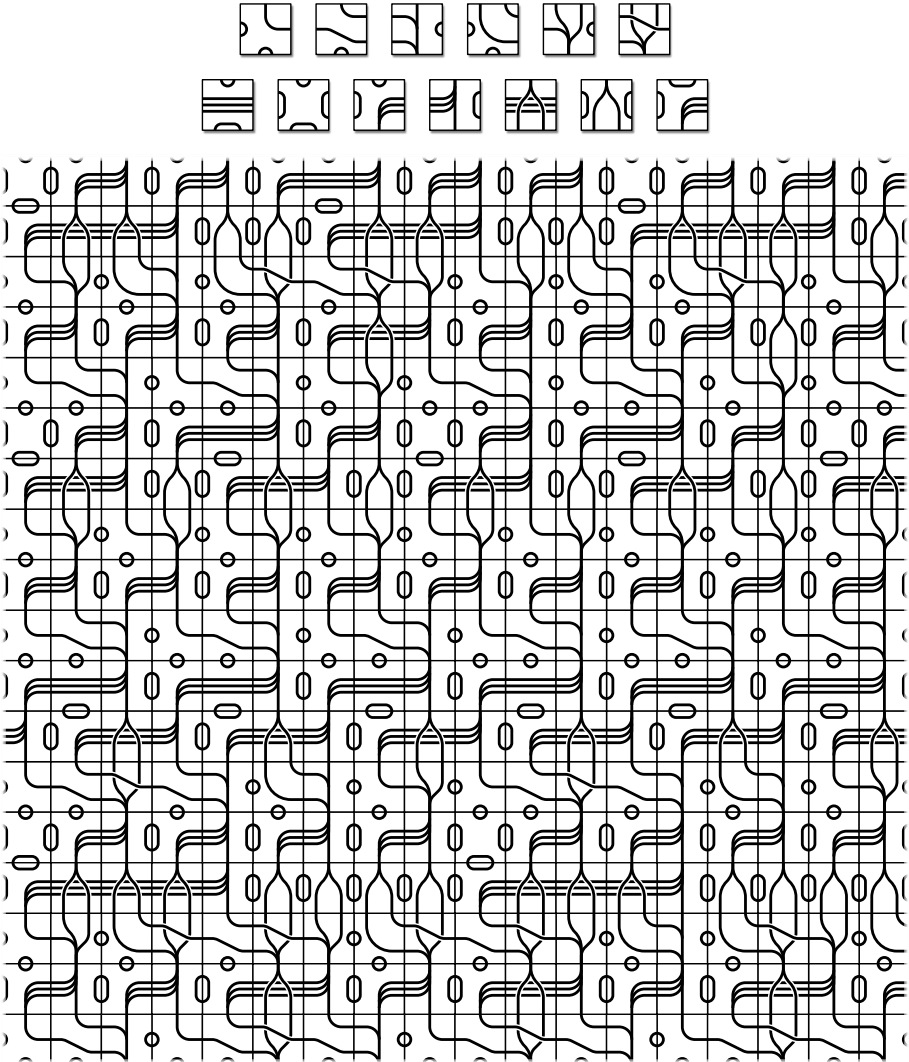


# *Models of Computation*

*Jeff Erickson*



December 28, 2018



<http://jeffe.cs.illinois.edu/teaching/algorithms/> • <http://algorithms.wtf/>

© Copyright 2014–2017 Jeff Erickson. Last update December 28, 2018.

For the most recent edition, see <http://jeffe.cs.illinois.edu/teaching/algorithms/>.

*I'm writing a book.  
I've got the page numbers done,  
so now I just have to fill in the rest.*

— Stephen Wright

## About These Notes

These are lecture notes that I wrote for the course “Algorithms and Models of Computation” at the University of Illinois, Urbana-Champaign for the first time in Fall 2014, and revised in Fall 2016. This course is a broad introduction to theoretical computer science, aimed at third-year computer science and computer engineering majors, that covers both fundamental topics in algorithms, for which I already have copious notes, and fundamental topics on formal languages and automata, for which I wrote the notes you are reading now.

The most recent revision of these notes (or nearly so) is available online at <http://jeffe.cs.illinois.edu/teaching/algorithms/> (or at the shorter URLs <http://algorithms.fyi/> and <http://algorithms.wtf/>) along with my algorithms notes and a near-complete archive of past homeworks and exams from all my theoretical computer science classes. I plan to revise and reorganize these whenever I teach this material, so you may find more recent versions on the web page of whatever course I am currently teaching.

## About the Exercises

Each note ends with several exercises, many of which I used in homeworks, discussion sections, or exams. \*Stars indicate more challenging problems (which I have *not* used in homeworks, discussion sections, or exams). Many of these exercises were contributed by my amazing teaching assistants:

Alex Steiger, Chao Xu, Charles Carlson, Connor Clark, Gail Steitz, Grant Czajkowski, Hsien-Chih Chang, Junqing Deng, Konstantinos Koiliaris, Nick Bachmair, Spencer Gordon, Tana Wattanawaroon, and Yipu Wang

**Please do not ask me for solutions to the exercises.** If you are a student, seeing the solution will rob you of the experience of solving the problem yourself, which is the only way to learn the material. If you are an instructor, you shouldn't ask your students to solve problems that you can't solve yourself. (I don't always follow my own advice, so some of the problems are buggy.)

## Caveat Lector!

These notes are best viewed as an unfinished first draft. You should assume the notes contain several major errors, in addition to the usual unending supply of typos, fencepost errors, off-by-one errors, and brain farts. Before Fall 2014, I had not taught this material in more than two decades. Moreover, the course itself is still very new—Lenny Pitt and I developed the course and offered the first pilot in Spring 2014 (with Lenny presenting the formal language material)—so even the choice of which material to emphasize, sketch, or exclude is still very much in flux.

I would sincerely appreciate feedback of any kind, especially bug reports.

Thanks, and enjoy!

— Jeff

# Contents

<b>1</b>	<b>Strings</b>	<b>1</b>
<b>2</b>	<b>Regular Languages</b>	<b>19</b>
<b>3</b>	<b>Finite-State Machines</b>	<b>31</b>
<b>4</b>	<b>Nondeterministic Automata</b>	<b>55</b>
<b>5</b>	<b>Context-Free Languages</b>	<b>77</b>
<b>6</b>	<b>Turing Machines</b>	<b>97</b>
<b>7</b>	<b>Undecidability</b>	<b>117</b>
<b>8</b>	<b>Universal Models</b>	<b>137</b>
<b>9</b>	<b>Nondeterministic Turing Machines</b>	<b>145</b>

# 1 Strings

Throughout this course, we will discuss dozens of algorithms and computational models that manipulate sequences: one-dimensional arrays, linked lists, blocks of text, walks in graphs, sequences of executed instructions, and so on. Ultimately the input and output of any algorithm must be representable as a finite string of symbols—the raw contents of some contiguous portion of the computer’s memory. Reasoning about computation requires reasoning about strings.

This note lists several formal definitions and formal induction proofs related to strings. These definitions and proofs are *intentionally* much more detailed than normally used in practice—most people’s intuition about strings is fairly accurate—but the extra precision is necessary for any sort of formal proof. It may be helpful to think of this material as part of the “assembly language” of theoretical computer science. We normally think about computation at a *much* higher level of abstraction, but ultimately every argument must “compile” down to these (and similar) definitions.

But the actual definitions and theorems are *not* the point. The point of playing with this material is to get some experience working with formal/mechanical definitions and proofs, especially inductive definitions and recursive proofs. Or should I say recursive definitions and inductive proofs? Whatever, they’re the same thing. Strings are a particularly simple and convenient playground for  $\left\{ \begin{smallmatrix} \text{induct} \\ \text{recurs} \end{smallmatrix} \right\}$ ion; we’ll see *many* more examples throughout the course. When you read this note, don’t just look at the *content* of the definitions and proofs; pay close attention to their *structure* and the *process* for creating them.

## 1.1 Strings

Fix an arbitrary finite set  $\Sigma$  called the **alphabet**; the individual elements of  $\Sigma$  are called **symbols** or **characters**. As a notational convention, I will always use lower-case letters near the start of the English alphabet ( $a, b, c, \dots$ ) as symbol variables, and *never* as explicit symbols. For explicit symbols, I will always use fixed-width upper-case letters ( $A, B, C, \dots$ ), digits ( $0, 1, 2, \dots$ ),

or other symbols ( $\diamond$ ,  $\$$ ,  $\#$ ,  $\bullet$ , ...) that are clearly distinguishable from variables. For further emphasis, I will almost always typeset explicit symbols in **RED**.

A **string** (or **word**) over  $\Sigma$  is a finite sequence of zero or more symbols from  $\Sigma$ . Formally, a string  $w$  over  $\Sigma$  is defined recursively as either

- the empty string, denoted by the Greek letter  $\varepsilon$  (epsilon), or
- an ordered pair  $(a, x)$ , where  $a$  is a symbol in  $\Sigma$  and  $x$  is a string over  $\Sigma$ .

We normally write either  $a \cdot x$  or simply  $ax$  to denote the ordered pair  $(a, x)$ . Similarly, we normally write explicit strings as sequences of symbols instead of nested ordered pairs; for example, **STRING** is convenient shorthand for the formal expression  $(S, (T, (R, (I, (N, (G, \varepsilon)))))$ ). As a notational convention, I will always use lower-case letters near the end of the English alphabet ( $\dots, w, x, y, z$ ) for string variables, and **SHOUTY** $\diamond$ **RED** $\diamond$ **MONOSPACED** $\diamond$ **TEXT** to typeset explicit (non-empty) strings.

The set of all strings over  $\Sigma$  is denoted  $\Sigma^*$  (pronounced “sigma star”). It is very important to remember that every element of  $\Sigma^*$  is a *finite* string, although  $\Sigma^*$  itself is an infinite set containing strings of every possible *finite* length.

## 1.2 Two Recursive Functions

Our first several proofs about strings will involve two natural functions, one giving the length of a string, the other gluing two strings together into a larger string. These two function behave exactly the way you think they do, but if we actually want to *prove* anything about them, we first have to *define* them in a way that supports formal proofs. Because the objects on which these functions act—strings—are defined recursively, the functions must also be defined recursively.

The **length**  $|w|$  of a string  $w$  is the number of symbols in  $w$ , defined formally as follows:

$$|w| := \begin{cases} 0 & \text{if } w = \varepsilon, \\ 1 + |x| & \text{if } w = ax. \end{cases}$$

For example, the string **FIFTEEN** has length 7, the string **SEVEN** has length 5, and the string **5** has length 1. Although they are formally different objects, we do not normally distinguish between symbols and strings of length 1.

The **concatenation** of two strings  $x$  and  $y$ , denoted either  $x \cdot y$  or simply  $xy$ , is the unique string containing the characters of  $x$  in order, followed by the characters in  $y$  in order. For example, the string **NOWHERE** is the concatenation of the strings **NOW** and **HERE**; that is, **NOW**  $\cdot$  **HERE** = **NOWHERE**. (On the other hand, **HERE**  $\cdot$  **NOW** = **HERENOW**.) Formally, concatenation is defined recursively as follows:

$$w \cdot z := \begin{cases} z & \text{if } w = \varepsilon, \\ a \cdot (x \cdot z) & \text{if } w = ax. \end{cases}$$

(Here I’m using a larger dot  $\cdot$  to formally distinguish the operator that concatenates two arbitrary strings from the syntactic sugar  $\cdot$  that builds a string from a single character and a string.)

When we describe the concatenation of more than two strings, we normally omit all dots and parentheses, writing  $wxyz$  instead of  $(w \cdot (x \cdot y)) \cdot z$ , for example. This simplification is justified by the fact (which we will prove shortly) that the function  $\cdot$  is associative.

### 1.3 Induction on Strings

Induction is *the* standard technique for proving statements about recursively defined objects. Hopefully you are already comfortable proving statements about *natural numbers* via induction, but induction is actually a far more general technique. Several different variants of induction can be used to prove statements about more general structures; here I describe the variant that I recommend (and actually use in practice). This variant follows two primary design considerations:

- **The case structure of the proof should mirror the case structure of the recursive definition.** For example, if you are proving something about all *strings*, your proof should have two cases: Either  $w = \varepsilon$ , or  $w = ax$  for some symbol  $a$  and string  $x$ .
- **The inductive hypothesis should be as strong as possible.** The (strong) inductive hypothesis for statements about natural numbers is *always* “Assume there is no counterexample  $k$  such that  $k < n$ .” I recommend adopting a similar inductive hypothesis for strings: “Assume there is no counterexample  $x$  such that  $|x| < |w|$ .” Then for the case  $w = ax$ , we have  $|x| = |w| - 1 < |w|$  by definition of  $|w|$ , so the inductive hypothesis applies to  $x$ .

Thus, string-induction proofs have the following boilerplate structure. Suppose we want to prove that every string is **perfectly cromulent**, whatever that means. The white boxes hide additional proof details that, among other things, depend on the precise definition of “**perfectly cromulent**”.

**Proof:** Let  $w$  be an arbitrary string.

Assume, for every string  $x$  such that  $|x| < |w|$ , that  $x$  is **perfectly cromulent**.

There are two cases to consider.

- Suppose  $w = \varepsilon$ .

Therefore,  $w$  is **perfectly cromulent**.

- Suppose  $w = ax$  for some symbol  $a$  and string  $x$ .

The induction hypothesis implies that  $x$  is **perfectly cromulent**.

Therefore,  $w$  is **perfectly cromulent**.

In both cases, we conclude that  $w$  is **perfectly cromulent**.

□

Here are three canonical examples of this proof structure. When developing proofs in this style, I strongly recommend first **mindlessly** writing **the green text (the boilerplate)** with lots of space for each case, then filling in **the red text (the actual theorem and the induction hypothesis)**, and only then starting to actually think.

Many students are confused (or at least bored and distracted) by the fact that we are proving *mind-bogglingly* obvious facts. If you're one of these students, try to remember that **the lemmas themselves are not the point**. Pay close attention to the *structure* of the proofs. Notice how each proof follows the boilerplate described above. Notice how every sentence of the proof follows *mechanically* from earlier sentences, definitions, and the rules of standard logic and arithmetic.

**Lemma 1.1.** *Adding nothing does nothing: For every string  $w$ , we have  $w \bullet \varepsilon = w$ .*

**Proof:** Let  $w$  be an arbitrary string. Assume that  $x \bullet \varepsilon = x$  for every string  $x$  such that  $|x| < |w|$ . There are two cases to consider:

- Suppose  $w = \varepsilon$ .

$$\begin{aligned} w \bullet \varepsilon &= \varepsilon \bullet \varepsilon && \text{because } w = \varepsilon, \\ &= \varepsilon && \text{by definition of concatenation,} \\ &= w && \text{because } w = \varepsilon. \end{aligned}$$

- Suppose  $w = ax$  for some symbol  $a$  and string  $x$ .

$$\begin{aligned} w \bullet \varepsilon &= (a \cdot x) \bullet \varepsilon && \text{because } w = ax, \\ &= a \cdot (x \bullet \varepsilon) && \text{by definition of concatenation,} \\ &= a \cdot x && \text{by the inductive hypothesis,} \\ &= w && \text{because } w = ax. \end{aligned}$$

In both cases, we conclude that  $w \bullet \varepsilon = w$ . □

**Lemma 1.2.** Concatenation adds length:  $|w \bullet x| = |w| + |x|$  for all strings  $w$  and  $x$ .

**Proof:** Let  $w$  and  $x$  be arbitrary strings. Assume that  $|y \bullet x| = |y| + |x|$  for every string  $y$  such that  $|y| < |w|$ . (Notice that we are using induction only on  $w$ , not on  $x$ .) There are two cases to consider:

- Suppose  $w = \varepsilon$ .

$$\begin{aligned} |w \bullet x| &= |\varepsilon \bullet x| && \text{because } w = \varepsilon \\ &= |x| && \text{by definition of } \bullet \\ &= |\varepsilon| + |x| && |\varepsilon| = 0 \text{ by definition of } |\cdot| \\ &= |w| + |x| && \text{because } w = \varepsilon \end{aligned}$$

- Suppose  $w = ay$  for some symbol  $a$  and string  $y$ .

$$\begin{aligned} |w \bullet x| &= |ay \bullet x| && \text{because } w = ay \\ &= |a \cdot (y \bullet x)| && \text{by definition of } \bullet \\ &= 1 + |y \bullet x| && \text{by definition of } |\cdot| \\ &= 1 + |y| + |x| && \text{by the inductive hypothesis} \\ &= |ay| + |x| && \text{by definition of } |\cdot| \\ &= |w| + |x| && \text{because } w = ay \end{aligned}$$

In both cases, we conclude that  $|w \bullet x| = |w| + |x|$ . □

**Lemma 1.3.** Concatenation is associative:  $(w \bullet x) \bullet y = w \bullet (x \bullet y)$  for all strings  $w$ ,  $x$ , and  $y$ .

**Proof:** Let  $w$ ,  $x$ , and  $y$  be arbitrary strings. Assume that  $(z \bullet x) \bullet y = z \bullet (x \bullet y)$  for every string  $z$  such that  $|z| < |w|$ . (Again, we are using induction only on  $w$ .) There are two cases to consider.



- Suppose  $w = \varepsilon$ .

$$\begin{aligned}
 (w \bullet x) \bullet y &= (\varepsilon \bullet x) \bullet y && \text{because } w = \varepsilon \\
 &= x \bullet y && \text{by definition of } \bullet \\
 &= \varepsilon \bullet (x \bullet y) && \text{by definition of } \bullet \\
 &= w \bullet (x \bullet y) && \text{because } w = \varepsilon
 \end{aligned}$$

- Suppose  $w = az$  for some symbol  $a$  and some string  $z$ .

$$\begin{aligned}
 (w \bullet x) \bullet y &= (az \bullet x) \bullet y && \text{because } w = az \\
 &= (a \bullet (z \bullet x)) \bullet y && \text{by definition of } \bullet \\
 &= a \bullet ((z \bullet x) \bullet y) && \text{by definition of } \bullet \\
 &= a \bullet (z \bullet (x \bullet y)) && \text{by the inductive hypothesis} \\
 &= az \bullet (x \bullet y) && \text{by definition of } \bullet \\
 &= w \bullet (x \bullet y) && \text{because } w = az
 \end{aligned}$$

In both cases, we conclude that  $(w \bullet x) \bullet y = w \bullet (x \bullet y)$ . □

## 1.4 More Than One Path up the Mountain

This is not the *only* boilerplate that one can use for induction proofs on strings. For example, we can model our case analysis on the following observation, whose easy proof we leave as an exercise (hint, hint): A string  $w \in \Sigma^*$  is **non-empty** if and only if either

- $w = a$  for some symbol  $a \in \Sigma$ , or
- $w = x \bullet y$  for some *non-empty* strings  $x$  and  $y$ .

In the latter case, Lemma 1.2 implies that  $|x| < |w|$  and  $|y| < |w|$ , so in an inductive proof, we can apply the inductive hypothesis to either  $x$  or  $y$  (or even both).

Here is a proof of Lemma 1.3 that uses this alternative recursive structure:

**Proof:** Let  $w$ ,  $x$ , and  $y$  be arbitrary strings. Assume that  $(z \bullet x') \bullet y' = z \bullet (x' \bullet y')$  for all strings  $x'$ ,  $y'$ , and  $z$  such that  $|z| < |w|$ . (We need a stronger induction hypothesis here than in the previous proofs!) There are **three** cases to consider.

- Suppose  $w = \varepsilon$ .

$$\begin{aligned}
 (w \bullet x) \bullet y &= (\varepsilon \bullet x) \bullet y && \text{because } w = \varepsilon \\
 &= x \bullet y && \text{by definition of } \bullet \\
 &= \varepsilon \bullet (x \bullet y) && \text{by definition of } \bullet \\
 &= w \bullet (x \bullet y) && \text{because } w = \varepsilon
 \end{aligned}$$

- Suppose  $w$  is equal to some symbol  $a$ .

$$\begin{aligned}
 (w \bullet x) \bullet y &= (a \bullet x) \bullet y && \text{because } w = a \\
 &= (a \bullet x) \bullet y && \text{because } a \bullet z = a \cdot z \text{ by definition of } \bullet \\
 &= a \bullet (x \bullet y) && \text{by definition of } \bullet \\
 &= a \bullet (x \bullet y) && \text{because } a \bullet z = a \cdot z \text{ by definition of } \bullet \\
 &= w \bullet (x \bullet y) && \text{because } w = a
 \end{aligned}$$

- Suppose  $w = u \cdot v$  for some nonempty strings  $u$  and  $v$ .

$$\begin{aligned}
 (w \cdot x) \cdot y &= ((u \cdot v) \cdot x) \cdot y && \text{because } w = u \cdot v \\
 &= (u \cdot (v \cdot x)) \cdot y && \text{by the inductive hypothesis, because } |u| < |w| \\
 &= u \cdot ((v \cdot x) \cdot y) && \text{by the inductive hypothesis, because } |u| < |w| \\
 &= u \cdot (v \cdot (x \cdot y)) && \text{by the inductive hypothesis, because } |v| < |w| \\
 &= (u \cdot v) \cdot (x \cdot y) && \text{by the inductive hypothesis, because } |u| < |w| \\
 &= w \cdot (x \cdot y) && \text{because } w = u \cdot v
 \end{aligned}$$

In all three cases, we conclude that  $(w \cdot x) \cdot y = w \cdot (x \cdot y)$ . □

## 1.5 Indices, Substrings, and Subsequences

Finally, I'll conclude this note by formally defining several other common functions and terms related to strings.

For any string  $w$  and any integer  $1 \leq i \leq |w|$ , the expression  $w_i$  denotes the  $i$ th symbol in  $w$ , counting from left to right. More formally,  $w_i$  is recursively defined as follows:

$$w_i := \begin{cases} a & \text{if } w = ax \text{ and } i = 1, \\ x_{i-1} & \text{if } w = ax \text{ and } i > 1. \end{cases}$$

As one might reasonably expect,  $w_i$  is formally undefined if  $i < 1$  or  $w = \varepsilon$ , and therefore (by induction) if  $i > |w|$ . The integer  $i$  is called the **index** of  $w_i$ .

We sometimes write strings as a concatenation of their constituent symbols using this subscript notation:  $w = w_1 w_2 \cdots w_{|w|}$ . While completely standard, this notation is slightly misleading, because it *incorrectly* suggests that the string  $w$  contains at least three symbols, when in fact  $w$  could be a single symbol or even the empty string.

In actual code, subscripts are usually expressed using the bracket notation  $w[i]$ . Brackets were introduced as a typographical convention over a hundred years ago because subscripts and superscripts<sup>1</sup> were difficult or impossible to type.<sup>2</sup> We sometimes write strings as explicit arrays  $w[1..n]$ , with the understanding that  $n = |w|$ . Again, this notation is potentially misleading; always remember that  $n$  might be zero; the string/array could be empty.

A **substring** of a string  $w$  is another string obtained from  $w$  by deleting zero or more symbols from the beginning and from the end. Formally, a string  $y$  is a substring of  $w$  if and only if there are strings  $x$  and  $z$  such that  $w = xyz$ . Extending the array notation for strings, we write  $w[i..j]$  to denote the substring of  $w$  starting at  $w_i$  and ending at  $w_j$ . More formally, we define

$$w[i..j] := \begin{cases} \varepsilon & \text{if } j < i, \\ w_i \cdot w[i+1..j] & \text{otherwise.} \end{cases}$$

---

<sup>1</sup>The same bracket notation is also used for bibliographic references, instead of the traditional footnote/endnote superscripts, for exactly the same reasons.

<sup>2</sup>A **typewriter** is an obsolete mechanical device loosely resembling a computer keyboard. Pressing a key on a typewriter moves a lever (called a “typebar”) that strikes a cloth ribbon full of ink against a piece of paper, leaving the image of a single character. Many historians believe that the ordering of letters on modern keyboards (**QWERTYUIOP**) evolved in the late 1800s, reaching its modern form on the 1874 Sholes & Glidden Type-Writer™, in part to separate many common letter pairs, to prevent typebars from jamming against each other; this is also why the keys on most modern keyboards are arranged in a slanted grid. (The common folk theory that the ordering was deliberately intended to slow down typists doesn’t withstand careful scrutiny.) A more recent theory suggests that the ordering was influenced by telegraph<sup>3</sup> operators, who found older alphabetic arrangements confusing, in part because of ambiguities in American Morse Code.

A **proper substring** of  $w$  is any substring other than  $w$  itself. For example, **LAUGH** is a proper substring of **SLAUGHTER**. Whenever  $y$  is a (proper) substring of  $w$ , we also call  $w$  a (proper) **superstring** of  $y$ .

A **prefix** of  $w[1..n]$  is any substring of the form  $w[1..j]$ . Equivalently, a string  $p$  is a **prefix** of another string  $w$  if and only if there is a string  $x$  such that  $px = w$ . A **proper prefix** of  $w$  is any prefix except  $w$  itself. For example, **DIE** is a proper prefix of **DIET**.

Similarly, a **suffix** of  $w[1..n]$  is any substring of the form  $w[i..n]$ . Equivalently, a string  $s$  is a **suffix** of a string  $w$  if and only if there is a string  $x$  such that  $xs = w$ . A **proper suffix** of  $w$  is any suffix except  $w$  itself. For example, **YES** is a proper suffix of **EYES**, and **HE** is both a proper prefix and a proper suffix of **HEADACHE**.

A **subsequence** of a string  $w$  is a string obtained by deleting zero or more symbols from *anywhere* in  $w$ . More formally,  $z$  is a subsequence of  $w$  if and only if

- $z = \varepsilon$ , or
- $w = ax$  for some symbol  $a$  and some string  $x$  such that  $z$  is a subsequence of  $x$ .
- $w = ax$  and  $z = ay$  for some symbol  $a$  and some strings  $x$  and  $y$ , and  $y$  is a subsequence of  $x$ .

A **proper subsequence** of  $w$  is any subsequence of  $w$  other than  $w$  itself. Whenever  $z$  is a (proper) subsequence of  $w$ , we also call  $w$  a (proper) **supersequence** of  $z$ .

Substrings and subsequences are not the same objects; don't confuse them! Every substring of  $w$  is also a subsequence of  $w$ , but not every subsequence is a substring. For example, **METAL** is a subsequence, but not a substring, of **MEATBALL**. To emphasize the distinction, we sometimes redundantly refer to substrings of  $w$  as **contiguous** substrings, meaning all their symbols appear together in  $w$ .

## Exercises

Most of the following exercises ask for proofs of various claims about strings. Here “prove” means give a complete, self-contained, formal proof by inductive definition-chasing, using the boilerplate structure recommended in Section 1.3. Feel free to use Lemmas 1.1, 1.2, and 1.3 without proof, but don't assume any other facts about strings that you have not actually proved. (Some later exercises rely on results proved in earlier exercises.) Do not appeal to intuition, and do not use the words “obvious” or “clearly” or “just”. Most of these claims *are* in fact obvious; the real exercise is understanding and formally expressing *why* they're obvious.

---

<sup>3</sup>A **telegraph** is an obsolete electromechanical communication device consisting of an electrical circuit with a switch at one end and an electromagnet at the other. The sending operator would press and release a key, closing and opening the circuit, originally causing the electromagnet to push a stylus onto a moving paper tape, leaving marks that could be decoded by the receiving operator. (Operators quickly discovered that they could directly decode the clicking sounds made by the electromagnet, and so the paper tape became obsolete almost immediately.) The most common scheme within the US to encode symbols, developed by Alfred Vail and Samuel Morse in 1837, used (mostly) short (·) and long (—) marks—now called “dots” and “dashes”, or “dits” and “dahs”—separated by gaps of various lengths. American Morse code (as it became known) was ambiguous; for example, the letter **Z** and the string **SE** were both encoded by the sequence · · · · (“di-di-dit, dit”). This ambiguity has been blamed for the **S** key's position on the typewriter keyboard near **E** and **Z**.

Vail and Morse were of course not the first people to propose encoding symbols as strings of bits. That honor apparently falls to Francis Bacon, who devised a five-bit binary encoding of the alphabet (except for the letters **J** and **U**) in 1605 as the basis for a steganographic code—a method of hiding secret message in otherwise normal text.

**Note to instructors:** *Do not assign any of these problems before solving them yourself, especially on exams.* It's very easy to underestimate the difficulty of these problems, or at least the lengths of their solutions, which for exams is a reasonable proxy for difficulty. Also, several later exercises rely implicitly on identities like  $\#(a, x \cdot y) = \#(a, x) + \#(a, y)$  that are only proved in earlier exercises. It's unfair to assign these problems to students without telling them which earlier results they can use.

## Useful Facts About Strings

1. Let  $w$  be an arbitrary string, and let  $n = |w|$ . Prove each of the following statements.
  - (a)  $w$  has exactly  $n + 1$  prefixes.
  - (b)  $w$  has exactly  $n$  proper suffixes.
  - (c)  $w$  has at most  $n(n + 1)/2$  distinct substrings. (Why “at most”?)
  - (d)  $w$  has at most  $2^n - 1$  distinct proper subsequences. (Why “at most”?)
  
2. Prove the following useful identities for all strings  $w, x, y$ , and  $z$  directly from the definition of  $\cdot$ , *without* referring to the length of any string.
  - (a) If  $x \cdot y = x$ , then  $y = \varepsilon$ .
  - (b) If  $x \cdot y = y$ , then  $x = \varepsilon$ .
  - (c) If  $x \cdot z = y \cdot z$ , then  $x = y$ .
  - (d) If  $x \cdot y = x \cdot z$ , then  $y = z$ .
  
3. Prove the following useful fact about substrings. An arbitrary string  $x$  is a substring of another arbitrary string  $w = u \cdot v$  if and only if at least one of the following conditions holds:
  - $x$  is a substring of  $u$ .
  - $x$  is a substring of  $v$ .
  - $x = yz$  where  $y$  is a suffix of  $u$  and  $z$  is a prefix of  $v$ .
  
4. Let  $w$  be an arbitrary string, and let  $n = |w|$ . Prove the following statements for all indices  $1 \leq i \leq j \leq k \leq n$ .
  - (a)  $|w[i..j]| = j - i + 1$
  - (b)  $w[i..j] \cdot w[j + 1..k] = w[i..k]$
  - (c)  $w^R[i..j] = (w[i'..j'])^R$  where  $i' + j = j' + i = |w| + 1$ .

## Recursive Functions

5. For any symbol  $a$  and any string  $w$ , let  $\#(a, w)$  denote the number of occurrences of  $a$  in  $w$ . For example,  $\#(A, BANANA) = 3$  and  $\#(X, FLIBBERTIGIBBET) = 0$ .
  - (a) Give a formal recursive definition of the function  $\# : \Sigma \times \Sigma^* \rightarrow \mathbb{N}$ .

- (b) Prove that  $\#(a, xy) = \#(a, x) + \#(a, y)$  for every symbol  $a$  and all strings  $x$  and  $y$ . Your proof must rely on both your answer to part (a) and the formal recursive definition of string concatenation.
- (c) Prove the following identity for all alphabets  $\Sigma$  and all strings  $w \in \Sigma^*$ :

$$|w| = \sum_{a \in \Sigma} \#(a, w)$$

[Hint: Don't try to use induction on  $\Sigma$ .]

6. The **reversal**  $w^R$  of a string  $w$  is defined recursively as follows:

$$w^R := \begin{cases} \varepsilon & \text{if } w = \varepsilon \\ x^R \cdot a & \text{if } w = a \cdot x \end{cases}$$

- (a) Prove that  $(w^R)^R = w$  for every string  $w$ .
- (b) Prove that  $|w^R| = |w|$  for every string  $w$ .
- (c) Prove that  $(w \cdot x)^R = x^R \cdot w^R$  for all strings  $w$  and  $x$ .
- (d) Prove that  $\#(a, w^R) = \#(a, w)$  for every string  $w$  and every symbol  $a$ . (See Exercise 5.)
7. For any string  $w$  and any non-negative integer  $n$ , let  $w^n$  denote the string obtained by concatenating  $n$  copies of  $w$ ; more formally, we define

$$w^n := \begin{cases} \varepsilon & \text{if } n = 0 \\ w \cdot w^{n-1} & \text{otherwise} \end{cases}$$

For example,  $(\text{BLAH})^5 = \text{BLAHBLAHBLAHBLAHBLAH}$  and  $\varepsilon^{374} = \varepsilon$ .

- (a) Prove that  $w^m \cdot w^n = w^{m+n}$  for every string  $w$  and all non-negative integers  $n$  and  $m$ .
- (b) Prove that  $\#(a, w^n) = n \cdot \#(a, w)$  for every string  $w$ , every symbol  $a$ , and every non-negative integer  $n$ . (See Exercise 5.)
- (c) Prove that  $(w^R)^n = (w^n)^R$  for every string  $w$  and every non-negative integer  $n$ .
- (d) Prove that for all strings  $x$  and  $y$  that if  $x \cdot y = y \cdot x$ , then  $x = w^m$  and  $y = w^n$  for some string  $w$  and some non-negative integers  $m$  and  $n$ . [Hint: Careful with  $\varepsilon$ !]
8. The **complement**  $w^c$  of a string  $w \in \{0, 1\}^*$  is obtained from  $w$  by replacing every 0 in  $w$  with a 1 and vice versa. The complement function can be defined recursively as follows:

$$w^c := \begin{cases} \varepsilon & \text{if } w = \varepsilon \\ 1 \cdot x^c & \text{if } w = 0x \\ 0 \cdot x^c & \text{if } w = 1x \end{cases}$$

- (a) Prove that  $|w| = |w^c|$  for every string  $w$ .

- (b) Prove that  $(x \cdot y)^c = x^c \cdot y^c$  for all strings  $x$  and  $y$ .
- (c) Prove that  $\#(\mathbf{1}, w) = \#(\mathbf{0}, w^c)$  for every string  $w$ .
- (d) Prove that  $(w^R)^c = (w^c)^R$  for every string  $w$ .
- (e) Prove that  $(w^n)^c = (w^c)^n$  for every string  $w$  and every non-negative integer  $n$ .

9. For any string  $w \in \{\mathbf{0}, \mathbf{1}, \mathbf{2}\}^*$ , let  $w^+$  denote the string obtained from  $w$  by replacing each symbol  $a$  in  $w$  by the symbol corresponding to  $(a + 1) \bmod 3$ . for example,  $\mathbf{0102101}^+ = \mathbf{1210212}$ . This function can be defined more formally as follows:

$$w^+ := \begin{cases} \varepsilon & \text{if } w = \varepsilon \\ \mathbf{1} \cdot x^+ & \text{if } w = \mathbf{0}x \\ \mathbf{2} \cdot x^+ & \text{if } w = \mathbf{1}x \\ \mathbf{0} \cdot x^+ & \text{if } w = \mathbf{2}x \end{cases}$$

- (a) Prove that  $|w| = |w^+|$  for every string  $w \in \{\mathbf{0}, \mathbf{1}, \mathbf{2}\}^*$ .
  - (b) Prove that  $(x \cdot y)^+ = x^+ \cdot y^+$  for all strings  $x, y \in \{\mathbf{0}, \mathbf{1}, \mathbf{2}\}^*$ .
  - (c) Prove that  $\#(\mathbf{1}, w^+) = \#(\mathbf{0}, w)$  for every string  $w \in \{\mathbf{0}, \mathbf{1}, \mathbf{2}\}^*$ .
  - (d) Prove that  $(w^+)^R = (w^R)^+$  for every string  $w \in \{\mathbf{0}, \mathbf{1}, \mathbf{2}\}^*$ .
10. For any string  $w \in \{\mathbf{0}, \mathbf{1}\}^*$ , let  $\text{swap}(w)$  denote the string obtained from  $w$  by swapping the first and second symbols, the third and fourth symbols, and so on. For example:

$$\begin{aligned} \text{swap}(\mathbf{101}) &= \mathbf{011} \\ \text{swap}(\mathbf{100111}) &= \mathbf{011011} \\ \text{swap}(\mathbf{101100011011}) &= \mathbf{011100100111}. \end{aligned}$$

The  $\text{swap}$  function can be formally defined as follows:

$$\text{swap}(w) := \begin{cases} \varepsilon & \text{if } w = \varepsilon \\ w & \text{if } w = \mathbf{0} \text{ or } w = \mathbf{1} \\ ba \cdot \text{swap}(x) & \text{if } w = abx \text{ for some } a, b \in \{\mathbf{0}, \mathbf{1}\} \text{ and } x \in \{\mathbf{0}, \mathbf{1}\}^* \end{cases}$$

- (a) Prove that  $|\text{swap}(w)| = |w|$  for every string  $w$ .
  - (b) Prove that  $\text{swap}(\text{swap}(w)) = w$  for every string  $w$ .
  - (c) Prove that  $\text{swap}(w^R) = (\text{swap}(w))^R$  for every string  $w$  such that  $|w|$  is even. [Hint: Your proof must invoke **four** different recursive definitions: reversal  $w^R$ , concatenation  $\cdot$ , length  $|w|$ , and the  $\text{swap}$  function!]
11. For any string  $w \in \{\mathbf{0}, \mathbf{1}\}^*$ , let  $\text{sort}(w)$  denote the string obtained by sorting the characters in  $w$ . For example,  $\text{sort}(\mathbf{010101}) = \mathbf{000111}$ . The sort function can be defined recursively as follows:

$$\text{sort}(w) := \begin{cases} \varepsilon & \text{if } w = \varepsilon \\ \mathbf{0} \cdot \text{sort}(x) & \text{if } w = \mathbf{0}x \\ \text{sort}(x) \cdot \mathbf{1} & \text{if } w = \mathbf{1}x \end{cases}$$

Prove the following for all strings  $w, x, y \in \{\mathbf{0}, \mathbf{1}\}^*$ :

- (a) Prove that  $\text{sort}(w) \in 0^*1^*$  for every string  $w \in \{0, 1\}^*$ .
- (b) Prove that  $\#(0, w) = \#(0, \text{sort}(w))$  for every string  $w \in \{0, 1\}^*$ .
- (c) Prove that  $|w| = |\text{sort}(w)|$ , for every string  $w \in \{0, 1\}^*$ .
- (d) Prove that  $\text{sort}(w) = 0^{\#(0,w)}1^{\#(1,w)}$ , for every string  $w \in \{0, 1\}^*$ . (In other words, prove that our recursive definition is correct.)
- (e) Prove that  $\text{sort}(w^R) = \text{sort}(w)$ , for every string  $w \in \{0, 1\}^*$ .

12. Consider the following recursively defined function:

$$\text{merge}(x, y) := \begin{cases} y & \text{if } x = \varepsilon \\ x & \text{if } y = \varepsilon \\ 0 \cdot \text{merge}(w, y) & \text{if } x = 0w \\ 0 \cdot \text{merge}(x, z) & \text{if } y = 0z \\ 1 \cdot \text{merge}(w, y) & \text{if } x = 1w \text{ and } y = 1z \end{cases}$$

For example:

$$\begin{aligned} \text{merge}(10, 10) &= 1010 \\ \text{merge}(10, 010) &= 01010 \\ \text{merge}(010, 0001100) &= 0000101100 \end{aligned}$$

- (a) Prove that  $\text{merge}(x, y) \in 0^*1^*$  for all strings  $x, y \in 0^*1^*$ . (The regular expression  $0^*1^*$  is shorthand for the language  $\{0^a1^b \mid a, b \geq 0\}$ .)
- (b) Prove that  $\text{sort}(x \cdot y) = \text{merge}(\text{sort}(x), \text{sort}(y))$  for all strings  $x, y \in \{0, 1\}^*$ . (The  $\text{sort}$  function is defined in the previous exercise.)

13. Consider the following pair of mutually recursive functions on strings:

$$\text{evens}(w) := \begin{cases} \varepsilon & \text{if } w = \varepsilon \\ \text{odds}(x) & \text{if } w = ax \end{cases} \quad \text{odds}(w) := \begin{cases} \varepsilon & \text{if } w = \varepsilon \\ a \cdot \text{evens}(x) & \text{if } w = ax \end{cases}$$

For example,  $\text{evens}(\text{MISSISSIPPI}) = \text{ISSIP}$  and  $\text{odds}(\text{MISSISSIPPI}) = \text{MSISPI}$ .

- (a) Prove the following identity for all strings  $w$  and  $x$ :

$$\text{evens}(w \cdot x) = \begin{cases} \text{evens}(w) \cdot \text{evens}(x) & \text{if } |w| \text{ is even,} \\ \text{evens}(w) \cdot \text{odds}(x) & \text{if } |w| \text{ is odd.} \end{cases}$$

- (b) State and prove a similar identity for  $\text{odds}(w \cdot x)$ .
- (c) Prove that every string  $w$  is a shuffle of  $\text{evens}(w)$  and  $\text{odds}(w)$ .

14. Consider the following recursively defined function:

$$\text{stutter}(w) := \begin{cases} \varepsilon & \text{if } w = \varepsilon \\ aa \cdot \text{stutter}(x) & \text{if } w = ax \end{cases}$$

For example,  $\text{stutter}(\text{MISSISSIPPI}) = \text{MMIISSSSIISSSSIIPPPPII}$ .

- (a) Prove that  $|stutter(w)| = 2|w|$  for every string  $w$ .
- (b) Prove that  $evens(stutter(w)) = w$  for every string  $w$ .
- (c) Prove that  $odds(stutter(w)) = w$  for every string  $w$ .
- (d) Prove that  $stutter(w)$  is a shuffle of  $w$  and  $w$ , for every string  $w$ .
- (e) Prove that  $w$  is a palindrome if and only if  $stutter(w)$  is a palindrome, for every string  $w$ .

15. Consider the following recursive function:

$$faro(w, z) := \begin{cases} z & \text{if } w = \varepsilon \\ a \cdot faro(z, x) & \text{if } w = ax \end{cases}$$

For example,  $faro(\mathbf{0011}, \mathbf{0101}) = \mathbf{00011011}$ . (A "faro shuffle" splits a deck of cards into two equal piles and then perfectly interleaves them.)

- (a) Prove that  $|faro(x, y)| = |x| + |y|$  for all strings  $x$  and  $y$ .
  - (b) Prove that  $faro(w, w) = stutter(w)$  for every string  $w$ .
  - (c) Prove that  $faro(odds(w), evens(w)) = w$  for every string  $w$ .
16. For any string  $w$ , let  $declutter(w)$  denote the string obtained from  $w$  by deleting any symbol that equals its immediate successor. For example,  $declutter(\mathbf{MISSISSIPPI}) = \mathbf{MISISPI}$ , and  $declutter(\mathbf{ABBCCCAAACCCBBA}) = \mathbf{ABCACBA}$ .
- (a) Given a recursive definition for the function  $declutter$ .
  - (b) Using your recursive definition, prove that  $declutter(stutter(w)) = declutter(w)$  for every string  $w$ .
  - (c) Using your recursive definition, prove that  $declutter(w^R) = (declutter(w))^R$  for every string  $w$ .
  - (d) Using your recursive definition, prove that  $w$  is a palindrome if and only if  $declutter(w)$  is a palindrome, for every string  $w$ .

17. Consider the following recursively defined function

$$hanoi(w) = \begin{cases} \varepsilon & \text{if } w = \varepsilon \\ hanoi(x) \cdot a \cdot hanoi(x) & \text{if } w = ax \end{cases}$$

Prove that  $|hanoi(w)| = 2^{|w|} - 1$  for every string  $w$ .

18. Consider the following recursively defined function

$$slog(w) = \begin{cases} \varepsilon & \text{if } w = \varepsilon \\ a \cdot slog(evens(w)) & \text{if } w = ax \end{cases}$$

Prove that  $|slog(w)| = \lceil \log_2(|w| + 1) \rceil$  for every string  $w$ .



19. Consider the following recursively defined function

$$\text{bitrev}(w) = \begin{cases} w & \text{if } |w| \leq 1 \\ \text{bitrev}(\text{odds}(w)) \cdot \text{bitrev}(\text{evens}(w)) & \text{otherwise} \end{cases}$$

(a) Prove that  $|\text{bitrev}(w)| = |w|$  for every string  $w$ .

\* (b) Prove that  $\text{bitrev}(\text{bitrev}(w)) = w$  for every string  $w$  such that  $|w|$  is a power of 2.

20. The **binary value** of any string  $w \in \{0, 1\}^*$  is the integer whose binary representation (possibly with leading 0s) is  $w$ . The value function can be defined recursively as follows:

$$\text{value}(w) := \begin{cases} 0 & \text{if } w = \varepsilon \\ 2 \cdot \text{value}(x) & \text{if } w = x \cdot 0 \\ 2 \cdot \text{value}(x) + 1 & \text{if } w = x \cdot 1 \end{cases}$$

(a) Prove that  $\text{value}(w) + \text{value}(w^c) = 2^{|w|} - 1$  for every string  $w \in \{0, 1\}^*$ .

(b) Prove that  $\text{value}(x \cdot y) = \text{value}(x) \cdot 2^{|y|} + \text{value}(y)$  for all strings  $x, y \in \{0, 1\}^*$ .

\* (c) Prove that  $\text{value}(x)$  is divisible by 3 if and only if  $\text{value}(x^R)$  is divisible by 3.

### Recursively Defined Sets

20. Recursively define a set  $L$  of strings over the alphabet  $\{0, 1\}$  as follows:

- The empty string  $\varepsilon$  is in  $L$ .
- For any two strings  $x$  and  $y$  in  $L$ , the string  $0x1y0$  is also in  $L$ .
- These are the only strings in  $L$ .

(a) Prove that the string  $000010101010010100$  is in  $L$ .

(b) Prove by induction that every string in  $L$  has exactly twice as many 0s as 1s. (You may assume the identity  $\#(a, xy) = \#(a, x) + \#(a, y)$  for any symbol  $a$  and any strings  $x$  and  $y$ ; see Exercise 5(b).)

(c) Give an example of a string with exactly twice as many 0s as 1s that is *not* in  $L$ .

21. Recursively define a set  $L$  of strings over the alphabet  $\{0, 1\}$  as follows:

- The empty string  $\varepsilon$  is in  $L$ .
- For any two strings  $x$  and  $y$  in  $L$ , the string  $0x1y$  is also in  $L$ .
- For any two strings  $x$  and  $y$  in  $L$ , the string  $1x0y$  is also in  $L$ .
- These are the only strings in  $L$ .

(a) Prove that the string  $01000110111001$  is in  $L$ .

(b) Prove by induction that every string in  $L$  has exactly the same number of 0s and 1s. (You may assume the identity  $\#(a, xy) = \#(a, x) + \#(a, y)$  for any symbol  $a$  and any strings  $x$  and  $y$ ; see Exercise 5(b).)

(c) Prove by induction that  $L$  contains every string with the same number of 0s and 1s.

22. Recursively define a set  $L$  of strings over the alphabet  $\{0, 1\}$  as follows:

- The empty string  $\varepsilon$  is in  $L$ .
- For any strings  $x$  in  $L$ , the strings  $0x1$  and  $1x0$  are also in  $L$ .
- For any two strings  $x$  and  $y$  in  $L$ , the string  $x \cdot y$  is also in  $L$ .
- These are the only strings in  $L$ .

(a) Prove that the string  $01000110111001$  is in  $L$ .

(b) Prove by induction that every string in  $L$  has exactly the same number of 0s and 1s. (You may assume the identity  $\#(a, xy) = \#(a, x) + \#(a, y)$  for any symbol  $a$  and any strings  $x$  and  $y$ ; see Exercise 5(b).)

(c) Prove by induction that every string with the same number of 0s and 1s is in  $L$ .

23. Recursively define a set  $L$  of strings over the alphabet  $\{0, 1, 2\}$  as follows:

- The empty string  $\varepsilon$  is in  $L$ .
- For any string  $x$  in  $L$ , the string  $0x$  is also in  $L$ .
- For any strings  $x$  and  $y$  in  $L$ , the strings  $1x2y$  and  $2x1y$  are also in  $L$ .
- These are the only strings in  $L$ .

(a) Prove that the string  $001201110220121220$  is in  $L$ .

(b) For any string  $w \in \{0, 1, 2\}^*$ , let  $\text{ModThree}(w)$  denote the sum of the digits of  $w$  modulo 3. This function can be defined recursively as follows:

$$\text{ModThree}(w) = \begin{cases} 0 & \text{if } w = \varepsilon \\ \text{ModThree}(x) & \text{if } w = 0x \\ (\text{ModThree}(x) + 1) \bmod 3 & \text{if } w = 1x \\ (\text{ModThree}(x) + 2) \bmod 3 & \text{if } w = 2x \end{cases}$$

For example,  $\text{ModThree}(2211) = 0$  and  $\text{ModThree}(001201110220121220) = 0$ . Prove that  $\text{ModThree}(w) = 0$  for every string in  $w \in L$ .

(c) Find a string  $w \in \{0, 1, 2\}^*$  such that  $\text{ModThree}(w) = 0$  but  $w \notin L$ . Prove that your answer is correct.

(d) Prove that  $\#(1, w) = \#(2, w)$  for every string  $w \in L$ .

(e) Find a string  $w \in \{0, 1, 2\}^*$  such that  $\#(1, w) = \#(2, w)$  but  $w \notin L$ . Prove that your answer is correct.

(f) Prove that  $L = \{w \in \{0, 1, 2\}^* \mid \text{ModThree}(w) = 0 \text{ and } \#(1, w) = \#(2, w)\}$ .

24. A **palindrome** is a string that is equal to its reversal.

- Give a recursive definition of a palindrome over the alphabet  $\Sigma$ .
- Prove that any string  $p$  meets your recursive definition of a palindrome if and only if  $p = p^R$ .
- Using your recursive definition, prove that the strings  $w \cdot w^R$  and  $w \cdot a \cdot w^R$  are palindromes, for every string  $w$  and symbol  $a$ .
- Using your recursive definition, prove that  $p^n$  is a palindrome for every palindrome  $p$  and every natural number  $n$ . (See Exercise 7.)
- Using your recursive definition, prove that for every palindrome  $p$ , there is at most one symbol  $a$  such that  $\#(a, p)$  is odd. (See Exercise 5.)

25. A string  $w \in \Sigma^*$  is called a **shuffle** of two strings  $x, y \in \Sigma^*$  if at least one of the following recursive conditions is satisfied:

- $w = x = y = \varepsilon$ .
- $w = aw'$  and  $x = ax'$  and  $w'$  is a shuffle of  $x'$  and  $y$ , for some  $a \in \Sigma$  and some  $w', x' \in \Sigma^*$ .
- $w = aw'$  and  $y = ay'$  and  $w'$  is a shuffle of  $x$  and  $y'$ , for some  $a \in \Sigma$  and some  $w', y' \in \Sigma^*$ .

For example, the string **B**ANANANAN**A**SA is a shuffle of the strings **B**ANANA and ANAN**A**S**A**.

- Prove that if  $w$  is a shuffle of  $x$  and  $y$ , then  $|w| = |x| + |y|$ .
- Prove that  $w$  is a shuffle of  $x$  and  $y$  if and only if  $w^R$  is a shuffle of  $x^R$  and  $y^R$ .

26. For any positive integer  $n$ , the **Fibonacci string**  $F_n$  is defined recursively as follows:

$$F_n = \begin{cases} \mathbf{0} & \text{if } n = 1, \\ \mathbf{1} & \text{if } n = 2, \\ F_{n-2} \cdot F_{n-1} & \text{otherwise.} \end{cases}$$

For example,  $F_6 = \mathbf{10101101}$  and  $F_7 = \mathbf{0110110101101}$ .

- Prove that for every integer  $n \geq 2$ , the string  $F_n$  can also be obtained from  $F_{n-1}$  by replacing every occurrence of **0** with **1** and replacing every occurrence of **1** with **01**. More formally, prove that  $F_n = \text{Finc}(F_{n-1})$ , where

$$\text{Finc}(w) = \begin{cases} \varepsilon & \text{if } w = \varepsilon \\ \mathbf{1} \cdot \text{Finc}(x) & \text{if } w = \mathbf{0}x \\ \mathbf{01} \cdot \text{Finc}(x) & \text{if } w = \mathbf{1}x \end{cases}$$

[Hint: First prove that  $\text{Finc}(x \cdot y) = \text{Finc}(x) \cdot \text{Finc}(y)$ .]

- Prove that the Fibonacci string  $F_n$  begins with **1** if and only if  $n$  is even.
- Prove that **00** is not a substring of any Fibonacci string  $F_n$ .
- Prove that **111** is not a substring of any Fibonacci string  $F_n$ .

- \*(e) Prove that  $01010$  is not a substring of any Fibonacci string  $F_n$ .
- \*(f) Find another string  $w$  that is not a substring of any Fibonacci string  $F_n$ , such that  $00$  and  $111$  and  $01010$  are not substrings of  $w$ .
- ★(g) Find a set of strings  $F$  with the following properties:
  - No string in  $F$  is a substring of any Fibonacci string  $F_n$ .
  - No string in  $F$  is a proper substring of any other string in  $F$ .
  - For all strings  $x \in \{0, 1\}^*$ , if  $x$  has no substrings in  $F$ , then  $x$  is a substring of some Fibonacci string  $F_n$ .
- ★(h) Prove that the reversal of each Fibonacci string is a substring of another Fibonacci string. More formally, prove that for every integer  $n \geq 0$ , the string  $F_n^R$  is a substring of  $F_m$  for some integer  $m \geq n$ .

\*27. Prove that the following three properties of strings are in fact identical.

- A string  $w \in \{0, 1\}^*$  is **balanced** if it satisfies one of the following conditions:
  - $w = \varepsilon$ ,
  - $w = 0x1$  for some balanced string  $x$ , or
  - $w = xy$  for some balanced strings  $x$  and  $y$ .
- A string  $w \in \{0, 1\}^*$  is **erasable** if it satisfies one of the following conditions:
  - $w = \varepsilon$ , or
  - $w = x01y$  for some strings  $x$  and  $y$  such that  $xy$  is erasable. (The strings  $x$  and  $y$  are not necessarily erasable.)
- A string  $w \in \{0, 1\}^*$  is **conservative** if it satisfies **both** of the following conditions:
  - $w$  has an equal number of  $0$ s and  $1$ s, and
  - no prefix of  $w$  has more  $0$ s than  $1$ s.

- (a) Prove that every balanced string is erasable.
- (b) Prove that every erasable string is conservative.
- (c) Prove that every conservative string is balanced.

[Hint: To develop intuition, it may be helpful to think of  $0$ s as left brackets and  $1$ s as right brackets, but **don't** invoke this intuition in your proofs.]

\*28. A string  $w \in \{0, 1\}^*$  is **equitable** if it has an equal number of  $0$ s and  $1$ s.

- (a) Prove that a string  $w$  is equitable if and only if it satisfies one of the following conditions:
  - $w = \varepsilon$ ,
  - $w = 0x1$  for some equitable string  $x$ ,
  - $w = 1x0$  for some equitable string  $x$ , or
  - $w = xy$  for some equitable strings  $x$  and  $y$ .

(b) Prove that a string  $w$  is equitable if and only if it satisfies one of the following conditions:

- $w = \varepsilon$ ,
- $w = x01y$  for some strings  $x$  and  $y$  such that  $xy$  is equitable, or
- $w = x10y$  for some strings  $x$  and  $y$  such that  $xy$  is equitable.

In the last two cases, the individual strings  $x$  and  $y$  are not necessarily equitable.

(c) Prove that a string  $w$  is equitable if and only if it satisfies one of the following conditions:

- $w = \varepsilon$ ,
- $w = xy$  for some balanced string  $x$  and some equitable string  $y$ , or
- $w = x^Ry$  for some for some balanced string  $x$  and some equitable string  $y$ .

(See the previous exercise for the definition of “balanced”.)

## 2 Regular Languages

### 2.1 Languages

A *formal language* (or just a *language*) is a set of strings over some finite alphabet  $\Sigma$ , or equivalently, an arbitrary subset of  $\Sigma^*$ . For example, each of the following sets is a language:

- The empty set  $\emptyset$ .<sup>1</sup>
- The set  $\{\varepsilon\}$ .
- The set  $\{0, 1\}^*$ .
- The set  $\{\text{THE, OXFORD, ENGLISH, DICTIONARY}\}$ .
- The set of all subsequences of  $\text{THE}\diamond\text{OXFORD}\diamond\text{ENGLISH}\diamond\text{DICTIONARY}$ .
- The set of all words in *The Oxford English Dictionary*.
- The set of all strings in  $\{0, 1\}^*$  with an odd number of 1s.
- The set of all strings in  $\{0, 1\}^*$  that represent a prime number in base 13.
- The set of all sequences of turns that solve the Rubik’s cube (starting in some fixed configuration)
- The set of all python programs that print “Hello World!”

As a notational convention, I will always use italic upper-case letters (usually  $L$ , but also  $A$ ,  $B$ ,  $C$ , and so on) to represent languages.

Formal languages are not “languages” in the same sense that English, Klingon, and Python are “languages”. Strings in a formal language do not necessarily carry any “meaning”, nor are they necessarily assembled into larger units (“sentences” or “paragraphs” or “packages”) according to some “grammar”.

---

<sup>1</sup>The empty set symbol  $\emptyset$  was introduced in 1939 by André Weil, as a member of the pseudonymous mathematical collective Nikolai Bourbaki. The symbol derives from the Norwegian letter Ø, which pronounced like a German ö or a sound of disgust, and *not* from the Greek letter  $\phi$ . Calling the empty set “fie” or “fee” makes the baby Jesus cry.

It is **very** important to distinguish between three “empty” objects. Many beginning students have trouble keeping these straight.

- $\emptyset$  is the empty *language*, which is a set containing zero strings.  $\emptyset$  is not a string.
- $\{\varepsilon\}$  is a language containing exactly one string, which has length zero.  $\{\varepsilon\}$  is not empty, and it is not a string.
- $\varepsilon$  is the empty *string*, which is a sequence of length zero.  $\varepsilon$  is not a language.

## 2.2 Building Languages

Languages can be combined and manipulated just like any other sets. Thus, if  $A$  and  $B$  are languages over  $\Sigma$ , then their union  $A \cup B$ , intersection  $A \cap B$ , difference  $A \setminus B$ , and symmetric difference  $A \oplus B$  are also languages over  $\Sigma$ , as is the complement  $\bar{A} := \Sigma^* \setminus A$ . However, there are two more useful operators that are specific to sets of *strings*.

The **concatenation** of two languages  $A$  and  $B$ , again denoted  $A \bullet B$  or just  $AB$ , is the set of all strings obtained by concatenating an arbitrary string in  $A$  with an arbitrary string in  $B$ :

$$A \bullet B := \{xy \mid x \in A \text{ and } y \in B\}.$$

For example, if  $A = \{\text{HOCUS}, \text{ABRACA}\}$  and  $B = \{\text{POCUS}, \text{DABRA}\}$ , then

$$A \bullet B = \{\text{HOCUSPOCUS}, \text{ABRACAPOCUS}, \text{HOCUSDABRA}, \text{ABRACADABRA}\}.$$

In particular, for every language  $A$ , we have

$$\emptyset \bullet A = A \bullet \emptyset = \emptyset \quad \text{and} \quad \{\varepsilon\} \bullet A = A \bullet \{\varepsilon\} = A.$$

The **Kleene closure** or **Kleene star**<sup>2</sup> of a language  $L$ , denoted  $L^*$ , is the set of all strings obtained by concatenating a sequence of zero or more strings from  $L$ . For example,  $\{0, 11\}^* = \{\varepsilon, 0, 11, 000, 011, 110, 0000, 0011, 0110, 1100, 1111, 00000, 00011, 011110011011, \dots\}$ . More formally,  $L^*$  is defined recursively as the set of all strings  $w$  such that either

- $w = \varepsilon$ , or
- $w = xy$ , for some strings  $x \in L$  and  $y \in L^*$ .

This definition immediately implies that

$$\emptyset^* = \{\varepsilon\}^* = \{\varepsilon\}.$$

For any other language  $L$ , the Kleene closure  $L^*$  is infinite and contains arbitrarily long (but *finite*!) strings. Equivalently,  $L^*$  can also be defined as the smallest superset of  $L$  that contains the empty string  $\varepsilon$  and is closed under concatenation (hence “closure”). The set of all strings  $\Sigma^*$  is, just as the notation suggests, the Kleene closure of the alphabet  $\Sigma$  (where each symbol is viewed as a string of length 1).

A useful variant of the Kleene closure operator is the **Kleene plus**, defined as  $L^+ := L \bullet L^*$ . Thus,  $L^+$  is the set of all strings obtained by concatenating a sequence of **one** or more strings from  $L$ .

The following identities, which we state here without (easy) proofs, are useful for designing, simplifying, and understanding languages.

---

<sup>2</sup>named after logician Stephen Cole Kleene, who actually pronounced his last name “**clay**-knee”, not “clean” or “cleanie” or “claynuh” or “dimaggio”.

**Lemma 2.1.** *The following identities hold for all languages  $A$ ,  $B$ , and  $C$ :*

- (a)  $A \cup B = B \cup A$ .
- (b)  $(A \cup B) \cup C = A \cup (B \cup C)$ .
- (c)  $\emptyset \cdot A = A \cdot \emptyset = \emptyset$ .
- (d)  $\{\varepsilon\} \cdot A = A \cdot \{\varepsilon\} = A$ .
- (e)  $(A \cdot B) \cdot C = A \cdot (B \cdot C)$ .
- (f)  $A \cdot (B \cup C) = (A \cdot B) \cup (A \cdot C)$ .
- (g)  $(A \cup B) \cdot C = (A \cdot C) \cup (B \cdot C)$ .

**Lemma 2.2.** *The following identities hold for every language  $L$ :*

- (a)  $L^* = \{\varepsilon\} \cup L^+ = L^* \cdot L^* = (L \cup \{\varepsilon\})^* = (L \setminus \{\varepsilon\})^* = \{\varepsilon\} \cup L \cup (L^+ \cdot L^+)$ .
- (b)  $L^+ = L \cdot L^* = L^* \cdot L = L^+ \cdot L^* = L^* \cdot L^+ = L \cup (L^+ \cdot L^+)$ .
- (c)  $L^+ = L^*$  if and only if  $\varepsilon \in L$ .

**Lemma 2.3 (Arden's Rule).** *For any languages  $A$ ,  $B$ , and  $L$  such that  $L = A \cdot L \cup B$ , we have  $A^* \cdot B \subseteq L$ . Moreover, if  $A$  does not contain the empty string, then  $L = A \cdot L \cup B$  if and only if  $L = A^* \cdot B$ .*

## 2.3 Regular Languages and Regular Expressions

Intuitively, a language is **regular** if it can be constructed from individual strings using any combination of union, concatenation, and unbounded repetition. More formally, a language  $L$  is regular if and only if it satisfies one of the following (recursive) conditions:

- $L$  is empty;
- $L$  contains exactly one string (which could be the empty string  $\varepsilon$ );
- $L$  is the union of two regular languages;
- $L$  is the concatenation of two regular languages; or
- $L$  is the Kleene closure of a regular language.

Regular languages are normally described using a compact notation called **regular expressions**, which omit braces around one-string sets, use  $+$  to represent union instead of  $\cup$ , and juxtapose subexpressions to represent concatenation instead of using an explicit operator  $\cdot$ . By convention, in the absence of parentheses, the  $*$  operator has highest precedence, followed by the (implicit) concatenation operator, followed by  $+$ .

For example, the regular expression  $10^*$  is shorthand for the language  $\{1\} \cdot \{0\}^*$  (containing all strings consisting of a **1** followed by zero or more **0**s), and *not* the language  $\{10\}^*$  (containing all strings of even length that start with **1** and alternate between **1** and **0**). As a larger example, the regular expression

$$0^*0 + 0^*1(10^*1 + 01^*0)^*10^*$$

represents the language

$$(\{0\}^* \cdot \{0\}) \cup (\{0\}^* \cdot \{1\} \cdot ((\{1\} \cdot \{0\}^* \cdot \{1\}) \cup (\{0\} \cdot \{1\}^* \cdot \{0\}))^* \cdot \{1\} \cdot \{0\}^*).$$



Most of the time we do not distinguish between regular expressions and the languages they represent, for the same reason that we do not normally distinguish between the arithmetic expression “ $2+2$ ” and the integer 4, or the symbol  $\pi$  and the area of the unit circle. However, we sometimes need to refer to regular expressions themselves *as strings*. In those circumstances, we write  $L(R)$  to denote the language represented by the regular expression  $R$ . String  $w$  *matches* regular expression  $R$  if and only if  $w \in L(R)$ .

Here are several more examples of regular expressions and the languages they represent.

- $0^*$  — the set of all strings of 0s, including the empty string.
- $00000^*$  — the set of all strings consisting of at least four 0s.
- $(00000)^*$  — the set of all strings of 0s whose length is a multiple of 5.
- $(0 + 1)^*$  — the set of all binary strings.
- $(\epsilon + 1)(01)^*(\epsilon + 0)$  — the set of all strings of alternating 0s and 1s, or equivalently, the set of all binary strings that do not contain the substrings 00 or 11.
- $(0 + 1)^*0000(0 + 1)^*$  — the set of all binary strings that contain the substring 0000.
- $((\epsilon + 0 + 00 + 000)1)^*(\epsilon + 0 + 00 + 000)$  — the set of all binary strings that do not contain the substring 0000.
- $((0 + 1)(0 + 1))^*$  — the set of all binary strings whose length is even.
- $1^*(01^*01^*)^*$  — the set of all binary strings with an even number of 0s.
- $0 + 1(0 + 1)^*00$  — the set of all non-negative binary numerals divisible by 4 and with no redundant leading 0s.
- $0^*0 + 0^*1(10^*1 + 01^*0)^*10^*$  — the set of all non-negative binary numerals divisible by 3, possibly with redundant leading 0s.

The last example should *not* be obvious. It is straightforward, but *really* tedious, to prove by induction that every string in  $0^*0 + 0^*1(10^*1 + 01^*0)^*10^*$  is the binary representation of a non-negative multiple of 3. It is similarly straightforward, but *even more* tedious, to prove that the binary representation of *every* non-negative multiple of 3 matches this regular expression. In a later note, we will see a systematic method for deriving regular expressions for some languages that avoids (or more accurately, *automates*) this tedium.

Two regular expressions  $R$  and  $R'$  are *equivalent* if they describe the same language. For example, the regular expressions  $(0 + 1)^*$  and  $(1 + 0)^*$  are equivalent, because the union operator is commutative. More subtly, the regular expressions  $(0 + 1)^*$  and  $(0^*1^*)^*$  and  $(00 + 01 + 10 + 11)^*(0 + 1 + \epsilon)$  are all equivalent; intuitively, these three expressions represent *different ways of thinking about* the language  $\{0, 1\}^*$ . In fact, almost every regular language can be represented by infinitely many distinct but equivalent regular expressions, even if we ignore ultimately trivial equivalences like  $L = (L\emptyset)^*L\epsilon + \emptyset$ .

Give some examples of designing regular expressions via Arden's rule. For example: All strings that don't contain the substring 0011.

## 2.4 Things What Ain't Regular Expressions

Many computing environments and programming languages support patterns called *regexen* (singular *regex*, pluralized like *ox*) that are considerably more general and powerful than regular expressions. Regexen include special symbols representing negation, character classes (for example, upper-case letters, or digits), contiguous ranges of characters, line and word boundaries, limited repetition (as opposed to the unlimited repetition allowed by *\**), back-references to earlier subexpressions, and even local variables. Despite its obvious etymology, a regex is *not* necessarily a regular expression, and it does *not* necessarily describe a regular language!<sup>3</sup>

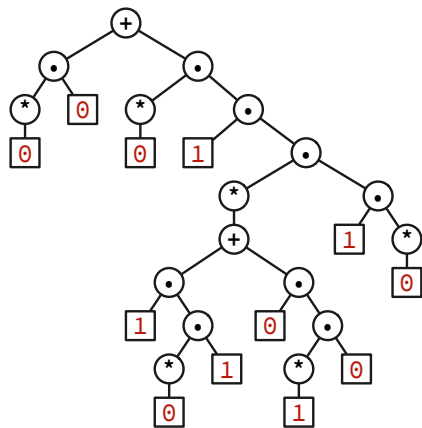
Another type of pattern that is often confused with regular expression are *globs*, which are patterns used in most Unix shells and some scripting languages to represent sets file names. Globbs include symbols for arbitrary single characters (*?*), single characters from a specified range (*[a-z]*), arbitrary substrings (*\**), and substrings from a specified finite set (*{foo,ba{r,z}}*). Globbs are significantly *less* powerful than regular expressions.

## 2.5 Regular Expression Trees

Regular expressions are convenient notational shorthand for a more explicit representation of regular languages called *regular expression trees*. A regular expression tree is formally defined as one of the following:

- A leaf node labeled  $\emptyset$ .
- A leaf node labeled with a string in  $\Sigma^*$ .
- A node labeled  $+$  with two children, each of which is the root of a regular expression tree.
- A node labeled  $\cdot$  with two children, each of which is the root of a regular expression tree.
- A node labeled  $*$  with one child, which is the root of a regular expression tree.

These cases mirror the definition of regular language exactly. A leaf labeled  $\emptyset$  represents the empty language; a leaf labeled with a string represents the language containing only that string; a node labeled  $+$  represents the union of the languages represented by its two children; a node labeled  $\cdot$  represents the concatenation of the languages represented by its two children; and a node labeled  $*$  represents the Kleene closure of the languages represented by its child.



A regular expression tree for  $0^*0 + 0^*1(10^*1 + 01^*0)^*10^*$

<sup>3</sup>However, regexen are not all-powerful, either; see <http://stackoverflow.com/a/1732454/775369>.

The **size** of a regular expression is the number of nodes in its regular expression tree. The size of a regular expression could be either larger or smaller than its length as a raw string. On the one hand, concatenation nodes in the tree are not represented by symbols in the string; on the other hand, parentheses in the string are not represented by nodes in the tree. For example, the regular expression  $0^*0 + 0^*1(10^*1 + 01^*0)^*10^*$  has size 29, but the corresponding raw string  $0^*0+0^*1(10^*1+01^*0)^*10^*$  has length 22.

A **subexpression** of a regular expression  $R$  is another regular expression  $S$  whose regular expression tree is a subtree of some regular expression tree for  $R$ . A **proper subexpression** of  $R$  is any subexpression except  $R$  itself. Every subexpression of  $R$  is also a substring of  $R$ , but not every substring is a subexpression. For example, the substring  $10^*1$  is a proper subexpression of  $0^*0 + 0^*1(10^*1 + 01^*0)^*10^*$ . However, the substrings  $0^*0 + 0^*1$  and  $0^*1 + 01^*$  are *not* subexpressions of  $0^*0 + 0^*1(10^*1 + 01^*0)^*10^*$ , even though they are well-formed regular expressions.

## 2.6 Proofs about Regular Expressions

The standard strategy for proving properties of regular expressions, just as for any other recursively defined structure, to argue inductively on the recursive *structure* of the expression, rather than considering the regular expression as a raw string.

Induction proofs about regular expressions follow a standard boilerplate that mirrors the recursive definition of regular languages. Suppose we want to prove that every regular expression is **perfectly cromulent**, whatever that means. The white boxes hide additional proof details that, among other things, depend on the precise definition of “**perfectly cromulent**”. The boilerplate structure is *longer* than the boilerplate for string induction proofs, but don’t be fooled into thinking it’s *harder*. The five cases in the proof mirror the five cases in the recursive definition of regular language.

**Proof:** Let  $R$  be an arbitrary regular expression.

Assume that **every proper subexpression of  $R$  is perfectly cromulent**.

There are five cases to consider.

- Suppose  $R = \emptyset$ .

Therefore,  $R$  is perfectly cromulent.

- Suppose  $R$  is a single string.

Therefore,  $R$  is perfectly cromulent.

- Suppose  $R = S + T$  for some regular expressions  $S$  and  $T$ .

The induction hypothesis implies that  $S$  and  $T$  are perfectly cromulent.

Therefore,  $R$  is perfectly cromulent.

- Suppose  $R = S \cdot T$  for some regular expressions  $S$  and  $T$ .

The induction hypothesis implies that  $S$  and  $T$  are perfectly cromulent.

Therefore,  $R$  is perfectly cromulent.

- Suppose  $R = S^*$  for some regular expression.  $S$ .

The induction hypothesis implies that  $S$  is perfectly cromulent.

Therefore,  $R$  is perfectly cromulent.

In all cases, we conclude that  $w$  is perfectly cromulent.

□

Students uncomfortable with structural induction can instead induct on the size of the regular expression (defined as the number of nodes in the corresponding regular expression tree). This variant changes only the statement of inductive hypothesis, not the structure of the proof itself; the rest of the boilerplate is utterly identical.

**Proof:** Let  $R$  be an arbitrary regular expression.

Assume that **every regular expression smaller than  $R$**  is perfectly cromulent.

There are five cases to consider.

⋮

In all cases, we conclude that  $w$  is perfectly cromulent.

□

Here is an example of the structural induction boilerplate in action. Again, this proof is *longer* than a typical induction proof about strings or integers, but each individual case is still just a short exercise in definition-chasing.

**Lemma 2.4.** *Every regular expression that does not use the symbol  $\emptyset$  represents a non-empty language.*

**Proof:** Let  $R$  be an arbitrary regular expression that does not use the symbol  $\emptyset$ . Assume that every proper subexpression of  $R$  that does not use the symbol  $\emptyset$  represents a non-empty language. There are five cases to consider, mirroring the definition of  $R$ .

- If  $R = \emptyset$ , we have a contradiction; we can ignore this case.
- If  $R$  is a single string  $w$ , then  $L(R)$  contains the string  $w$ . (In fact,  $L(R) = \{w\}$ .)
- Suppose  $R = S + T$  for some regular expressions  $S$  and  $T$ .  
 $S$  does not use the symbol  $\emptyset$ , because otherwise  $R$  would.  
Thus, the inductive hypothesis implies that  $L(S)$  is non-empty.  
Choose an arbitrary string  $x \in L(S)$ .  
Then  $L(R) = L(S + T) = L(S) \cup L(T)$  also contains the string  $x$ .
- Suppose  $R = S \cdot T$  for some regular expressions  $S$  and  $T$ .  
Neither  $S$  nor  $T$  uses the symbol  $\emptyset$ , because otherwise  $R$  would.  
Thus, the inductive hypothesis implies that both  $L(S)$  and  $L(T)$  are non-empty.  
Choose arbitrary strings  $x \in L(S)$  and  $y \in L(T)$ .  
Then  $L(R) = L(S \cdot T) = L(S) \cdot L(T)$  contains the string  $xy$ .
- Suppose  $R = S^*$  for some regular expression  $S$ .  
Then  $L(R)$  contains the empty string  $\varepsilon$ .

In every case, the language  $L(R)$  is non-empty.

□

Similarly, most algorithms that accept regular expressions as input actually require regular expression *trees*, rather than regular expressions as raw *strings*. Fortunately, it is possible to **parse** any regular expression of length  $n$  into an equivalent regular expression tree in  $O(n)$  time. (The details of the parsing algorithm are beyond the scope of this chapter.) Thus, when we see an algorithmic problem that starts “Given a regular expression. . .”, we can assume without loss of generality that we are actually given a regular expression *tree*.

## 2.7 Proofs about Regular Languages

The same boilerplate also applies to inductive arguments about properties of regular *languages*. Languages themselves are just unadorned sets; they don’t have any recursive structure that we build an inductive proof around. In particular, proper subsets of an *infinite* language  $L$  are not necessarily “smaller” than  $L$ ! Rather than trying to argue directly about an arbitrary regular language  $L$ , we choose an arbitrary regular expression that represents  $L$ , and then build our inductive argument around *the recursive structure of that regular expression*.

**Lemma 2.5.** *Every non-empty regular language is represented by a regular expression that does not use the symbol  $\emptyset$ .*

**Proof:** Let  $R$  be an arbitrary regular expression; we need to prove that either  $L(R) = \emptyset$  or  $L(R) = L(R')$  for some  $\emptyset$ -free regular expression  $R'$ . For every proper subexpression  $S$  of  $R$ , assume that either  $L(S) = \emptyset$  or  $L(S) = L(S')$  for some  $\emptyset$ -free regular expression  $S'$ . There are five cases to consider, mirroring the definition of  $R$ .

- If  $R = \emptyset$ , then  $L(R) = \emptyset$ .
- If  $R$  is a single string  $w$ , then  $R$  is already  $\emptyset$ -free.
- Suppose  $R = S + T$  for some regular expressions  $S$  and  $T$ . There are four subcases to consider:
  - If  $L(S) = L(T) = \emptyset$ , then  $L(R) = L(S) \cup L(T) = \emptyset$ .
  - Suppose  $L(S) \neq \emptyset$  and  $L(T) = \emptyset$ . The inductive hypothesis implies that there is a  $\emptyset$ -free regular expression  $S'$  such that  $L(S') = L(S) = L(S) \cup L(T) = L(R)$ .
  - Suppose  $L(S) = \emptyset$  and  $L(T) \neq \emptyset$ . The inductive hypothesis implies that there is a  $\emptyset$ -free regular expression  $T'$  such that  $L(T') = L(T) = L(S) \cup L(T) = L(R)$ .
  - Finally, suppose  $L(S) \neq \emptyset$  and  $L(T) \neq \emptyset$ . The inductive hypothesis implies that there are  $\emptyset$ -free regular expressions  $S'$  and  $T'$  such that  $L(S') = L(S)$  and  $L(T') = L(T)$ . The regular expression  $S' + T'$  is  $\emptyset$ -free and  $L(S' + T') = L(S') \cup L(T') = L(S) \cup L(T) = L(S + T) = L(R)$ .
- Suppose  $R = S \cdot T$  for some regular expressions  $S$  and  $T$ . There are two subcases to consider.
  - If either  $L(S) = \emptyset$  or  $L(T) = \emptyset$  then  $L(R) = L(S) \cdot L(T) = \emptyset$ .
  - Otherwise, the inductive hypothesis implies that there are  $\emptyset$ -free regular expressions  $S'$  and  $T'$  such that  $L(S') = L(S)$  and  $L(T') = L(T)$ . The regular expression  $S' \cdot T'$  is  $\emptyset$ -free and equivalent to  $R$ .
- Suppose  $R = S^*$  for some regular expression  $S$ . There are two subcases to consider.

- If  $L(S) = \emptyset$ , then  $L(R) = \emptyset^* = \{\varepsilon\}$ , so  $R$  is represented by the  $\emptyset$ -free regular expression  $\varepsilon$ .
- Otherwise, The inductive hypothesis implies that there is a  $\emptyset$ -free regular expression  $S'$  such that  $L(S') = L(S)$ . The regular expression  $(S')^*$  is  $\emptyset$ -free and equivalent to  $R$ .

In all cases, either  $L(R) = \emptyset$  or  $R$  is equivalent to some  $\emptyset$ -free regular expression  $R'$ .  $\square$

## 2.8 Not Every Language is Regular

You may be tempted to conjecture that *all* languages are regular, but in fact, the following cardinality argument *almost all* languages are *not* regular. To make the argument concrete, let's consider languages over the single-symbol alphabet  $\{\diamond\}$ .

- Every regular expression over the one-symbol alphabet  $\{\diamond\}$  is itself a string over the seven-symbol alphabet  $\{\diamond, +, (, ), *, \varepsilon, \emptyset\}$ . By interpreting these symbols as the digits 1 through 7, we can interpret any string over this larger alphabet as the base-8 representation of some unique integer. Thus, the set of all regular expressions over  $\{\diamond\}$  is *at most* as large as the set of integers, and is therefore countably infinite. It follows that the set of all regular languages over  $\{\diamond\}$  is also countably infinite.
- On the other hand, for any real number  $0 \leq \alpha < 1$ , we can define a corresponding language

$$L_\alpha = \{\diamond^n \mid \alpha 2^n \bmod 1 \geq 1/2\}.$$

In other words,  $L_\alpha$  contains the string  $\diamond^n$  if and only if the  $(n+1)$ th bit in the binary representation of  $\alpha$  is equal to 1. For any distinct real numbers  $\alpha \neq \beta$ , the binary representations of  $\alpha$  and  $\beta$  must differ in some bit, so  $L_\alpha \neq L_\beta$ . We conclude that the set of *all* languages over  $\{\diamond\}$  is *at least* as large as the set of real numbers between 0 and 1, and is therefore uncountably infinite.

We will see several explicit examples of non-regular languages in later lectures. In particular, the set of all regular expressions over the alphabet  $\{0, 1\}$  is itself a non-regular language over the alphabet  $\{0, 1, +, (, ), *, \varepsilon, \emptyset\}$ !

## Exercises

- (a) Prove that  $\emptyset \cdot L = L \cdot \emptyset = \emptyset$ , for every language  $L$ .
- (b) Prove that  $\{\varepsilon\} \cdot L = L \cdot \{\varepsilon\} = L$ , for every language  $L$ .
- (c) Prove that  $(A \cdot B) \cdot C = A \cdot (B \cdot C)$ , for all languages  $A$ ,  $B$ , and  $C$ .
- (d) Prove that  $|A \cdot B| \leq |A| \cdot |B|$ , for all languages  $A$  and  $B$ . (The second  $\cdot$  is multiplication!)
  - Describe two languages  $A$  and  $B$  such that  $|A \cdot B| < |A| \cdot |B|$ .
  - Describe two languages  $A$  and  $B$  such that  $|A \cdot B| = |A| \cdot |B|$ .
- (e) Prove that  $L^*$  is finite if and only if  $L = \emptyset$  or  $L = \{\varepsilon\}$ .
- (f) Prove that if  $A \cdot B = B \cdot C$ , then  $A^* \cdot B = B \cdot C^* = A^* \cdot B \cdot C^*$ , for all languages  $A$ ,  $B$ , and  $C$ .
- (g) Prove that  $(A \cup B)^* = (A^* \cdot B^*)^*$ , for all languages  $A$  and  $B$ .

2. Recall that the reversal  $w^R$  of a string  $w$  is defined recursively as follows:

$$w^R := \begin{cases} \varepsilon & \text{if } w = \varepsilon \\ x^R \cdot a & \text{if } w = a \cdot x \end{cases}$$

The reversal  $L^R$  of any language  $L$  is the set of reversals of all strings in  $L$ :

$$L^R := \{w^R \mid w \in L\}.$$

- (a) Prove that  $(A \cdot B)^R = B^R \cdot A^R$  for all languages  $A$  and  $B$ .
  - (b) Prove that  $(L^R)^R = L$  for every language  $L$ .
  - (c) Prove that  $(L^*)^R = (L^R)^*$  for every language  $L$ .
3. Prove that each of the following regular expressions is equivalent to  $(0 + 1)^*$ .
- (a)  $\varepsilon + 0(0 + 1)^* + 1(1 + 0)^*$
  - (b)  $0^* + 0^*1(0 + 1)^*$
  - (c)  $((\varepsilon + 0)(\varepsilon + 1))^*$
  - (d)  $0^*(10^*)^*$
  - (e)  $(1^*0)^*(0^*1)^*$
4. For each of the following languages in  $\{0, 1\}^*$ , describe an equivalent regular expression. There are infinitely many correct answers for each language. (This problem will become significantly simpler after we've seen finite-state machines.)
- (a) Strings that end with the suffix  $0^9 = 000000000$ .
  - (b) All strings except  $010$ .
  - (c) Strings that contain the substring  $010$ .
  - (d) Strings that contain the subsequence  $010$ .
  - (e) Strings that do not contain the substring  $010$ .
  - (f) Strings that do not contain the subsequence  $010$ .
  - (g) Strings that contain an even number of occurrences of the substring  $010$ .
  - \* (h) Strings that contain an even number of occurrences of the substring  $000$ .
  - (i) Strings in which every occurrence of the substring  $00$  appears before every occurrence of the substring  $11$ .
  - (j) Strings  $w$  such that *in every prefix of  $w$* , the number of  $0$ s and the number of  $1$ s differ by at most 1.
  - \* (k) Strings  $w$  such that *in every prefix of  $w$* , the number of  $0$ s and the number of  $1$ s differ by at most 2.
  - \* (l) Strings in which the number of  $0$ s and the number of  $1$ s differ by a multiple of 3.
  - \* (m) Strings that contain an even number of  $1$ s and an odd number of  $0$ s.
  - ★ (n) Strings that represent a number divisible by 5 in binary.

5. For any string  $w$ , define  $\text{stutter}(w)$  as follows:

$$\text{stutter}(w) := \begin{cases} \varepsilon & \text{if } w = \varepsilon \\ aa \cdot \text{stutter}(x) & \text{if } w = ax \text{ for some symbol } a \text{ and string } x \end{cases}$$

Let  $L$  be an arbitrary regular language.

(a) Prove that the language  $\text{stutter}(L) = \{\text{stutter}(w) \mid w \in L\}$  is also regular.

\*(b) Prove that the language  $\text{stutter}^{-1}(L) = \{w \mid \text{stutter}(w) \in L\}$  is also regular.

6. Recall that the **reversal**  $w^R$  of a string  $w$  is defined recursively as follows:

$$w^R = \begin{cases} \varepsilon & \text{if } w = \varepsilon \\ x \cdot a & \text{if } w = ax \text{ for some symbol } a \text{ and some string } x \end{cases}$$

The reversal  $L^R$  of a language  $L$  is defined as the set of reversals of all strings in  $L$ :

$$L^R := \{w^R \mid w \in L\}$$

(a) Prove that  $(L^*)^R = (L^R)^*$  for every language  $L$ .

(b) Prove that the reversal of any regular language is also a regular language. (You may assume part (a) even if you haven't proved it yet.)

You may assume the following facts without proof:

- $L^* \cdot L^* = L^*$  for every language  $L$ .
- $(w^R)^R = w$  for every string  $w$ .
- $(x \cdot y)^R = y^R \cdot x^R$  for all strings  $x$  and  $y$ .

[Hint: Yes, all three proofs use induction, but induction on what? And yes, all **three** proofs.]

7. This problem considers two special classes of regular expressions.

- A regular expression  $R$  is **plus-free** if and only if it never uses the  $+$  operator.
- A regular expression  $R$  is **top-plus** if and only if either
  - $R$  is plus-free, or
  - $R = S + T$ , where  $S$  and  $T$  are top-plus.

For example,  $1((0^*10)^*1)^*0$  is plus-free and (therefore) top-plus;  $01^*0 + 10^*1 + \varepsilon$  is top-plus but not plus-free, and  $0(0 + 1)^*(1 + \varepsilon)$  is neither top-plus nor plus-free.

Recall that two regular expressions  $R$  and  $S$  are **equivalent** if they describe exactly the same language:  $L(R) = L(S)$ .

(a) Prove that for any top-plus regular expressions  $R$  and  $S$ , there is a top-plus regular expression that is equivalent to  $RS$ .

(b) Prove that for any top-plus regular expression  $R$ , there is a **plus-free** regular expression  $S$  such that  $R^*$  and  $S^*$  are equivalent.



- (c) Prove that for any regular expression, there is an equivalent top-plus regular expression.

You may assume the following facts without proof, for all regular expressions  $A$ ,  $B$ , and  $C$ :

- $A(B + C)$  is equivalent to  $AB + AC$ .
- $(A + B)C$  is equivalent to  $AC + BC$ .
- $(A + B)^*$  is equivalent to  $(A^*B^*)^*$ .

8. (a) Describe and analyze an efficient algorithm to determine, given a regular expression  $R$ , whether  $L(R) = \emptyset$ .
- (b) Describe and analyze an efficient algorithm to determine, given a regular expression  $R$ , whether  $L(R) = \{\varepsilon\}$ . [Hint: Use part (a).]
- (c) Describe and analyze an efficient algorithm to determine, given a regular expression  $R$ , whether  $L(R)$  is finite. [Hint: Use parts (a) and (b).]

In each problem, assume you are given  $R$  as a regular expression *tree*, not just a raw string.

## 3 Finite-State Machines

### 3.1 Intuition

Suppose we want to determine whether a given string  $w[1..n]$  of bits represents a multiple of 5 in binary. After a bit of thought, you might realize that you can read the bits in  $w$  one at a time, from left to right, keeping track of the value modulo 5 of the prefix you have read so far.

```
MULTIPLEOF5( $w[1..n]$ ):
   $rem \leftarrow 0$ 
  for  $i \leftarrow 1$  to  $n$ 
     $rem \leftarrow (2 \cdot rem + w[i]) \bmod 5$ 
  if  $rem = 0$ 
    return TRUE
  else
    return FALSE
```

Aside from the loop index  $i$ , which we need just to read the entire input string, this algorithm has a single local variable  $rem$ , which has only four different values: 0, 1, 2, 3, or 4.

This algorithm already runs in  $O(n)$  time, which is the best we can hope for—after all, we have to read every bit in the input—but we can speed up the algorithm *in practice*. Let's define a **change** or **transition** function  $\delta: \{0, 1, 2, 3, 4\} \times \{\textcolor{red}{0}, \textcolor{red}{1}\} \rightarrow \{0, 1, 2, 3, 4\}$  as follows:

$$\delta(q, a) = (2q + a) \bmod 5.$$

(Here I'm implicitly converting the symbols  $\textcolor{red}{0}$  and  $\textcolor{red}{1}$  to the corresponding integers 0 and 1.) Since we already know all values of the transition function, we can store them in a precomputed table, and then replace the computation in the main loop of MULTIPLEOF5 with a simple array lookup.

We can also modify the return condition to check for different values modulo 5. To be completely general, we replace the final if-then-else lines with another array lookup, using an array  $A[0..4]$  of booleans describing which final mod-5 values are “acceptable”.

After both of these modifications, our algorithm looks like one of the following, depending on whether we want something iterative or recursive (with  $q = 0$  in the initial call):

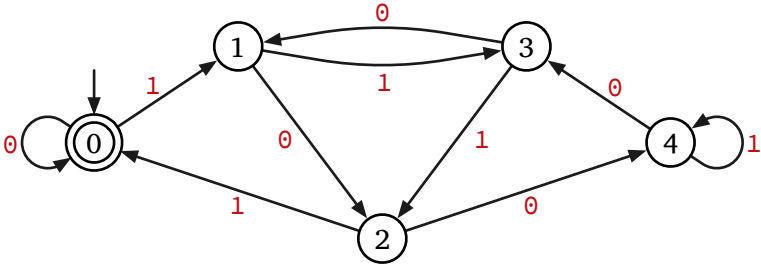
```
DOSOMETHINGCOOL( $w[1..n]$ ):
   $q \leftarrow 0$ 
  for  $i \leftarrow 1$  to  $n$ 
     $q \leftarrow \delta[q, w[i]]$ 
  return  $A[q]$ 
```

```
DOSOMETHINGCOOL( $q, w$ ):
  if  $w = \epsilon$ 
    return  $A[q]$ 
  else
    decompose  $w = a \cdot x$ 
    return DOSOMETHINGCOOL( $\delta(q, a), x$ )
```

If we want to use our new DoSOMETHINGCOOL algorithm to implement MULTIPLEOF5, we simply give the arrays  $\delta$  and  $A$  the following hard-coded values:

$q$	$\delta[q, 0]$	$\delta[q, 1]$	$A[q]$
0	0	1	TRUE
1	2	3	FALSE
2	4	0	FALSE
3	1	2	FALSE
4	3	4	FALSE

We can also visualize the behavior of DoSOMETHINGCOOL by drawing a directed graph, whose vertices represent possible values of the variable  $q$ —the possible *states* of the algorithm—and whose edges are labeled with input symbols to represent transitions between states. Specifically, the graph includes the labeled directed edge  $p \xrightarrow{a} q$  if and only if  $\delta(p, a) = q$ . To indicate the proper return value, we draw the “acceptable” final states using doubled circles. Here is the resulting graph for MULTIPLEOF5:



State-transition graph for MULTIPLEOF5

If we run the MULTIPLEOF5 algorithm on the string 00101110110 (representing the number 374 in binary), the algorithm performs the following sequence of transitions:

$$0 \xrightarrow{0} 0 \xrightarrow{0} 0 \xrightarrow{1} 1 \xrightarrow{0} 2 \xrightarrow{1} 0 \xrightarrow{1} 1 \xrightarrow{1} 3 \xrightarrow{0} 1 \xrightarrow{1} 3 \xrightarrow{1} 2 \xrightarrow{0} 4$$

Because the final state is not the “acceptable” state 0, the algorithm correctly returns FALSE. We can also think of this sequence of transitions as a walk in the graph, which is completely determined by the start state 0 and the sequence of edge labels; the algorithm returns TRUE if and only if this walk ends at an “acceptable” state.

### 3.2 Formal Definitions

The object we have just described is an example of a *finite-state machine*. A finite-state machine is a formal model of any system/machine/algorithm that can exist in a finite number of *states* and that transitions among those states based on sequence of *input* symbols.

Finite-state machines are also known as *deterministic finite-state automata*, abbreviated *DFAs*. The word “deterministic” means that the behavior of the machine is completely *determined* by the input string; we’ll discuss nondeterministic automata in the next lecture. The word “automaton” (the singular of “automata”) comes from ancient Greek αὐτόματος meaning “self-acting”, from the roots αὐτο- (“self”) and -ματος (“thinking, willing”, the root of Latin *mentus*).

Formally, every finite-state machine consists of five components:

- An arbitrary finite set  $\Sigma$ , called the *input alphabet*.

- Another arbitrary finite set  $Q$ , whose elements are called **states**.<sup>1</sup>
- An arbitrary **transition** function  $\delta : Q \times \Sigma \rightarrow Q$ .
- A **start state**  $s \in Q$ .
- A subset  $A \subseteq Q$  of **accepting states**.

The behavior of a finite-state machine is governed by an **input string**  $w$ , which is a finite sequence of symbols from the input alphabet  $\Sigma$ . The machine **reads** the symbols in  $w$  one at a time in order (from left to right). At all times, the machine has a **current state**  $q$ ; initially  $q$  is the machine's start state  $s$ . Each time the machine reads a symbol  $a$  from the input string, its current state **transitions** from  $q$  to  $\delta(q, a)$ . After all the characters have been read, the machine **accepts**  $w$  if the current state is in  $A$  and **rejects**  $w$  otherwise. In other words, every finite state machine runs the algorithm `DoSomethingCool!`

More formally, we extend the transition function  $\delta : Q \times \Sigma \rightarrow Q$  of any finite-state machine to a function  $\delta^* : Q \times \Sigma^* \rightarrow Q$  that transitions on *strings* as follows:

$$\delta^*(q, w) := \begin{cases} q & \text{if } w = \varepsilon, \\ \delta^*(\delta(q, a), x) & \text{if } w = ax. \end{cases}$$

Finally, a finite-state machine **accepts** a string  $w$  if and only if  $\delta^*(s, w) \in A$ , and **rejects**  $w$  otherwise. (Compare this definition with the recursive formulation of `DoSomethingCool!`)

For example, our final `MULTIPLEOF5` algorithm is a DFA with the following components:

- input alphabet:  $\Sigma = \{0, 1\}$
- state set:  $Q = \{0, 1, 2, 3, 4\}$
- transition function:  $\delta(q, a) = (2q + a) \bmod 5$
- start state:  $s = 0$
- accepting states:  $A = \{0\}$

This machine rejects the string `00101110110`, because

$$\begin{aligned} \delta^*(0, 00101110110) &= \delta^*(\delta(0, 0), 0101110110) \\ &= \delta^*(0, 0101110110) = \delta^*(\delta(0, 0), 101110110) \\ &= \delta^*(0, 101110110) = \delta^*(\delta(0, 1), 01110110) = \dots \\ &\quad \vdots \\ \dots &= \delta^*(1, 110) = \delta^*(\delta(1, 1), 10) \\ &= \delta^*(3, 10) = \delta^*(\delta(3, 1), 0) \\ &= \delta^*(2, 0) = \delta^*(\delta(3, 0), \varepsilon) \\ &= \delta^*(4, \varepsilon) = 4 \notin A. \end{aligned}$$

---

<sup>1</sup>It's unclear why we use the letter  $Q$  to refer to the state set, and lower-case  $q$  to refer to a generic state, but that is now the firmly-established notational standard. Although the formal study of finite-state automata began much earlier, its modern formulation was established in a 1959 paper by Michael Rabin and Dana Scott, for which they won the Turing award. Rabin and Scott called the set of states  $S$ , used lower-case  $s$  for a generic state, and called the start state  $s_0$ . On the other hand, in the 1936 paper for which the Turing award was named, Alan Turing used  $q_1, q_2, \dots, q_R$  to refer to states (or “ $m$ -configurations”) of a generic Turing machine. Turing may have been mirroring the standard notation  $Q$  for configuration spaces in classical mechanics, also of uncertain origin.

We have already seen a more graphical representation of this entire sequence of transitions:

$$0 \xrightarrow{0} 0 \xrightarrow{0} 0 \xrightarrow{1} 1 \xrightarrow{0} 2 \xrightarrow{1} 0 \xrightarrow{1} 1 \xrightarrow{1} 3 \xrightarrow{0} 1 \xrightarrow{1} 3 \xrightarrow{1} 2 \xrightarrow{0} 4$$

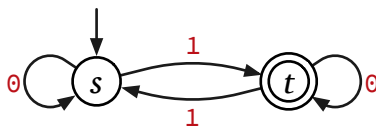
The arrow notation is easier to read and write for specific examples, but surprisingly, most people actually find the more formal functional notation easier to use in formal proofs. Try them both!

We can equivalently define a DFA as a directed graph whose vertices are the states  $Q$ , whose edges are labeled with symbols from  $\Sigma$ , such that every vertex has exactly one outgoing edge with each label. In our drawings of finite state machines, the start state  $s$  is always indicated by an incoming arrow, and the accepting states  $A$  are always indicated by doubled circles. By induction, for any string  $w \in \Sigma^*$ , this graph contains a unique walk that starts at  $s$  and whose edges are labeled with the symbols in  $w$  in order. The machine accepts  $w$  if this walk ends at an accepting state. This graphical formulation of DFAs is incredibly useful for developing intuition and even designing DFAs. For proofs, it's largely a matter of taste whether to write in terms of extended transition functions or labeled graphs, but (as much as I wish otherwise) I actually find it easier to write **correct** proofs using the functional formulation.

### 3.3 Another Example

The following drawing shows a finite-state machine with input alphabet  $\Sigma = \{0, 1\}$ , state set  $Q = \{s, t\}$ , start state  $s$ , a single accepting state  $t$ , and the transition function

$$\delta(s, 0) = s, \quad \delta(s, 1) = t, \quad \delta(t, 0) = t, \quad \delta(t, 1) = s.$$



A simple finite-state machine.

For example, the two-state machine  $M$  at the top of this page accepts the string  $00101110100$  after the following sequence of transitions:

$$s \xrightarrow{0} s \xrightarrow{0} s \xrightarrow{1} t \xrightarrow{0} t \xrightarrow{1} s \xrightarrow{1} t \xrightarrow{1} s \xrightarrow{0} s \xrightarrow{1} t \xrightarrow{0} t \xrightarrow{0} t.$$

The same machine  $M$  rejects the string  $11101101$  after the following sequence of transitions:

$$s \xrightarrow{1} t \xrightarrow{1} s \xrightarrow{1} t \xrightarrow{0} t \xrightarrow{1} s \xrightarrow{1} t \xrightarrow{0} t \xrightarrow{1} s.$$

Finally,  $M$  rejects the empty string, because the start state  $s$  is not an accepting state.

From these examples and others, it is easy to conjecture that the language of  $M$  is the set of all strings of  $0$ s and  $1$ s with an odd number of  $1$ s. So let's prove it!

**Proof (tedious case analysis):** Let  $\#(a, w)$  denote the number of times symbol  $a$  appears in string  $w$ . We will prove the following stronger claims by induction, for any string  $w$ .

$$\delta^*(s, w) = \begin{cases} s & \text{if } \#(1, w) \text{ is even} \\ t & \text{if } \#(1, w) \text{ is odd} \end{cases} \quad \text{and} \quad \delta^*(t, w) = \begin{cases} t & \text{if } \#(1, w) \text{ is even} \\ s & \text{if } \#(1, w) \text{ is odd} \end{cases}$$

Let's begin. Let  $w$  be an arbitrary string. Assume that for any string  $x$  that is shorter than  $w$ , we have  $\delta^*(s, x) = s$  and  $\delta^*(t, x) = t$  if  $x$  has an even number of  $1$ s, and  $\delta^*(s, x) = t$  and  $\delta^*(t, x) = s$  if  $x$  has an odd number of  $1$ s. There are five cases to consider.

- If  $w = \varepsilon$ , then  $w$  contains an even number of  $1$ s and  $\delta^*(s, w) = s$  and  $\delta^*(t, w) = t$  by definition.
- Suppose  $w = 1x$  and  $\#(1, w)$  is even. Then  $\#(1, x)$  is odd, which implies

$$\begin{aligned}
 \delta^*(s, w) &= \delta^*(\delta(s, 1), x) && \text{by definition of } \delta^* \\
 &= \delta^*(t, x) && \text{by definition of } \delta \\
 &= s && \text{by the inductive hypothesis} \\
 \delta^*(t, w) &= \delta^*(\delta(t, 1), x) && \text{by definition of } \delta^* \\
 &= \delta^*(s, x) && \text{by definition of } \delta \\
 &= T && \text{by the inductive hypothesis}
 \end{aligned}$$

Since the remaining cases are similar, I'll omit the line-by-line justification.

- If  $w = 1x$  and  $\#(1, w)$  is odd, then  $\#(1, x)$  is even, so the inductive hypothesis implies

$$\begin{aligned}
 \delta^*(s, w) &= \delta^*(\delta(s, 1), x) = \delta^*(t, x) = t \\
 \delta^*(t, w) &= \delta^*(\delta(t, 1), x) = \delta^*(s, x) = s
 \end{aligned}$$

- If  $w = 0x$  and  $\#(1, w)$  is even, then  $\#(1, x)$  is even, so the inductive hypothesis implies

$$\begin{aligned}
 \delta^*(s, w) &= \delta^*(\delta(s, 0), x) = \delta^*(s, x) = s \\
 \delta^*(t, w) &= \delta^*(\delta(t, 0), x) = \delta^*(t, x) = t
 \end{aligned}$$

- Finally, if  $w = 0x$  and  $\#(1, w)$  is odd, then  $\#(1, x)$  is odd, so the inductive hypothesis implies

$$\begin{aligned}
 \delta^*(s, w) &= \delta^*(\delta(s, 0), x) = \delta^*(s, x) = t \\
 \delta^*(t, w) &= \delta^*(\delta(t, 0), x) = \delta^*(t, x) = s
 \end{aligned}$$

□

Notice that this proof contains  $|Q|^2 \cdot |\Sigma| + |Q|$  separate inductive arguments. For every pair of states  $p$  and  $q$ , we must argue about the language of all strings  $w$  such that  $\delta^*(p, w) = q$ , and we must consider every possible first symbol in  $w$ . We must also argue about  $\delta(p, \varepsilon)$  for every state  $p$ . Each of those arguments is typically straightforward, but it's easy to get lost in the deluge of cases.

For this particular proof, however, we can reduce the number of cases by switching from tail recursion to *head* recursion. The following identity holds for all strings  $x \in \Sigma^*$  and symbols  $a \in \Sigma$ :

$$\delta^*(q, xa) = \delta(\delta^*(q, x), a)$$

We leave the inductive proof of this identity as a straightforward exercise (hint, hint).

**Proof (clever renaming, head induction):** Let's rename the states with the integers 0 and 1 instead of  $s$  and  $t$ . Then the transition function can be described concisely as  $\delta(q, a) = (q + a) \bmod 2$ . We claim that for every string  $w$ , we have  $\delta^*(0, w) = \#(1, w) \bmod 2$ .

Let  $w$  be an arbitrary string, and assume that for any string  $x$  that is shorter than  $w$  that  $\delta^*(0, x) = \#(1, x) \bmod 2$ . There are only two cases to consider: either  $w$  is empty or it isn't.

- If  $w = \varepsilon$ , then  $\delta^*(0, w) = 0 = \#(1, w) \bmod 2$  by definition.

- Otherwise,  $w = xa$  for some string  $x$  and some symbol  $a$ , and we have

$$\begin{aligned}
\delta^*(0, w) &= \delta(\delta^*(0, x), a) && \text{by definition of } \delta^* \\
&= \delta(\#(\mathbf{1}, x) \bmod 2, a) && \text{by the inductive hypothesis} \\
&= (\#(\mathbf{1}, x) \bmod 2 + a) \bmod 2 && \text{by definition of } \delta \\
&= (\#(\mathbf{1}, x) + a) \bmod 2 && \text{by definition of } \bmod 2 \\
&= (\#(\mathbf{1}, x) + \#(\mathbf{1}, a)) \bmod 2 && \text{because } \#(\mathbf{1}, \mathbf{0}) = 0 \text{ and } \#(\mathbf{1}, \mathbf{1}) = 1 \\
&= (\#(\mathbf{1}, xa)) \bmod 2 && \text{by definition of } \# \\
&= (\#(\mathbf{1}, w)) \bmod 2 && \text{because } w = xa \quad \square
\end{aligned}$$

Hmmm. This “clever” proof is certainly shorter than the earlier brute-force proof, but is it actually *better*? Simpler? More intuitive? Easier to understand? I’m skeptical. Sometimes brute force really is more effective.

### 3.4 Real-World Examples

Finite-state machines were first formally defined in the mid-20th century, but people have been building automata for centuries, if not millennia. Many of the earliest records about automata are clearly mythological—for example, the brass giant Talus created by Hephaestus to guard Crete against intruders—but others are more believable, such as King-Shu’s construction of a flying magpie from wood and bamboo in China around 500BCE.

Perhaps the most common examples of finite-state automata are **clocks**. For example, the Swiss railway clock designed by Hans Hilfiker in 1944 has hour and minute hands that can indicate any time between 1:00 and 12:59. The minute hands advance discretely once per minute when they receive an electrical signal from a central master clock.<sup>2</sup> Thus, a Swiss railway clock is a finite-state machine with 720 states, one input symbol, and a simple transition function:

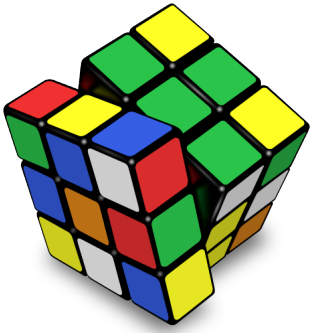
$$\begin{aligned}
Q &= \{(h, m) \mid 0 \leq h < 12 \text{ and } 0 \leq m \leq 59\} \\
\Sigma &= \{\text{tick}\} \\
\delta((h, m), \text{tick}) &= \begin{cases} (h, m + 1) & \text{if } m < 59 \\ (h + 1, 0) & \text{if } h < 11 \text{ and } m = 59 \\ (0, 0) & \text{if } h = 11 \text{ and } m = 59 \end{cases}
\end{aligned}$$

This clock doesn’t *quite* match our abstraction, because there’s no “start” state or “accepting” states, unless perhaps you consider the “accepting” state to be the time when your train arrives.

A more playful example of a finite-state machine is the **Rubik’s cube**, a well-known mechanical puzzle invented independently by Ernő Rubik in Hungary and Terutoshi Ishigi in Japan in the mid-1970s. This puzzle has precisely 519,024,039,293,878,272,000 distinct configurations. In the unique *solved* configuration, each of the six faces of the cube shows exactly one color. We can change the configuration of the cube by rotating one of the six faces of the cube by 90 degrees, either clockwise or counterclockwise. The cube has six faces (front, back, left, right, up, and down), so there are exactly twelve possible turns, typically represented by the symbols **R, L, F, B, U, D,  $\bar{R}$ ,  $\bar{L}$ ,  $\bar{F}$ ,  $\bar{B}$ ,  $\bar{U}$ ,  $\bar{D}$** , where the letter indicates which face to turn and the presence or absence of a bar over the letter

<sup>2</sup>A second hand was added to the Swiss Railway clocks in the mid-1950s, which sweeps continuously around the clock in approximately 58½ seconds and then pauses at 12:00 until the next minute signal “to bring calm in the last moment and ease punctual train departure”. Let’s ignore that.

indicates turning counterclockwise or clockwise, respectively. Thus, we can represent a Rubik’s cube as a finite-state machine with 519,024,039,293,878,272,000 states and an input alphabet with 12 symbols; or equivalently, as a directed graph with 519,024,039,293,878,272,000 vertices, each with 12 outgoing edges. In practice, the number of states is *far* too large for us to actually draw the machine or explicitly specify its transition function; nevertheless, the number of states is still finite. If we let the start state  $s$  and the sole accepting state be the solved state, then the language of this finite state machine is the set of all move sequences that leave the cube unchanged.



Three finite-state machines.

### 3.5 A Brute-Force Design Example

As usual in algorithm design, there is no purely mechanical recipe—no *automatic* method—no *algorithm*—for building DFAs in general. Here I’ll describe one systematic approach that works reasonably well, although it tends to produce DFAs with many more states than necessary.

#### 3.5.1 DFAs are Algorithms

The basic approach is to try to *construct an algorithm* that looks like MULTIPLEOF5: A simple for-loop through the symbols, using a *constant* number of variables, where each variable (except the loop index) has only a *constant* number of possible values. Here, “constant” means an actual number that is not a function of the input size  $n$ . You should be able to compute the number of possible values for each variable *at compile time*.

For example, the following algorithm determines whether a given string in  $\Sigma = \{0, 1\}$  contains the substring **11**.

```

CONTAINS11(w[1..n]):
  found ← FALSE
  for i ← 1 to n
    if i = 1
      last2 ← w[1]
    else
      last2 ← w[i − 1] · w[i]
    if last2 = 11
      found ← TRUE
  return found

```

Aside from the loop index, this algorithm has exactly two variables.



- A boolean flag *found* indicating whether we have seen the substring **11**. This variable has exactly two possible values: **TRUE** and **FALSE**.
- A string *last2* containing the last (up to) three symbols we have read so far. This variable has exactly 7 possible values:  $\varepsilon$ , **0**, **1**, **00**, **01**, **10**, and **11**.

Thus, altogether, the algorithm can be in at most  $2 \times 7 = 14$  possible states, one for each possible pair (*found*, *last2*). Thus, we can encode the behavior of **CONTAINS11** as a DFA with fourteen states, where the start state is (**FALSE**,  $\varepsilon$ ) and the accepting states are all seven states of the form (**TRUE**, \*). The transition function is described in the following table (split into two parts to save space):

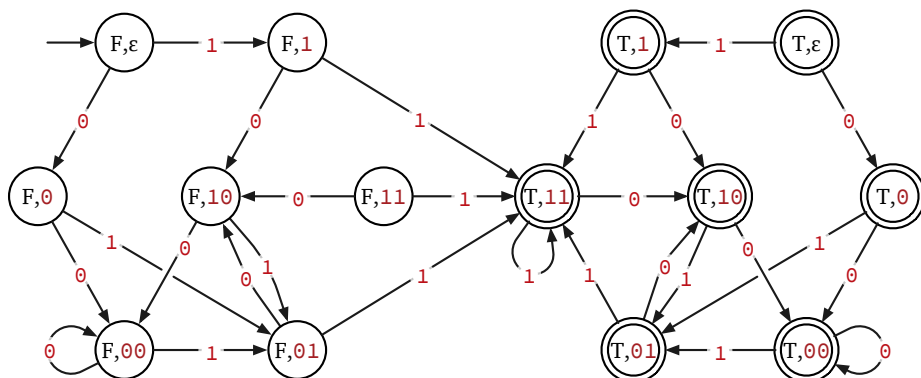
<i>q</i>	$\delta[q, 0]$	$\delta[q, 1]$	<i>q</i>	$\delta[q, 0]$	$\delta[q, 1]$
( <b>FALSE</b> , $\varepsilon$ )	( <b>FALSE</b> , <b>0</b> )	( <b>FALSE</b> , <b>1</b> )	( <b>TRUE</b> , $\varepsilon$ )	( <b>TRUE</b> , <b>0</b> )	( <b>TRUE</b> , <b>1</b> )
( <b>FALSE</b> , <b>0</b> )	( <b>FALSE</b> , <b>00</b> )	( <b>FALSE</b> , <b>01</b> )	( <b>TRUE</b> , <b>0</b> )	( <b>TRUE</b> , <b>00</b> )	( <b>TRUE</b> , <b>01</b> )
( <b>FALSE</b> , <b>1</b> )	( <b>FALSE</b> , <b>10</b> )	( <b>TRUE</b> , <b>11</b> )	( <b>TRUE</b> , <b>1</b> )	( <b>TRUE</b> , <b>10</b> )	( <b>TRUE</b> , <b>11</b> )
( <b>FALSE</b> , <b>00</b> )	( <b>FALSE</b> , <b>00</b> )	( <b>FALSE</b> , <b>01</b> )	( <b>TRUE</b> , <b>00</b> )	( <b>TRUE</b> , <b>00</b> )	( <b>TRUE</b> , <b>01</b> )
( <b>FALSE</b> , <b>01</b> )	( <b>FALSE</b> , <b>10</b> )	( <b>TRUE</b> , <b>11</b> )	( <b>TRUE</b> , <b>01</b> )	( <b>TRUE</b> , <b>10</b> )	( <b>TRUE</b> , <b>11</b> )
( <b>FALSE</b> , <b>10</b> )	( <b>FALSE</b> , <b>00</b> )	( <b>FALSE</b> , <b>01</b> )	( <b>TRUE</b> , <b>10</b> )	( <b>TRUE</b> , <b>00</b> )	( <b>TRUE</b> , <b>01</b> )
( <b>FALSE</b> , <b>11</b> )	( <b>FALSE</b> , <b>10</b> )	( <b>TRUE</b> , <b>11</b> )	( <b>TRUE</b> , <b>11</b> )	( <b>TRUE</b> , <b>10</b> )	( <b>TRUE</b> , <b>11</b> )

For example, given the input string **1001011100**, this DFA performs the following sequence of transitions and then accepts.

(**FALSE**,  $\varepsilon$ )  $\xrightarrow{1}$  (**FALSE**, **1**)  $\xrightarrow{0}$  (**FALSE**, **10**)  $\xrightarrow{0}$  (**FALSE**, **00**)  $\xrightarrow{1}$  (**FALSE**, **01**)  $\xrightarrow{0}$  (**FALSE**, **10**)  $\xrightarrow{1}$  (**FALSE**, **01**)  $\xrightarrow{1}$  (**TRUE**, **11**)  $\xrightarrow{1}$  (**TRUE**, **11**)  $\xrightarrow{0}$  (**TRUE**, **10**)  $\xrightarrow{0}$  (**TRUE**, **00**)

### 3.5.2 ... but Algorithms can be Wasteful

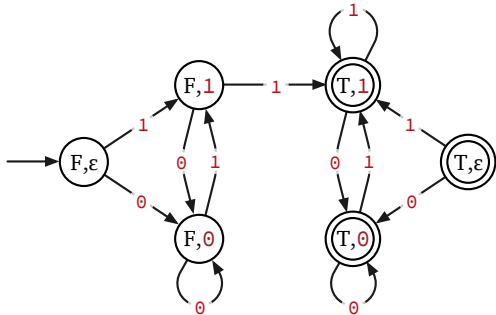
You can probably guess that the brute-force DFA we just constructed has considerably more states than necessary, especially after seeing its transition graph:



Our brute-force DFA for strings containing the substring **11**

For example, the state (**FALSE**, **11**) has no incoming transitions, so we can just delete it. (This state would indicate that we've never read **11**, but the last two symbols we read were **11**, which

is impossible!) More significantly, we don't need actually to remember both of the last two symbols, but only the penultimate symbol, because the last symbol is the one we're currently reading. This observation allows us to reduce the number of states from fourteen to only six.

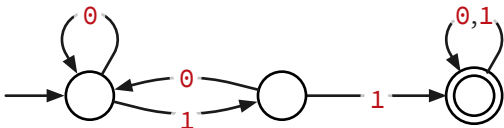


A less brute-force DFA for strings containing the substring 11

But even this DFA has more states than necessary. Once the flag part of the state is set to TRUE, we know the machine will eventually accept, so we might as well merge all the accepting states together. More subtly, because both transitions out of (FALSE, 0) and (FALSE, ε) lead to the same states, we can merge those two states together as well. After all these optimizations, we obtain the following DFA with just three states:

- The start state, which indicates that the machine has not read the substring 11 and did not just read the symbol 1.
- An intermediate state, which indicates that the machine has not read the substring 11 but just read the symbol 1.
- A unique accept state, which indicates that the machine has read the substring 11.

This is the smallest possible DFA for this language.



A minimal DFA for superstrings of 11

While it is important not to use an excessive number of states when we design DFAs—too many states makes a DFA hard to understand—there is really no point in trying to reduce DFAs *by hand* to the absolute minimum number of states. Clarity is much more important than brevity (especially in this class), and DFAs with too *few* states can *also* be hard to understand. At the end of this note, I'll describe an efficient algorithm that automatically transforms any given DFA into an equivalent DFA with the fewest possible states.

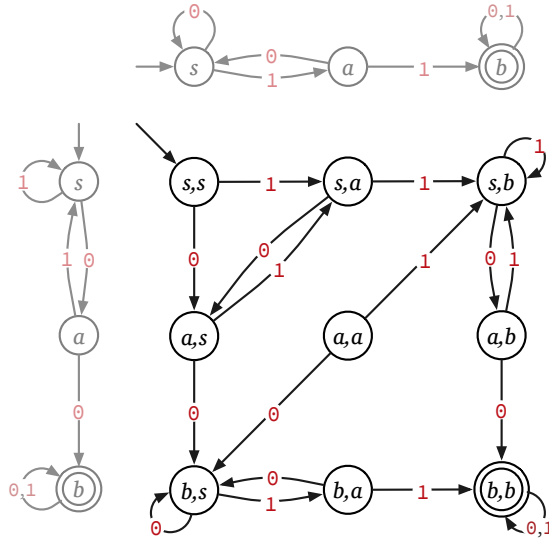
### 3.6 Combining DFAs: The Product Construction

Now suppose we want to accept all strings that contain both 00 and 11 as substrings, in either order. Intuitively, we'd like to run two DFAs in parallel—the DFA  $M_{00}$  to detect superstrings of 00 and a similar DFA  $M_{11}$  obtained from  $M_{00}$  by swapping  $0 \leftrightarrow 1$  everywhere—and then accept the input string if and only if *both* of these DFAs accept.

In fact, we can encode precisely this “parallel computation” into a single DFA using the following *product construction* first proposed by Edward Moore in 1956:

- The states of the new DFA are all ordered pairs  $(p, q)$ , where  $p$  is a state in  $M_{00}$  and  $q$  is a state in  $M_{11}$ .
- The start state of the new DFA is the pair  $(s, s')$ , where  $s$  is the start state of  $M_{00}$  and  $s'$  is the start state of  $M_{11}$ .
- The new DFA includes the transition  $(p, q) \xrightarrow{a} (p', q')$  if and only if  $M_{00}$  contains the transition  $p \xrightarrow{a} p'$  and  $M_{11}$  contains the transition  $q \xrightarrow{a} q'$ .
- Finally,  $(p, q)$  is an accepting state of the new DFA if and only if  $p$  is an accepting state in  $M_{00}$  and  $q$  is an accepting state in  $M_{11}$ .

The resulting nine-state DFA is shown on the next page, with the two factor DFAs  $M_{00}$  and  $M_{11}$  shown in gray for reference. (The state  $(a, a)$  can be removed, because it has no incoming transition, but let's not worry about that now.)



Building a DFA for the language of strings containing both  $00$  and  $11$ .

More generally, let  $M_1 = (\Sigma, Q_1, \delta_1, s_1, A_1)$  be an arbitrary DFA that accepts some language  $L_1$ , and let  $M_2 = (\Sigma, Q_2, \delta_2, s_2, A_2)$  be an arbitrary DFA that accepts some language  $L_2$  (over the same alphabet  $\Sigma$ ). We can construct a third DFA  $M = (\Sigma, Q, \delta, s, A)$  that accepts the intersection language  $L_1 \cap L_2$  as follows.

$$\begin{aligned}
 Q &:= Q_1 \times Q_2 = \{(p, q) \mid p \in Q_1 \text{ and } q \in Q_2\} \\
 \delta((p, q), a) &:= (\delta_1(p, a), \delta_2(q, a)) \\
 s &:= (s_1, s_2) \\
 A &:= A_1 \times A_2 = \{(p, q) \mid p \in A_1 \text{ and } q \in A_2\}
 \end{aligned}$$

To convince ourselves that this product construction is actually correct, let's consider the extended transition function  $\delta^*: (Q \times Q') \times \Sigma^* \rightarrow (Q \times Q')$ , which acts on strings instead of individual symbols. Recall that this function is defined recursively as follows:

$$\delta^*((p, q), w) := \begin{cases} (p, q) & \text{if } w = \varepsilon, \\ \delta^*(\delta((p, q), a), x) & \text{if } w = ax. \end{cases}$$

This function behaves exactly as we should expect:

**Lemma 3.1.**  $\delta^*((p, q), w) = (\delta_1^*(p, w), \delta_2^*(q, w))$  for any string  $w$ .

**Proof:** Let  $w$  be an arbitrary string. Assume  $\delta^*((p, q), x) = (\delta_1^*(p, x), \delta_2^*(q, x))$  for every string  $x$  that is shorter than  $w$ . As usual, there are two cases to consider.

- First suppose  $w = \varepsilon$ :

$$\begin{aligned} \delta^*((p, q), \varepsilon) &= (p, q) && \text{by the definition of } \delta^* \\ &= (\delta_1^*(p, \varepsilon), q) && \text{by the definition of } \delta_1^* \\ &= (\delta_1^*(p, \varepsilon), \delta_2^*(q, \varepsilon)) && \text{by the definition of } \delta_2^* \end{aligned}$$

- Now suppose  $w = ax$  for some symbol  $a$  and some string  $x$ :

$$\begin{aligned} \delta^*((p, q), ax) &= \delta^*(\delta((p, q), a), x) && \text{by the definition of } \delta^* \\ &= \delta^*((\delta_1(p, a), \delta_2(q, a)), x) && \text{by the definition of } \delta \\ &= (\delta_1^*((\delta_1(p, a), x), \delta_2^*(\delta_2(q, a), x))) && \text{by the induction hypothesis} \\ &= (\delta_1^*(p, ax), \delta_2^*(q, ax)) && \text{by the definitions of } \delta_1^* \text{ and } \delta_2^*. \end{aligned}$$

In both cases, we conclude that  $\delta^*((p, q), w) = (\delta_1^*(p, w), \delta_2^*(q, w))$ . □

An immediate consequence of this lemma is that for every string  $w$ , we have  $\delta^*(s, w) \in A$  if and only if both  $\delta_1^*(s_1, w) \in A_1$  and  $\delta_2^*(s_2, w) \in A_2$ . In other words,  $M$  accepts  $w$  if and only if **both**  $M_1$  accepts  $w$  **and**  $M_2$  accept  $w$ , as required.

As usual, this construction technique does not necessarily yield *minimal* DFAs. For example, in our first example of a product DFA, illustrated above, the central state  $(a, a)$  cannot be reached by any other state and is therefore redundant. Whatever.

Similar product constructions can be used to build DFAs that accept any other boolean combination of languages; in fact, the only part of the construction that changes is the choice of accepting states. For example:

- To accept the union  $L_1 \cup L_2$ , define  $A = \{(p, q) \mid p \in A_1 \text{ **or** } q \in A_2\}$ .
- To accept the difference  $L_1 \setminus L_2$ , define  $A = \{(p, q) \mid p \in A_1 \text{ **but** } q \notin A_2\}$ .
- To accept the symmetric difference  $L_1 \oplus L_2$ , define  $A = \{(p, q) \mid p \in A_1 \text{ **xor** } q \in A_2\}$ .

Examples of these constructions are shown on the next page.

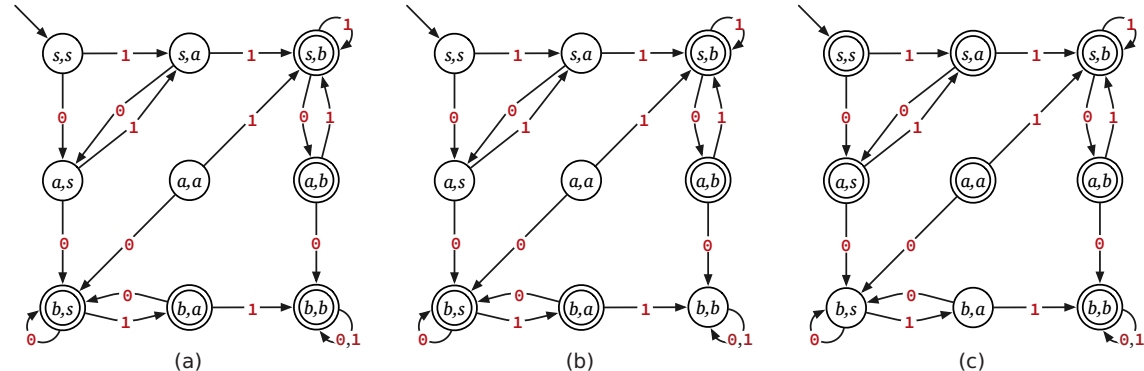
Moreover, by cascading this product construction, we can construct DFAs that accept arbitrary boolean combinations of arbitrary finite collections of regular languages.

### 3.7 Automatic Languages and Closure Properties

The **language** of a finite state machine  $M$ , denoted  $L(M)$ , is the set of all strings in  $\Sigma^*$  that  $M$  accepts. More formally, if  $M = (\Sigma, Q, \delta, s, A)$ , then

$$L(M) := \{w \in \Sigma^* \mid \delta^*(s, w) \in A\}.$$

We call a language **automatic** if it is the language of some finite state machine. Our product construction examples let us prove that the set of automatic languages is **closed** under simple boolean operations.



DFA for (a) strings that contain 00 or 11, (b) strings that contain either 00 or 11 but not both, and (c) strings that contain 11 if they contain 00. These DFA are identical except for their choices of accepting states.

**Theorem 3.2.** Let  $L$  and  $L'$  be arbitrary automatic languages over an arbitrary alphabet  $\Sigma$ .

- $\bar{L} = \Sigma^* \setminus L$  is automatic.
- $L \cup L'$  is automatic.
- $L \cap L'$  is automatic.
- $L \setminus L'$  is automatic.
- $L \oplus L'$  is automatic.

Eager students may have noticed that a Google search for the phrase “automatic language” turns up **no** results that are relevant for this class, except perhaps this lecture note. That’s because “automatic” is just a synonym for “regular”! This equivalence was first observed by Stephen Kleene (the inventor of regular expressions) in 1956.

**Theorem 3.3 (Kleene).** For any regular expression  $R$ , there is a DFA  $M$  such that  $L(R) = L(M)$ . For any DFA  $M$ , there is a regular expression  $R$  such that  $L(M) = L(R)$ .

Unfortunately, we don’t yet have all the tools we need to prove Kleene’s theorem; we’ll return to the proof in the next lecture note, after we have introduced *nondeterministic* finite-state machines. The proof is actually constructive—there are explicit algorithms that transform arbitrary DFA into equivalent regular expressions and vice versa.<sup>3</sup>

This equivalence between regular and automatic languages implies that the set of **regular** languages is also closed under simple boolean operations. The union of two regular languages is regular *by definition*, but it’s much less obvious that *every* boolean combination of regular languages can also be described by regular expressions.

**Corollary 3.4.** Let  $L$  and  $L'$  be arbitrary **regular** languages over an arbitrary alphabet  $\Sigma$ .

- $\bar{L} = \Sigma^* \setminus L$  is regular.
- $L \cap L'$  is regular.
- $L \setminus L'$  is regular.
- $L \oplus L'$  is regular.

Conversely, because concatenations and Kleene closures of regular languages are regular *by definition*, we can immediately conclude that concatenations and Kleene closures of automatic languages are automatic.

<sup>3</sup>These conversion algorithms run in exponential time in the worst case, but that’s unavoidable. There are regular languages whose smallest accepting DFA is exponentially larger than their smallest regular expression, and there are regular languages whose smallest regular expression is exponentially larger than their smallest accepting DFA.

**Corollary 3.5.** Let  $L$  and  $L'$  be arbitrary automatic languages.

- $L \cdot L'$  is automatic.
- $L^*$  is automatic.

These results give us several options to prove that a given languages is regular or automatic. We can either (1) build a regular expression that describes the language, (2) build a DFA that accepts the language, or (3) build the language from simpler pieces from other regular/automatic languages. (Later we'll see a fourth option, and possibly even a fifth.)

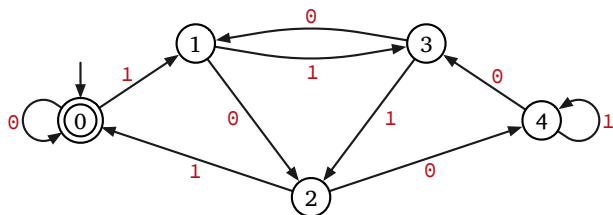
### 3.8 Proving a Language is Not Regular

But now suppose we're faced with a language  $L$  where none of these techniques seem to work. How would we prove  $L$  is *not* regular? By Theorem ??, it suffices to prove that there is no finite-state automaton that accepts  $L$ . Equivalently, we need to prove that any automaton that accepts  $L$  requires infinitely many states. That may sound tricky, what with the “infinitely many”, but there's actually a fairly simple technique to prove exactly that.

#### 3.8.1 Distinguishing Suffixes

Perhaps the single most important feature of DFAs is that they have no memory other than the current state. Once a DFA enters a particular state, all *future* transitions depend only on that state and *future* input symbols; *past* input symbols are simply forgotten.

For example, consider our very first DFA, which accepts the binary representations of integers divisible by 5.



DFA accepting binary multiples of 5.

The strings **0010** and **11011** both lead this DFA to state 2, although they follow different transitions to get there. Thus, for any string  $z$ , the strings **0010** $z$  and **11011** $z$  also lead to the same state in this DFA. In particular, **0010** $z$  leads to the accepting state if and only if **11011** $z$  leads to the accepting state. It follows that **0010** $z$  is divisible by 5 if and only if **11011** $z$  is divisible by 5.

More generally, any DFA  $M = (\Sigma, Q, s, A, \delta)$  defines an equivalence relation over  $\Sigma^*$ , where two strings  $x$  and  $y$  are equivalent if and only if they lead to the same state, or more formally, if  $\delta^*(s, x) = \delta^*(s, y)$ . If  $x$  and  $y$  are equivalent strings, then for any string  $z$ , the strings  $xz$  and  $yz$  are also equivalent. In particular,  $M$  accepts  $xz$  if and only if  $M$  accepts  $yz$ . Thus, if  $L$  is the language accepted by  $M$ , then  $xz \in L$  if and only if  $yz \in L$ . In short, if the *machine* can't distinguish between  $x$  and  $y$ , then the *language* can't distinguish between  $xz$  and  $yz$  for any suffix  $z$ .

Now let's turn the previous argument on its head. Let  $L$  be an arbitrary language, and let  $x$  and  $y$  be arbitrary strings. A **distinguishing suffix** for  $x$  and  $y$  (with respect to  $L$ ) is a third string  $z$  such that *exactly one* of the strings  $xz$  and  $yz$  is in  $L$ . If  $x$  and  $y$  have a distinguishing

suffix  $z$ , then in *any* DFA that accepts  $L$ , the strings  $xz$  and  $yz$  must lead to different states, and therefore the strings  $x$  and  $y$  must lead to different states!

For example, let  $L_5$  denote the the set of all strings over  $\{0, 1\}$  that represent multiples of 5 in binary. Then the strings  $x = 01$  and  $y = 0011$  are distinguished by the suffix  $z = 01$ :

$$\begin{aligned}xz &= 01 \cdot 01 = 0101 \in L_5 && \text{(because } 0101_2 = 5\text{)} \\yz &= 0011 \cdot 01 = 001101 \notin L_5 && \text{(because } 001101_2 = 13\text{)}\end{aligned}$$

It follows that in *every* DFA that accepts  $L_5$ , the strings  $01$  and  $0011$  lead to different states. Moreover, since neither  $01$  nor  $0011$  belong to  $L_5$ , every DFA that accepts  $L_5$  must have at least two *non-accepting* states, and therefore at least three states overall.

### 3.8.2 Fooling Sets

A **fooling set** for a language  $L$  is a set  $F$  of strings such that *every* pair of strings in  $F$  has a distinguishing suffix. For example,  $F = \{0, 1, 10, 11, 100\}$  is a fooling set for the language  $L_5$  of binary multiples of 5, because each pair of strings in  $F$  has a distinguishing suffix:

- $0$  distinguishes  $0$  and  $1$ ;
- $0$  distinguishes  $0$  and  $10$ ;
- $0$  distinguishes  $0$  and  $11$ ;
- $0$  distinguishes  $0$  and  $100$ ;
- $1$  distinguishes  $1$  and  $10$ ;
- $01$  distinguishes  $1$  and  $11$ ;
- $01$  distinguishes  $1$  and  $100$ ;
- $1$  distinguishes  $10$  and  $11$ ;
- $1$  distinguishes  $10$  and  $100$ ;
- $11$  distinguishes  $11$  and  $100$ .

Each of these five strings leads to a different state, for *any* DFA  $M$  that accepts  $L_5$ . Thus, *every* DFA that accepts the language  $L_5$  has at least five states. And hey, look, we already have a DFA for  $L_5$  with five states, so that's the best we can do!

More generally, for *any* language  $L$ , and *any* fooling set  $F$  for  $L$ , *every* DFA that accepts  $L$  must have at least  $|F|$  states. In particular, if the fooling set  $F$  is *infinite*, then every DFA that accepts  $L$  must have an *infinite* number of states. But there's no such thing as a **finite**-state machine with an **infinite** number of states!

*If  $L$  has an infinite fooling set, then  $L$  is not regular.*

This is arguably both the simplest and most powerful method for proving that a language is non-regular. Here are a few canonical examples of the fooling-set technique in action.

**Lemma 3.6.** *The language  $L = \{0^n 1^n \mid n \geq 0\}$  is not regular.*

**Proof:** Consider the infinite set  $F = \{0^n \mid n \geq 0\}$ , or more simply  $F = 0^*$ .

Let  $x$  and  $y$  be arbitrary distinct strings in  $F$ .

The definition of  $F$  implies  $x = 0^i$  and  $y = 0^j$  for some integers  $i \neq j$ .

The suffix  $z = 1^i$  distinguishes  $x$  and  $y$ , because  $xz = 0^i 1^i \in L$ , but  $yz = 0^j 1^i \notin L$ .

Thus, every pair of distinct strings in  $F$  has a distinguishing suffix.

In other words,  $F$  is a fooling set for  $L$ .

Because  $F$  is infinite,  $L$  cannot be regular. □

**Lemma 3.7.** *The language  $L = \{ww^R \mid w \in \Sigma^*\}$  of even-length palindromes is not regular.*

**Proof:** Let  $F$  denote the set  $0^*1$ , and let  $x$  and  $y$  be arbitrary distinct strings in  $F$ . Then we must have  $x = 0^i 1$  and  $y = 0^j 1$  for some integers  $i \neq j$ . The suffix  $z = 10^i$  distinguishes  $x$  and  $y$ , because  $xz = 0^i 110^i \in L$ , but  $yz = 0^j 110^i \notin L$ . We conclude that  $F$  is a fooling set for  $L$ . Because  $F$  is infinite,  $L$  cannot be regular. □

**Lemma 3.8.** *The language  $L = \{0^{2^n} \mid n \geq 0\}$  is not regular.*

**Proof ( $F = L$ ):** Let  $x$  and  $y$  be arbitrary distinct strings in  $L$ . Then we must have  $x = 0^{2^i}$  and  $y = 0^{2^j}$  for some integers  $i \neq j$ . The suffix  $z = 0^{2^i}$  distinguishes  $x$  and  $y$ , because  $xz = 0^{2^i+2^i} = 0^{2^{i+1}} \in L$ , but  $yz = 0^{2^i+2^j} \notin L$ . We conclude that  $L$  itself is a fooling set for  $L$ . Because  $L$  is infinite,  $L$  cannot be regular. □

**Proof ( $F = 0^*$ ):** Let  $x$  and  $y$  be arbitrary distinct strings in  $0^*$ . Then we must have  $x = 0^i$  and  $y = 0^j$  for some integers  $i \neq j$ ; without loss of generality, assume  $i < j$ . Let  $k$  be any positive integer such that  $2^k > j$ . Consider the suffix  $z = 0^{2^k-i}$ . We have  $xz = 0^{i+(2^k-i)} = 0^{2^k} \in L$ , but  $yz = 0^{j+(2^k-i)} = 0^{2^k-i+j} \notin L$ , because

$$2^k < 2^k - i + j < 2^k + j < 2^k + 2^k = 2^{k+1}.$$

Thus,  $z$  is a distinguishing suffix for  $x$  and  $y$ . We conclude that  $0^*$  is a fooling set for  $L$ . Because  $L$  is infinite,  $L$  cannot be regular. □

**Proof ( $F = 0^*$  again):** Let  $x$  and  $y$  be arbitrary distinct strings in  $0^*$ . Then we must have  $x = 0^i$  and  $y = 0^j$  for some integers  $i \neq j$ ; without loss of generality, assume  $i < j$ . Let  $k$  be any positive integer such that  $2^{k-1} > j$ . Consider the suffix  $z = 0^{2^k-j}$ . We have  $xz = 0^{i+(2^k-j)} = 0^{2^k-j+i} \notin L$ , because

$$2^{k-1} < 2^k - 2^{k-1} + i < 2^k - j + i < 2^k.$$

On the other hand,  $yz = 0^{j+(2^k-j)} = 0^{2^k} \in L$ . Thus,  $z$  is a distinguishing suffix for  $x$  and  $y$ . We conclude that  $0^*$  is a fooling set for  $L$ . Because  $L$  is infinite,  $L$  cannot be regular. □

The previous examples show the flexibility of this proof technique; a single non-regular language can have many different infinite fooling sets,<sup>4</sup> and each pair of strings in any fooling set can have many different distinguishing suffixes. Fortunately, we only have to find *one* infinite set  $F$  and *one* distinguishing suffix for each pair of strings in  $F$ .

**Lemma 3.9.** *The language  $L = \{0^p \mid p \text{ is prime}\}$  is not regular.*

---

<sup>4</sup>At some level, this observation is trivial. If  $F$  is an infinite fooling set for  $L$ , then every infinite subset of  $F$  is also an infinite fooling set for  $L$ !



**Proof ( $F = \emptyset^*$ ):** Again, we use  $\emptyset^*$  as our fooling set, but the actual argument is somewhat more complicated than in our earlier examples.

Let  $x$  and  $y$  be arbitrary distinct strings in  $\emptyset^*$ . Then we must have  $x = \emptyset^i$  and  $y = \emptyset^j$  for some integers  $i \neq j$ ; without loss of generality, assume that  $i < j$ . Let  $p$  be any prime number larger than  $i$ . Because  $p + 0(j - i)$  is prime and  $p + p(j - i) > p$  is not, there must be a positive integer  $k \leq p$  such that  $p + (k - 1)(j - i)$  is prime but  $p + k(j - i)$  is not. Then I claim that the suffix  $z = \emptyset^{p+(k-1)j-ki}$  distinguishes  $x$  and  $y$ :

$$\begin{aligned} xz &= \emptyset^i \emptyset^{p+(k-1)j-ki} = \emptyset^{p+(k-1)(j-i)} \in L && \text{because } p + (k - 1)(j - i) \text{ is prime;} \\ yz &= \emptyset^j \emptyset^{p+(k-1)j-ki} = \emptyset^{p+k(j-i)} \notin L && \text{because } p + k(j - i) \text{ is not prime.} \end{aligned}$$

(Because  $i < j$  and  $i < p$ , the suffix  $\emptyset^{p+(k-1)j-ki} = \emptyset^{(p-i)+(k-1)(j-i)}$  has positive length and therefore *actually exists!*) We conclude that  $\emptyset^*$  is indeed a fooling set for  $L$ , which implies that  $L$  is not regular.  $\square$

**Proof ( $F = L$ ):** Let  $x$  and  $y$  be arbitrary distinct strings in  $L$ . Then we must have  $x = \emptyset^p$  and  $y = \emptyset^q$  for some primes  $p \neq q$ ; without loss of generality, assume  $p < q$ .

Now consider strings of the form  $\emptyset^{p+k(q-p)}$ . Because  $p + 0(q - p)$  is prime and  $p + p(q - p) > p$  is not prime, there must be a non-negative integer  $k < p$  such that  $p + k(p - q)$  is prime but  $p + (k + 1)(p - q)$  is not prime. I claim that the suffix  $z = \emptyset^{k(q-p)}$  distinguishes  $x$  and  $y$ :

$$\begin{aligned} xz &= \emptyset^p \emptyset^{k(q-p)} = \emptyset^{p+k(p-q)} \in L && \text{because } p + k(p - q) \text{ is prime;} \\ yz &= \emptyset^q \emptyset^{k(q-p)} = \emptyset^{p+(k+1)(q-p)} \notin L && \text{because } p + (k + 1)(p - q) \text{ is not prime.} \end{aligned}$$

We conclude that  $L$  is a fooling set for itself!! Because  $L$  is infinite,  $L$  cannot be regular!  $\square$

Obviously the most difficult part of this technique is coming up with an appropriate fooling set. Fortunately, *most* languages  $L$ —in particular, almost all languages that students are asked to prove non-regular on homeworks or exams—fall into one of two categories:

- Some simple regular language like  $\emptyset^*$  or  $1\emptyset^*1$  or  $(\emptyset 1)^*$  is a fooling set for  $L$ . In particular, the fooling set is a regular language with one Kleene star and no  $+$ .
- The language  $L$  itself is a fooling set for  $L$ .

The most important point to remember is that **you choose** the fooling set  $F$ , and you can use that fooling set to effectively impose additional structure on the language  $L$ .

I'm not sure yet how to express this effectively, but here is some more intuition about choosing fooling sets and distinguishing suffixes.

As a sanity check, try to write an *algorithm* to recognize strings in  $L$ , as described at the start of this note, where the only variable that can take on an unbounded number of values is the loop index  $i$ . (I should probably rewrite that template as a while-loop or tail recursion, but anyway. . . ) If you succeed, the language is regular. But if you fail, it's probably because there are counters of string variables that you can't get rid of. **One of those unavoidable counters is the basis for your fooling set.**

For example, any algorithm that recognizes the language  $\{0^n 1^n 2^n \mid n \geq 0\}$  “obviously” has to count 0s and 1s in the input string. (We can avoid counting 2s by decrementing the 0 counter.) Because the 0s come first in the string, this intuition suggests using strings of the form  $0^n$  as our fooling set and matching strings of the form  $1^n 2^n$  as distinguishing suffixes. (This is a rare example of an “obvious” fact that is actually true.)

It's also important to remember that when you choose the fooling set, you can effectively impose additional structure that isn't present in the language already. For example, to prove that the language  $L = \{w \in (0+1)^* \mid \#(0, w) = (1, w)\}$  is not regular, we can use strings of the form  $0^n$  as our fooling set and matching strings of the form  $1^n$  as distinguishing suffixes, **exactly** as we did for  $\{0^n 1^n \mid n \geq 0\}$ . The fact that  $L$  contains strings that start with 1 is irrelevant. There may be more equivalence classes that our proof doesn't find, but since we found an infinite set of equivalence class, we don't care.

At some level, this fooling set proof is implicitly considering the simpler language  $L \cap 0^* 1^* = \{0^n 1^n \mid n \geq 0\}$ . If  $L$  were regular, then  $L \cap 0^* 1^*$  would also be regular, because regular languages are closed under intersection.

### \*3.9 The Myhill-Nerode Theorem

The fooling set technique implies a *necessary* condition for a language to be accepted by a DFA—the language must have no infinite fooling sets. In fact, this condition is also *sufficient*. The following powerful theorem was first proved by Anil Nerode in 1958, strengthening a 1957 result of John Myhill.<sup>5</sup> We write  $x \equiv_L y$  if  $xz \in L \iff yz \in L$  for all strings  $z$ .

**The Myhill-Nerode Theorem.** *For any language  $L$ , the following are equal:*

- (a) *the minimum number of states in a DFA that accepts  $L$ ,*
- (b) *the maximum size of a fooling set for  $L$ , and*
- (c) *the number of equivalence classes of  $\equiv_L$ .*

*In particular,  $L$  is accepted by a DFA if and only if every fooling set for  $L$  is finite.*

**Proof:** Let  $L$  be an arbitrary language.

We have already proved that the size of any fooling set for  $L$  is at most the number of states in any DFA that accepts  $L$ , so (a)  $\geq$  (b). It also follows directly from the definitions that  $F \subseteq \Sigma^*$  is a fooling set for  $L$  if and only if  $F$  contains at most one string in each equivalence class of  $\equiv_L$ ; thus, (b) = (c). We complete the proof by showing that (a)  $\leq$  (c).

We have already proved that if  $\equiv_L$  has an infinite number of equivalence classes, there is no DFA that accepts  $L$ , so assume that the number of equivalence classes is finite. For any string  $w$ ,

---

<sup>5</sup>Myhill considered the finer equivalence relation  $x \sim_L y$ , meaning  $wxz \in L$  if and only if  $wyz \in L$  for all strings  $w$  and  $z$ , and proved that  $L$  is regular if and only if  $\sim_L$  defines a finite number of equivalence classes. Like most of Myhill's early automata research, this result appears in an unpublished Air Force technical report. The modern Myhill-Nerode theorem appears (in an even more general form) as a minor lemma in Nerode's 1958 paper, which (not surprisingly) does not cite Myhill.

let  $[w]$  denote its equivalence class. We define a DFA  $M_{\equiv} = (\Sigma, Q, s, A, \delta)$  as follows:

$$\begin{aligned} Q &:= \{[w] \mid w \in \Sigma^*\} \\ s &:= [\varepsilon] \\ A &:= \{[w] \mid w \in L\} \\ \delta([w], a) &:= [w \bullet a] \end{aligned}$$

We claim that this DFA accepts the language  $L$ ; this claim completes the proof of the theorem.

But before we can prove anything about this DFA, we first need to verify that it is actually well-defined. Let  $x$  and  $y$  be two strings such that  $[x] = [y]$ . By definition of  $L$ -equivalence, for any string  $z$ , we have  $xz \in L$  if and only if  $yz \in L$ . It immediately follows that for any symbol  $a \in \Sigma$  and any string  $z'$ , we have  $xaz' \in L$  if and only if  $yaz' \in L$ . Thus, by definition of  $L$ -equivalence, we have  $[xa] = [ya]$  for every symbol  $a \in \Sigma$ . We conclude that the function  $\delta$  is indeed well-defined.

An easy inductive proof implies that  $\delta^*([\varepsilon], x) = [x]$  for every string  $x$ . Thus,  $M$  accepts string  $x$  if and only if  $[x] = [w]$  for some string  $w \in L$ . But if  $[x] = [w]$ , then by definition (setting  $z = \varepsilon$ ), we have  $x \in L$  if and only if  $w \in L$ . So  $M$  accepts  $x$  if and only if  $x \in L$ . In other words,  $M$  accepts  $L$ , as claimed, so the proof is complete.  $\square$

### \*3.10 Minimal Automata

Given a DFA  $M = (\Sigma, Q, s, A, \delta)$ , suppose we want to find another DFA  $M' = (\Sigma, Q', s', A', \delta')$  with the fewest possible states that accepts the same language. In this final section, we describe an efficient algorithm to minimize DFAs, first described (in slightly different form) by Edward Moore in 1956. We analyze the running time of Moore's in terms of two parameters:  $n = |Q|$  and  $\sigma = |\Sigma|$ .

In the preprocessing phase, we find and remove any states that cannot be reached from the start state  $s$ ; this filtering can be performed in  $O(n\sigma)$  time using any graph traversal algorithm. So from now on we assume that all states are reachable from  $s$ .

Now we recursively define two states  $p$  and  $q$  in the remaining DFA to be **distinguishable**, written  $p \not\sim q$ , if at least one of the following conditions holds:

- $p \in A$  and  $q \notin A$ ,
- $p \notin A$  and  $q \in A$ , or
- $\delta(p, a) \not\sim \delta(q, a)$  for some  $a \in \Sigma$ .

Equivalently,  $p \not\sim q$  if and only if there is a string  $z$  such that exactly one of the states  $\delta^*(p, z)$  and  $\delta^*(q, z)$  is accepting. (Sound familiar?) Intuitively, the main algorithm assumes that all states are equivalent until proven otherwise, and then repeatedly looks for state pairs that can be proved distinguishable.

The main algorithm maintains a two-dimensional table, indexed by the states, where  $\text{Dist}[p, q] = \text{TRUE}$  indicates that we have proved states  $p$  and  $q$  are distinguishable. Initially, for all states  $p$  and  $q$ , we set  $\text{Dist}[p, q] \leftarrow \text{TRUE}$  if  $p \in A$  and  $q \notin A$  or vice versa, and  $\text{Dist}[p, q] = \text{FALSE}$  otherwise. Then we repeatedly consider each pair of states and each symbol to find more distinguishable pairs, until we make a complete pass through the table without modifying it. The table-filling algorithm can be summarized as follows:

```

MINDFATABLE( $\Sigma, Q, s, A, \delta$ ):
  for all  $p \in Q$ 
    for all  $q \in Q$ 
      if  $(p \in A \text{ and } q \notin A) \text{ or } (p \notin A \text{ and } q \in A)$ 
         $\text{Dist}[p, q] \leftarrow \text{TRUE}$ 
      else
         $\text{Dist}[p, q] \leftarrow \text{FALSE}$ 
   $\text{notdone} \leftarrow \text{TRUE}$ 
  while  $\text{notdone}$ 
     $\text{notdone} \leftarrow \text{FALSE}$ 
    for all  $p \in Q$ 
      for all  $q \in Q$ 
        if  $\text{Dist}[p, q] = \text{FALSE}$ 
          for all  $a \in \Sigma$ 
            if  $\text{Dist}[\delta(p, a), \delta(q, a)]$ 
               $\text{Dist}[p, q] \leftarrow \text{TRUE}$ 
               $\text{notdone} \leftarrow \text{TRUE}$ 
  return  $\text{Dist}$ 

```

The algorithm must eventually halt, because there are only a finite number of entries in the table that can be marked. In fact, the main loop is guaranteed to terminate after at most  $n$  iterations, which implies that the entire algorithm runs in  $O(\sigma n^3)$  time. Once the table is filled,<sup>6</sup> any two states  $p$  and  $q$  such that  $\text{Dist}(p, q) = \text{FALSE}$  are equivalent and can be merged into a single state. The remaining details of constructing the minimized DFA are straightforward.

Need to prove that the main loop terminates in at most  $n$  iterations.

With more care, Moore's minimization algorithm can be modified to run in  $O(\sigma n^2)$  time. A faster DFA minimization algorithm, due to John Hopcroft, runs in  $O(\sigma n \log n)$  time.

## Example

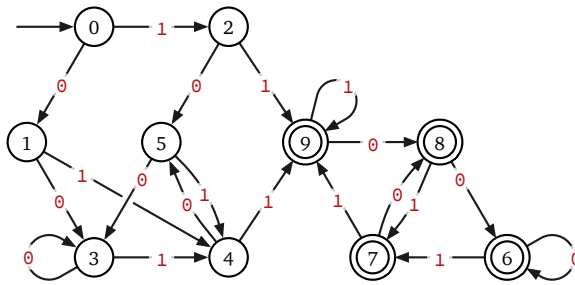
To get a better idea how this algorithm works, let's visualize its execution on our earlier brute-force DFA for strings containing the substring **11**. This DFA has four unreachable states:  $(\text{FALSE}, \mathbf{11})$ ,  $(\text{TRUE}, \epsilon)$ ,  $(\text{TRUE}, \mathbf{0})$ , and  $(\text{TRUE}, \mathbf{1})$ . We remove these states, and relabel the remaining states for easier reference. (In an actual implementation, the states would almost certainly be represented by indices into an array anyway, not by mnemonic labels.)

The main algorithm initializes (the bottom half of) a  $10 \times 10$  table as follows. (In the following figures, cells marked  $\times$  have value **TRUE** and blank cells have value **FALSE**.)

<sup>6</sup>More experienced readers should be enraged by the mere suggestion that any algorithm merely *fills in a table*, as opposed to *evaluating a recurrence*. This algorithm is no exception. Consider the boolean function  $\text{Dist}(p, q, k)$ , which equals **TRUE** if and only if  $p$  and  $q$  can be distinguished by some string of length at most  $k$ . This function obeys the following recurrence:

$$\text{Dist}(p, q, k) = \begin{cases} (p \in A) \oplus (q \in A) & \text{if } k = 0, \\ \text{Dist}(p, q, k-1) \vee \bigvee_{a \in \Sigma} \text{Dist}(\delta(p, a), \delta(q, a), k-1) & \text{otherwise.} \end{cases}$$

Moore's "table-filling" algorithm is just a space-efficient dynamic programming algorithm to evaluate this recurrence.



Our brute-force DFA for strings containing the substring **11**, after removing all four unreachable states

	0	1	2	3	4	5	6	7	8
1									
2									
3									
4									
5									
6	x	x	x	x	x	x			
7	x	x	x	x	x	x			
8	x	x	x	x	x	x			
9	x	x	x	x	x	x			

In the first iteration of the main loop, the algorithm discovers several distinguishable pairs of states. For example, the algorithm sets  $Dist[0, 2] \leftarrow \text{TRUE}$  because  $Dist[\delta(0, 1), \delta(2, 1)] = Dist[2, 9] = \text{TRUE}$ . After the iteration ends, the table looks like this:

	0	1	2	3	4	5	6	7	8
1									
2	x	x							
3			x						
4	x	x		x					
5			x		x				
6	x	x	x	x	x	x			
7	x	x	x	x	x	x			
8	x	x	x	x	x	x			
9	x	x	x	x	x	x			

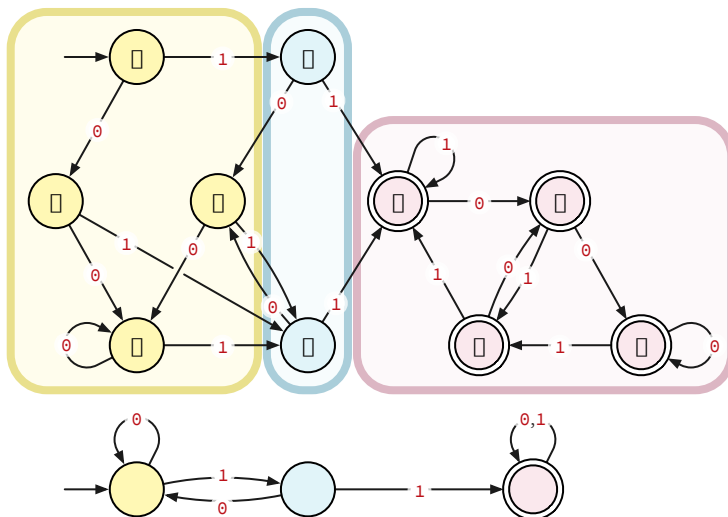
The second iteration of the while loop makes no further changes to the table—We got lucky!—so the algorithm terminates.

The final table implies that the 10 states of our DFA fall into exactly three equivalence classes:  $\{0, 1, 3, 5\}$ ,  $\{2, 4\}$ , and  $\{6, 7, 8, 9\}$ . Replacing each equivalence class with a single state gives us the three-state DFA that we already discovered.

### Exercises

- For each of the following languages in  $\{0, 1\}^*$ , describe a deterministic finite-state machine that accepts that language. There are infinitely many correct answers for each language. “Describe” does not necessarily mean “draw”.

(a) Only the string **0110**.



Equivalence classes of states in our DFA, and the resulting minimal equivalent DFA.

- (b) Every string except  $0110$ .
- (c) Strings that contain the substring  $0110$ .
- (d) Strings that do not contain the substring  $0110$ .
- \* (e) Strings that contain an even number of occurrences of the substring  $0110$ . (For example, this language contains the strings  $0110110$  and  $01011$ .)
- (f) Strings that contain the *subsequence*  $0110$ .
- (g) Strings that do not contain the *subsequence*  $0110$ .
- \* (h) Strings that contain an even number of occurrences of the *subsequence*  $0110$ .
- (i) Strings that contain an even number of  $1$ s and an odd number of  $0$ s.
- (j) Every string that represents a number divisible by 7 in binary.
- (k) Every string whose reversal represents a number divisible by 7 in binary.
- (l) Strings in which the substrings  $01$  and  $10$  appear the same number of times.
- (m) Strings such that in every prefix, the number of  $0$ s and the number of  $1$ s differ by at most 1.
- (n) Strings such that in every prefix, the number of  $0$ s and the number of  $1$ s differ by at most 4.
- (o) Strings that end with  $0^{10} = 0000000000$ .
- (p) All strings in which the number of  $0$ s is even if and only if the number of  $1$ s is *not* divisible by 3.
- (q) All strings that are both the binary representation of an integer divisible by 3 and the ternary (base-3) representation of an integer divisible by 4.
- (r) Strings in which the number of  $1$ s is even, the number of  $0$ s is divisible by 3, the overall length is divisible by 5, the binary value is divisible by 7, the binary value of the reversal is divisible by 11, and does not contain thirteen  $1$ s in a row. [Hint: This is more tedious than difficult.]
- \* (s) Strings  $w$  such that  $\binom{|w|}{2} \bmod 6 = 4$ .

- \* (t) Strings  $w$  such that  $F_{\#(\mathbf{10}, w)} \bmod 10 = 4$ , where  $\#(\mathbf{10}, w)$  denotes the number of times  $\mathbf{10}$  appears as a substring of  $w$ , and as usual  $F_n$  is the  $n$ th Fibonacci number:

$$F_n = \begin{cases} 0 & \text{if } n = 0 \\ 1 & \text{if } n = 1 \\ F_{n-1} + F_{n-2} & \text{otherwise} \end{cases}$$

- ★ (u) Strings  $w$  such that  $F_{\#(\mathbf{1}\cdots\mathbf{0}, w)} \bmod 10 = 4$ , where  $\#(\mathbf{1}\cdots\mathbf{0}, w)$  denotes the number of times  $\mathbf{10}$  appears as a *subsequence* of  $w$ , and as usual  $F_n$  is the  $n$ th Fibonacci number:

$$F_n = \begin{cases} 0 & \text{if } n = 0 \\ 1 & \text{if } n = 1 \\ F_{n-1} + F_{n-2} & \text{otherwise} \end{cases}$$

2. (a) Let  $L \subseteq \mathbf{0}^*$  be an arbitrary *unary* language. Prove that  $L^*$  is regular.  
 (b) Prove that there is a binary language  $L \subseteq (\mathbf{0} + \mathbf{1})^*$  such that  $L^*$  is not regular.

3. Prove that none of the following languages is automatic.

- (a)  $\{\mathbf{0}^{n^2} \mid n \geq 0\}$   
 (b)  $\{\mathbf{0}^{n^3} \mid n \geq 0\}$   
 (c)  $\{\mathbf{0}^{f(n)} \mid n \geq 0\}$ , where  $f(n)$  is *any* fixed polynomial in  $n$  with degree at least 2.  
 (d)  $\{\mathbf{0}^n \mid n \text{ is composite}\}$   
 (e)  $\{\mathbf{0}^n \mathbf{10}^n \mid n \geq 0\}$   
 (f)  $\{\mathbf{0}^m \mathbf{1}^n \mid m \neq n\}$   
 (g)  $\{\mathbf{0}^m \mathbf{1}^n \mid m < 3n\}$   
 (h)  $\{\mathbf{0}^{2n} \mathbf{1}^n \mid n \geq 0\}$   
 (i)  $\{w \in (\mathbf{0} + \mathbf{1})^* \mid \#(\mathbf{0}, w) = \#(\mathbf{1}, w)\}$   
 (j)  $\{w \in (\mathbf{0} + \mathbf{1})^* \mid \#(\mathbf{0}, w) < \#(\mathbf{1}, w)\}$   
 (k)  $\{\mathbf{0}^m \mathbf{1}^n \mid m/n \text{ is an integer}\}$   
 (l)  $\{\mathbf{0}^m \mathbf{1}^n \mid m \text{ and } n \text{ are relatively prime}\}$   
 (m)  $\{\mathbf{0}^m \mathbf{1}^n \mid n - m \text{ is a perfect square}\}$   
 (n)  $\{w\#w \mid w \in (\mathbf{0} + \mathbf{1})^*\}$   
 (o)  $\{ww \mid w \in (\mathbf{0} + \mathbf{1})^*\}$   
 (p)  $\{w\#\mathbf{0}^{|w|} \mid w \in (\mathbf{0} + \mathbf{1})^*\}$   
 (q)  $\{w\mathbf{0}^{|w|} \mid w \in (\mathbf{0} + \mathbf{1})^*\}$   
 (r)  $\{xy \mid x, y \in (\mathbf{0} + \mathbf{1})^* \text{ and } |x| = |y| \text{ but } x \neq y\}$   
 (s)  $\{\mathbf{0}^m \mathbf{1}^n \mathbf{0}^{m+n} \mid m, n \geq 0\}$   
 (t)  $\{\mathbf{0}^m \mathbf{1}^n \mathbf{0}^{mn} \mid m, n \geq 0\}$   
 (u) Strings in which the substrings  $\mathbf{00}$  and  $\mathbf{11}$  appear the same number of times.

- (v) Strings of the form  $w_1 \# w_2 \# \cdots \# w_n$  for some  $n \geq 2$ , where  $w_i \in \{0, 1\}^*$  for every index  $i$ , and  $w_i = w_j$  for some indices  $i \neq j$ .
  - (w) The set of all palindromes in  $(0 + 1)^*$  whose length is divisible by 7.
  - (x)  $\{w \in (0 + 1)^* \mid w \text{ is the binary representation of a perfect square}\}$
  - ★(y)  $\{w \in (0 + 1)^* \mid w \text{ is the binary representation of a prime number}\}$
4. For each of the following languages over the alphabet  $\Sigma = \{0, 1\}$ , either prove that the language is regular (by constructing an appropriate DFA or regular expression) or prove that the language is not regular (using fooling sets). Recall that  $\Sigma^+$  denotes the set of all nonempty strings over  $\Sigma$ . [Hint: Believe it or not, most of these languages **are** actually regular.]
- (a)  $\{0^n w 1^n \mid w \in \Sigma^* \text{ and } n \geq 0\}$
  - (b)  $\{0^n 1^n w \mid w \in \Sigma^* \text{ and } n \geq 0\}$
  - (c)  $\{w 0^n 1^n x \mid w, x \in \Sigma^* \text{ and } n \geq 0\}$
  - (d)  $\{0^n w 1^n x \mid w, x \in \Sigma^* \text{ and } n \geq 0\}$
  - (e)  $\{0^n w 1 x 0^n \mid w, x \in \Sigma^* \text{ and } n \geq 0\}$
  - (f)  $\{0^n w 0^n \mid w \in \Sigma^+ \text{ and } n > 0\}$
  - (g)  $\{w 0^n w \mid w \in \Sigma^+ \text{ and } n > 0\}$
  - (h)  $\{w x w \mid w, x \in \Sigma^*\}$
  - (i)  $\{w x w \mid w, x \in \Sigma^+\}$
  - (j)  $\{w x w^R \mid w, x \in \Sigma^+\}$
  - (k)  $\{w w x \mid w, x \in \Sigma^+\}$
  - (l)  $\{w w^R x \mid w, x \in \Sigma^+\}$
  - (m)  $\{w x w y \mid w, x, y \in \Sigma^+\}$
  - (n)  $\{w x w^R y \mid w, x, y \in \Sigma^+\}$
  - (o)  $\{x w w y \mid w, x, y \in \Sigma^+\}$
  - (p)  $\{x w w^R y \mid w, x, y \in \Sigma^+\}$
  - (q)  $\{w x x w \mid w, x \in \Sigma^+\}$
  - ★(r)  $\{w x w^R x \mid w, x \in \Sigma^+\}$
  - (s) All strings  $w$  such that no prefix of  $w$  is a palindrome.
  - (t) All strings  $w$  such that no prefix of  $w$  with length at least 3 is a palindrome.
  - (u) All strings  $w$  such that no substring of  $w$  with length at least 3 is a palindrome.
  - (v) All strings  $w$  such that no prefix of  $w$  with positive even length is a palindrome.
  - (w) All strings  $w$  such that no substring of  $w$  with positive even length is a palindrome.
  - (x) Strings in which the substrings  $00$  and  $11$  appear the same number of times.
  - (y) Strings in which the substrings  $01$  and  $10$  appear the same number of times.



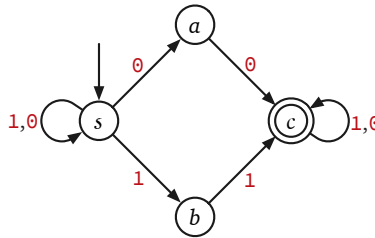
5. Let  $F$  and  $L$  be arbitrary infinite languages in  $\{0, 1\}^*$ .

- (a) Suppose for any two distinct strings  $x, y \in F$ , there is a string  $w \in \Sigma^*$  such that  $wx \in L$  and  $wy \notin L$ . (We can reasonably call  $w$  a *distinguishing prefix* for  $x$  and  $y$ .) Prove that  $L$  cannot be regular. [Hint: The reversal of a regular language is regular.]
- \* (b) Suppose for any two distinct strings  $x, y \in F$ , there are two (possibly equal) strings  $w, z \in \Sigma^*$  such that  $wxz \in L$  and  $wyz \notin L$ . Prove that  $L$  cannot be regular.

## 4 Nondeterminism

### 4.1 Nondeterministic State Machines

The following diagram shows something that looks like a finite-state machine over the alphabet  $\{0, 1\}$ , but on closer inspection, it is not consistent with our earlier definitions. On one hand, there are two transitions out of  $s$  for each input symbol. On the other hand, states  $a$  and  $b$  are each missing an outgoing transition.



A nondeterministic finite-state automaton

Nevertheless, there is a sense in which this machine “accepts” the set of all strings that contain either  $00$  or  $11$  as a substring. Imagine that when the machine reads a symbol in state  $s$ , it makes a **choice** about which transition to follow. If the input string contains the substring  $00$ , then it is *possible* for the machine to end in the accepting state  $c$ , by *choosing* to move into state  $a$  when it reads a  $0$  immediately before another  $0$ . Similarly, if the input string contains the substring  $11$ , it is *possible* for the machine to end in the accepting state  $c$ . On the other hand, if the input string does not contain either  $00$  or  $11$ —or in other words, if the input alternates between  $0$  and  $1$ —there are no choices that lead the machine to the accepting state. If the machine incorrectly chooses to transition to state  $a$  and then reads a  $1$ , or transitions to  $b$  and then reads  $0$ , it explodes; the only way to avoid an explosion is to stay in state  $s$ .

This object is an example of a **nondeterministic finite-state automaton**, or **NFA**, so named because its behavior is not uniquely *determined* by the input string. Formally, every NFA has five components:

- An arbitrary finite set  $\Sigma$ , called the **input alphabet**.
- Another arbitrary finite set  $Q$ , whose elements are called **states**.
- An arbitrary **transition** function  $\delta : Q \times \Sigma \rightarrow 2^Q$ .
- A **start state**  $s \in Q$ .

- A subset  $A \subseteq Q$  of **accepting states**.

The only difference from the formal definition of *deterministic* finite-state automata is the domain of the transition function. In a DFA, the transition function always returns a single state; in an NFA, the transition function returns a *set* of states, which could be empty, or all of  $Q$ , or anything in between.

Just like DFAs, the behavior of an NFA is governed by an **input string**  $w \in \Sigma^*$ , which the machine reads one symbol at a time, from left to right. Unlike DFAs, however, an NFA does not maintain a single current state, but rather a *set* of current states. Whenever the NFA reads a symbol  $a$ , its set of current states changes from  $C$  to  $\delta(C, a) := \bigcup_{q \in C} \delta(q, a)$ . After all symbols have been read, the NFA **accepts**  $w$  if its current state set contains *at least one* accepting state and **rejects**  $w$  otherwise. In particular, if the set of current states ever becomes empty, it will stay empty forever, and the NFA will reject.

More formally, we define the function  $\delta^*: Q \times \Sigma^* \rightarrow 2^Q$  that transitions on *strings* as follows:

$$\delta^*(q, w) := \begin{cases} \{q\} & \text{if } w = \varepsilon, \\ \bigcup_{r \in \delta(q, a)} \delta^*(r, x) & \text{if } w = ax. \end{cases}$$

The NFA  $(Q, \Sigma, \delta, s, A)$  **accepts**  $w \in \Sigma^*$  if and only if  $\delta^*(s, w) \cap A \neq \emptyset$ .

We can equivalently define an NFA as a directed graph whose vertices are the states  $Q$ , whose edges are labeled with symbols from  $\Sigma$ . We no longer require that every vertex has exactly one outgoing edge with each label; it may have several such edges or none. An NFA accepts a string  $w$  if the graph contains *at least one* walk from the start state to an accepting state whose label is  $w$ .

It's arguably more natural to an arbitrary set of start states  $S \subseteq Q$  instead of just one. Then an NFA accepts a string  $w$  if and only if there is a sequence of transitions consistent with  $w$  from *some* start state to *some* accepting state, or more formally if  $\delta^*(S, q) \cap A \neq \emptyset$ . Change the definition and chase through all the theorems? Or prove equivalence and bounce back and forth, like we already do for  $\varepsilon$ -transitions?

## 4.2 Intuition

There are at least three useful ways to think about non-determinism.

**Clairvoyance.** Whenever an NFA reads symbol  $a$  in state  $q$ , it *chooses* the next state from the set  $\delta(q, a)$ , always *magically* choosing a state that leads to the NFA accepting the input string, unless no such choice is possible. As the BSD fortune file put it, “Nondeterminism means never having to say you’re wrong.”<sup>1</sup> Of course real machines can’t actually look into the future; that’s why I used the word “magic”.

**Parallel threads.** An arguably more “realistic” view is that when an NFA reads symbol  $a$  in state  $q$ , it spawns an independent execution thread for each state in  $\delta(q, a)$ . In particular, if  $\delta(q, a)$  is empty, the current thread simply dies. The NFA accepts if *at least one* thread is in an accepting state after it reads the last input symbol.

<sup>1</sup>This sentence is a riff on a horrible aphorism that was (sadly) popular in the US in the 70s and 80s. Fortunately, everyone seems to have forgotten the original saying, except maybe for that one time it was mocked on *The Simpsons*. Ah, who am I kidding? Nobody remembers *The Simpsons* either.

Equivalently, we can imagine that when an NFA reads symbol  $a$  in state  $q$ , it branches into several parallel universes, one for each state in  $\delta(q, a)$ . If  $\delta(q, a)$  is empty, the NFA destroys the universe (including itself). Similarly, if the NFA finds itself in a non-accepting state when the input ends, the NFA destroys the universe. Thus, when the input is gone, only universes in which the NFA somehow chose a path to an accept state still exist. One slight disadvantage of this metaphor is that if an NFA reads a string that is not in its language, it destroys *all* universes.

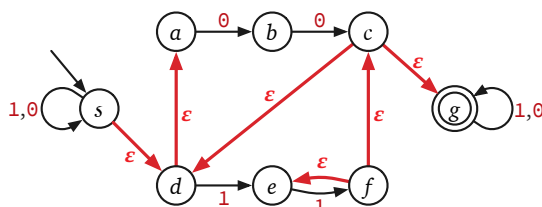
**Proofs/oracles.** Finally, we can treat NFAs not as a mechanism for *computing* something, but as a mechanism for *verifying proofs*. If we want to *prove* that a string  $w$  contains one of the suffixes **00** or **11**, it suffices to demonstrate a single walk in our example NFA that starts at  $s$  and ends at  $c$ , and whose edges are labeled with the symbols in  $w$ . Equivalently, whenever the NFA faces a nontrivial choice, the prover can simply tell the NFA which state to move to next.

This intuition can be formalized as follows. Consider a *deterministic* finite state machine whose input alphabet is the product  $\Sigma \times \Omega$  of an *input* alphabet  $\Sigma$  and an *oracle* alphabet  $\Omega$ . Equivalently, we can imagine that this DFA reads simultaneously from two strings of the same length: the *input* string  $w$  and the *oracle* string  $\omega$ . In either formulation, the transition function has the form  $\delta: Q \times (\Sigma \times \Omega) \rightarrow Q$ . As usual, this DFA accepts the pair  $(w, \omega) \in (\Sigma \times \Omega)^*$  if and only if  $\delta^*(s, (w, \omega)) \in A$ . Finally,  $M$  **nondeterministically accepts** the string  $w \in \Sigma^*$  if there is an oracle string  $\omega \in \Omega^*$  with  $|\omega| = |w|$  such that  $(w, \omega) \in L(M)$ .

### 4.3 $\epsilon$ -Transitions

It is fairly common for NFAs to include so-called  **$\epsilon$ -transitions**, which allow the machine to change state without reading an input symbol. An NFA with  $\epsilon$ -transitions accepts a string  $w$  if and only if there is a sequence of transitions  $s \xrightarrow{a_1} q_1 \xrightarrow{a_2} q_2 \xrightarrow{a_3} \dots \xrightarrow{a_\ell} q_\ell$  where the final state  $q_\ell$  is accepting, each  $a_i$  is either  $\epsilon$  or a symbol in  $\Sigma$ , and  $a_1 a_2 \dots a_\ell = w$ .

For example, consider the following NFA with  $\epsilon$ -transitions. (For this example, we indicate the  $\epsilon$ -transitions using large red arrows; we won't normally do that.) This NFA deliberately has more  $\epsilon$ -transitions than necessary.



A (rather silly) NFA with  $\epsilon$ -transitions

The NFA starts as usual in state  $s$ . If the input string is **100111**, the the machine might non-deterministically choose the following transitions and then accept.

$$s \xrightarrow{\epsilon} s \xrightarrow{1} s \xrightarrow{\epsilon} d \xrightarrow{\epsilon} a \xrightarrow{0} b \xrightarrow{0} c \xrightarrow{\epsilon} d \xrightarrow{1} e \xrightarrow{1} f \xrightarrow{\epsilon} e \xrightarrow{1} f \xrightarrow{\epsilon} c \xrightarrow{\epsilon} g$$

More formally, the transition function in an NFA with  $\epsilon$ -transitions has a slightly larger domain  $\delta: Q \times (\Sigma \cup \{\epsilon\}) \rightarrow 2^Q$ . The  **$\epsilon$ -reach** of a state  $q \in Q$  consists of all states  $r$  that satisfy one of the following conditions:

- either  $r = q$ ,

- or  $r \in \delta(q', \varepsilon)$  for some state  $q'$  in the  $\varepsilon$ -reach of  $q$ .

In other words,  $r$  is in the  $\varepsilon$ -reach of  $q$  if there is a (possibly empty) sequence of  $\varepsilon$ -transitions leading from  $q$  to  $r$ . For example, in the example NFA above, the  $\varepsilon$ -reach of state  $f$  is  $\{a, c, d, f, g\}$ .

Now we redefine the extended transition function  $\delta^*: Q \times \Sigma^* \rightarrow 2^Q$ , which transitions on arbitrary strings, as follows:

$$\delta^*(p, w) := \begin{cases} \varepsilon\text{-reach}(p) & \text{if } w = \varepsilon, \\ \bigcup_{r \in \varepsilon\text{-reach}(p)} \bigcup_{q \in \delta(r, a)} \delta^*(q, x) & \text{if } w = ax. \end{cases}$$

If we abuse notation by writing  $\delta(S, a) = \bigcup_{q \in S} \delta(q, a)$  and  $\delta^*(S, w) = \bigcup_{q \in S} \delta^*(q, w)$  and  $\varepsilon\text{-reach}(S) = \bigcup_{q \in S} \varepsilon\text{-reach}(q)$  for any subset of states  $S \subseteq Q$ , this definition simplifies as follows:

$$\delta^*(p, w) := \begin{cases} \varepsilon\text{-reach}(p) & \text{if } w = \varepsilon, \\ \delta^*(\delta(\varepsilon\text{-reach}(p), a), x) & \text{if } w = ax. \end{cases}$$

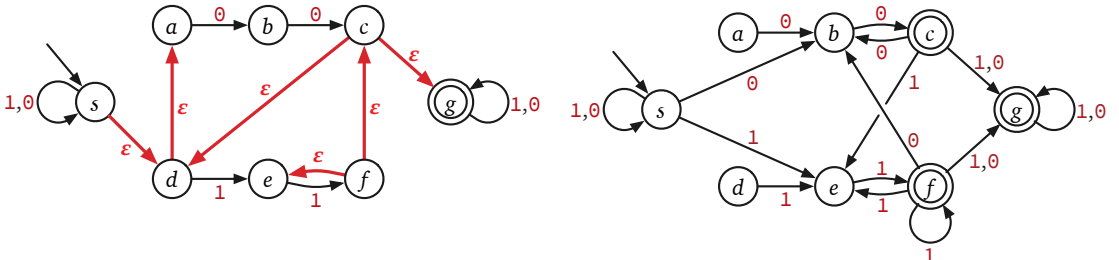
Finally, as usual, an NFA with  $\varepsilon$ -transitions accepts a string  $w$  if and only if  $\delta^*(s, w)$  contains at least one accepting state.

Although it may appear at first that  $\varepsilon$ -transitions give us a more powerful set of machines, NFAs with and without  $\varepsilon$ -transitions are actually equivalent. Given an NFA  $M = (\Sigma, Q, s, A, \delta)$  with  $\varepsilon$ -transitions, we can construct an equivalent NFA  $M' = (\Sigma, Q', s', A', \delta')$  without  $\varepsilon$ -transitions as follows:

$$\begin{aligned} Q' &:= Q \\ s' &= s \\ A' &= \{q \in Q \mid \varepsilon\text{-reach}(q) \cap A \neq \emptyset\} \\ \delta'(q, a) &= \delta(\varepsilon\text{-reach}(q), a) \end{aligned}$$

Straightforward definition-chasing now implies that  $M$  and  $M'$  accept exactly the same language. Thus, whenever we reason about or design NFAs, we are free to either allow or forbid  $\varepsilon$ -transitions, whichever is more convenient for the task at hand.

For example, our previous NFA with  $\varepsilon$ -transitions can be transformed into an equivalent NFA without  $\varepsilon$ -transitions, as shown in the figure below. The NFA on the right has two unreachable states  $a$  and  $d$ , but whatever.



A (rather silly) NFA with  $\varepsilon$ -transitions, and an equivalent NFA without  $\varepsilon$ -transitions

This reduction might be easier to understand incrementally.

- For every transition pair  $p \xrightarrow{\varepsilon} q \xrightarrow{a} r$ , add a direct transition  $p \xrightarrow{a} r$ . This addition does not change the accepted language.
- For each transition  $p \xrightarrow{\varepsilon} q$  where  $q$  is an accepting state, make  $p$  an accepting state. This modification does not change the accepted language.
- When no more of the previous modifications are possible, delete all  $\varepsilon$ -transitions. This modification does not change the accepted language.

## 4.4 Kleene's Theorem

We are now finally in a position to prove the following fundamental fact, first observed by Steven Kleene in 1951:

**Theorem 4.1.** *A language  $L$  can be described by a regular expression if and only if  $L$  is the language accepted by a DFA.*

We will prove Kleene's fundamental theorem in four stages:

- Every DFA can be transformed into an equivalent NFA.
- Every NFA can be transformed into an equivalent DFA.
- Every regular expression can be transformed into an equivalent NFA.
- Every NFA can be transformed into an equivalent regular expression.

The first of these four transformations is completely trivial; a DFA is just a special type of NFA where the transition function always returns a single state. Unfortunately, the other three transformations require a bit more work.

## 4.5 NFA to DFA: The Subset Construction

In the parallel-thread model of NFA execution, an NFA does not have a single current state, but rather a *set* of current states. The evolution of this set of states is *determined* by a modified transition function  $\delta': 2^Q \times \Sigma \rightarrow 2^Q$ , defined by setting  $\delta'(P, a) := \bigcup_{p \in P} \delta(p, a)$  for any set of states  $P \subseteq Q$  and any symbol  $a \in \Sigma$ . When the NFA finishes reading its input string, it accepts if and only if the current set of states intersects the set  $A$  of accepting states.

This formulation makes the NFA completely deterministic! We have just shown that any NFA  $M = (\Sigma, Q, s, A, \delta)$  is equivalent to a DFA  $M' = (\Sigma, Q', s', A', \delta')$  defined as follows:

$$\begin{aligned} Q' &:= 2^Q \\ s' &:= \{s\} \\ A' &:= \{S \subseteq Q \mid S \cap A \neq \emptyset\} \\ \delta'(q', a) &:= \bigcup_{p \in q'} \delta(p, a) \quad \text{for all } q' \subseteq Q \text{ and } a \in \Sigma. \end{aligned}$$

Similarly, any NFA with  $\varepsilon$ -transitions is equivalent to a DFA defined as follows:

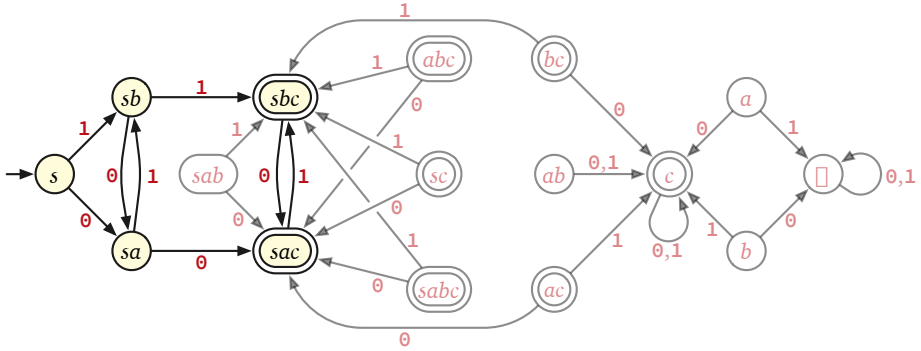
$$\begin{aligned} Q' &:= 2^Q \\ s' &:= \{s\} \end{aligned}$$

$$A' := \{S \subseteq Q \mid \varepsilon\text{-reach}(S) \cap A \neq \emptyset\}$$

$$\delta'(q', a) := \bigcup_{p \in q'} \bigcup_{r \in \varepsilon\text{-reach}(p)} \delta(r, a) \quad \text{for all } q' \subseteq Q \text{ and } a \in \Sigma.$$

This conversion from NFA to DFA is often called the **subset construction**, but that name is somewhat misleading; it's not a "construction" so much as a change in perspective.

For example, the subset construction converts the 4-state NFA on the first page of this note into the following 16-state DFA. To simplify notation, I've named each DFA state using a simple string, omitting the braces and commas from the corresponding subset of NFA states; for example, DFA state *sbc* corresponds to the subset  $\{s, b, c\}$  of NFA states.



The 16-state DFA obtained from our first 4-state NFA by the subset construction.  
Only the five yellow states are reachable from the start state.

An obvious disadvantage of this "construction" is that it (usually) leads to DFAs with far more states than necessary, in part because many states cannot even be reached from the start state. In the example above, there are eleven unreachable states; only five states are reachable from *s*.

## Incremental Subset Construction

Instead of building the entire subset DFA and then discarding the unreachable states, we can avoid the unreachable states from the beginning by constructing the DFA incrementally, essentially by performing a breadth-first search of the DFA graph.

To execute this algorithm by hand, we prepare a table with  $|\Sigma| + 3$  columns, with one row for each DFA state we discover. In order, these columns record the following information:

- The DFA state (as a subset of NFA states)
- The  $\varepsilon$ -reach of the corresponding subset of NFA states
- Whether the DFA state is accepting (that is, whether the  $\varepsilon$ -reach intersects *A*)
- The output of the transition function for each symbol in  $\Sigma$ .

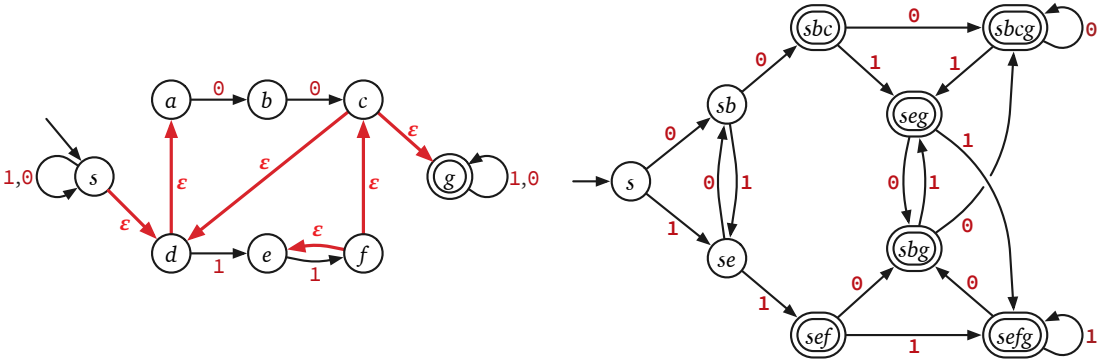
We start with DFA-state  $\{s\}$  in the first row and first column. Whenever we discover an unexplored state in one of the last  $|\Sigma|$  columns, we copy it to the left column in a new row. To reduce notational clutter, we write all subsets of NFA states without braces or commas.

For example, given the NFA with  $\varepsilon$ -transitions from Section 4.3, the standard subset construction would produce a DFA with 256 states, but the incremental subset construction produces a nine-state DFA, described by the following table and illustrated on the next page. We would fill in the first row, for the starting DFA state *s*, as follows:

- The  $\varepsilon$ -reach of NFA state  $s$  is  $\{s, a, d\}$ , so we write  $sad$  in the first column.
- None of the NFA states  $\{s, a, d\}$  is an accepting state, so  $\{s\}$  is not an accepting state of the DFA, so we do *not* check the second column.
- Next,  $\delta'(\{s, a, d\}, 0) = \delta(s, 0) \cup \delta(a, 0) \cup \delta(d, 0) = \{s\} \cup \{b\} \cup \emptyset = \{s, b\}$ , so we write  $sb$  in the third column. Because  $sb$  does not already appear in the first column in any existing row, we have discovered a new DFA state! We start a new row for DFA state  $sb$ .
- Finally,  $\delta'(\{s, a, d\}, 1) = \delta(s, 1) \cup \delta(a, 1) \cup \delta(d, 1) = \{s\} \cup \emptyset \cup \{e\} = \{s, e\}$ , so we write  $se$  in the fourth column, and we start a new row for the new DFA state  $se$ .

We now have two new rows to fill in, corresponding to states  $sb$  and  $se$ . The algorithm continues filling in rows (and discovering new rows) until all rows are filled, ending with the following table:

$q'$	$\varepsilon\text{-reach}(q')$	$q' \in A'?$	$\delta'(q', 0)$	$\delta'(q', 1)$
$s$	$sad$		$sb$	$se$
$sb$	$sabd$		$sbc$	$se$
$se$	$sade$		$sb$	$sef$
$sbc$	$sabcdg$	✓	$sbcg$	$seg$
$sef$	$sacdefg$	✓	$sbg$	$sefg$
$sbcg$	$sabcdg$	✓	$sbcg$	$seg$
$seg$	$sadeg$	✓	$sbg$	$sefg$
$sbg$	$sabdg$	✓	$sbcg$	$seg$
$sefg$	$sacdefg$	✓	$sbg$	$sefg$



An eight-state NFA with  $\varepsilon$ -transitions, and the output of the incremental subset construction for that NFA.

Although it avoids unreachable states, the incremental subset algorithm still gives us a DFA with far more states than necessary, intuitively because it keeps looking for  $00$  and  $11$  substrings even after it's already found one. After all, after the NFA finds both  $00$  and  $11$  as substrings, it doesn't kill all the other parallel execution threads, because it *can't*. NFAs often have significantly fewer states than equivalent DFAs, but that efficiency also makes them kind of stupid.

## 4.6 Regular Expression to NFA: Thompson's Algorithm

We now turn to the core of Kleene's theorem, which claims that regular languages (described by regular expressions) and automatic languages (accepted by finite-state automata) are the same.



**Lemma 4.2.** Every regular language is accepted by a nondeterministic finite-state automaton.

**Proof:** In fact, we will prove the following stronger claim: Every regular language is accepted by an NFA with exactly one accepting state, which is different from its start state. The following construction was first described by Ken Thompson in 1968. Thompson’s algorithm actually proves a stronger statement: For any regular language  $L$ , there is an NFA that accepts  $L$  that has exactly one accepting state  $t$ , which is distinct from the starting state  $s$ .

Let  $R$  be an arbitrary regular expression over an arbitrary finite alphabet  $\Sigma$ . Assume that for any sub-expression  $S$  of  $R$ , the language described by  $S$  is accepted by an NFA with one accepting state distinct from its start state, which we denote pictorially by  $\textcircled{S}$ . There are six cases to consider—three base cases and three recursive cases—mirroring the recursive definition of a regular expression.

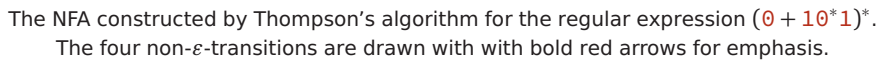
- If  $R = \emptyset$ , then  $L(R) = \emptyset$  is accepted by the trivial NFA  $\textcircled{\emptyset}$ .
- If  $R = \varepsilon$ , then  $L(R) = \{\varepsilon\}$  is accepted by a different trivial NFA  $\textcircled{\varepsilon}$ .
- If  $R = a$  for some symbol  $a \in \Sigma$ , then  $L(R) = \{a\}$  is accepted by the NFA  $\textcircled{a}$ . (The case where  $R$  is a single string with length greater than 1 reduces to the single-symbol case by concatenation, as described in the next case.)
- Suppose  $R = ST$  for some regular expressions  $S$  and  $T$ . The inductive hypothesis implies that the languages  $L(S)$  and  $L(T)$  are accepted by NFAs  $\textcircled{S}$  and  $\textcircled{T}$ , respectively. Then  $L(R) = L(ST) = L(S) \cdot L(T)$  is accepted by the NFA  $\textcircled{S} \xrightarrow{\varepsilon} \textcircled{T}$ , built by connecting the two component NFAs in series.
- Suppose  $R = S + T$  for some regular expressions  $S$  and  $T$ . The inductive hypothesis implies that the language  $L(S)$  and  $L(T)$  are accepted by NFAs  $\textcircled{S}$  and  $\textcircled{T}$ , respectively. Then  $L(R) = L(S + T) = L(S) \cup L(T)$  is accepted by the NFA  $\textcircled{S} \text{ and } \textcircled{T}$ , built by connecting the two component NFAs in parallel with new start and accept states.
- Finally, suppose  $R = S^*$  for some regular expression  $S$ . The inductive hypothesis implies that the language  $L(S)$  is accepted by an NFA  $\textcircled{S}$ . Then the language  $L(R) = L(S^*) = L(S)^*$

is accepted by the NFA

In all cases, the language  $L(R)$  is accepted by an NFA with one accepting state, which is different from its start state, as claimed.  $\square$

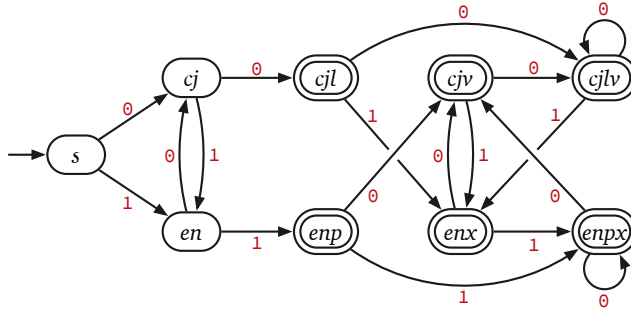
As an example, given the regular expression  $(0 + 10^*1)^*$  of strings containing an even number of 1s, Thompson’s algorithm produces a 14-state NFA shown on the next page. As this example shows, Thompson’s algorithm tends to produce NFAs with many redundant states. Fortunately, just as there are for DFAs, there are algorithms that can reduce any NFA to an equivalent NFA with the smallest possible number of states.

Interestingly, applying the incremental subset algorithm to Thompson’s NFA tends to yield a DFA with relatively few states, in part because the states in Thompson’s NFA tend to have large  $\varepsilon$ -reach, and in part because relatively few of those states are the targets of non- $\varepsilon$ -transitions. Starting with the Thompson’s NFA for  $(0 + 10^*1)^*$ , for example, the incremental subset construction yields a DFA with just five states.



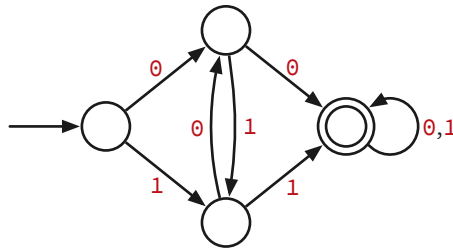
Given this NFA as input, the incremental subset construction computes the following table, leading to a DFA with just nine states. Yeah, the  $\varepsilon$ -reaches get a bit ridiculous; unfortunately, this is typical for Thompson's NFA. As usual, the resulting DFA has far more states than necessary.

$q'$	$\varepsilon\text{-reach}(q')$	$q' \in A'$ ?	$\delta'(q', 0)$	$\delta'(q', 1)$
$s$	$sabdghim$		$cj$	$en$
$cj$	$sabdfghijkm$		$cjl$	$en$
$en$	$sabdfghmno$		$cj$	$enp$
$cjl$	$sabdfghijklmqrtuwz$	✓	$cjlv$	$enx$
$enp$	$sabdfghmnopqrtuwz$	✓	$cjv$	$enpx$
$cjlv$	$sabdfghijklmqrtuvwyz$	✓	$cjlv$	$enx$
$enx$	$sabdfghmnopqrtuvwxyz$	✓	$cjv$	$enpx$
$cjv$	$sabdfghijkmrtuvwyz$	✓	$cjlv$	$enx$
$enpx$	$sabdfghmnopqrtuvwxyz$	✓	$cjv$	$enpx$



The DFA computed by the incremental subset algorithm from Thompson's NFA for  $(0 + 1)^*(00 + 11)(0 + 1)^*$ .

Finally, the DFA-minimization algorithm from the previous lecture note correctly discovers that all six accepting states of the incremental-subset DFA are equivalent, and thus reduces the DFA to just four states.



The minimal DFA that accepts the language  $(0 + 1)^*(00 + 11)(0 + 1)^*$ .

#### \*4.8 NFA to Regular Expression: Han and Wood's Algorithm

The only component of Kleene's theorem left to prove is that every language accepted by an NFA is regular. I'll describe a relatively recent proof, due Yo-Sub Han and Derick Wood in 2005<sup>2</sup>, that is morally equivalent to Kleene's 1951 argument, but uses more modern standard notation.

Recall that a standard NFA can be represented by a state-transition graph, whose vertices are the states and whose edges represent possible transitions. Each edge is labeled with a single symbol in  $\Sigma$ . A string  $w \in \Sigma^*$  is accepted if and only if there is a sequence of transitions

$$s \xrightarrow{a_1} q_1 \xrightarrow{a_2} q_2 \xrightarrow{a_3} \cdots \xrightarrow{a_\ell} q_\ell$$

<sup>2</sup>Yo-Sub Han\* and Derick Wood. The generalization of generalized automata: Expression automata. *International Journal of Foundations of Computer Science* 16(3):499–510, 2005.

where the final state  $q_\ell$  is accepting and  $a_1a_2\cdots a_\ell = w$ .

We've already seen that NFAs can be generalized to include  $\varepsilon$ -transitions; we can push this generalization further. A **string NFA** allows each transition  $p \rightarrow q$  to be labeled with an *arbitrary string*  $x(p \rightarrow q) \in \Sigma^*$ . We are allowed to transition from state  $p$  to state  $q$  if the label  $x(p \rightarrow q)$  is a *prefix* of the remaining input. Thus, a string  $w \in \Sigma^*$  is accepted if and only if there is a sequence of transitions

$$s \xrightarrow{x_1} q_1 \xrightarrow{x_2} q_2 \xrightarrow{x_3} \cdots \xrightarrow{x_\ell} q_\ell$$

where the final state  $q_\ell$  is accepting, and  $x_1 \cdot x_2 \cdot \cdots \cdot x_\ell = w$ . Thus, an NFA with  $\varepsilon$ -transitions is just a string NFA where every label has length 0 or 1. Any string NFA can be converted into an equivalent standard NFA, by subdividing each edge  $p \rightarrow q$  into a path of length  $|x(p \rightarrow q)|$  (unless  $x(p \rightarrow q) = \varepsilon$ ).

Finally, Han and Wood define an **expression automaton** as a finite-state machine where each transition  $p \rightarrow q$  is labeled with an arbitrary *regular expression*  $R(p \rightarrow q)$ . We can transition from state  $p$  to state  $q$  if *any* prefix of the remaining input matches the regular expression  $R(p \rightarrow q)$ . Thus, a string  $w \in \Sigma^*$  is accepted by an expression automaton if and only if there is a sequence of transitions

$$s \xrightarrow{R_1} q_1 \xrightarrow{R_2} q_2 \xrightarrow{R_3} \cdots \xrightarrow{R_\ell} q_\ell$$

where the final state  $q_\ell$  is accepting, and we can write  $w = x_1 \cdot x_2 \cdot \cdots \cdot x_\ell = w$ , where each substring  $x_i$  matches the corresponding regular expression  $R_i$ .

More formally, an expression automaton consists of the following components:

- A finite set  $\Sigma$  called the **input alphabet**
- Another finite set  $Q$  whose elements are called **states**
- A unique **start state**  $s \in Q$
- A unique **target state**  $t \in Q \setminus \{s\}$
- A **transition function**  $R: (Q \setminus \{t\}) \times (Q \setminus \{s\}) \rightarrow \text{Reg}(\Sigma)$ , where  $\text{Reg}(\Sigma)$  is the set of regular expressions over  $\Sigma$ .

The requirement that the start and target states are unique and disatinct is not essential to the model. We impose this requirement for convenience of the equivalence proof; it can be easily enforced using  $\varepsilon$ -transitions.

Expression automata are even more nondeterministic than NFAs. A single string could match several (even infinitely many) transition sequences from  $s$  to  $t$ , and it could match each of those sequences in several (even infinitely many) different ways. A string  $w$  is accepted if *any* decomposition of  $w$  into a sequence of substrings matches *any* sequence of transitions from  $s$  to  $t$ . Conversely, a string might match *no* state sequences, in which case the string is rejected.

Two extreme special cases of expression automata are already familiar. First, every regular language is clearly the language of an expression automaton with exactly two states. Second, with only minor modifications, any DFA or NFA can be converted into an expression automaton with trivial transition expressions. Thompson's algorithm can be used to transform any expression automaton into a standard NFA (with  $\varepsilon$ -transitions), by recursively expanding any nontrivial transition expression. To complete the proof of Kleene's theorem, we show how to convert an arbitrary expression automaton into a regular expression, by repeatedly deleting vertices.

**Lemma 4.3.** *Every expression automaton accepts a regular language.*

**Proof:** Let  $E = (Q, \Sigma, R, s, t)$  be an arbitrary expression automaton. Assume that any expression automaton with fewer states than  $E$  accepts a regular language. There are two cases to consider, depending on the number of states in  $Q$ :

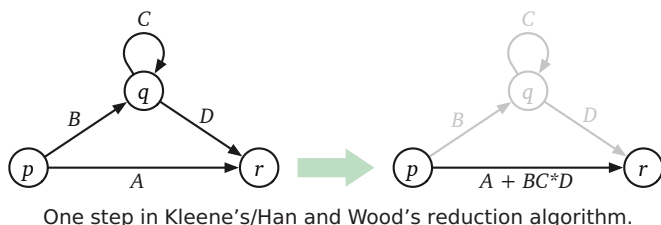
- If  $Q = \{s, t\}$ , then trivially,  $E$  accepts the regular language  $R(s \rightarrow t)$ .
- On the other hand, suppose  $Q$  has more than two states; fix an arbitrary state  $q \in Q \setminus \{s, t\}$ . We modify the automaton, without changing its language, so that state  $q$  is redundant and can be removed. Define a new transition function  $R' : Q \times Q \rightarrow \text{Reg}(\Sigma)$  by setting

$$R'(p \rightarrow r) := R(p \rightarrow r) + R(p \rightarrow q)R(q \rightarrow q)^*R(q \rightarrow r).$$

With this modified transition function in place, any string  $w$  that matches the sequence  $p \rightarrow q \rightarrow q \rightarrow \dots \rightarrow q \rightarrow r$  with any number of  $q$ 's also matches the single transition  $p \rightarrow r$ . Thus, by induction, if  $w$  matches a sequence of states, it also matches the subsequence obtained by removing all  $q$ 's. Let  $E'$  be the expression automaton with states  $Q' = Q \setminus \{q\}$  that uses this modified transition function  $R'$ . This new automaton accepts exactly the same strings as the original automaton  $E$ . Because  $E'$  has fewer states than  $E$ , the inductive hypothesis implies  $E'$  accepts a regular language.

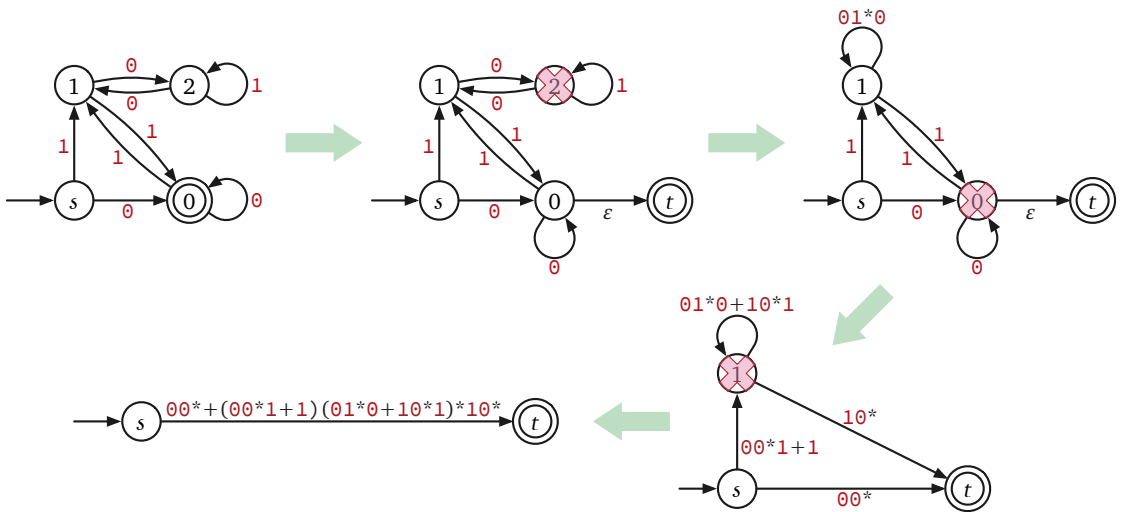
In both cases, we conclude that  $E$  accepts a regular language. □

This proof can be mechanically translated into an algorithm to convert any DFA or NFA into an equivalent regular expression, via a sequence of expression automata with fewer and fewer states, but increasingly complex transition expressions.



The figure on the next page shows Han and Wood's algorithm in action, starting with a DFA that accepts the binary representations of non-negative integers divisible by 3, possibly with extra leading 0s. (State  $i$  means the binary number we've read so far is congruent to  $i \pmod 3$ .) First we convert the DFA into an expression automaton by adding a new accept state. (We don't need to add a new start state, because there are no transitions the original start state  $s$ .) Then we remove state 2, then state 0, and finally state 1, updating the transition expressions between any remaining states at each iteration. For the sake of clarity, edges  $p \rightarrow q$  with  $R(p \rightarrow q) = \emptyset$  are omitted from the figures. The final regular expression  $00^* + (00^*1 + 1)(10^*1 + 01^*0)^*10^*$  can be slightly simplified to  $0^*0 + 0^*1(10^*1 + 01^*0)^*10^*$ , which is precisely the regular expression we gave for this language back in Lecture Note 2!

Given an NFA with  $n$  states (including  $s$  and  $t$ ), Han and Wood's algorithm iteratively removes  $n - 2$  states, updating  $O(n^2)$  transition expressions in each iteration. If the concatenation and Kleene star operations could be performed in constant time, the resulting algorithm would run in  $O(n^3)$  time. However, in the worst case, the transition expressions grows in length by roughly a factor of 4 in each iteration, so the final expression has length  $\Theta(4^n)$ . If we insist on representing the expressions as explicit strings, the worst-case running time is actually  $\Theta(4^n)$ .



Converting a DFA into an equivalent regular expression using Han and Wood's algorithm.

## 4.9 Regular Language Transformations

We have already seen that many functions of regular languages are themselves regular: unions, concatenations, and Kleene closure by definition. and intersections and differences by product construction on DFA. However, the set of regular languages is closed under a much richer class of functions.

Suppose we wanted to prove that regular languages are closed under some function  $f$ ; that is, for every regular language  $L$ , we want to prove that the language  $f(L)$  is also regular. There are two general techniques to prove such a statement:

- Describe an algorithm that transforms an arbitrary regular expression  $R$  into a new regular expression  $R'$  such that  $L(R') = f(L(R))$ .
- Describe an algorithm that transforms an arbitrary DFA  $M$  into a new NFA  $M'$  such that  $L(M') = f(L(M))$ .

The equivalence between regular expressions and finite automata implies that *in principle* we can always use either technique, but in practice, the second one is far more powerful and usually simpler. The asymmetry in the second technique is important. We start with a DFA for  $L$  to impose as much structure as possible in the input; we aim for an NFA with  $\epsilon$ -transitions to give ourselves as much freedom as possible in the output.<sup>3</sup>

For our first example, I'll describe proofs using both techniques.

**Lemma 4.4.** *For any regular language  $L$ , the language  $L^R = \{w^R \mid w \in L\}$  is also regular.*

**Proof (regular expression to regular expression):** Let  $R$  be an arbitrary regular expression such that  $L = L(R)$ . Assume for any proper subexpression  $S$  of  $R$  that  $L(S)^R$  is regular. There are five cases to consider, mirroring the recursive definition of regular expressions:

- If  $R = \emptyset$ , then  $L^R = L = \emptyset$ , so  $L(R) = L^R$ .
- Suppose  $R$  consists of a single word  $w$ . Let  $R' = w^R$ . Then  $L(R') = \{w^R\} = L^R$ .

<sup>3</sup>We could give ourselves even more freedom by constructing an *expression* automaton, but creativity thrives on constraint.

- Suppose  $R = A + B$ . The inductive hypothesis implies that there are regular expressions  $A'$  and  $B'$  such that  $L(A') = L(A)^R$  and  $L(B') = L(B)^R$ . Let  $R' = A' + B'$ . Then  $L(R') = L(A') \cup L(B') = L(A)^R \cup L(B)^R = (L(A) \cup L(B))^R = L^R$ .
- Suppose  $R = A \cdot B$ . The inductive hypothesis implies that there are regular expressions  $A'$  and  $B'$  such that  $L(A') = L(A)^R$  and  $L(B') = L(B)^R$ . Let  $R' = B' \cdot A'$ . Then  $L(R') = L(B') \cdot L(A') = L(B)^R \cdot L(A)^R = (L(A) \cdot L(B))^R = L^R$ .
- Finally, suppose  $R = A^*$ . The inductive hypothesis implies that there is a regular expression  $A'$  such that  $L(A') = L(A)^R$ . Let  $R' = (A')^*$ . Then  $L(R') = L(A')^* = (L(A)^R)^* = (L(A)^*)^R = L^R$ .

In all cases, we have constructed a regular expression  $R'$  such that  $L(R') = L^R$ . We conclude that  $L^R$  is regular.  $\square$

Careful readers may be unsatisfied with the previous proof, because it assumes several “obvious” properties of string and language reversal. Specifically, for all strings  $x$  and  $y$  and all languages  $L$  and  $L'$ , we assumed the following:

- $(x \cdot y)^R = y^R \cdot x^R$
- $(L \cdot L')^R = (L')^R \cdot L^R$ .
- $(L \cup L')^R = L^R \cup (L')^R$ .
- $(L^*)^R = (L^R)^*$ .

All of these claims are all easy to prove by inductive definition-chasing.

**Proof (DFA to NFA):** Let  $M = (\Sigma, Q, s, A, \delta)$  be an arbitrary DFA that accepts  $L$ . We construct an NFA  $M^R = (\Sigma, Q^R, s^R, A^R, \delta^R)$  with  $\varepsilon$ -transitions that accepts  $L^R$ , intuitively by reversing every transition in  $M$ , and swapping the roles of the start state and the accepting states. Because  $M$  does not have a unique accepting state, we need to introduce a special start state  $s^R$ , with  $\varepsilon$ -transitions to each accepting state in  $M$ . These are the only  $\varepsilon$ -transitions in  $M^R$ .

$$Q^R = Q \cup \{s^R\}$$

$$A^R = \{s\}$$

$$\delta^R(s^R, \varepsilon) = A$$

$$\delta^R(s^R, a) = \emptyset \quad \text{for all } a \in \Sigma$$

$$\delta^R(q, \varepsilon) = \emptyset \quad \text{for all } q \in Q$$

$$\delta^R(q, a) = \{p \mid q \in \delta(p, a)\} \quad \text{for all } q \in Q \text{ and } a \in \Sigma$$

Routine inductive definition-chasing now implies that the reversal of any sequence  $q_0 \rightarrow q_1 \rightarrow \dots \rightarrow q_\ell$  of transitions in  $M$  is a valid sequence  $q_\ell \rightarrow q_{\ell-1} \rightarrow \dots \rightarrow q_0$  of transitions in  $M^R$ . Because the transitions retain their labels (but reverse directions), it follows that  $M$  accepts any string  $w$  if and only if  $M^R$  accepts  $w^R$ .

We conclude that the NFA  $M^R$  accepts  $L^R$ , so  $L^R$  must be regular.  $\square$

**Lemma 4.5.** *For any regular language  $L$ , the language  $\text{half}(L) := \{w \mid ww \in L\}$  is also regular.*

**Proof:** Let  $M = (\Sigma, Q, s, A, \delta)$  be an arbitrary DFA that accepts  $L$ .

Intuitively, we construct an NFA  $M'$  that reads its input string  $w$  and simulates the original DFA  $M$  reading the input string  $ww$ . Our overall strategy has three parts:



- First  $M'$  non-deterministically *guesses* the state  $h = \delta^*(s, w)$  that  $M$  will reach after reading input  $w$ . (We can't just run  $M$  on input  $w$  to compute the correct state  $h$ , because that would consume the input string!)
- Then  $M'$  runs two copies of  $M$  in parallel (using a product construction): a “left” copy starting at  $s$  and a “right” copy starting at the (guessed) halfway state  $h$ .
- Finally, when  $M'$  is done reading  $w$ , it accepts if and only if the first copy of  $M$  actually stopped in state  $h$  (so our initial guess was correct) and the second copy of  $M$  stopped in an accepting state. That is,  $M'$  accepts if and only if  $\delta^*(s, w) = h$  and  $\delta^*(h, w) \in A$ .

To implement this strategy,  $M'$  needs to maintain *three* states of  $M$ : the state of the left copy of  $M$ , the guess  $h$  for the halfway state, and the state of the right copy of  $M$ . The first and third states evolve according to the transition function  $\delta$ , but the second state never changes. Finally, to implement the non-deterministic guessing,  $M'$  includes a special start state  $s'$  with  $\varepsilon$ -transitions to every triple of the form  $(s, h, h)$ .

Summing up, our new NFA  $M' = (\Sigma, Q', s', A', \delta')$  is formally defined as follows.

$$\begin{aligned}
 Q' &= (Q \times Q \times Q) \cup \{s'\} \\
 A' &= \{(h, h, q) \mid h \in Q \text{ and } q \in A\} \\
 \delta'(s', \varepsilon) &= \{(s, h, h) \mid h \in Q\} \\
 \delta'(s', a) &= \emptyset && \text{for all } a \in \Sigma \\
 \delta'((p, h, q), \varepsilon) &= \emptyset && \text{for all } p, h, q \in Q \\
 \delta'((p, h, q), a) &= \{(\delta(p, a), h, \delta(q, a))\} && \text{for all } p, h, q \in Q \text{ and } a \in \Sigma
 \end{aligned}$$

□

## Exercises

1. For each of the following regular expressions, describe or draw two finite-state machines:

- An NFA that accepts the same language, constructed using Thompson's algorithm.
- An equivalent DFA, built from the previous NFA using the incremental subset construction. For each state in your DFA, identify the corresponding subset of states in your NFA. Your DFA should have no unreachable states.

- $(01 + 10)^*(0 + 1 + \varepsilon)$
- $(\varepsilon + 1)(01)^*(\varepsilon + 0)$
- $1^* + (10)^* + (100)^*$
- $(\varepsilon + 0 + 00)(1 + 10 + 100)^*$
- $((0 + 1)(0 + 1))^*$
- $\varepsilon + 0(0 + 1)^* + 1(1 + 0)^*$

2. The accepting language of an NFA  $M = (\Sigma, Q, s, A, \delta)$  is defined as follows:

$$L(M) := \{w \in \Sigma^* \mid \delta^*(s, w) \cap A \neq \emptyset\}.$$

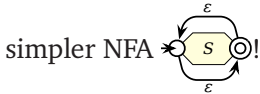
Kleene's theorem (described here as Han and Wood's algorithm) implies that  $L(M)$  is regular. Prove that the following languages associated with  $M$  are also regular:



- (a)  $L^\forall(M) := \{w \in \Sigma^* \mid A \subseteq \delta^*(s, w)\}$ . That is, a string  $w$  is in the language  $L^\forall(M)$  if and only if  $\delta^*(s, w)$  contains *every* accepting states.
- (b)  $L^\subseteq(M) := \{w \in \Sigma^* \mid \delta^*(s, w) \subseteq A\}$ . That is, a string  $w$  is in the language  $L^\subseteq(M)$  if and only if  $\delta^*(s, w)$  contains *only* accepting states.
- (c)  $L^\equiv(M) := \{w \in \Sigma^* \mid \delta^*(s, w) = A\}$ . That is, a string  $w$  is in the language  $L^\equiv(M)$  if and only if  $\delta^*(s, w)$  is exactly the set of accepting states.

3. A certain professor who really should know better once woke up in the middle of the night with a startling revelation—Thompson’s algorithm doesn’t need all those  $\varepsilon$ -transitions! Filled with the certainty that only sleep deprivation can bring, he ran to his laptop and quickly changed two cases in his description of Thompson’s algorithm.

- When  $R = S \cdot T$ , instead of connecting the accept state of  $\textcircled{S}$  to the start state of  $\textcircled{T}$  with an  $\varepsilon$ -transition, we can **just** identify those two states to build the simpler NFA  $\textcircled{S} \textcircled{T}$ !
- When  $R = S^*$ , instead of introducing two new states and four  $\varepsilon$ -transitions, we can **just** add two  $\varepsilon$ -transitions between the start and accept states of  $\textcircled{S}$  to build the



Satisfied with his simplification, he thanked the penguin who gave him the idea, and then flew his hat back into the ocean marshmallows, where a giant man with the head of a dog gave him the power of bread. The next morning, while he was proudly teaching his new simplified proof for the first time, he realized his horrible mistake.

Prove that *neither* of the professor’s optimizations is actually correct.

- (a) Find a regular expression  $R$ , such that the NFA constructed from  $R$  by Thompson’s algorithm **with only the first modification** accepts strings that are not in  $L(R)$ .
- (b) Find a regular expression  $R$ , such that the NFA constructed from  $R$  by Thompson’s algorithm **with only the second modification** accepts strings that are not in  $L(R)$ .

4. A **Moore machine** is a variant of a finite-state automaton that produces output; Moore machines are sometimes called finite-state *transducers*. For purposes of this problem, a Moore machine formally consists of six components:

- A finite set  $\Sigma$  called the input alphabet
- A finite set  $\Gamma$  called the output alphabet
- A finite set  $Q$  whose elements are called states
- A start state  $s \in Q$
- A transition function  $\delta: Q \times \Sigma \rightarrow Q$
- An output function  $\omega: Q \rightarrow \Gamma$

More intuitively, a Moore machine is a graph with a special start vertex, where every node (state) has one outgoing edge labeled with each symbol from the input alphabet, and each node (state) is additionally labeled with a symbol from the output alphabet.

The Moore machine reads an input string  $w \in \Sigma^*$  one symbol at a time. For each symbol, the machine changes its state according to the transition function  $\delta$ , and then outputs the symbol  $\omega(q)$ , where  $q$  is the new state. Formally, we recursively define a *transducer* function  $\omega^*: \Sigma^* \times Q \rightarrow \Gamma^*$  as follows:

$$\omega^*(w, q) = \begin{cases} \varepsilon & \text{if } w = \varepsilon \\ \omega(\delta(a, q)) \cdot \omega^*(x, \delta(a, q)) & \text{if } w = ax \end{cases}$$

Given the input string  $w \in \Sigma^*$ , the machine outputs the string  $\omega^*(w, s) \in \Gamma^*$ . To simplify notation, we define  $M(w) = \omega^*(w, s)$ .

Finally, the **output language**  $L^\circ(M)$  of a Moore machine  $M$  is the set of all strings that the machine can output:

$$L^\circ(M) := \{M(w) \mid w \in \Sigma^*\}$$

- (a) Let  $M$  be an arbitrary Moore machine. Prove that  $L^\circ(M)$  is a regular language.
- (b) Let  $M$  be an arbitrary Moore machine whose input alphabet  $\Sigma$  and output alphabet  $\Gamma$  are identical. Prove that the language

$$L^=(M) = \{w \in \Sigma^* \mid M(w) = w\}$$

is regular. Strings in  $L^=(M)$  are also called *fixed points* of the function  $M: \Sigma^* \rightarrow \Sigma^*$ .

- \* (c) As in part (b), let  $M$  be an arbitrary Moore machine whose input and output alphabets are identical. Prove that the language  $\{w \in \Sigma^* \mid M(M(w)) = w\}$  is regular.

[Hint: Parts (a) and (b) are easier than they look!]

5. A **Mealy machine** is a variant of a finite-state automaton that produces output; Mealy machines are sometimes called finite-state *transducers*. For purposes of this problem, a Mealy machine formally consists of six components:

- A finite set  $\Sigma$  called the input alphabet
- A finite set  $\Gamma$  called the output alphabet
- A finite set  $Q$  whose elements are called states
- A start state  $s \in Q$
- A transition function  $\delta: Q \times \Sigma \rightarrow Q$
- An output function  $\omega: Q \times \Sigma \rightarrow \Gamma$

More intuitively, a Mealy machine is a graph with a special start vertex, where every node (state) has one outgoing edge labeled with each symbol from the input alphabet, and each edge (transition) is additionally labeled with a symbol from the output alphabet. (Mealy machines are closely related to *Moore* machines, which produce output at each *state* instead of at each transition.)

The Mealy machine reads an input string  $w \in \Sigma^*$  one symbol at a time. For each symbol, the machine changes its state according to the transition function  $\delta$ , and simultaneously outputs a symbol according to the output function  $\omega$ . Formally, we recursively define a *transducer* function  $\omega^*: Q \times \Sigma^* \rightarrow \Gamma^*$  as follows:

$$\omega^*(q, w) = \begin{cases} \varepsilon & \text{if } w = \varepsilon \\ \omega(q, a) \cdot \omega^*(\delta(q, a), x) & \text{if } w = ax \end{cases}$$

Given any input string  $w \in \Sigma^*$ , the machine outputs the string  $\omega^*(w, s) \in \Gamma^*$ . To simplify notation, we define  $M(w) = \omega^*(w, s)$ .

Finally, the **output language**  $L^\circ(M)$  of a Mealy machine  $M$  is the set of all strings that the machine can output:

$$L^\circ(M) := \{M(w) \mid w \in \Sigma^*\}$$

- (a) Let  $M$  be an arbitrary Mealy machine. Prove that  $L^\circ(M)$  is a regular language.
- (b) Let  $M$  be an arbitrary Mealy machine whose input alphabet  $\Sigma$  and output alphabet  $\Gamma$  are identical. Prove that the language

$$L^=(M) = \{w \in \Sigma^* \mid w = \omega^*(s, w)\}$$

is regular.  $L^=(M)$  consists of all strings  $w$  such that  $M$  outputs  $w$  when given input  $w$ ; these are also called *fixed points* for the transducer function  $\omega^*$ .

- \* (c) As in part (b), let  $M$  be an arbitrary Mealy machine whose input and output alphabets are identical. Prove that the language  $\{w \in \Sigma^* \mid M(M(w)) = w\}$  is regular.

[Hint: Parts (a) and (b) are easier than they look!]

6. Let  $L \subseteq \Sigma^*$  be an arbitrary regular language. Prove that the following languages are regular. Assume  $\# \in \Sigma$ .

- (a)  $\text{censor}(L) := \{\#^{|w|} \mid w \in L\}$
- (b)  $\text{dehash}(L) = \{\text{dehash}(w) \mid w \in L\}$ , where  $\text{dehash}(w)$  is the subsequence of  $w$  obtained by deleting every  $\#$ .
- (c)  $\text{insert}\#(L) := \{x\#y \mid xy \in L\}$ .
- (d)  $\text{delete}\#(L) := \{xy \mid x\#y \in L\}$ .
- (e)  $\text{prefix}(L) := \{x \in \Sigma^* \mid xy \in L \text{ for some } y \in \Sigma^*\}$
- (f)  $\text{suffix}(L) := \{y \in \Sigma^* \mid xy \in L \text{ for some } x \in \Sigma^*\}$
- (g)  $\text{substring}(L) := \{y \in \Sigma^* \mid xyz \in L \text{ for some } x, z \in \Sigma^*\}$
- (h)  $\text{superstring}(L) := \{xyz \mid y \in L \text{ and } x, z \in \Sigma^*\}$
- (i)  $\text{cycle}(L) := \{xy \mid x, y \in \Sigma^* \text{ and } yx \in L\}$
- (j)  $\text{prefmax}(L) := \{x \in L \mid xy \in L \iff y = \varepsilon\}$ .
- (k)  $\text{sufmin}(L) := \{xy \in L \mid y \in L \iff x = \varepsilon\}$ .

- (l)  $\text{minimal}(L) := \{w \in L \mid \text{no proper substring of } w \text{ is in } L\}$ .
- (m)  $\text{maximal}(L) := \{w \in L \mid \text{no proper superstring of } w \text{ is in } L\}$ .
- (n)  $\text{evens}(L) := \{\text{evens}(w) \mid w \in L\}$ , where  $\text{even}(w)$  is the subsequence of  $w$  containing every even-indexed symbol. For example,  $\text{evens}(\text{EVENINDEX}) = \text{VNNE}$ .
- (o)  $\text{evens}^{-1}(L) := \{w \in \Sigma^* \mid \text{evens}(w) \in L\}$ .
- (p)  $\text{subseq}(L) := \{x \in \Sigma^* \mid x \text{ is a subsequence of some } y \in L\}$
- (q)  $\text{superseq}(L) := \{x \in \Sigma^* \mid \text{some } y \in L \text{ is a subsequence of } x\}$
- (r)  $\text{swap}(L) := \{\text{swap}(w) \mid w \in L\}$ , where  $\text{swap}(w)$  is defined recursively as follows:

$$\text{swap}(w) = \begin{cases} w & \text{if } |w| \leq 1 \\ ba \cdot \text{swap}(x) & \text{if } w = abx \text{ for some } a, b \in \Sigma \text{ and } x \in \Sigma^* \end{cases}$$

- (s)  $\text{oneswap}(L) := \{xbay \mid xaby \in L \text{ where } a, b \in \Sigma \text{ and } x, y \in \Sigma^*\}$ .
- (t)  $\text{left}(L) := \{x \in \Sigma^* \mid xy \in L \text{ for some } y \in \Sigma^* \text{ where } |x| = |y|\}$
- (u)  $\text{right}(L) := \{y \in \Sigma^* \mid xy \in L \text{ for some } x \in \Sigma^* \text{ where } |x| = |y|\}$
- (v)  $\text{middle}(L) := \{y \in \Sigma^* \mid xyz \in L \text{ for some } x, z \in \Sigma^* \text{ where } |x| = |y| = |z|\}$
- (w)  $\text{halfseq}(L) := \{w \in \Sigma^* \mid w \text{ is a subsequence of some string } x \in L \text{ where } |x| = 2 \cdot |w|\}$
- (x)  $\text{third}(L) := \{w \in \Sigma^* \mid www \in L\}$
- (y)  $\text{palin}(L) := \{w \in \Sigma^* \mid ww^R \in L\}$
- (z)  $\text{drome}(L) := \{w \in \Sigma^* \mid w^R w \in L\}$

7. Let  $L$  and  $L'$  be arbitrary regular languages over the alphabet  $\{0, 1\}$ . Prove that the following languages are also regular:

- (a)  $L \sqcap L' := \{x \sqcap y \mid x \in L \text{ and } y \in L' \text{ and } |x| = |y|\}$ , where  $x \sqcap y$  denotes bitwise-and. For example,  $0011 \sqcap 0101 = 0001$ .
- (b)  $L \sqcup L' := \{x \sqcup y \mid x \in L \text{ and } y \in L' \text{ with } |x| = |y|\}$ , where  $x \sqcup y$  denotes bitwise-or. For example,  $0011 \sqcup 0101 = 0111$ .
- (c)  $L \boxplus L' := \{x \boxplus y \mid x \in L \text{ and } y \in L' \text{ with } |x| = |y|\}$ , where  $x \boxplus y$  denotes bitwise-exclusive-or. For example,  $0011 \boxplus 0101 = 0110$ .
- (d)  $\text{faro}(L, L') := \{\text{faro}(x, z) \mid x \in L \text{ and } z \in L' \text{ with } |x| = |z|\}$ , where

$$\text{faro}(x, z) := \begin{cases} z & \text{if } x = \varepsilon \\ a \cdot \text{faro}(z, y) & \text{if } x = ay \end{cases}$$

For example,  $\text{faro}(0011, 0101) = 00011011$ .

- (e)  $\text{shuffles}(L, L') := \bigcup_{w \in L, y \in L'} \text{shuffles}(w, y)$ , where  $\text{shuffles}(w, y)$  is the set of all strings obtained by shuffling  $w$  and  $y$ , or equivalently, all strings in which  $w$  and  $y$  are

complementary subsequences. Formally:

$$\text{shuffles}(w, y) = \begin{cases} \{y\} & \text{if } w = \varepsilon \\ \{w\} & \text{if } y = \varepsilon \\ \{a\} \cdot \text{shuffles}(x, y) \cup \{b\} \cdot \text{shuffles}(w, z) & \text{if } w = ax \text{ and } y = bz \end{cases}$$

For example,  $\text{shuffles}(\mathbf{01}, \mathbf{10}) = \{\mathbf{0101}, \mathbf{0110}, \mathbf{1001}, \mathbf{1010}\}$  and  $\text{shuffles}(\mathbf{00}, \mathbf{11}) = \{\mathbf{0011}, \mathbf{0101}, \mathbf{1001}, \mathbf{0110}, \mathbf{1010}, \mathbf{1100}\}$ .

8. (a) Let  $\text{inc}: \{0, 1\}^* \rightarrow \{0, 1\}^*$  denote the *increment* function, which transforms the binary representation of an arbitrary integer  $n$  into the binary representation of  $n + 1$ , truncated to the same number of bits. For example:

$$\text{inc}(\mathbf{0010}) = \mathbf{0011} \quad \text{inc}(\mathbf{0111}) = \mathbf{1000} \quad \text{inc}(\mathbf{1111}) = \mathbf{0000} \quad \text{inc}(\varepsilon) = \varepsilon$$

Let  $L \subseteq \{0, 1\}^*$  be an arbitrary regular language. Prove that  $\text{inc}(L) = \{\text{inc}(w) \mid w \in L\}$  is also regular.

- (b) Let  $\text{dbl}: \{0, 1\}^* \rightarrow \{0, 1\}^*$  denote the *doubling* function, which transforms the binary representation of an arbitrary integer  $n$  into the binary representation of  $2n$ , truncated to the same number of bits. For example:

$$\text{dbl}(\mathbf{0010}) = \mathbf{0100} \quad \text{dbl}(\mathbf{0111}) = \mathbf{1110} \quad \text{dbl}(\mathbf{1111}) = \mathbf{1110} \quad \text{dbl}(\varepsilon) = \varepsilon$$

Let  $L \subseteq \{0, 1\}^*$  be an arbitrary regular language. Prove that  $\text{dbl}(L) = \{\text{dbl}(w) \mid w \in L\}$  is also regular.

- \* (c) Let  $\text{tpl}: \{0, 1\}^* \rightarrow \{0, 1\}^*$  denote the *tripling* function, which transforms the binary representation of an arbitrary integer  $n$  into the binary representation of  $3n$ , truncated to the same number of bits. For example:

$$\text{tpl}(\mathbf{0010}) = \mathbf{0110} \quad \text{tpl}(\mathbf{0111}) = \mathbf{0101} \quad \text{tpl}(\mathbf{1111}) = \mathbf{1101} \quad \text{tpl}(\varepsilon) = \varepsilon$$

Let  $L \subseteq \{0, 1\}^*$  be an arbitrary regular language. Prove that  $\text{tpl}(L) = \{\text{tpl}(w) \mid w \in L\}$  is also regular. [Hint: It may be easier to consider the language  $\text{tpl}(L^R)^R$  first.]

- \*9. Let  $L \subseteq \Sigma^*$  be an arbitrary regular language. Prove that the following languages are regular.

(a)  $\text{sqrt}(L) := \{x \in \Sigma^* \mid xy \in L \text{ for some } y \in \Sigma^* \text{ such that } |y| = |x|^2\}$

(b)  $\text{log}(L) := \{x \in \Sigma^* \mid xy \in L \text{ for some } y \in \Sigma^* \text{ such that } |y| = 2^{|x|}\}$

(c)  $\text{flog}(L) := \{x \in \Sigma^* \mid xy \in L \text{ for some } y \in \Sigma^* \text{ such that } |y| = F_{|x|}\}$ , where  $F_n$  is the  $n$ th Fibonacci number.

- \*10. Let  $L \subseteq \Sigma^*$  be an arbitrary regular language. Prove that the following languages are regular.

(a)  $\text{somerep}(L) := \{w \in \Sigma^* \mid w^n \in L \text{ for some } n \geq 0\}$

- (b)  $\text{allreps}(L) := \{w \in \Sigma^* \mid w^n \in L \text{ for every } n \geq 0\}$
- (c)  $\text{manyreps}(L) := \{w \in \Sigma^* \mid w^n \in L \text{ for infinitely many } n \geq 0\}$
- (d)  $\text{fewreps}(L) := \{w \in \Sigma^* \mid w^n \in L \text{ for finitely many } n \geq 0\}$
- (e)  $\text{powers}(L) := \{w \in \Sigma^* \mid w^{2^n} \in L \text{ for some } n \geq 0\}$
- ★(f)  $\text{whatthe}_N(L) := \{w \in \Sigma^* \mid w^n \in L \text{ for some } n \in N\}$ , where  $N$  is an **arbitrary** fixed set of non-negative integers. [Hint: You only have to prove that an accepting NFA exists; you don't have to describe how to construct it.]

[Hint: For each of these languages, there is an accepting NFA with at most  $q^q$  states, where  $q$  is the number of states in some DFA that accepts  $L$ .]

- ★11. For any string  $w \in (\mathbf{0} + \mathbf{1})^*$ , let  $\langle w \rangle_2$  denote the integer represented by  $w$  in binary. For example:

$$\langle \varepsilon \rangle_2 = 0 \quad \langle \mathbf{0010} \rangle_2 = 2 \quad \langle \mathbf{0111} \rangle_2 = 7 \quad \langle \mathbf{1111} \rangle_2 = 15$$

Let  $L$  and  $L'$  be arbitrary regular languages over the alphabet  $\{\mathbf{0}, \mathbf{1}\}$ . Prove that the following language is also regular:

$$\{w \in (\mathbf{0} + \mathbf{1})^* \mid \langle w \rangle_2 = \langle x \rangle_2 + \langle y \rangle_2 \text{ for some strings } x \in L \text{ and } y \in L'\}$$

- ★12. Let  $L \subseteq \Sigma^*$  be an arbitrary regular language. Prove that the following languages are regular.

- (a)  $\text{repsqrt}(L) = \{w \in \Sigma^* \mid w^{|w|} \in L\}$ .
- (b)  $\text{replog}(L) = \{w \in \Sigma^* \mid w^{2^{|w|}} \in L\}$ .
- (c)  $\text{repflog}(L) = \{w \in \Sigma^* \mid w^{F_{|w|}} \in L\}$ , where  $F_n$  is the  $n$ th Fibonacci number.

[Hint: The NFAs for these languages use a **LOT** of states. Let  $M = (\Sigma, Q, s, A, \delta)$  be a DFA that accepts  $L$ . Imagine that you somehow know  $\delta^*(q, w)$  in advance, for every state  $q \in Q$ . Ha, ha, ha! Mine is an evil laugh!]

## 5 Context-Free Languages and Grammars

### 5.1 Definitions

Intuitively, a language is regular if it can be built from individual strings by concatenation, union, and repetition. In this note, we consider a wider class of **context-free** languages, which are languages that can be built from individual strings by concatenation, union, and *recursion*.

Formally, a language is context-free if and only if it has a certain type of recursive description known as a **context-free grammar**, which is a structure with the following components:

- A finite set  $\Sigma$ , whose elements are called **symbols** or **terminals**.
- A finite set  $\Gamma$  disjoint from  $\Sigma$ , whose elements are called **non-terminals**.
- A finite set  $R$  of **production rules** of the form  $A \rightarrow w$ , where  $A \in \Gamma$  is a non-terminal and  $w \in (\Sigma \cup \Gamma)^*$  is a string of symbols and variables.
- A **starting** non-terminal, typically denoted  $S$ .

For example, the following eight production rules describe a context free grammar with terminals  $\Sigma = \{0, 1\}$  and non-terminals  $\Gamma = \{S, A, B, C\}$ :

$$\begin{array}{llll} S \rightarrow A & A \rightarrow 0A & B \rightarrow B1 & C \rightarrow \varepsilon \\ S \rightarrow B & A \rightarrow 0C & B \rightarrow C1 & C \rightarrow 0C1 \end{array}$$

Normally we write grammars more compactly by combining the right sides of all rules for each non-terminal into one list, with alternatives separated by vertical bars.<sup>1</sup> For example, the previous grammar can be written more compactly as follows:

$$\begin{array}{l} S \rightarrow A \mid B \\ A \rightarrow 0A \mid 0C \\ B \rightarrow B1 \mid C1 \\ C \rightarrow \varepsilon \mid 0C1 \end{array}$$

---

<sup>1</sup>Yes, this means we now have *three* symbols  $\cup$ ,  $+$ , and  $|$  with exactly the same meaning. Sigh.

For the rest of this lecture, I will *almost* always use the following notational conventions.

- Monospaced digits (0, 1, 2, ...) and symbols ( $\diamond$ , \$, #, •, ...) are explicit terminals.
- Early lower-case Latin letters (a, b, c, ...) represent unknown or arbitrary terminals in  $\Sigma$ .
- Upper-case Latin letters (A, B, C, ...) and the letter S represent non-terminals in  $\Gamma$ .
- Late lower-case Latin letters (... , w, x, y, z) represent strings in  $(\Sigma \cup \Gamma)^*$ , whose characters could be either terminals or non-terminals.

We can **apply** a production rule to a string in  $(\Sigma \cup \Gamma)^*$  by replacing any instance of the non-terminal on the left of the rule with the string on the right. More formally, for any strings  $x, y, z \in (\Sigma \cup \Gamma)^*$  and any non-terminal  $A \in \Gamma$ , applying the production rule  $A \rightarrow y$  to the string  $xAz$  yields the string  $xyz$ . We use the notation  $xAz \rightsquigarrow xyz$  to describe this application. For example, we can apply the rule  $C \rightarrow 0C1$  to the string  $00C1BAC0$  in two different ways:

$$00\underline{C}1BAC0 \rightsquigarrow 00\underline{0C1}1BAC0 \quad 00C1BAC\underline{0} \rightsquigarrow 00C1BA\underline{0C1}0$$

More generally, for any strings  $x, z \in (\Sigma \cup \Gamma)^*$ , we say that  $z$  **derives from**  $x$ , written  $x \rightsquigarrow^* z$ , if we can transform  $x$  into  $z$  by applying a finite sequence of production rules, or more formally, if either

- $x = z$ , or
- $x \rightsquigarrow y$  and  $y \rightsquigarrow^* z$  for some string  $y \in (\Sigma \cup \Gamma)^*$ .

Straightforward definition-chasing implies that, for any strings  $w, x, y, z \in (\Sigma \cup \Gamma)^*$ , if  $x \rightsquigarrow^* y$ , then  $wxz \rightsquigarrow^* wyz$ .

The **language**  $L(w)$  of any string  $w \in (\Sigma \cup \Gamma)^*$  is the set of all strings in  $\Sigma^*$  that derive from  $w$ :

$$L(w) := \{x \in \Sigma^* \mid w \rightsquigarrow^* x\}.$$

The language **generated by** a context-free grammar  $G$ , denoted  $L(G)$ , is the language of its starting non-terminal. Finally, a language is **context-free** if it is generated by some context-free grammar.

Context-free grammars are sometimes used to model natural languages. In this context, the symbols are *words*, and the strings in the languages are *sentences*. For example, the following grammar describes a simple subset of English sentences. (Here I diverge from the usual notation conventions. Strings in  $\langle \text{angle brackets} \rangle$  are non-terminals, and regular strings are terminals.)

$$\langle \text{sentence} \rangle \rightarrow \langle \text{noun phrase} \rangle \langle \text{verb phrase} \rangle \langle \text{noun phrase} \rangle$$

$$\langle \text{noun phrase} \rangle \rightarrow \langle \text{adjective phrase} \rangle \langle \text{noun} \rangle$$

$$\langle \text{adj. phrase} \rangle \rightarrow \langle \text{article} \rangle \mid \langle \text{possessive} \rangle \mid \langle \text{adjective phrase} \rangle \langle \text{adjective} \rangle$$

$$\langle \text{verb phrase} \rangle \rightarrow \langle \text{verb} \rangle \mid \langle \text{adverb} \rangle \langle \text{verb phrase} \rangle$$

$$\langle \text{noun} \rangle \rightarrow \text{dog} \mid \text{trousers} \mid \text{daughter} \mid \text{nose} \mid \text{homework} \mid \text{time lord} \mid \text{pony} \mid \dots$$

$$\langle \text{article} \rangle \rightarrow \text{the} \mid \text{a} \mid \text{some} \mid \text{every} \mid \text{that} \mid \dots$$

$$\langle \text{possessive} \rangle \rightarrow \langle \text{noun phrase} \rangle \text{'s} \mid \text{my} \mid \text{your} \mid \text{his} \mid \text{her} \mid \dots$$

$$\langle \text{adjective} \rangle \rightarrow \text{friendly} \mid \text{furious} \mid \text{moist} \mid \text{green} \mid \text{severed} \mid \text{timey-wimey} \mid \text{little} \mid \dots$$

$$\langle \text{verb} \rangle \rightarrow \text{ate} \mid \text{found} \mid \text{wrote} \mid \text{killed} \mid \text{mangled} \mid \text{saved} \mid \text{invented} \mid \text{broke} \mid \dots$$

$$\langle \text{adverb} \rangle \rightarrow \text{squarely} \mid \text{incompetently} \mid \text{barely} \mid \text{sort of} \mid \text{awkwardly} \mid \text{totally} \mid \dots$$

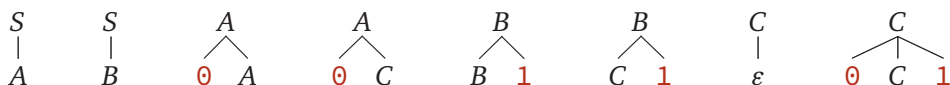


## 5.2 Parse Trees

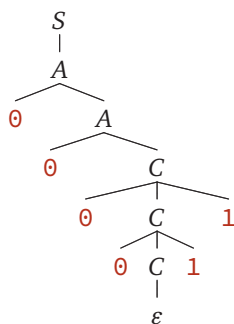
It is often useful to visualize derivations of strings in  $L(G)$  using a **parse tree**. The parse tree for a string  $w \in L(G)$  is a rooted ordered tree where

- Each leaf is labeled with a terminal or the empty string  $\varepsilon$ . Concatenating these in order from left to right yields the string  $w$ .
- Each internal node is labeled with a non-terminal. In particular, the root is labeled with the start non-terminal  $S$ .
- For each internal node  $v$ , there is a production rule  $A \rightarrow \omega$  where  $A$  is the label of  $v$  and the symbols in  $\omega$  are the labels of the children of  $v$  in order from left to right.

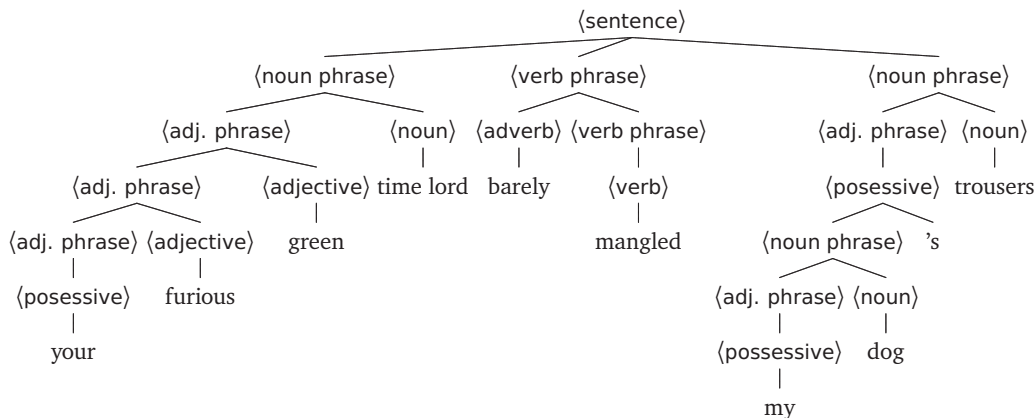
In other words, the production rules of the grammar describe *template trees* that can be assembled into larger parse trees. For example, the simple grammar on the previous page has the following templates, one for each production rule:



The same grammar gives us the following parse tree for the string 000011:

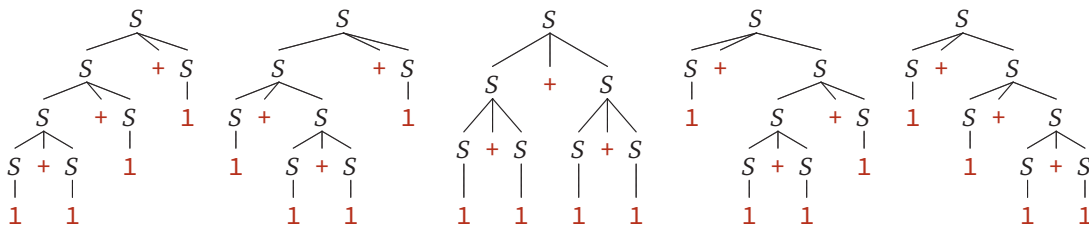


Our more complicated “English” grammar gives us parse trees like the following:



Any parse tree that contains at least one node with more than one non-terminal child corresponds to several different derivations. For example, when deriving an “English” sentence, we have a choice of whether to expand the first ⟨noun phrase⟩ (“your furious green time lord”) before or after the second (“my dog’s trousers”).

A string  $w$  is **ambiguous** with respect to a grammar if there is more than one parse tree for  $w$ , and a grammar  $G$  is **ambiguous** if some string is ambiguous with respect to  $G$ . Neither of the previous example grammars is ambiguous. However, the grammar  $S \rightarrow 1 \mid S+S$  is ambiguous, because the string  $1+1+1+1$  has five different parse trees:



A context-free language  $L$  is *inherently ambiguous* if every context-free grammar that generates  $L$  is ambiguous. The language generated by the previous grammar (the regular language  $(\mathbf{1}+)^*\mathbf{1}$ ) is *not* inherently ambiguous, because the unambiguous grammar  $S \rightarrow \mathbf{1} \mid \mathbf{1}+S$  generates the same language.

### 5.3 From Grammar to Language

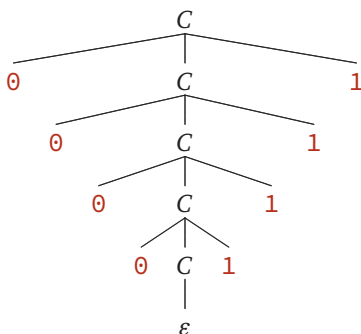
Let's figure out the language generated by our first example grammar

$$S \rightarrow A \mid B \quad A \rightarrow \textcircled{0}A \mid \textcircled{0}C \quad B \rightarrow B\textcolor{red}{1} \mid C\textcolor{red}{1} \quad C \rightarrow \varepsilon \mid \textcircled{0}C\textcolor{red}{1}.$$

Since the production rules for non-terminal  $C$  do not refer to any other non-terminal, let's begin by figuring out  $L(C)$ . After playing around with the smaller grammar  $C \rightarrow \varepsilon \mid \textcircled{0}C\textcircled{1}$  for a few seconds, you can probably guess that its language is  $\{\varepsilon, \textcircled{0}\textcircled{1}, \textcircled{0}\textcircled{0}\textcircled{1}\textcircled{1}, \textcircled{0}\textcircled{0}\textcircled{0}\textcircled{1}\textcircled{1}\textcircled{1}, \dots\}$ , that is, the set all of strings of the form  $\textcircled{0}^n\textcircled{1}^n$  for some integer  $n$ . For example, we can derive the string  $\textcircled{0}\textcircled{0}\textcircled{0}\textcircled{1}\textcircled{1}\textcircled{1}\textcircled{1}$  from the start non-terminal  $S$  using the following derivation:

$$C \rightsquigarrow \bullet C1 \rightsquigarrow \bullet\bullet C11 \rightsquigarrow \bullet\bullet\bullet C111 \rightsquigarrow \bullet\bullet\bullet\bullet C1111 \rightsquigarrow \bullet\bullet\bullet\bullet\epsilon 1111 = \bullet\bullet\bullet\bullet 1111$$

The same derivation can be viewed as the following parse tree:



In fact, it is not hard to *prove* by induction that  $L(C) = \{\mathbf{0}^n \mathbf{1}^n \mid n \geq 0\}$  as follows. As usual when we prove that two sets  $X$  and  $Y$  are equal, the proof has two stages: one stage to prove  $X \subseteq Y$ , the other to prove  $Y \subseteq X$ .

**Lemma 5.1.**  $C \rightsquigarrow^* \mathbf{0}^n \mathbf{1}^n$  for every non-negative integer  $n$ .

**Proof:** Fix an arbitrary non-negative integer  $n$ . Assume that  $C \rightsquigarrow^* 0^k 1^k$  for every non-negative integer  $k < n$ . There are two cases to consider.

- If  $n = 0$ , then  $\mathbf{0}^n \mathbf{1}^n = \varepsilon$ . The rule  $C \rightarrow \varepsilon$  implies that  $C \rightsquigarrow \varepsilon$  and therefore  $C \rightsquigarrow^* \varepsilon$ .
- Suppose  $n > 0$ . The inductive hypothesis implies that  $C \rightsquigarrow^* \mathbf{0}^{n-1} \mathbf{1}^{n-1}$ . Thus, the rule  $C \rightarrow \mathbf{0}C\mathbf{1}$  implies that  $C \rightsquigarrow \mathbf{0}C\mathbf{1} \rightsquigarrow^* \mathbf{0}(\mathbf{0}^{n-1} \mathbf{1}^{n-1})\mathbf{1} = \mathbf{0}^n \mathbf{1}^n$ .

In both cases, we conclude that  $C \rightsquigarrow^* \mathbf{0}^n \mathbf{1}^n$ , as claimed.  $\square$

**Lemma 5.2.** *For every string  $w \in L(C)$ , we have  $w = \mathbf{0}^n \mathbf{1}^n$  for some non-negative integer  $n$ .*

**Proof:** Fix an arbitrary string  $w \in L(C)$ . Assume that for any string  $x \in L(C)$  such that  $|x| < |w|$ , we have  $x = \mathbf{0}^k \mathbf{1}^k$  for some non-negative integer  $k$ . There are two cases to consider, one for each production rule.

- If  $w = \varepsilon$ , then  $w = \mathbf{0}^0 \mathbf{1}^0$ .
- Otherwise,  $w = \mathbf{0}x\mathbf{1}$  for some string  $x \in L(C)$ . Because  $|x| = |w| - 2 < |w|$ , the inductive hypothesis implies that  $x = \mathbf{0}^k \mathbf{1}^k$  for some integer  $k$ . Then we have  $w = \mathbf{0}^{k+1} \mathbf{1}^{k+1}$ .

In both cases, we conclude that  $w = \mathbf{0}^n \mathbf{1}^n$  for some non-negative integer  $n$ , as claimed.  $\square$

The first proof uses induction on strings, following the boilerplate proposed in the very first lecture; in particular, the case analysis mirrors the recursive definition of “string”. The second proof uses *structural induction* on the parse tree of the string  $\mathbf{0}^n \mathbf{1}^n$ ; the case analysis mirrors the recursive definition of the language of  $S$ , as described by the production rules. In both proofs, as in every proof by induction, the inductive hypothesis is “Assume there is no smaller counterexample.”

Similar analysis implies that  $L(A) = \{\mathbf{0}^m \mathbf{1}^n \mid m > n\}$  and  $L(B) = \{\mathbf{0}^m \mathbf{1}^n \mid m < n\}$ , and therefore  $L(S) = \{\mathbf{0}^m \mathbf{1}^n \mid m \neq n\}$ .

### 5.3.1 Careful With Those Epsilons

There is an important subtlety in the proof of Lemma 5.2. The proof is written as induction on the length of the string  $w$ ; unfortunately, this induction pattern does not work for all context-free grammars. Consider the following ambiguous grammar

$$S \rightarrow \varepsilon \mid SS \mid \mathbf{0}S\mathbf{1} \mid \mathbf{1}S\mathbf{0}.$$

A bit of experimentation should convince you that  $L(S)$  is the language of all binary strings with the same number of  $\mathbf{0}$ s and  $\mathbf{1}$ s. We cannot use the string-induction boilerplate for this grammar, because there are arbitrarily long<sup>2</sup> derivations of the form

$$S \rightsquigarrow SS \rightsquigarrow S \rightsquigarrow SS \rightsquigarrow S \rightsquigarrow SS \rightsquigarrow SS \rightsquigarrow \dots \rightsquigarrow w,$$

which alternately apply the productions  $S \rightarrow SS$  and  $S \rightarrow \varepsilon$ . Specifically, even if we knew that our arbitrary string  $w$  can be written as  $xy$  for some strings  $x, y \in L(S)$ , we cannot guarantee that  $|x| < |w|$  and  $|y| < |w|$ , so we cannot apply the standard string-induction hypothesis.

However, we can still argue inductively about this grammar, by considering a *minimum-length* derivation of  $w$ , and basing the case analysis on the first production in this derivation. Here’s an example of this induction boilerplate in action, with the modified boilerplate language highlighted.

<sup>2</sup>but not infinite; derivations are finite *by definition*!

**Lemma 5.3.** For every string  $w \in L(S)$ , we have  $\#(0, w) = \#(1, w)$ .

**Proof:** Let  $w$  be an arbitrary string in  $L(S)$ . Fix an minimum-length derivation of  $w$ .

Assume that for any string  $x \in L(S)$  that is shorter than  $w$ , we have  $x = 0^k 1^k$  for some non-negative integer  $k$ . There are four cases to consider, depending on the first production in our fixed derivation.

- Suppose the first production is  $S \rightarrow \varepsilon$ . Then  $w = \varepsilon$  and therefore  $\#(0, w) = \#(1, w) = 0$  by definition.
- Suppose the first production is  $S \rightarrow SS$ . Then  $w = xy$  for some strings  $x, y \in L(S)$ . Both  $x$  and  $y$  must be non-empty; otherwise, we could shorten our derivation of  $w$ . Thus, both  $x$  and  $y$  are shorter than  $w$ . The inductive hypothesis implies  $\#(0, x) = \#(1, x)$  and  $\#(0, y) = \#(1, y)$ , so  $\#(0, w) = \#(0, x) + \#(0, y) = \#(1, x) + \#(1, y) = \#(1, w)$ .
- Suppose the first production is  $S \rightarrow 0S1$ . Then  $w = 0x1$  for some string  $x \in L(S)$ . The inductive hypothesis implies  $\#(0, x) = \#(1, x)$  so  $\#(0, w) = \#(0, x) + 1 = \#(1, x) + 1 = \#(1, w)$ .
- Finally, suppose the first production is  $S \rightarrow 1S0$ . Then  $w = 1x0$  for some string  $x \in L(S)$ . The inductive hypothesis implies  $\#(0, x) = \#(1, x)$  so  $\#(0, w) = \#(0, x) + 1 = \#(1, w) + 1 = \#(1, w)$ .

In all cases, we conclude that  $\#(0, w) = \#(1, w)$ , as claimed.  $\square$

Another (more traditional) way to handle this issue is to fix an *arbitrary* derivation, and then induct on the length of the derivation, rather than the length of the string itself. The case analysis is still based on the first production in the chosen derivation.

In fact, this subtlety only matters for grammars that either contain a *nullable* non-terminal  $A$  such that  $A \rightsquigarrow^* \varepsilon$  or *equivalent* nonterminals  $A$  and  $B$  such that  $A \rightsquigarrow^* B$  and  $B \rightsquigarrow^* A$ . We describe algorithms to identify these pathologies and remove them from the grammar (without changing its language) in Section 5.9 below.

### 5.3.2 Mutual Induction

Another pitfall in induction proofs for context-free languages is that non-terminals may invoke each other. Consider, for example, the grammar

$$S \rightarrow 0A1 \mid \varepsilon \quad A \rightarrow 1S0 \mid \varepsilon$$

Because each non-terminal appears on the right side of a production rule for the other, we must argue about  $L(S)$  and  $L(A)$  simultaneously.

**Lemma 5.4.**  $L(S) = (01)^*$ .

**Proof:** We actually prove simultaneously that  $L(S) = (01)^*$  and  $L(A) = (10)^*$ .

First, we claim that for any non-negative integer  $n$ , we have  $(01)^n \in L(S)$  and  $(10)^n \in L(A)$ . Let  $n$  be an arbitrary non-negative integer, and assume, for all non-negative integers  $m < n$ , that  $(01)^m \in L(S)$  and  $(10)^m \in L(A)$ . There are two cases to consider.

- If  $n = 0$ , the production rules  $S \rightarrow \varepsilon$  and  $A \rightarrow \varepsilon$  immediately imply  $S \rightsquigarrow \varepsilon = (01)^n$  and  $A \rightsquigarrow \varepsilon = (10)^n$ .

- Suppose  $n > 0$ . We easily observe that  $(01)^n = 0(10)^{n-1}1$ , so the production rule  $S \rightarrow 0A1$  and inductive hypothesis imply  $S \rightsquigarrow 0A1 \rightsquigarrow^* (01)^n$ . Symmetrically,  $(10)^n = 1(01)^{n-1}0$ , so the production rule  $A \rightarrow 1S0$  and the inductive hypothesis implies  $A \rightsquigarrow 1S0 \rightsquigarrow^* (10)^n$ .

Next we claim that for every string  $w \in L(S)$ , we have  $w = (01)^n$  for some non-negative integer  $n$ , and for every string  $w \in L(A)$ , we have  $w = (10)^n$  for some non-negative integer  $n$ . The proof requires two stages.

- Let  $w$  be an arbitrary string in  $L(S)$ , and assume for all  $x \in L(A)$  such that  $|x| < |w|$  that  $x = (10)^n$  for some non-negative integer  $n$ . There are two cases to consider.
  - If  $w = \varepsilon$ , then  $w = (01)^0$ .
  - Suppose  $w = 0x1$  for some string  $x \in L(A)$ . The inductive hypothesis implies  $x = (10)^n$  for some non-negative integer  $n$ . It follows that  $w = 0(10)^n1 = (01)^{n+1}$ .
- Let  $w$  be an arbitrary string in  $L(A)$ , and assume for all  $x \in L(S)$  such that  $|x| < |w|$  that  $x = (01)^n$  for some non-negative integer  $n$ . There are two cases to consider.
  - If  $w = \varepsilon$ , then  $w = (10)^0$ .
  - Suppose  $w = 1x0$  for some string  $x \in L(S)$ . The inductive hypothesis implies  $x = (01)^n$  for some non-negative integer  $n$ . It follows that  $w = 1(01)^n0 = (10)^{n+1}$ .

Together these two claims imply  $L(S) = (01)^*$  and  $L(A) = (10)^*$ , as required.  $\square$

## 5.4 More Examples

Here are some more examples of context-free languages and grammars that generate them, along with brief sketches of correctness proofs.

- Palindromes in  $\{0, 1\}^*$ :

$$S \rightarrow 0S0 \mid 1S1 \mid 0 \mid 1 \mid \varepsilon$$

This grammar is a straightforward translation of the recursive definition of palindrome.

- Strings in  $(0 + 1)^*$  that are *not* palindromes.

$$\begin{aligned} S &\rightarrow 0S0 \mid 1S1 \mid 0Z1 \mid 1Z0 \\ Z &\rightarrow \varepsilon \mid 0Z \mid 1Z \end{aligned}$$

A string  $w$  is a non-palindrome if and only if  $w = x0z1x^R$  or  $w = x1z0x^R$  for some (possibly empty) strings  $x$  and  $y$ .

- Strings in  $\{0, 1\}^*$  with the same number of 0s and 1s:

$$S \rightarrow 0S1 \mid 1S0 \mid SS \mid \varepsilon$$

A non-empty string  $w$  has the same number of 0s and 1s if and only one of the following conditions holds:

- We can write  $w = xy$  for some non-empty strings  $x$  and  $y$  such that  $\#(0, x) = \#(1, x)$  and  $\#(0, y) = \#(1, y)$ .

- $\#(0, x) > \#(1, x)$  for every non-empty proper prefix  $x$  of  $w$ . In this case,  $w = 0z1$  for some string  $z$  with  $\#(0, z) = \#(1, z)$ .
- $\#(0, x) < \#(1, x)$  for every non-empty proper prefix  $x$  of  $w$ . In this case,  $w = 1z0$  for some string  $z$  with  $\#(0, z) = \#(1, z)$ .

- Strings in  $\{0, 1\}^*$  with the same number of 0s and 1s, again:

$$S \rightarrow 0S1S \mid 1S0S \mid \varepsilon$$

Let  $w$  be any non-empty string such that  $\#(0, w) = \#(1, w)$ , let  $x$  be the shortest non-empty prefix of  $w$  such that  $\#(0, x) = \#(1, x)$ , and let  $y$  be the complementary suffix of  $w$ , so  $w = xy$ . It is not hard to prove that  $x$  begins and ends with different symbols, so either  $w = 0z1y$  or  $w = 1z0y$ , where  $\#(0, y) = \#(1, y)$  and  $\#(0, z) = \#(1, z)$ .

- Strings in  $\{0, 1\}^*$  in which the number of 0s is greater than or equal to the number of 1s:

$$S \rightarrow 0S1 \mid 0S \mid 1S0 \mid S0 \mid SS \mid \varepsilon \qquad S \rightarrow 0S1S \mid 0SS \mid 1S0S \mid S0S \mid \varepsilon$$

We have to different grammars, each constructed from a grammar for strings with equal 0s and 1s by either dropping the 1 or keeping the 1 from the right side of each production rule containing a 1. For example, we split the production rule  $S \rightarrow 0S1$  in the first grammar into two production rules  $S \rightarrow 0S1$  and  $S \rightarrow 0S$ .

If we add the trivial production  $S \rightarrow 0$  to the first grammar, we can remove two redundant productions to get the simpler grammar

$$S \rightarrow 0S1 \mid 1S0 \mid SS \mid 0 \mid \varepsilon$$

- Strings in  $\{0, 1\}^*$  with *different* numbers of 0s and 1s:

$$\begin{array}{ll} S \rightarrow O \mid I & \text{(different)} \\ O \rightarrow E0O \mid E0E & \text{(more 0s)} \\ I \rightarrow E1I \mid E1E & \text{(more 1s)} \\ E \rightarrow 0E1E \mid 1E0E \mid \varepsilon & \text{(equal)} \end{array}$$

We can argue correctness by considering each non-terminal in turn, in reverse order.

- $E$  generates all strings with the same number of 0s and 1s, as in the previous example.
- $I$  generates all strings with more 1s than 0s. Any such string can be decomposed into its longest prefix with the same number of 0s and 1s ( $E$ ), followed by a 0, followed by a suffix with at least as many 0s as 1s ( $I$  or  $E$ ).
- Symmetrically,  $O$  generates all strings with more 0s than 1s.
- Finally,  $S$  generates all strings with different numbers of 0s and 1s. Any such string either has more 0s ( $O$ ) or more 1s ( $I$ )

- Balanced strings of parentheses:

$$S \rightarrow (S) \mid SS \mid \varepsilon \qquad \text{or} \qquad S \rightarrow (S)S \mid \varepsilon$$

Here we have two grammars for the same language. The first one uses simpler productions, and is a bit closer to the natural recursive definition. However, the first grammar is

ambiguous — consider the string  $()()()$  — while the second grammar is not. The second grammar decomposes any balanced string of parentheses into its shortest non-empty balanced prefix, which must start with  $($  and end with  $)$ , and the remaining suffix, which must be balanced.

- Unbalanced strings of parentheses—the complement of the previous language:

$$\begin{array}{ll}
 S \rightarrow L \mid RX & \text{(unbalanced)} \\
 L \rightarrow E(L \mid E(E & \text{(more left parens)} \\
 R \rightarrow E)R \mid E)E & \text{(more right parens)} \\
 E \rightarrow \varepsilon \mid (E)E \mid )E(E & \text{(equal left and right)} \\
 X \rightarrow \varepsilon \mid (X \mid )X & \text{(anything)}
 \end{array}$$

A string  $w$  of parens is balanced if and only if both (a)  $w$  has the same number of left and right parens and (b) no prefix of  $w$  has more right parens than left parens. (Proving this fact is a good homework exercise.) Thus, a string  $w$  of parens is *unbalanced* if and only if *either*  $w$  has more left parens than right parens *or* some prefix of  $w$  has more right parens than left parens.

- Arithmetic expressions, possibly with redundant parentheses, over the variables  $X$  and  $Y$ :

$$\begin{array}{ll}
 E \rightarrow E+T \mid T & \text{(expressions)} \\
 T \rightarrow T \times F \mid F & \text{(terms)} \\
 F \rightarrow (E) \mid X \mid Y & \text{(factors)}
 \end{array}$$

Every  $E$ expression is a sum of  $T$ terms, every  $T$ erm is a product of  $F$ actors, and every  $F$ actor is either a variable or a parenthesized  $E$ expression.

- Regular expressions over the alphabet  $\{0, 1\}$  *without* redundant parentheses

$$\begin{array}{ll}
 S \rightarrow T \mid T+S & \text{(Regular expressions)} \\
 T \rightarrow F \mid FT & \text{(Terms = summable expressions)} \\
 F \rightarrow \emptyset \mid W \mid (T+S) \mid X^* \mid (Y)^* & \text{(Factors = concatenable expressions)} \\
 X \rightarrow \emptyset \mid \varepsilon \mid 0 \mid 1 & \text{(Directly starrable expressions)} \\
 Y \rightarrow T+S \mid F \bullet T \mid X^* \mid (Y)^* \mid ZZ & \text{(Starrable expressions needing parens)} \\
 W \rightarrow \varepsilon \mid Z & \text{(Words = strings)} \\
 Z \rightarrow 0 \mid 1 \mid ZZ & \text{(Non-empty strings)}
 \end{array}$$

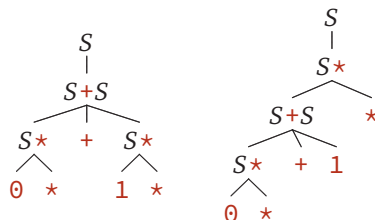
Every regular expression is a sum of terms; every term is a concatenation of factors. Every factor is either the empty-set symbol, a string, a nontrivial sum of terms in parens, or a starred expression. The expressions  $\emptyset^*$ ,  $\varepsilon^*$ ,  $0^*$ , and  $1^*$  require no parentheses; otherwise, the starred subexpression is either a nontrivial sum of terms, a nontrivial concatenation of factors, a starred expression, or a string of length 2 or more.

The “epsilon” symbol  $\varepsilon$  in the production rules for  $W$  and  $Z$  does *not* represent the empty string *per se*, but rather an actual symbol that might appear in a regular expression. The empty string is not a regular expression, but the one-symbol string  $\varepsilon$  is a regular expression that represents the set containing only the empty string!

The final grammar illustrates an important subtlety for certain applications of context-free grammars. This grammar is considerably more complicated than one might initially expect from the definition of regular languages. It's tempting to suggest a much simpler grammar like

$$S \rightarrow \emptyset \mid \varepsilon \mid 0 \mid 1 \mid S+S \mid SS \mid S^* \mid (S)$$

**but this is incorrect!** This grammar does correctly generate all regular expressions *as raw strings*, but it allows parse trees that do not respect the *meaning* of the regular expression. For example, this “simpler” grammar can parse the regular expression string  $0^{*}+1^{*}$  in two different ways:



The first tree correctly parses the regular expression string  $1^{*}+0^{*}$  as the expression  $(1^{*}) + (0^{*})$  but without the redundant parentheses. The second tree incorrectly parses the same string as  $(1^{*} + 0)^{*}$ , which describes a *very* different regular language!

## 5.5 Regular Languages are Context-Free

The following inductive argument proves that every regular language is also a context-free language. Let  $L$  be an arbitrary regular language, encoded by some regular expression  $R$ . Assume that any regular expression simpler than  $R$  represents a context-free language. (“Assume no smaller counterexample.”) We construct a context-free grammar for  $L$  as follows. There are several cases to consider.

- Suppose  $L$  is empty. Then  $L$  is generated by the trivial grammar  $S \rightarrow S$ .
- Suppose  $L = \{w\}$  for some string  $w \in \Sigma^{*}$ . Then  $L$  is generated by the grammar  $S \rightarrow w$ .
- Suppose  $L$  is the union of some regular languages  $L_1$  and  $L_2$ . The inductive hypothesis implies that  $L_1$  and  $L_2$  are context-free. Let  $G_1$  be a context-free language for  $L_1$  with starting non-terminal  $S_1$ , and let  $G_2$  be a context-free language for  $L_2$  with starting non-terminal  $S_2$ , where the non-terminal sets in  $G_1$  and  $G_2$  are disjoint. Then  $L = L_1 \cup L_2$  is generated by the production rule  $S \rightarrow S_1 \mid S_2$ .
- Suppose  $L$  is the concatenation of some regular languages  $L_1$  and  $L_2$ . The inductive hypothesis implies that  $L_1$  and  $L_2$  are context-free. Let  $G_1$  be a context-free language for  $L_1$  with starting non-terminal  $S_1$ , and let  $G_2$  be a context-free language for  $L_2$  with starting non-terminal  $S_2$ , where the non-terminal sets in  $G_1$  and  $G_2$  are disjoint. Then  $L = L_1 L_2$  is generated by the production rule  $S \rightarrow S_1 S_2$ .
- Suppose  $L$  is the Kleene closure of some regular language  $L_1$ . The inductive hypothesis implies that  $L_1$  is context-free. Let  $G_1$  be a context-free language for  $L_1$  with starting non-terminal  $S_1$ . Then  $L = L_1^{*}$  is generated by the production rule  $S \rightarrow \varepsilon \mid S_1 S$ .

In every case, we have found a context-free grammar that generates  $L$ , which means  $L$  is context-free.

In the previous lecture note, we proved that the context-free language  $\{0^n 1^n \mid n \geq 0\}$  is not regular. (In fact, this is the *canonical example* of a non-regular language.) Thus, context-free grammars are strictly more expressive than regular expressions.



## 5.6 Not Every Language is Context-Free

Again, you may be tempted to conjecture that *every* language is context-free, but a variant of our earlier cardinality argument implies that this is not the case.

Any context-free grammar over the alphabet  $\Sigma$  can be encoded as a string over the alphabet  $\Sigma \cup \Gamma \cup \{\epsilon, \rightarrow, |, \$\}$ , where  $\epsilon$  indicates the end of the production rules for each non-terminal. For example, our example grammar

$$S \rightarrow A | B \qquad A \rightarrow 0A | 0C \qquad B \rightarrow B1 | C1 \qquad C \rightarrow \epsilon | 0C1$$

can be encoded as the string

$$S \rightarrow A | B \$ A \rightarrow 0A | 0C \$ B \rightarrow B1 | C1 \$ C \rightarrow \epsilon | 0C1 \$$$

We can further encode any such string as a *binary* string by associating each symbol in the set  $\Sigma \cup \Gamma \cup \{\epsilon, \rightarrow, |, \$\}$  with a different binary substring. Specifically, if we encode each of the grammar symbols  $\epsilon, \rightarrow, |, \$$  as a string of the form  $11^*0$ , each terminal in  $\Sigma$  as a string of the form  $011^*0$ , and each non-terminal as a string of the form  $0011^*0$ , we can unambiguously recover the grammar from the encoding. For example, applying the code

$$\begin{array}{lll} \epsilon \mapsto 10 & 0 \mapsto 010 & S \mapsto 0010 \\ \rightarrow \mapsto 110 & 1 \mapsto 0110 & A \mapsto 00110 \\ | \mapsto 1110 & & B \mapsto 001110 \\ \$ \mapsto 11110 & & C \mapsto 0011110 \end{array}$$

transforms our example grammar into the 136-bit string

$$\begin{array}{l} 00101100011011100011101111000110 \\ 11001000110111001000111101111000 \\ 11101100011100110111000111100110 \\ 11110001111011010111001000111100 \\ 11011110. \end{array}$$

Adding a **1** to the start of this bit string gives us the binary encoding of the integer

$$102\,231\,235\,533\,163\,527\,515\,344\,124\,802\,467\,059\,875\,038.$$

Our construction guarantees that two different context-free grammars over the same alphabet (ignoring changing the names of the non-terminals) yield different positive integers. Thus, the set of context-free grammars over any alphabet is *at most* as large as the set of integers, and is therefore countably infinite. (Most integers are not encodings of context-free grammars, but that only helps us.) It follows that the set of all context-free *languages* over any fixed alphabet is also countably infinite. But we already showed that the set of *all* languages over any alphabet is uncountably infinite. So almost all languages are non-context-free!

There are techniques for proving that specific languages are not context-free, just as there are for proving certain languages are not regular; unfortunately, they are beyond the scope of this course. In particular, the  $\{0^n 1^n 0^n \mid n \geq 0\}$  is not context-free. (In fact, this is the *canonical example* of a non-context-free language.)

## \*5.7 Recursive Automata

All the flavors of finite-state automata we have seen so far describe/encode/accept/compute *regular* languages; these are precisely the languages that can be constructed from individual strings by union, concatenation, and unbounded repetition. Just as context-free grammars are recursive generalizations of regular expressions, we can define a class of machines called *recursive automata*, which generalize (nondeterministic) finite-state automata. Recursive automata were introduced by Walter Woods in 1970 for natural language parsing; Wodds' terminology *recursive transition networks* is more common among computational linguists.

Formally, a **recursive automaton** consists of the following components:

- A non-empty finite set  $\Sigma$ , called the **input alphabet**
- Another non-empty finite set  $N$ , disjoint from  $\Sigma$ , whose elements are called **module names**
- A **start name**  $S \in N$
- A set  $M = \{M_A \mid A \in N\}$  of NFAs, called **modules**, over the alphabet  $\Sigma \cup N$ . Each module  $M_A$  has the following components:
  - A finite set  $Q_A$  of **states**, such that  $Q_A \cap Q_B = \emptyset$  for all  $A \neq B$
  - A **start** state  $s_A \in Q_A$
  - A unique **terminal** or **accepting** state  $t_A \in Q_A$
  - A nondeterministic **transition function**  $\delta_A: Q_A \times (\Sigma \cup \{\varepsilon\} \cup N) \rightarrow 2^{Q_A}$ .

Equivalently, we have a single global transition function  $\delta: Q \times (\Sigma \cup \{\varepsilon\} \cup N) \rightarrow 2^Q$ , where  $Q = \bigcup_{A \in N} Q_A$ , such that for any name  $A$  and any state  $q \in Q_A$  we have  $\delta(q) \subseteq Q_A$ . Machine  $M_S$  is called the **main module**.

A **configuration** of a recursive automaton is a triple  $(w, q, s)$ , where  $w$  is a string in  $\Sigma^*$  called the **input**,  $q$  is a state in  $Q$  called the **local state**, and  $s$  is a string in  $Q^*$  called the **stack**. The module containing the local state  $q$  is called the **active module**. A configuration can be changed by three types of transitions.

- A **read** transition consumes the first symbol in the input and changes the local state within the active module, just like a standard NFA.
- An **epsilon** transition changes the local state within the active module, without consuming any input characters, just like a standard NFA.
- A **call** transition chooses a module name  $A$ , pushes some state in  $\delta(q, A)$  onto the stack, and then changes the local state to  $s_A$  (thereby changing the active module to  $M_A$ ), without consuming any input characters.
- Finally, if the current state is the terminal state of the active module *and* the stack is non-empty, a **return** transition pops the top state off the stack and makes it the new local state (thereby possibly changing the active module), without consuming any input characters.

Symbolically, we can describe these transitions as follows:

<b>read:</b>	$(ax, q, \sigma) \mapsto (x, q', \sigma)$	for some $q' \in \delta(q, a)$
<b>epsilon:</b>	$(w, q, \sigma) \mapsto (w, q', \sigma)$	for some $q' \in \delta(q, \varepsilon)$
<b>call:</b>	$(w, q, \sigma) \mapsto (w, s_A, q' \cdot \sigma)$	for some $A \in N$ and some $q' \in \delta(q, A)$
<b>return:</b>	$(w, t_A, q \cdot \sigma) \mapsto (w, q, \sigma)$	

A recursive automaton **accepts** a string  $w$  if there is a *finite* sequence of transitions starting at the start configuration  $(w, s_S, \varepsilon)$  and ending at the terminal configuration  $(\varepsilon, t_S, \varepsilon)$ .

**Reformulate recursive automata using recursion-as-magic**, analogously to the non-determinism-as-magic in string NFAs or Han and Wood’s expression automata. A recursive automaton module  $M_A$  accepts a string  $w$  if and only if there is a finite sequence of transitions

$$s_A = q_0 \xrightarrow{\alpha_1} q_1 \xrightarrow{\alpha_2} q_2 \xrightarrow{\alpha_3} \dots \xrightarrow{\alpha_\ell} q_\ell = t_A$$

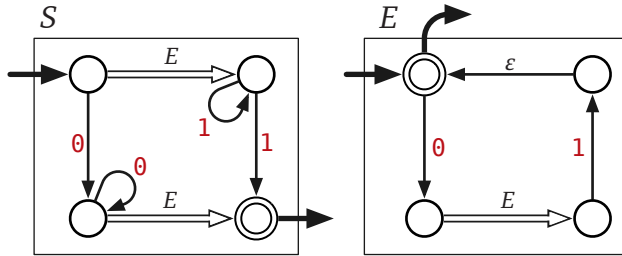
among states in  $Q_A$ , and a decomposition of  $w$  into substrings  $x_1 \cdot x_2 \cdot \dots \cdot x_\ell$ , where one of the following conditions holds for each index  $i$ :

- $\alpha_i = \varepsilon$  and  $x_i = \varepsilon$
- $\alpha_i \in \Sigma$  and  $x_i = \alpha_i$
- $\alpha_i \in N$  and module  $M_{\alpha_i}$  accepts  $x_i$ .

This model is then easily extended to more general transition labels:

- Recursive *string*-automata allow transitions to be labeled by either strings in  $\Sigma^*$  or module names, and the first two bullets in the previous list become “ $\alpha_i = x_i$ ”.
- Recursive *expression*-automata allow transitions to be labeled by either regular expressions over  $\Sigma$  or module names, and the first two bullets in the previous list become “ $\alpha_i$  is a regular expression and  $x_i$  matches  $\alpha_i$ ”.
- We could even consider *recursive-expression* automata, which allow transitions to be labeled by arbitrary regular expressions over  $\Sigma \cup N$ .
- We could even even consider recursive-*grammar* automata, which allow transitions to be labeled by arbitrary *context-free grammars* over  $\Sigma \cup N$ . WEEEE/EEEEW!

For example, the following recursive automaton accepts the language  $\{0^m 1^n \mid m \neq n\}$ . The recursive automaton has two component modules; the start machine named  $S$  and a “subroutine” named  $E$  (for “equal”) that accepts the language  $\{0^n 1^n \mid n \geq 0\}$ . White arrows indicate recursive transitions. The large arrow into each module indicates that module’s start state; the large arrow leading out of each module indicates that module’s terminal state.



A recursive automaton for the language  $\{0^m 1^n \mid m \neq n\}$

**Lemma 5.5.** *Every context-free language is accepted by a recursive automaton.*

**Proof:**

Direct construction from the CFG, with one module per nonterminal.

For example, the context-free grammar

$$S \rightarrow 0A \mid B1$$

$$A \rightarrow 0A \mid E$$

$$B \rightarrow B1 \mid E$$

$$E \rightarrow \varepsilon \mid 0E0$$

leads to the following recursive automaton with four modules:

Figure!

**Lemma 5.6.** *Every recursive automaton accepts a context-free language.*

**Proof (sketch):** Let  $R = (\Sigma, N, S, \delta, M)$  be an arbitrary recursive automaton. We define a context-free grammar  $G$  that describes the language accepted by  $R$  as follows.

The set of nonterminals in  $G$  is isomorphic to the state set  $Q$ ; that is, for each state  $q \in Q$ , the grammar contains a corresponding nonterminal  $[q]$ . The language of  $[q]$  will be the set of strings  $w$  such that there is a finite sequence of transitions starting at the start configuration  $(w, q, \varepsilon)$  and ending at the terminal configuration  $(\varepsilon, t, \varepsilon)$ , where  $t$  is the terminal state of the module containing  $q$ .

The grammar has four types of production rules, corresponding to the four types of transitions:

- **read:** For each symbol  $a$  and each pair of states  $p$  and  $q$  such that  $p \in \delta(q, a)$ , the grammar contains the production rule  $[q] \rightarrow a[p]$ .
- **epsilon:** For any two states  $p$  and  $q$  such that  $p \in \delta(q, \varepsilon)$ , the grammar contains the production rule  $[q] \rightarrow [p]$ .
- **call:** Each name  $A$  and each pair of states  $p$  and  $q$  such that  $p \in \delta(q, A)$ , the grammar contains the production rule  $[q] \rightarrow [s_A][p]$ .
- **return:** Each name  $A$ , the grammar contains the production rule  $[t_A] \rightarrow \varepsilon$ .

Finally, the starting nonterminal of  $G$  is  $[s_S]$ , which corresponds to the start state of the main module.

We can now argue inductively that the grammar  $G$  and the recursive automaton  $R$  describe the same language. Specifically, any sequence of transitions in  $R$  from  $(w, s_S, \varepsilon)$  to  $(\varepsilon, t_S, \varepsilon)$  can be transformed mechanically into a derivation of  $w$  from the nonterminal  $[s_S]$  in  $G$ . Symmetrically, the **leftmost** derivation of any string  $w$  in  $G$  can be mechanically transformed into an accepting sequence of transitions in  $R$ . We omit the straightforward but tedious details.  $\square$

For example, the recursive automaton on the previous page gives us the following context-free grammar. To make the grammar more readable, I've renamed the nonterminals corresponding to start and terminal states:  $S = [s_S]$ ,  $T = [t_S]$ , and  $E = [s_E] = [t_E]$ :

$$S \rightarrow EA \mid 0B \quad E \rightarrow \varepsilon \mid 0X$$

$$A \rightarrow 1A \mid 1T \quad X \rightarrow EY$$

$$B \rightarrow 0B \mid ET \quad Y \rightarrow 1Z$$

$$T \rightarrow \varepsilon \quad Z \rightarrow E$$

Our earlier proofs imply that we can forbid  $\varepsilon$ -transitions or even allow regular-expression transitions in our recursive automata without changing the set of languages they accept.

## 5.8 Chomsky Normal Form

For many algorithmic problems involving context-free grammars, it is helpful to consider grammars with a particular special structure called **Chomsky normal form**, abbreviated **CNF**:

- The starting non-terminal  $S$  does not appear on the right side of any production rule.
- The starting non-terminal  $S$  *may* have the production rule  $S \rightarrow \varepsilon$ .
- The right side of every other production rule is either a single terminal symbol or a string of exactly two non-terminals—that is, every other production rule has the form  $A \rightarrow BC$  or  $A \rightarrow a$ .

A particularly attractive feature of CNF grammars is that they yield *full binary* parse trees; in particular, every parse tree for a string of length  $n > 0$  has exactly  $2n - 1$  non-terminal nodes. Consequently, any string of length  $n$  in the language of a CNF grammar can be derived in exactly  $2n - 1$  production steps. It follows that we can actually determine whether a string belongs to the language of a CNF grammar by brute-force consideration of all possible derivations of the appropriate length.

For arbitrary context-free grammars, there is no similar upper bound on the length of a derivation, and therefore no similar brute-force membership algorithm, because the grammar may contain additional  **$\varepsilon$ -productions** of the form  $A \rightarrow \varepsilon$  and/or **unit productions** of the form  $A \rightarrow B$ , where both  $A$  and  $B$  are non-terminals. Unit productions introduce nodes of degree 1 into any parse tree, and  $\varepsilon$ -productions introduce leaves that do not contribute to the word being parsed.

Fortunately, it is possible to determine membership in the language of an arbitrary context-free grammar, thanks to the following theorem. Two context-free grammars are **equivalent** if they define the same language.

**Every context-free grammar is equivalent to a grammar in Chomsky normal form.**

Moreover, there are algorithms to automatically convert any context-free grammar into Chomsky normal form. Unfortunately, these conversion algorithms are quite complex, but for most applications of context-free grammars, the details of the conversion are unimportant—it's enough to know that the algorithms exist. For the sake of completeness, however, I will describe one such conversion algorithm in the next section.

## \*5.9 CNF Conversion Algorithm

I'll actually prove a stronger statement: Not only can we convert any context-free grammar into Chomsky normal form, but we can do so *quickly*. We analyze the running time of our conversion algorithm in terms of the **total length** of the input grammar, which is just the number of symbols needed to write down the grammar. Up to constant factors, the total length is the sum of the lengths of the production rules.

**Theorem 5.7.** *Given an arbitrary context-free grammar with total length  $L$ , we can compute an equivalent grammar in Chomsky normal form with total length  $O(L^2)$  in  $O(L^2)$  time.*

Our algorithm consists of several relatively straightforward stages. Efficient implementation of some of these stages requires standard graph-traversal algorithms, which are described in a different part of the course.

**o. Add a new starting non-terminal.** Add a new non-terminal  $S'$  and a production rule  $S' \rightarrow S$ , where  $S$  is the starting non-terminal for the given grammar.  $S'$  will be the starting non-terminal for the resulting CNF grammar. (In fact, this step is necessary only when  $S \rightsquigarrow^* \varepsilon$ , but at this point in the conversion process, we don't yet know whether that's true.)

**1. Decompose long production rules.** For each production rule  $A \rightarrow \omega$  whose right side  $w$  has length greater than two, add new production rules of length two that still permit the derivation  $A \rightsquigarrow^* \omega$ . Specifically, suppose  $\omega = \alpha\chi$  for some symbol  $\alpha \in \Sigma \cup \Gamma$  and string  $\chi \in (\Sigma \cup \Gamma)^*$ . The algorithm replaces  $A \rightarrow \omega$  with two new production rules  $A \rightarrow \alpha B$  and  $B \rightarrow \chi$ , where  $B$  is a new non-terminal, and then (if necessary) recursively decomposes the production rule  $B \rightarrow \chi$ . For example, we would replace the long production rule  $A \rightarrow \textcolor{red}{0}BC\textcolor{red}{1}CB$  with the following sequence of short production rules, where each  $A_i$  is a new non-terminal:

$$A \rightarrow \textcolor{red}{0}A_1 \quad A_1 \rightarrow BA_2 \quad A_2 \rightarrow CA_3 \quad A_3 \rightarrow \textcolor{red}{1}A_4 \quad A_4 \rightarrow CB$$

This stage can significantly increase the number of non-terminals and production rules, but it increases the *total length* of all production rules by at most a small constant factor.<sup>3</sup> Moreover, for the remainder of the conversion algorithm, every production rule has length at most two. The running time of this stage is  $O(L)$ .

**2. Identify nullable non-terminals.** A non-terminal  $A$  is **nullable** if and only if  $A \rightsquigarrow^* \varepsilon$ . The recursive definition of  $\rightsquigarrow^*$  implies that  $A$  is nullable if and only if the grammar contains a production rule  $A \rightarrow \omega$  where  $\omega$  consists entirely of nullable non-terminals (in particular, if  $\omega = \varepsilon$ ). You may be tempted to transform this recursive characterization directly into a recursive algorithm, but this is a bad idea; the resulting algorithm would fall into an infinite loop if (for example) the same non-terminal appeared on both sides of the same production rule. Instead, we apply the following **fixed-point** algorithm, which repeatedly scans through the entire grammar until a complete scan discovers no new nullable non-terminals.

```

NULLABLES( $\Sigma, \Gamma, R, S$ ):
   $\Gamma_\varepsilon \leftarrow \emptyset$       ⟨⟨known nullable non-terminals⟩⟩
  done  $\leftarrow$  FALSE
  while  $\neg$ done
    done  $\leftarrow$  TRUE
    for each non-terminal  $A \in \Gamma \setminus \Gamma_\varepsilon$ 
      for each production rule  $A \rightarrow \omega$ 
        if  $\omega \in \Gamma_\varepsilon^*$ 
          add  $A$  to  $\Gamma_\varepsilon$ 
          done  $\leftarrow$  FALSE
  return  $\Gamma_\varepsilon$ 

```

At this point in the conversion algorithm, if  $S'$  is **not** identified as nullable, then we can safely remove it from the grammar and use the original starting nonterminal  $S$  instead.

As written, NULLABLES runs in  $O(nL) = O(L^2)$  time, where  $n$  is the number of non-terminals in  $\Gamma$ . Each iteration of the main loop except the last adds at least one non-terminal to  $\Gamma_\varepsilon$ , so the

---

<sup>3</sup>In most textbook descriptions of the CFG conversion algorithm, this stage is performed *last*, after removing  $\varepsilon$ -productions and unit productions. But with the stages in that traditional order, removing  $\varepsilon$ -productions could *exponentially* increase the length of the grammar in the worst case! Consider the production rule  $A \rightarrow (BC)^k$ , where  $B$  is nullable but  $C$  is not. Decomposing this rule first and then removing  $\varepsilon$ -productions introduces about  $3k$  new rules; whereas, removing  $\varepsilon$ -productions first introduces  $2^k$  new rules, most of which then must then be further decomposed!

algorithm halts after at most  $n + 1 \leq L$  iterations, and in each iteration, we examine at most  $L$  production rules. There is a faster implementation of NULLABLES that runs in  $O(n + L) = O(L)$  time,<sup>4</sup> but since other parts of the conversion algorithm already require  $O(L^2)$  time, we needn't bother.

**3. Eliminate  $\varepsilon$ -productions.** First, remove every production rule of the form  $A \rightarrow \varepsilon$ . Then for each production rule  $A \rightarrow w$ , add all possible new production rules of the form  $A \rightarrow w'$ , where  $w'$  is a **non-empty** string obtained from  $w$  by removing one nullable non-terminal. For example, if the grammar contained the production rule  $A \rightarrow BC$ , where  $B$  and  $C$  are both nullable, we would add two new production rules  $A \rightarrow B \mid C$ . Finally, if the starting nonterminal  $S'$  was identified as nullable in the previous stage, add the production rule  $S' \rightarrow \varepsilon$ ; this will be the *only*  $\varepsilon$ -production in the final grammar. This phase of the conversion runs in  $O(L)$  time and at most triples the number of production rules.

**4. Merge equivalent non-terminals.** We say that two non-terminals  $A$  and  $B$  are **equivalent** if they can be derived from each other:  $A \rightsquigarrow^* B$  and  $B \rightsquigarrow^* A$ . Because we have already removed  $\varepsilon$ -productions, any such derivation must consist entirely of unit productions. For example, in the grammar

$$S \rightarrow B \mid C, \quad A \rightarrow B \mid D \mid CC \mid \textcolor{red}{0}, \quad B \rightarrow C \mid AD \mid \textcolor{red}{1}, \quad C \rightarrow A \mid DA, \quad D \rightarrow BA \mid CS,$$

non-terminals  $A, B, C$  are all equivalent, but  $S$  is not in that equivalence class (because we cannot derive  $S$  from  $A$ ) and neither is  $D$  (because we cannot derive  $A$  from  $D$ ).

Construct a directed graph  $G$  whose vertices are the non-terminals and whose edges correspond to unit productions, in  $O(L)$  time. Then two non-terminals are equivalent if and only if they are in the same strong component of  $G$ . Compute the strong components of  $G$  in  $O(L)$  time using, for example, the algorithm of Kosaraju and Sharir. Then merge all the non-terminals in each equivalence class into a single non-terminal. Finally, remove any unit productions of the form  $A \rightarrow A$ . The total running time for this phase is  $O(L)$ . Starting with our example grammar above, merging  $B$  and  $C$  with  $A$  and removing the production  $A \rightarrow A$  gives us the simpler grammar

$$S \rightarrow A, \quad A \rightarrow AA \mid D \mid DA \mid \textcolor{red}{0} \mid \textcolor{red}{1}, \quad D \rightarrow AA \mid AS.$$

We could further simplify the grammar by merging all non-terminals reachable from  $S$  using only unit productions (in this case, merging non-terminals  $S$  and  $S$ ), but this further simplification is unnecessary.

**5. Remove unit productions.** Once again, we construct a directed graph  $G$  whose vertices are the non-terminals and whose edges correspond to unit productions, in  $O(L)$  time. Because no two non-terminals are equivalent,  $G$  is acyclic. Thus, using topological sort, we can index the non-terminals  $A_1, A_2, \dots, A_n$  such that for every unit production  $A_i \rightarrow A_j$  we have  $i < j$ , again in  $O(L)$  time; moreover, we can assume that the starting non-terminal is  $A_1$ . (In fact, both the dag  $G$  and the linear ordering of non-terminals was already computed in the previous phase!!)

Then for each index  $j$  in decreasing order, for each unit production  $A_i \rightarrow A_j$  and each production  $A_j \rightarrow \omega$ , we add a new production rule  $A_i \rightarrow \omega$ . At this point, all unit productions are

---

<sup>4</sup>Consider the bipartite graph whose vertices correspond to non-terminals and the right sides of production rules, with one edge per rule. The faster algorithm is a modified breadth-first search of this graph, starting at the vertex representing  $\varepsilon$ .



redundant and can be removed. Applying this algorithm to our example grammar above gives us the grammar

$$S \rightarrow AA \mid AS \mid DA \mid \textcolor{red}{0} \mid \textcolor{red}{1}, \quad A \rightarrow AA \mid AS \mid DA \mid \textcolor{red}{0} \mid \textcolor{red}{1}, \quad D \rightarrow AA \mid AS.$$

In the worst case, each production rule for  $A_n$  is copied to each of the other  $n - 1$  non-terminals. Thus, this phase runs in  $\Theta(nL) = O(L^2)$  time and increases the length of the grammar to  $\Theta(nL) = O(L^2)$  in the worst case.

**This phase dominates the running time of the CNF conversion algorithm.** Unlike previous phases, no faster algorithm for removing unit transformations is known! There are grammars of length  $L$  with unit productions such that any equivalent grammar without unit productions has length  $\Omega(L^{1.499999})$  (for any desired number of 9s), but this lower bound does not rule out the possibility of an algorithm that runs in only  $O(L^{3/2})$  time. Closing the gap between  $\Omega(L^{3/2-\epsilon})$  and  $O(L^2)$  has been an open problem since the early 1980s!

**6. Protect terminals.** Finally, for each terminal  $a \in \Sigma$ , we introduce a new non-terminal  $A_a$  and a new production rule  $A_a \rightarrow a$ , and then replace  $a$  with  $A_a$  in every production rule of length two.

**This completes the conversion to Chomsky normal form!** As claimed, the total running time of the algorithm is  $O(L^2)$ , and the total length of the output grammar is also  $O(L^2)$ .

To see the conversion algorithm in action, let's apply these stages one at a time to our very first example grammar for the language  $\{0^m 1^n \mid m \neq n\}$ :

$$S \rightarrow A \mid B \quad A \rightarrow \textcolor{red}{0}A \mid \textcolor{red}{0}C \quad B \rightarrow B\textcolor{red}{1} \mid C\textcolor{red}{1} \quad C \rightarrow \varepsilon \mid \textcolor{red}{0}C\textcolor{red}{1}$$

o. Add a new starting non-terminal  $S'$ .

$$\underline{S' \rightarrow S} \quad S \rightarrow A \mid B \quad A \rightarrow \textcolor{red}{0}A \mid \textcolor{red}{0}C \quad B \rightarrow B\textcolor{red}{1} \mid C\textcolor{red}{1} \quad C \rightarrow \varepsilon \mid \textcolor{red}{0}C\textcolor{red}{1}$$

1. Decompose the long production rule  $C \rightarrow \textcolor{red}{0}C\textcolor{red}{1}$ .

$$S' \rightarrow S \quad S \rightarrow A \mid B \quad A \rightarrow \textcolor{red}{0}A \mid \textcolor{red}{0}C \quad B \rightarrow B\textcolor{red}{1} \mid C\textcolor{red}{1} \quad \underline{C \rightarrow \varepsilon \mid \textcolor{red}{0}D} \quad \underline{D \rightarrow C\textcolor{red}{1}}$$

2. Identify  $C$  as the only nullable non-terminal. Because  $S'$  is not nullable, remove the production rule  $S' \rightarrow S$ .

3. Eliminate the  $\varepsilon$ -production  $C \rightarrow \varepsilon$ .

$$S \rightarrow A \mid B \quad A \rightarrow \textcolor{red}{0}A \mid \textcolor{red}{0}C \mid \underline{\textcolor{red}{0}} \quad B \rightarrow B\textcolor{red}{1} \mid C\textcolor{red}{1} \mid \underline{\textcolor{red}{1}} \quad C \rightarrow \textcolor{red}{0}D \quad D \rightarrow C\textcolor{red}{1} \mid \underline{\textcolor{red}{1}}$$

4. No two non-terminals are equivalent, so there's nothing to merge.

5. Remove the unit productions  $S' \rightarrow S$ ,  $S \rightarrow A$ , and  $S \rightarrow B$ .

$$\underline{S \rightarrow \textcolor{red}{0}A \mid \textcolor{red}{0}C \mid B\textcolor{red}{1} \mid C\textcolor{red}{1} \mid \textcolor{red}{0} \mid \textcolor{red}{1}} \\ A \rightarrow \textcolor{red}{0}A \mid \textcolor{red}{0}C \mid \textcolor{red}{0} \quad B \rightarrow B\textcolor{red}{1} \mid C\textcolor{red}{1} \mid \textcolor{red}{1} \quad C \rightarrow \textcolor{red}{0}D \quad D \rightarrow C\textcolor{red}{1} \mid \textcolor{red}{1}.$$

6. Finally, protect the terminals  $\textcolor{red}{0}$  and  $\textcolor{red}{1}$  to obtain the final CNF grammar.

$$\underline{S \rightarrow EA \mid EC \mid BF \mid CF \mid \textcolor{red}{0} \mid \textcolor{red}{1}} \\ A \rightarrow \underline{EA \mid EC} \mid \textcolor{red}{0} \quad B \rightarrow \underline{BF \mid CF} \mid \textcolor{red}{1} \\ C \rightarrow \underline{ED} \quad D \rightarrow \underline{CF} \mid \textcolor{red}{1} \\ \underline{E \rightarrow \textcolor{red}{0}} \quad \underline{F \rightarrow \textcolor{red}{1}}$$



## Exercises

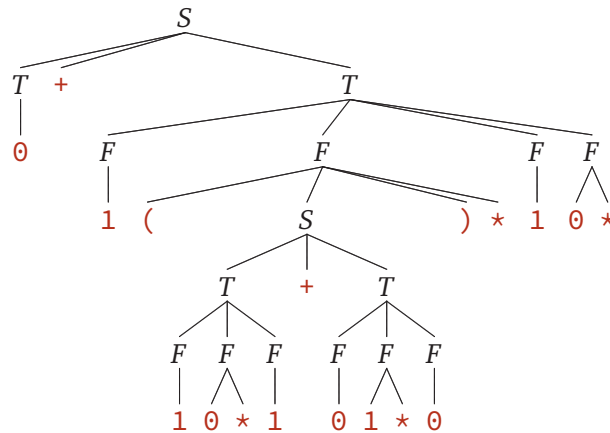
1. Describe context-free grammars that generate each of the following languages. The function  $\#(x, w)$  returns the number of occurrences of the **substring**  $x$  in the string  $w$ . For example,  $\#(0, 101001) = 3$  and  $\#(010, 1010100011) = 2$ . These are *not* listed in order of increasing difficulty.
  - (a) All strings in  $\{0, 1\}^*$  whose length is divisible by 5.
  - (b) All strings in  $\{0, 1\}^*$  representing a non-negative multiple of 5 in binary.
  - (c)  $\{w \in \{0, 1\}^* \mid \#(0, w) = 2 \cdot \#(1, w)\}$
  - (d)  $\{w \in \{0, 1\}^* \mid \#(0, w) \neq 2 \cdot \#(1, w)\}$
  - (e)  $\{w \in \{0, 1\}^* \mid \#(00, w) = \#(11, w)\}$
  - (f)  $\{w \in \{0, 1\}^* \mid \#(01, w) = \#(10, w)\}$
  - (g)  $\{w \in \{0, 1\}^* \mid \#(0, w) = \#(1, w) \text{ and } |w| \text{ is a multiple of } 3\}$
  - (h)  $\{0, 1\}^* \setminus \{0^n 1^n \mid n \geq 0\}$
  - (i)  $\{0^n 1^{2n} \mid n \geq 0\}$
  - (j)  $\{0, 1\}^* \setminus \{0^n 1^{2n} \mid n \geq 0\}$
  - (k)  $\{0^n 1^m \mid 0 \leq 2m \leq n < 3m\}$
  - (l)  $\{0^i 1^j 2^{i+j} \mid i, j \geq 0\}$
  - (m)  $\{0^i 1^j 2^k \mid i = j \text{ or } j = k\}$
  - (n)  $\{0^i 1^j 2^k \mid i \neq j \text{ or } j \neq k\}$
  - (o)  $\{0^i 1^j 0^j 1^i \mid i, j \geq 0\}$
  - (p)  $\{w \# 0^{\#(0, w)} \mid w \in \{0, 1\}^*\}$
  - (q)  $\{xy \mid x, y \in \{0, 1\}^* \text{ and } x \neq y \text{ and } |x| = |y|\}$
  - (r)  $\{x \# y^R \mid x, y \in \{0, 1\}^* \text{ and } x \neq y\}$
  - (s)  $\{x \# y \mid x, y \in \{0, 1\}^* \text{ and } \#(0, x) = \#(1, y)\}$
  - (t)  $\{0, 1\}^* \setminus \{ww \mid w \in \{0, 1\}^*\}$
  - (u) All strings in  $\{0, 1\}^*$  that are *not* palindromes.
  - (v) All strings in  $\{ (, ), \diamond \}^*$  in which the parentheses are balanced and the symbol  $\diamond$  appears at most four times. For example,  $(( ( ( ) ) ) )$  and  $(\diamond ( ( ( ) ) ) \diamond ( ) ( ) ) \diamond$  and  $\diamond \diamond \diamond$  are strings in this language, but  $) ( ( )$  and  $(\diamond \diamond \diamond) \diamond$  are not.
2. Describe recursive automata for each of the languages in problem 1. (“Describe” does not necessarily mean “draw”!)
3. Prove that if  $L$  is a context-free language, then  $L^R$  is also a context-free language. [Hint: How do you reverse a context-free grammar?]
4. Consider a generalization of context-free grammars that allows any *regular expression* over  $\Sigma \cup \Gamma$  to appear on the right side of a production rule. Without loss of generality, for each non-terminal  $A \in \Gamma$ , the generalized grammar contains a single regular expression  $R(A)$ . To

apply a production rule to a string, we replace any non-terminal  $A$  with an arbitrary word in the language described by  $R(A)$ . As usual, the language of the generalized grammar is the set of all strings that can be derived from its start non-terminal.

For example, the following generalized context-free grammar describes the language of all regular expressions over the alphabet  $\{0, 1\}$ :

$S \rightarrow (T^+)^*T + \emptyset$	(Regular expressions)
$T \rightarrow \mathfrak{E} + F^*F$	(Terms = summable expressions)
$F \rightarrow (\mathbf{0} + \mathbf{1} + (S))(\star + \varepsilon)$	(Factors = concatenable expressions)

Here is a parse tree for the regular expression  $0+1(10^*1+01^*0)^*10^*$  (which represents the set of all binary numbers divisible by 3):



Prove that every *generalized* context-free grammar describes a context-free language. In other words, show that allowing regular expressions to appear in production rules does not increase the expressive power of context-free grammars.

## 6 Turing Machines

In 1936, a few months before his 24th birthday, Alan Turing launched computer science as a modern intellectual discipline. In a single remarkable paper, Turing provided the following results:

- A simple formal model of mechanical computation now known as *Turing machines*.
- A description of a single *universal* machine that can be used to compute *any* function computable by *any* other Turing machine.
- A proof that no Turing machine can solve the *halting problem*—Given the formal description of an arbitrary Turing machine  $M$ , does  $M$  halt or run forever?
- A proof that no Turing machine can determine whether an arbitrary given proposition is provable from the axioms of first-order logic. This is Hilbert and Ackermann’s famous *Entscheidungsproblem* (“decision problem”).
- Compelling arguments<sup>1</sup> that his machines can execute *arbitrary* “calculation by finite means”.

Although Turing did not know it at the time, he was not the first to prove that the *Entscheidungsproblem* had no algorithmic solution. The first such proof is implicit in the work of Kurt Gödel; in lectures on his incompleteness theorems<sup>2</sup> at Princeton in 1934, Gödel described a model of **general recursive functions**, which he largely credited to Jacques Herbrand.<sup>3</sup> Gödel’s incompleteness theorem can be viewed as a proof that some propositions cannot be proved using general recursive functions. However, Gödel viewed his definition as a “heuristic principle” rather than an accurate model of effective computation, or even a complete definition of recursion.

The first *published* proof was written by Alonzo Church and published just a few months before Turing’s paper, using another different model of computation, now called the **untyped  $\lambda$ -calculus**. Turing and Church developed their results independently; indeed, Turing rushed

---

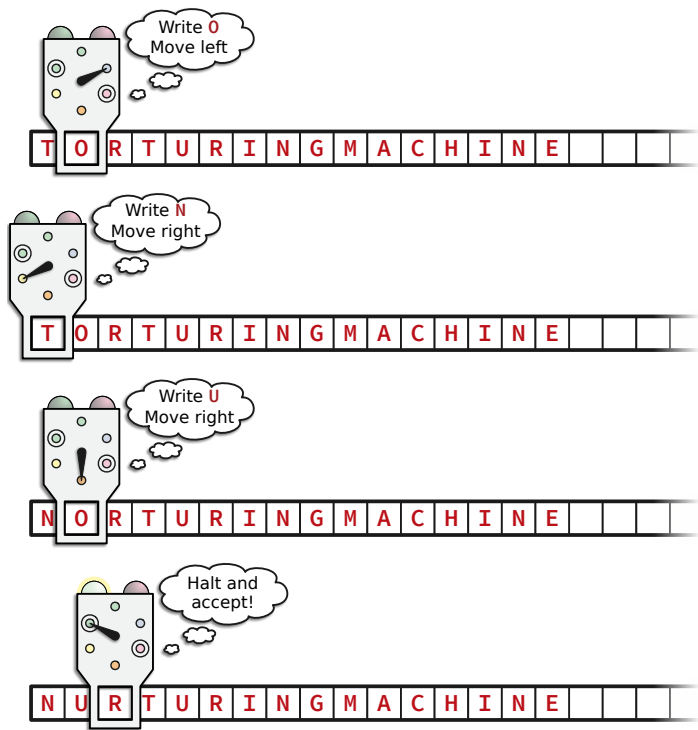
<sup>1</sup>As Turing put it, “All arguments which can be given are bound to be, fundamentally, appeals to intuition, and for this reason rather unsatisfactory mathematically.” The claim that anything that can be computed can be computed using Turing machines is now known as the **Church-Turing thesis**.

<sup>2</sup>which he published at the ripe old age of 25

<sup>3</sup>Herbrand was a brilliant French mathematician who was killed in a mountain-climbing accident at the age of 23.

the submission of his own paper immediately after receiving a copy of Church’s paper, pausing only long enough to prove that any function computable via  $\lambda$ -calculus can also be computed by a Turing machine and vice versa.<sup>4</sup> Church was the referee for Turing’s paper; between the paper’s submission and its acceptance, Turing was admitted to Princeton as a PhD student, where Church became his advisor. He finished his PhD two years later.

Informally, Turing described an abstract machine with a finite number of *internal states* that has access to “memory” in the form of a *tape*. The tape consists of a semi-infinite sequence of *cells*, each containing a single symbol from some arbitrary finite alphabet. The Turing machine can access the tape only through its *head*, which is positioned over a single cell. Initially, the tape contains an arbitrary finite *input string* followed by an infinite sequence of *blanks*, and the head is positioned over the first cell on the tape. In a single iteration, the machine reads the symbol in that cell, possibly write a new symbol into that cell, possibly changes its internal state, possibly moves the head to a neighboring cell, and possibly halts. The precise behavior of the machine at each iteration is entirely determined by its internal state and the symbol that it reads. When the machine halts, it indicates whether it has *accepted* or *rejected* the original input string.



A few iterations of a six-state Turing machine.

### 6.1 Why Bother?

Students used to thinking of computation in terms of higher-level operations like random memory accesses, function calls, and recursion may wonder why we should even consider a model as simple and constrained as Turing machines. Admittedly, Turing machines are a terrible model for thinking about *fast* computation; simple operations that take constant time in the standard

<sup>4</sup>At roughly the same time as Turing’s paper, Church and his more recent PhD graduate Steven Kleene independently proved that all general recursive functions are definable in the  $\lambda$ -calculus and vice versa. At the time, Kleene was 27, and Church was 33.

random-access model can require *arbitrarily* many steps on a Turing machine. Worse, seemingly minor variations in the precise definition of “Turing machine” can have significant impact on problem complexity. As a simple example (which will make more sense later), we can reverse a string of  $n$  bits in  $O(n)$  time using a two-tape Turing machine, but the same task provably requires  $\Omega(n^2)$  time on a single-tape machine.

But here we are not interested in finding *fast* algorithms, or indeed in finding algorithms at all, but rather in proving that some problems cannot be solved by *any* computational means. Such a bold claim requires a formal definition of “computation” that is simple enough to support formal argument, but still powerful enough to describe arbitrary algorithms. Turing machines are ideal for this purpose. In particular, Turing machines are powerful enough to *simulate other Turing machines*, while still simple enough to let us build up this self-simulation from scratch, unlike more complex but efficient models like the standard random-access machine

(Arguably, self-simulation is even simpler in Church’s  $\lambda$ -calculus, or in Schönfinkel and Curry’s combinator calculus, which is one of many reasons those models are more common in the design and analysis of programming languages than Turing machines. Those models are much more abstract; in particular, they are harder to show equivalent to standard iterative models of computation.)

## 6.2 Formal Definitions

Formally, a Turing machine consists of the following components. (Hang on; it’s a long list.)

- An arbitrary finite set  $\Gamma$  with at least two elements, called the **tape alphabet**.
- An arbitrary symbol  $\square \in \Gamma$ , called the **blank symbol** or just the **blank**.
- An arbitrary nonempty subset  $\Sigma \subseteq (\Gamma \setminus \{\square\})$ , called the **input alphabet**.
- Another arbitrary finite set  $Q$  whose elements are called **states**.
- Three distinct special states **start**, **accept**, **reject**  $\in Q$ .
- A **transition** function  $\delta: (Q \setminus \{\text{accept}, \text{reject}\}) \times \Gamma \rightarrow Q \times \Gamma \times \{-1, +1\}$ .

A **configuration** or **global state** of a Turing machine is represented by a triple  $(q, x, i) \in Q \times \Gamma^* \times \mathbb{N}$ , indicating that the machine’s internal state is  $q$ , the tape contains the string  $x$  followed by an infinite sequence of blanks, and the head is located at position  $i$ . Trailing blanks in the tape string are ignored; the triples  $(q, x, i)$  and  $(q, x\square, i)$  describe exactly the same configuration.

The transition function  $\delta$  describes the evolution of the machine. For example,  $\delta(q, a) = (p, b, -1)$  means that when the machine reads symbol  $a$  in state  $q$ , it changes its internal state to  $p$ , writes symbol  $b$  onto the tape at its current location (replacing  $a$ ), and then decreases its position by 1 (or more intuitively, moves one step to the left). If the position of the head becomes negative, no further transitions are possible, and the machine **crashes**.

We write  $(p, x, i) \Rightarrow_M (q, y, j)$  to indicate that Turing machine  $M$  transitions from the first configuration to the second in one step. (The symbol  $\Rightarrow$  is often pronounced “yields”; I will omit the subscript  $M$  if the machine is clear from context.) For example,  $\delta(p, a) = (q, b, \pm 1)$  means that

$$(p, xay, i) \Rightarrow (q, xby, i \pm 1)$$

for any non-negative integer  $i$ , any string  $x$  of length  $i$ , and any string  $y$ . The evolution of any Turing machine is **deterministic**; each configuration  $C$  yields a *unique* configuration  $C'$ . We write  $C \Rightarrow^* C'$  to indicate that there is a (possibly empty) sequence of transitions from configuration  $C$  to configuration  $C'$ . (The symbol  $\Rightarrow^*$  can be pronounced “eventually yields”.)

The initial configuration is (**start**,  $w$ , 0) for some arbitrary (and possibly empty) **input string**  $w \in \Sigma^*$ . If  $M$  eventually reaches the **accept** state—more formally, if  $(\text{start}, w, 0) \Rightarrow^* (\text{accept}, x, i)$  for some string  $x \in \Gamma^*$  and some integer  $i$ —we say that  $M$  **accepts** the original input string  $w$ . Similarly, if  $M$  eventually reaches the **reject** state, we say that  $M$  **rejects**  $w$ . We must emphasize that “rejects” and “does not accept” are *not* synonyms; if  $M$  crashes or runs forever, then  $M$  neither accepts nor rejects  $w$ .

We distinguish between two different senses in which a Turing machine can “accept” a language. Let  $M$  be a Turing machine with input alphabet  $\Sigma$ , and let  $L \subseteq \Sigma^*$  be an arbitrary language over  $\Sigma$ .

- $M$  **recognizes** or **accepts**  $L$  if and only if  $M$  accepts every string in  $L$  but nothing else. A language is **recognizable** (or *semi-computable* or *recursively enumerable*) if it is recognized by some Turing machine.
- $M$  **decides**  $L$  if and only if  $M$  accepts every string in  $L$  and rejects every string in  $\Sigma^* \setminus L$ . Equivalently,  $M$  decides  $L$  if and only if  $M$  recognizes  $L$  and halts (without crashing) on all inputs. A language is **decidable** (or *computable* or *recursive*) if it is decided by some Turing machine.

Trivially, every decidable language is recognizable, but (as we will see later), not every recognizable language is decidable.

## 6.3 A First Example

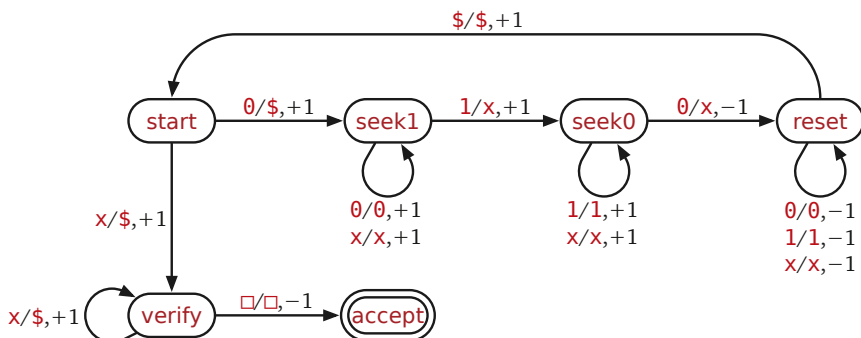
Consider the language  $L = \{0^n 1^n 0^n \mid n \geq 0\}$ . This language is neither regular nor context-free, but it can be decided by the following six-state Turing machine. The alphabets and states of the machine are defined as follows:

$$\Gamma = \{0, 1, \$, x, \square\}$$

$$\Sigma = \{0, 1\}$$

$$Q = \{\text{start}, \text{seek1}, \text{seek0}, \text{reset}, \text{verify}, \text{accept}, \text{reject}\}$$

The transition function is described in the table on the next page; all unspecified transitions lead to the **reject** state. The figure below shows a graphical representation of the same machine, which resembles a drawing of a DFA, but with output symbols and actions specified on each edge. For example, we indicate the transition  $\delta(p, 0) = (q, 1, +1)$  by writing  $0/1, +1$  next to the arrow from state  $p$  to state  $q$ .



A graphical representation of the example Turing machine

$\delta(p, a) = (q, b, \Delta)$	explanation
$\delta(\text{start}, 0) = (\text{seek1}, \$, +1)$	mark first 0 and scan right
$\delta(\text{start}, x) = (\text{verify}, \$, +1)$	looks like we're done, but let's make sure
$\delta(\text{seek1}, 0) = (\text{seek1}, 0, +1)$	scan rightward for 1
$\delta(\text{seek1}, x) = (\text{seek1}, x, +1)$	
$\delta(\text{seek1}, 1) = (\text{seek0}, x, +1)$	mark 1 and continue right
$\delta(\text{seek0}, 1) = (\text{seek0}, 1, +1)$	scan rightward for 0
$\delta(\text{seek0}, x) = (\text{seek0}, x, +1)$	
$\delta(\text{seek0}, 0) = (\text{reset}, x, +1)$	mark 0 and scan left
$\delta(\text{reset}, 0) = (\text{reset}, 0, -1)$	scan leftward for \$
$\delta(\text{reset}, 1) = (\text{reset}, 1, -1)$	
$\delta(\text{reset}, x) = (\text{reset}, x, -1)$	
$\delta(\text{reset}, \$) = (\text{start}, \$, +1)$	step right and start over
$\delta(\text{verify}, x) = (\text{verify}, \$, +1)$	scan right for any unmarked symbol
$\delta(\text{verify}, \square) = (\text{accept}, \square, -1)$	success!

The transition function for a Turing machine that decides the language  $\{0^n 1^n 0^n \mid n \geq 0\}$ .

Finally, we trace the execution of this machine on two input strings:  $001100 \in L$  and  $00100 \notin L$ . In each configuration, we indicate the position of the head using a small triangle instead of listing the position explicitly. Notice that we automatically add blanks to the tape string as necessary. Proving that this machine actually decides  $L$ —and in particular, that it never crashes or infinite-loops—is a straightforward but tedious exercise in induction.

## 6.4 Variations

There are actually several formal models that all fall under the name “Turing machine”, each with small variations on the definition we’ve given. Although we do need to be explicit about *which* variant we want to use for any particular problem, the differences between the variants are relatively unimportant. For any machine defined in one model, there is an equivalent machine in each of the other models; in particular, all of these variants recognize the same languages and decide the same languages. For example:

- **Halting conditions.** Some models allow multiple accept and reject states, which (depending on the precise model) trigger acceptance or rejection either when the machine enters the state, or when the machine has no valid transitions out of such a state. Others include only explicit accept states, and either equate crashing with rejection or do not define a rejection mechanism at all. Still other models include halting as one of the possible *actions* of the machine, in addition to moving left or moving right; in these models, the machine accepts/rejects its input if and only if it halts in an accepting/non-accepting state.
- **Actions.** Some Turing machine models allow transitions that do not move the head, or that move the head by more than one cell in a single step. Others insist that a single step of the machine *either* writes a new symbol onto the tape *or* moves the head one step. Finally, as mentioned above, some models include halting as one of the available actions.
- **Transition function.** Some models of Turing machines, including Turing’s original definition, allow the transition function to be undefined on some state-symbol pairs. In this

(start, 001100)  
 ⇒ (seek1, \$01100)  
 ⇒ (seek1, \$01100)  
 ⇒ (seek0, \$0x100)  
 ⇒ (seek0, \$0x100)  
 ⇒ (reset, \$0x1x0)  
 ⇒ (reset, \$0x1x0)  
 ⇒ (reset, \$0x1x0)  
 ⇒ (reset, \$0x1x0)  
 ⇒ (start, \$0x1x0)  
 ⇒ (seek1, \$\$x1x0)  
 ⇒ (seek1, \$\$x1x0)  
 ⇒ (seek0, \$\$\$xx0)  
 ⇒ (seek0, \$\$\$xx0)  
 ⇒ (reset, \$\$\$xxx)  
 ⇒ (reset, \$\$\$xxx)  
 ⇒ (reset, \$\$\$xxx)  
 ⇒ (reset, \$\$\$xxx)  
 ⇒ (start, \$\$\$xxx)  
 ⇒ (verify, \$\$\$xx)  
 ⇒ (verify, \$\$\$xx)  
 ⇒ (verify, \$\$\$\$x)  
 ⇒ (verify, \$\$\$\$x)  
 ⇒ (accept, \$\$\$\$x) ⇒ **accept!**

The evolution of the example Turing machine on the input string 001100 ∈ L

(start, 00100)  
 ⇒ (seek1, \$0100)  
 ⇒ (seek1, \$0100)  
 ⇒ (seek0, \$0x00)  
 ⇒ (reset, \$0xx0)  
 ⇒ (reset, \$0xx0)  
 ⇒ (reset, \$0xx0)  
 ⇒ (start, \$0xx0)  
 ⇒ (seek1, \$\$xx0)  
 ⇒ (seek1, \$\$xx0)  
 ⇒ (seek1, \$\$xx0) ⇒ **reject!**

The evolution of the example Turing machine on the input string 00100 ∉ L



formulation, the transition function is given by a set  $\delta \subset Q \times \Gamma \times Q \times \Gamma \times \{+1, -1\}$ , such that for each state  $q$  and symbol  $a$ , there is at most one transition  $(q, a, \cdot, \cdot, \cdot) \in \delta$ . If the machine enters a configuration from which there is no transition, it halts and (depending on the precise model) either crashes or rejects. Others define the transition function as  $\delta: Q \times \Gamma \rightarrow Q \times (\Gamma \cup \{-1, +1\})$ , allowing the machine to *either* write a symbol to the tape *or* move the head in each step.

- **Beginning of the tape.** Some models forbid the head to move past the beginning of the tape, either by starting the tape with a special symbol that cannot be overwritten and that forces a rightward transition, or by declaring that a leftward transition at position 0 leaves the head in position 0, or even by pure fiat—declaring any machine that performs a leftward move at position 0 to be invalid.

To prove that any two of these variant “species” of Turing machine are equivalent, we must show how to transform a machine of one species into a machine of the other species that accepts and rejects the same strings. For example, let  $M = (\Gamma, \sqcup, \Sigma, Q, s, \text{accept}, \text{reject}, \delta)$  be a Turing machine with explicit accept and reject states. We can define an equivalent Turing machine  $M'$  that halts only when it moves left from position 0, and accepts only by halting while in an accepting state, as follows. We define the set of accepting states for  $M'$  as  $A = \{\text{accept}\}$  and define a new transition function

$$\delta'(q, a) := \begin{cases} (\text{accept}, a, -1) & \text{if } q = \text{accept} \\ (\text{reject}, a, -1) & \text{if } q = \text{reject} \\ \delta(q, a) & \text{otherwise} \end{cases}$$

Similarly, suppose someone gives us a Turing machine  $M = (\Gamma, \sqcup, \Sigma, Q, s, \text{accept}, \text{reject}, \delta)$  whose transition function  $\delta: Q \times \Gamma \rightarrow Q \times \Gamma \times \{-1, 0, +1\}$  allows the machine to transition without moving its head. Then we can construct an equivalent Turing machine  $M' = (\Gamma, \sqcup, \Sigma, Q', s, \text{accept}, \text{reject}, \delta')$  that moves its head at every transition by defining  $Q' := Q \times \{0, 1\}$  and

$$\begin{aligned} \delta'((p, 0), a) &:= \begin{cases} ((q, 1), b, +1) & \text{if } \delta(p, a) = (q, b, 0), \\ ((q, 0), b, \Delta) & \text{if } \delta(p, a) = (q, b, \Delta) \text{ and } \Delta \neq 0, \end{cases} \\ \delta'((p, 1), a) &:= ((p, 0), a, -1). \end{aligned}$$

## 6.5 Computing Functions

Turing machines can also be used to compute functions from strings to strings, instead of just accepting or rejecting strings. Since we don’t care about acceptance or rejection, we replace the explicit **accept** and **reject** states with a single halt state, and we define the **output** of the Turing machine to be the contents of the tape when the machine halts, after removing the infinite sequence of trailing blanks. More formally, for any Turing machine  $M$ , any string  $w \in \Sigma^*$ , and any string  $x \in \Gamma^*$  that does not end with a blank, we write  $\mathbf{M}(w) = x$  if and only if  $(w, s, 0) \Rightarrow_M^* (x, \text{halt}, i)$  for some integer  $i$ . If  $M$  does not halt on input  $w$ , then we write  $\mathbf{M}(w) \nearrow$ , which can be read either “ $M$  diverges on  $w$ ” or “ $M(w)$  is undefined.” We say that  $M$  **computes** the function  $f: \Sigma^* \rightarrow \Sigma^*$  if and only if  $\mathbf{M}(w) = f(w)$  for every string  $w$ .

### 6.5.1 Shifting

One basic operation that is used in many Turing machine constructions is **shifting** the input string a constant number of steps to the right or to the left. For example, given any input

string  $w \in \{0, 1\}^*$ , we can compute the string  $0w$  using a Turing machine with tape alphabet  $\Gamma = \{0, 1, \square\}$ , state set  $Q = \{0, 1, \text{halt}\}$ , start state 0, and the following transition function:

$$\begin{array}{l} \delta(p, a) = (q, b, \Delta) \\ \delta(0, 0) = (0, 0, +1) \\ \delta(0, 1) = (1, 0, +1) \\ \delta(0, \square) = (\text{halt}, 0, +1) \\ \delta(1, 0) = (0, 1, +1) \\ \delta(1, 1) = (1, 1, +1) \\ \delta(1, \square) = (\text{halt}, 1, +1) \end{array}$$

By increasing the number of states, we can build a Turing machine that shifts the input string any fixed number of steps in either direction. For example, a machine that shifts its input to the left by five steps might read the string from right to left, storing the five most recently read symbols in its internal state. A typical transition for such a machine would be  $\delta(12345, 0) = (01234, 5, -1)$ .

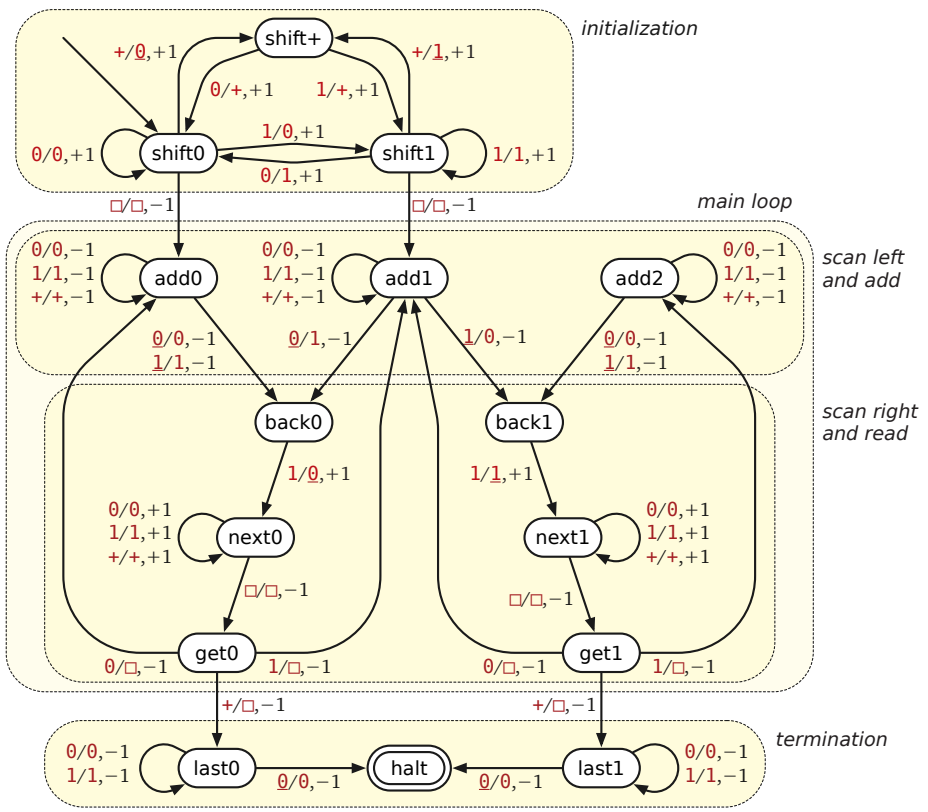
### 6.5.2 Binary Addition

With a more complex Turing machine, we can implement binary addition. The input is a string of the form  $w+x$ , where  $w, x \in \{0, 1\}^n$ , representing two numbers in binary; the output is the binary representation of  $w+x$ . To simplify our presentation, we assume that  $|w| = |x| > 0$ ; however, this restrictions can be removed with the addition of a few more states. The following figure shows the entire Turing machine at a glance. The machine uses the tape alphabet  $\Gamma = \{\square, 0, 1, +, \underline{0}, \underline{1}\}$ ; the start state is **shift0**. All missing transitions go to a **fail** state, indicating that the input was badly formed.

Execution of this Turing machine proceeds in several phases, each with its own subset of states, as indicated in the figure. The initialization phase scans the entire input, shifting it to the right to make room for the output string, marking the rightmost bit of  $w$ , and reading and erasing the last bit of  $x$ .

$$\begin{array}{l} \delta(p, a) = (q, b, \Delta) \\ \delta(\text{shift0}, 0) = (\text{shift0}, 0, +1) \\ \delta(\text{shift0}, 1) = (\text{shift1}, 0, +1) \\ \delta(\text{shift0}, +) = (\text{shift+}, \underline{0}, +1) \\ \delta(\text{shift0}, \square) = (\text{add0}, \square, -1) \\ \delta(\text{shift1}, 0) = (\text{shift0}, 1, +1) \\ \delta(\text{shift1}, 1) = (\text{shift1}, 1, +1) \\ \delta(\text{shift1}, +) = (\text{shift+}, \underline{1}, +1) \\ \delta(\text{shift1}, \square) = (\text{add1}, \square, -1) \\ \delta(\text{shift+}, 0) = (\text{shift0}, +, +1) \\ \delta(\text{shift+}, 1) = (\text{shift1}, +, +1) \end{array}$$

The first part of the main loop scans left to the marked bit of  $w$ , adds the bit of  $x$  that was just erased plus the carry bit from the previous iteration, and records the carry bit for the next iteration in the machines internal state.



A Turing machine that adds two binary numbers of the same length.

$\delta(p, a) = (q, b, \Delta)$	$\delta(p, a) = (q, b, \Delta)$	$\delta(p, a) = (q, b, \Delta)$
$\delta(\text{add0}, 0) = (\text{add0}, 0, -1)$	$\delta(\text{add1}, 0) = (\text{add1}, 0, -1)$	$\delta(\text{add2}, 0) = (\text{add2}, 0, -1)$
$\delta(\text{add0}, 1) = (\text{add0}, 0, -1)$	$\delta(\text{add1}, 1) = (\text{add1}, 0, -1)$	$\delta(\text{add2}, 1) = (\text{add2}, 0, -1)$
$\delta(\text{add0}, +) = (\text{add0}, 0, -1)$	$\delta(\text{add1}, +) = (\text{add1}, 0, -1)$	$\delta(\text{add2}, +) = (\text{add2}, 0, -1)$
$\delta(\text{add0}, \underline{0}) = (\text{back0}, 0, -1)$	$\delta(\text{add1}, \underline{0}) = (\text{back0}, 1, -1)$	$\delta(\text{add2}, \underline{0}) = (\text{back1}, 0, -1)$
$\delta(\text{add0}, \underline{1}) = (\text{back0}, 1, -1)$	$\delta(\text{add1}, \underline{1}) = (\text{back1}, 0, -1)$	$\delta(\text{add2}, \underline{1}) = (\text{back1}, 1, -1)$

The second part of the main loop marks the previous bit of  $w$ , scans right to the end of  $x$ , and then reads and erases the last bit of  $x$ , all while maintaining the carry bit.

$\delta(p, a) = (q, b, \Delta)$	$\delta(p, a) = (q, b, \Delta)$
$\delta(\text{back0}, 0) = (\text{next0}, \underline{0}, +1)$	$\delta(\text{back1}, 0) = (\text{next1}, \underline{0}, +1)$
$\delta(\text{back0}, 1) = (\text{next0}, \underline{1}, +1)$	$\delta(\text{back1}, 1) = (\text{next1}, \underline{1}, +1)$
$\delta(\text{next0}, 0) = (\text{next0}, 0, +1)$	$\delta(\text{next1}, 0) = (\text{next1}, 0, +1)$
$\delta(\text{next0}, 1) = (\text{next0}, 0, +1)$	$\delta(\text{next1}, 1) = (\text{next1}, 0, +1)$
$\delta(\text{next0}, +) = (\text{next0}, 0, +1)$	$\delta(\text{next1}, +) = (\text{next1}, 0, +1)$
$\delta(\text{next0}, \square) = (\text{get0}, \square, -1)$	$\delta(\text{next1}, \square) = (\text{get1}, \square, -1)$
$\delta(\text{get0}, 0) = (\text{add0}, \square, -1)$	$\delta(\text{get1}, 0) = (\text{add1}, \square, -1)$
$\delta(\text{get0}, 1) = (\text{add1}, \square, -1)$	$\delta(\text{get1}, 1) = (\text{add2}, \square, -1)$
$\delta(\text{get0}, +) = (\text{last0}, \square, -1)$	$\delta(\text{get1}, +) = (\text{last1}, \square, -1)$

Finally, after erasing the  $+$  in the last iteration of the main loop, the termination phase adds the last carry bit to the leftmost output bit and halts.

$$\begin{array}{l}
\delta(p, a) = (q, b, \Delta) \\
\delta(\text{last0}, \underline{0}) = (\text{last0}, 0, -1) \\
\delta(\text{last0}, 1) = (\text{last0}, 0, -1) \\
\delta(\text{last0}, \underline{0}) = (\text{halt}, 0, \quad) \\
\delta(\text{last1}, \underline{0}) = (\text{last1}, 0, -1) \\
\delta(\text{last1}, 1) = (\text{last1}, 0, -1) \\
\delta(\text{last1}, \underline{0}) = (\text{halt}, 1, \quad)
\end{array}$$

## 6.6 Variations on Tracks, Heads, and Tapes

### Multiple Tracks

It is sometimes convenient to endow the Turing machine tape with multiple *tracks*, each with its own tape alphabet, and allow the machine to read from and write to the same position on all tracks simultaneously. For example, to define a Turing machine with three tracks, we need three tape alphabets  $\Gamma_1$ ,  $\Gamma_2$ , and  $\Gamma_3$ , each with its own blank symbol, where (say)  $\Gamma_1$  contains the input alphabet  $\Sigma$  as a subset; we also need a transition function of the form

$$\delta : Q \times \Gamma_1 \times \Gamma_2 \times \Gamma_3 \rightarrow Q \times \Gamma_1 \times \Gamma_2 \times \Gamma_3 \times \{-1, +1\}$$

Describing a configuration of this machine requires a quintuple  $(q, x_1, x_2, x_3, i)$ , indicating that each track  $i$  contains the string  $x_i$  followed by an infinite sequence of blanks. The initial configuration is  $(\text{start}, w, \varepsilon, \varepsilon, 0)$ , with the input string written on the first track, and the other two tracks completely blank.

But any such machine is equivalent (if not *identical*) to a single-track Turing machine with the (still finite!) tape alphabet  $\Gamma := \Gamma_1 \times \Gamma_2 \times \Gamma_3$ . Instead of thinking of the tape as three infinite sequences of symbols, we think of it as a single infinite sequence of “records”, each containing three symbols. Moreover, there’s nothing special about the number 3 in this construction; a Turing machine with *any* constant number of tracks is equivalent to a single-track machine.

### Doubly-Infinite Tape

It is also sometimes convenient to allow the tape to be infinite in both directions, for example, to avoid boundary conditions. There are several ways to simulate a doubly-infinite tape on a machine with only a semi-infinite tape. Perhaps the simplest method is to use a semi-infinite tape with two tracks, one containing the cells with positive index and the other containing the cells with negative index in reverse order, with a special marker symbol at position zero to indicate the transition.

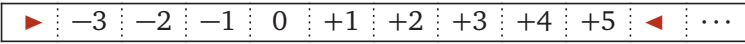
0	+1	+2	+3	+4	...
►	-1	-2	-3	-4	...

Another method is to shuffle the positive-index and negative-index cells onto a single track, and add additional states to allow the Turing machine to move two steps in a single transition. Again, we need a special symbol at the left end of the tape to indicate the transition:

►	0	-1	+1	-2	+2	-3	+3	...
---	---	----	----	----	----	----	----	-----

A third method maintains two sentinel symbols ► and ◄ that surround all other non-blank symbols on the tape. Whenever the machine reads the right sentinel ◄, we write a blank, move right, write ◄, move left, and then proceed as if we had just read a blank. On the other hand, when the machine reads the left sentinel ►, we shift the entire contents of the tape (up to and

including the right sentinel) one step to the right, then move back to the left sentinel, move right, write a blank, and finally proceed as if we had just read a blank. Since the Turing machine does not actually have access to the position of the head *as an integer*, shifting the head and the tape contents one step right has no effect on its future evolution.



Using either of the first two methods, we can simulate  $t$  steps of an arbitrary Turing machine with a doubly-infinite tape using only  $O(t)$  steps on a standard Turing machine. The third method, unfortunately, requires  $\Theta(t^2)$  steps in the worst case.

### Insertion and Deletion

We can also allow Turing machines to insert and delete cells on the tape, in addition to simply overwriting existing symbols. We’ve already seen how to insert a new cell: Leave a special mark on the tape (perhaps in a second track), shift everything to the right of this mark one cell to the right, scan left to the mark, erase the mark, and finally write the correct character into the new cell. Deletion is similar: Mark the cell to be deleted, shift everything to the right of the mark one step to the left, scan left to the mark, and erase the mark. We may also need to maintain a mark in some cell to the right every non-blank symbol, indicating that all cells further to the right are blank, so that we know when to stop shifting left or right.

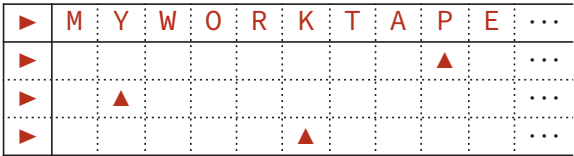
### Multiple Heads

Another convenient extension is to allow machines simultaneous access to more than one position on the tape. For example, to define a Turing machine with *three* heads, we need a transition function of the form

$$\delta: Q \times \Gamma^3 \rightarrow Q \times \Gamma^3 \times \{-1, +1\}^3.$$

Describing a configuration of such a machine requires a quintuple  $(q, x, i, j, k)$ , indicating that the machine is in state  $q$ , the tape contains string  $x$ , and the three heads are at positions  $i, j, k$ . The transition function tells us, given  $q$  and the three symbols  $x[i], x[j], x[k]$ , which three symbols to write on the tape and which direction to move each of the heads.

We can simulate this behavior with a single head by adding additional tracks to the tape that record the positions of each head. To simulate a machine  $M$  with three heads, we use a tape with four tracks: track 0 is the actual work tape; each of the remaining tracks has a single non-blank symbol recording the position of one of the heads. We also insert a special marker symbols at the left end of the tape.



We can simulate any single transition of  $M$ , starting with our single head at the left end of the tape, as follows. Throughout the simulation, we maintain the internal state of  $M$  as one of the components of our current state. First, for each  $i$ , we read the symbol under the  $i$ th head of  $M$  as follows:

Scan to the right to find the mark on track  $i$ , read the corresponding symbol from track 0 into our internal state, and then return to the left end of the tape.

At this point, our internal state records  $M$ 's current internal state and the three symbols under  $M$ 's heads. After one more transition (using  $M$ 's transition function), our internal state records  $M$ 's *next* state, the symbol to be written by each head, and the direction to move each head. Then, for each  $i$ , we write with and move the  $i$ th head of  $M$  as follows:

Scan to the right to find the mark on track  $i$ , write the correct symbol onto on track 0, move the mark on track  $i$  one step left or right, and then return to the left end of the tape.

Again, there is nothing special about the number 3 here; we can simulate machines with *any* fixed number of heads.

Careful analysis of this technique implies that for any integer  $k$ , we can simulate  $t$  steps of an arbitrary Turing machine with  $k$  independent heads in  $\Theta(t^2)$  time on a standard Turing machine with only one head. Unfortunately, this quadratic blowup is unavoidable. It is relatively easy to recognize the language of *marked palindromes*  $\{w \bullet w^R \mid w \in \{0, 1\}^*\}$  in  $O(n)$  time using a Turing machine with two heads, but recognizing this language provably requires  $\Omega(n^2)$  time on a standard machine with only one head. On the other hand, with much more sophisticated techniques, it is possible to simulate  $t$  steps of a Turing machine with  $k$  head, for any fixed integer  $k$ , using only  $O(t \log t)$  steps on a Turing machine with just *two* heads.

### Multiple Tapes

We can also allow machines with multiple independent tapes, each with its own head. To simulate such a machine with a single tape, we simply maintain each tape as an independent track with its own head. Equivalently, we can simulate a machine with  $k$  tapes using a single tape with  $2k$  tracks, half storing the contents of the  $k$  tapes and half storing the positions of the  $k$  heads.

▶	T	A	P	E	#	O	N	E		...
▶									▲	...
▶	T	A	P	E	#	T	W	O		...
▶			▲							...
▶	T	A	P	E	#	T	H	R	E	...
▶						▲				...

Just as for multiple tracks, for any constant  $k$ , we can simulate  $t$  steps of an arbitrary Turing machine with  $k$  independent tapes in  $\Theta(t^2)$  steps on a standard Turing machine with one tape, and this quadratic blowup is unavoidable. Moreover, it is possible to simulate  $t$  steps on a  $k$ -tape Turing machine using only  $O(t \log t)$  steps on a *two*-tape Turing machine using more sophisticated techniques. (This faster simulation is easier to obtain for multiple independent tapes than for multiple heads on the same tape.)

By combining these tricks, we can simulate a Turing machine with any fixed number of tapes, each of which may be infinite in one or both directions, each with any fixed number of heads and any fixed number of tracks, with at most a quadratic blowup in the running time.

## 6.7 Simulating a Real Computer

### 6.7.1 Subroutines and Recursion

Use a second tape/track as a “call stack”. Add save and restore actions. In the simplest formulation, subroutines do not have local memory. To call a subroutine, save the current state onto the call stack and jump to the first state of the subroutine. To return, restore (and remove) the return state from the call stack. We can simulate  $t$  steps of any recursive Turing machine with  $O(t)$  steps on a multitape standard Turing machine, or in  $O(t^2)$  steps on a standard Turing machine.

More complex versions of this simulation can adapt to

### 6.7.2 Random-Access Memory

Keep [address•data] pairs on a separate “memory” tape. Write address to an “address” tape; read data from or write data to a “data” tape. Add new or changed [address•data] pairs at the end of the memory tape. (Semantics of reading from an address that has never been written to?)

Suppose all memory accesses require at most  $\ell$  address and data bits. Then we can simulate the  $k$ th memory access in  $O(k\ell)$  steps on a multitape Turing machine or in  $O(k^2\ell^2)$  steps on a single-tape machine. Thus, simulating  $t$  memory accesses in a random-access machine with  $\ell$ -bit words requires  $O(t^2\ell)$  time on a multitape Turing machine, or  $O(t^3\ell^2)$  time on a single-tape machine.

## 6.8 Universal Turing Machines

With all these tools in hand, we can now describe the pinnacle of Turing machine constructions: the **universal** Turing machine. For modern computer scientists, it’s useful to think of a universal Turing machine as a “Turing machine *interpreter* written in Turing machine”. Just as the input to a Python interpreter is a string of Python source code, the input to our universal Turing machine  $U$  is a string  $\langle M, w \rangle$  that encodes both an arbitrary Turing machine  $M$  and a string  $w$  in the input alphabet of  $M$ . Given these encodings,  $U$  simulates the execution of  $M$  on input  $w$ ; in particular,

- $U$  accepts  $\langle M, w \rangle$  if and only if  $M$  accepts  $w$ .
- $U$  rejects  $\langle M, w \rangle$  if and only if  $M$  rejects  $w$ .

The next few pages, I will sketch a universal Turing machine  $U$  that uses the input alphabet  $\{0, 1, [, ], \bullet, | \}$  and a somewhat larger tape alphabet (via marks on additional tracks). However, I will *not* require that the Turing machines that  $U$  simulates have similarly small alphabets, so we first need a method to encode *arbitrary* input and tape alphabets.

### Encodings

Let  $M = (\Gamma, \square, \Sigma, Q, \text{start}, \text{accept}, \text{reject}, \delta)$  be an arbitrary Turing machine, with a single half-infinite tape and a single read-write head. (I will consistently indicate the states and tape symbols of  $M$  in *slanted green* to distinguish them from the *upright red* states and tape symbols of  $U$ .)

We encode each symbol  $a \in \Gamma$  as a unique string  $|a|$  of  $\lceil \lg(|\Gamma|) \rceil$  bits. Thus, if  $\Gamma = \{0, 1, \$, x, \square\}$ , we might use the following encoding:

$$\langle 0 \rangle = 001, \quad \langle 1 \rangle = 010, \quad \langle \$ \rangle = 011, \quad \langle x \rangle = 100, \quad \langle \square \rangle = 000.$$

The input string  $w$  is encoded by its sequence of symbol encodings, with separators  $\bullet$  between every pair of symbols and with brackets  $[$  and  $]$  around the whole string. For example, with this encoding, the input string  $001100$  would be encoded on the input tape as

$$\langle 001100 \rangle = [001 \bullet 001 \bullet 010 \bullet 010 \bullet 001 \bullet 001]$$

Similarly, we encode each state  $q \in Q$  as a distinct string  $\langle q \rangle$  of  $\lceil \lg|Q| \rceil$  bits. Without loss of generality, we encode the start state with all  $1$ s and the reject state with all  $0$ s. For example, if  $Q = \{\text{start}, \text{seek1}, \text{seek0}, \text{reset}, \text{verify}, \text{accept}, \text{reject}\}$ , we might use the following encoding:

$$\begin{array}{llll} \langle \text{start} \rangle = 111 & \langle \text{seek1} \rangle = 010 & \langle \text{seek0} \rangle = 011 & \langle \text{reset} \rangle = 100 \\ \langle \text{verify} \rangle = 101 & \langle \text{accept} \rangle = 110 & \langle \text{reject} \rangle = 000 & \end{array}$$

We encode the machine  $M$  itself as the string  $\langle M \rangle = [\langle \text{reject} \rangle \bullet \langle \square \rangle] \langle \delta \rangle$ , where  $\langle \delta \rangle$  is the concatenation of substrings  $[\langle p \rangle \bullet \langle a \rangle \mid \langle q \rangle \bullet \langle b \rangle \bullet \langle \Delta \rangle]$  encoding each transition  $\delta(p, a) = (q, b, \Delta)$  such that  $q \neq \text{reject}$ . We encode the actions  $\Delta = \pm 1$  by defining  $\langle -1 \rangle := 0$  and  $\langle +1 \rangle := 1$ . Conveniently, every transition string has exactly the same length. For example, with the symbol and state encodings described above, the transition  $\delta(\text{reset}, \$) = (\text{start}, \$, +1)$  would be encoded as

$$[100 \bullet 011 \mid 001 \bullet 011 \bullet 1].$$

Our first example Turing machine for recognizing  $\{0^n 1^n 0^n \mid n \geq 0\}$  would be represented by the following string (here broken into multiple lines for readability):

$$\begin{aligned} &[000 \bullet 000] [001 \bullet 001 \mid 010 \bullet 011 \bullet 1] [001 \bullet 100 \mid 101 \bullet 011 \bullet 1] \\ &\quad [010 \bullet 001 \mid 010 \bullet 001 \bullet 1] [010 \bullet 100 \mid 010 \bullet 100 \bullet 1] \\ &\quad [010 \bullet 010 \mid 011 \bullet 100 \bullet 1] [011 \bullet 010 \mid 011 \bullet 010 \bullet 1] \\ &\quad [011 \bullet 100 \mid 011 \bullet 100 \bullet 1] [011 \bullet 001 \mid 100 \bullet 100 \bullet 1] \\ &\quad [100 \bullet 001 \mid 100 \bullet 001 \bullet 0] [100 \bullet 010 \mid 100 \bullet 010 \bullet 0] \\ &\quad [100 \bullet 100 \mid 100 \bullet 100 \bullet 0] [100 \bullet 011 \mid 001 \bullet 011 \bullet 1] \\ &\quad [101 \bullet 100 \mid 101 \bullet 011 \bullet 1] [101 \bullet 000 \mid 110 \bullet 000 \bullet 0] \end{aligned}$$

Finally, we encode any *configuration* of  $M$  on  $U$ 's work tape by alternating between encodings of states and encodings of tape symbols. Thus, each tape cell is represented by the string  $[\langle q \rangle \bullet \langle a \rangle]$  indicating that (1) the cell contains symbol  $a$ ; (2) if  $q \neq \text{reject}$ , then  $M$ 's head is located at this cell, and  $M$  is in state  $q$ ; and (3) if  $q = \text{reject}$ , then  $M$ 's head is located somewhere else. Conveniently, each cell encoding uses exactly the same number of bits. We also surround the entire tape encoding with brackets  $[$  and  $]$ .

For example, with the encodings described above, the initial configuration  $(\text{start}, 001100, 0)$  for our first example Turing machine would be encoded on  $U$ 's tape as follows.

$$\underbrace{[111 \bullet 001]}_{\text{start } 0} \underbrace{[000 \bullet 001]}_{\text{reject } 0} \underbrace{[000 \bullet 010]}_{\text{reject } 1} \underbrace{[000 \bullet 010]}_{\text{reject } 1} \underbrace{[000 \bullet 001]}_{\text{reject } 0} \underbrace{[000 \bullet 001]}_{\text{reject } 0}$$

Similarly, the intermediate configuration  $(\text{reset}, \$0x1x0, 3)$  would be encoded as follows:

$$\underbrace{[000 \bullet 011]}_{\text{reject } \$} \underbrace{[000 \bullet 011]}_{\text{reject } 0} \underbrace{[000 \bullet 100]}_{\text{reject } x} \underbrace{[010 \bullet 010]}_{\text{reset } 1} \underbrace{[000 \bullet 100]}_{\text{reject } x} \underbrace{[000 \bullet 001]}_{\text{reject } 0}$$



## Input and Execution

Without loss of generality, we assume that the input to our universal Turing machine  $U$  is given on a separate read-only **input tape**, as the encoding of an arbitrary Turing machine  $M$  followed by an encoding of its input string  $x$ . Notice the substrings `[[` and `]]` each appear only once on the input tape, immediately before and after the encoded transition table, respectively.  $U$  also has a read-write **work tape**, which is initially blank.

We start by initializing the work tape with the encoding  $\langle \text{start}, x, 0 \rangle$  of the initial configuration of  $M$  with input  $x$ . First, we write `[[ $\langle \text{start} \rangle$ •`. Then we copy the encoded input string  $\langle x \rangle$  onto the work tape, but we change the punctuation as follows:

- Instead of copying the left bracket `[`, write `[[ $\langle \text{start} \rangle$ •`.
- Instead of copying each separator `•`, write `]] $\langle \text{reject} \rangle$ •`.
- Instead of copying the right bracket `]`, write two right brackets `]]`.

The state encodings  $\langle \text{start} \rangle$  and  $\langle \text{reject} \rangle$  can be copied directly from the beginning of  $\langle M \rangle$  (replacing `0`s for `1`s for  $\langle \text{start} \rangle$ ). Finally, we move the head back to the start of  $U$ 's tape.

At the start of each step of the simulation,  $U$ 's head is located at the start of the work tape. We scan through the work tape to the unique encoded cell `[[ $\langle p \rangle$ • $\langle a \rangle$ ]]` such that  $p \neq \text{reject}$ . Then we scan through the encoded transition function  $\langle \delta \rangle$  to find the unique encoded tuple `[[ $\langle p \rangle$ • $\langle a \rangle$  |  $\langle q \rangle$ • $\langle b \rangle$ • $\langle \Delta \rangle$ ]]` whose left half matches our the encoded tape cell. If there is no such tuple, then  $U$  immediately halts and rejects. Otherwise, we copy the right half  $\langle q \rangle$ • $\langle b \rangle$  of the tuple to the work tape. Now if  $q = \text{accept}$ , then  $U$  immediately halts and accepts. (We don't bother to encode  $\text{reject}$  transformations, so we know that  $q \neq \text{reject}$ .) Otherwise, we transfer the state encoding to either the next or previous encoded cell, as indicated by  $M$ 's transition function, and then continue with the next step of the simulation.

During the final state-copying phase, we ever read two right brackets `]]`, indicating that we have reached the right end of the tape encoding, we replace the second right bracket with `[[ $\langle \text{reject} \rangle$ • $\langle \square \rangle$ ]]` (mostly copied from the beginning of the machine encoding  $\langle M \rangle$ ) and then scan back to the left bracket we just wrote. This trick allows our universal machine to *pretend* that its tape contains an infinite sequence of *encoded* blanks `[[ $\langle \text{reject} \rangle$ • $\langle \square \rangle$ ]]` instead of *actual* blanks  $\square$ .

## Example

As an illustrative example, suppose  $U$  is simulating our first example Turing machine  $M$  on the input string `001100`. The execution of  $M$  on input  $w$  eventually reaches the configuration  $(\text{seek1}, \$\$x1x0, 3)$ . At the start of the corresponding step in  $U$ 's simulation,  $U$  is in the following configuration:

`[[000•011]] [000•011] [000•100] [010•010] [000•100] [000•001]]`  
▲

First  $U$  scans for the first encoded tape cell whose state is not  $\text{reject}$ . That is,  $U$  repeatedly compares the first half of each encoded state cell on the work tape with the prefix `[[ $\langle \text{reject} \rangle$ •` of the machine encoding  $\langle M \rangle$  on the input tape.  $U$  finds a match in the fourth encoded cell.

`[[000•011]] [000•011] [000•100] [010•010] [000•100] [000•001]]`  
▲

Next,  $U$  scans the machine encoding  $\langle M \rangle$  for the substring `[010•010` matching the current encoded cell.  $U$  eventually finds a match in the left side of the the encoded transition

$[010 \bullet 010 \mid 011 \bullet 100 \bullet 1]$ .  $U$  copies the state-symbol pair  $011 \bullet 100$  from the right half of this encoded transition into the current encoded cell. (The underline indicates which symbols are changed.)

$[[000 \bullet 011][000 \bullet 011][000 \bullet 100][\underline{011 \bullet 100}][000 \bullet 100][000 \bullet 001]]$

The encoded transition instructs  $U$  to move the current state encoding one cell to the right. (The underline indicates which symbols are changed.)

$[[000 \bullet 011][000 \bullet 011][000 \bullet 100][\underline{000 \bullet 100}][\underline{011 \bullet 100}][000 \bullet 001]]$

Finally,  $U$  scans left until it reads two left brackets  $[[$ ; this returns the head to the left end of the work tape to start the next step in the simulation.  $U$ 's tape now holds the encoding of  $M$ 's configuration ( $\text{seek}0, \$\$xx\text{seek}0, 4$ ), as required.

$\uparrow [[000 \bullet 011][000 \bullet 011][000 \bullet 100][000 \bullet 100][011 \bullet 100][000 \bullet 001]]$

## Exercises

In the following problems, a *standard* Turing machine has a single semi-infinite tape, one read-write head, and the input alphabet  $\Sigma = \{0, 1\}$ . For problems that ask you to *construct* a standard Turing machine, you may assume without loss of generality that the initial tape contains a special symbol  $\blacktriangleright$  just to the left of the input string, indicating the left end of the tape; the read-write head starts just to the right of this symbol. For problems that ask you to *simulate* a standard Turing machine, you may assume without loss of generality that the tape alphabet is  $\{0, 1, \square\}$ .

## Turing Machine Programming

1. Describe standard Turing machines that decide each of the following languages:
  - (a) Palindromes over the alphabet  $\{0, 1\}$
  - (b)  $\{ww \mid w \in \{0, 1\}^*\}$
  - (c)  $\{0^a 1^b 0^{ab} \mid a, b \in \mathbb{N}\}$
2. Let  $\langle n \rangle_2$  denote the binary representation of the non-negative integer  $n$ . For example,  $\langle 17 \rangle_2 = 10001$  and  $\langle 42 \rangle_2 = 101010$ . Describe standard Turing machines that compute the following functions from  $\{0, 1\}^*$  to  $\{0, 1\}^*$ :
  - (a)  $w \mapsto www$
  - (b)  $1^n 0 1^m \mapsto 1^{mn}$
  - (c)  $1^n \mapsto 1^{2^n}$
  - (d)  $1^n \mapsto \langle n \rangle_2$
  - (e)  $0^* \langle n \rangle_2 \mapsto 1^n$
  - (f)  $\langle n \rangle_2 \mapsto \langle n^2 \rangle_2$

3. Describe standard Turing machines that write each of the following infinite streams of bits onto their tape. Specifically, for each integer  $n$ , there must be a finite time after which the first  $n$  symbols on the tape always match the first  $n$  symbols in the target stream.

(a) An infinite stream of **1**s

(b) **0101101110111101111101111110** ..., where the  $n$ th block of **1**s has length  $n$ .

(c) The stream of bits whose  $n$ th bit is **1** if and only if  $n$  is prime.

(d) The *Thue-Morse sequence*  $T_0 \cdot T_1 \cdot T_2 \cdot T_3 \cdots$ , where

$$T_n := \begin{cases} \mathbf{0} & \text{if } n = 0 \\ \mathbf{1} & \text{if } n = 1 \\ T_{n-1} \cdot \overline{T_{n-1}} & \text{otherwise} \end{cases}$$

where  $\overline{w}$  indicates the binary string obtained from  $w$  by flipping every bit. Equivalently, the  $n$ th bit of the Thue Morse sequence is **0** if the binary representation of  $n$  has an even number of **1**s and **1** otherwise.

**01101001100101101001011001101001100101100110010110011010010110** ...

(e) The *Fibonacci sequence*  $F_0 \cdot F_1 \cdot F_2 \cdot F_3 \cdots$ , where

$$F_n := \begin{cases} \mathbf{0} & \text{if } n = 0 \\ \mathbf{1} & \text{if } n = 1 \\ F_{n-2} \cdot F_{n-1} & \text{otherwise} \end{cases}$$

**0101101011011010110101101101011011010110101101101011010101** ...

## Simulation by “Weaker” Machines

4. A *two-stack machine* is a Turing machine with two tapes with the following restricted behavior. At all times, on each tape, every cell to the right of the head is blank, and every cell at or to the left of the head is non-blank. Thus, a head can only move right by writing a non-blank symbol into a blank cell; symmetrically, a head can only move left by erasing the rightmost non-blank cell. Thus, each tape behaves like a stack. To avoid underflow, there is a special symbol at the start of each tape that cannot be overwritten. Initially, one tape contains the input string, with the head at its *last* symbol, and the other tape is empty (except for the start-of-tape symbol).

Prove that any standard Turing machine can be simulated by a two-stack machine. That is, given any standard Turing machine  $M$ , describe a two-stack machine  $M'$  that accepts and rejects exactly the same input strings as  $M$ .

5. A *k-register machine* is a finite-state automaton with  $k$  non-negative integer registers. Formally, a  $k$ -register machine consists of a finite set  $Q$  of states (which include **start**, **accept**, and **reject**) and a transition function

$$\delta : Q \times \{0, 1\}^k \rightarrow Q \times \{\text{halve}, \text{nop}, \text{double}, \text{double}+1\}^k$$

that takes the internal state and the *signs* of the registers as input, and produces a new internal state and *changes* to the registers as output. The instructions **halve**, **nop**, **double**, and **double+1** change any register with value  $n$  to  $n/2$ ,  $n$ ,  $2n$ , and  $2n + 1$ , respectively.

For example, if  $\delta(p, 0, 1, 0, 1) = (q, \text{halve}, \text{nop}, \text{halve}, \text{double}+1)$ , then from the configuration  $(p, 0, 2, 3, 1)$ , the machine would transition to  $(q, 1, 2, 1, 3)$ .

Prove that any standard Turing machine (with suitably encoded input and output) can be simulated by a two-register machine. The input to the register machine is encoded in reversed binary in one of the registers, so the parity of the register value is the first input bit; the other register is initially zero.

6. A ***k*-counter machine** (also known as a ***Minsky machine***) is a finite-state automaton with  $k$  non-negative integer registers. Formally, a  $k$ -counter machine consists of a finite set  $Q$  of states (which include **start**, **accept**, and **reject**) and a transition function

$$\delta: Q \times \{0, +\}^k \rightarrow Q \times \{\text{inc, nop, dec}\}^k$$

that takes the internal state and the *signs* of the registers as input, and produces a new internal state and changes to the registers as output. The instructions **inc**, **nop**, and **dec** change any register with value  $n$  to  $n + 1$ ,  $n$ , and  $n - 1$ , respectively. The transition function must forbid decrementing a register whose value is already zero.

For example, if  $\delta(p, 0, +, +, +) = (q, \text{inc}, \text{dec}, \text{nop}, \text{dec})$ , then from the configuration  $(p, 0, 2, 3, 1)$ , the machine would transition to  $(q, 1, 1, 3, 0)$ .

- (a) Prove that any standard Turing machine (with suitably encoded input and output) can be simulated by a three-counter machine. [*Hint: Simulate a two-**register** machine, using the third counter for scratch work.*]
- (b) Prove that any three-counter machine (with suitably encoded input and output) can be simulated by a two-counter machine. [*Hint: Store all three counters in a single integer of the form  $2^a 3^b 5^c$ , and use the other counter for scratch work.*]
- \* (c) Prove that a three-counter machine can compute a suitable encoding of any computable function. Specifically, for any computable function  $f : \mathbb{N} \rightarrow \mathbb{N}$ , prove there is a three-counter machine  $M$  that transforms any input  $(n, 0, 0)$  into  $(f(n), 0, 0)$ . [*Hint: First transform  $(n, 0, 0)$  to  $(2^n, 0, 0)$  using all three counters; then run a two- (or three-) counter TM simulation to obtain  $(2^{f(n)}, 0, 0)$ ; and finally transform  $(2^{f(n)}, 0, 0)$  to  $(f(n), 0, 0)$  using all three counters.*]
- ★ (d) Prove that not two-counter machine can transform  $(n, 0)$  to  $(2^n, 0)$ . This impossibility result was independently proved by Bärzdīņš in 1963, Yao in 1971, and Schroeppel in 1972. <sup>5</sup>

7. A **hole-punch** Turing machine is a standard Turing machine with two restrictions. First, the tape alphabet has only two symbols  $\square$  and  $\blacksquare$ , and thus the input alphabet is the singleton

<sup>5</sup>Ja. M. Barzdin' [Jānis Barzdiņš]. Ob odnom klasse mašin T'uringa (mašiny Minskogo) [On a class of Turing machines (Minsky machines)]. *Algebra i Logika* 1(6):42–51, 1963. In Russian. Sorry.

Oscar H. Ibarra, Nicholas Q. Tr  n. A note on simple programs with two variables. *Theoretical Computer Science* 112(2): 391–397, 1993.

Rich Schroepel. A two counter machine cannot calculate  $2^N$ . Artificial Intelligence Memo 257, MIT AI Lab, May 1972. [Schroepel claims that the same result was independently proved by Frances Yao in 1971.

set  $\{\blacksquare\}$ . Second, the machine can never write a blank ( $\square$ ) over a non-blank ( $\blacksquare$ ); intuitively, the machine can punch new holes ( $\blacksquare$ s) into the tape, but it cannot erase holes.

Prove that any standard Turing machine (with a unary input alphabet) can be simulated by a hole-punch Turing machine.

8. A **tag**-Turing machine has two heads: one can only read, the other can only write. Initially, the read head is located at the left end of the tape, and the write head is located at the first blank after the input string. At each transition, the read head can either move one cell to the right or stay put, but the write head *must* write a symbol to its current cell and move one cell to the right. Neither head can ever move to the left.

Prove that any standard Turing machine can be simulated by a tag-Turing machine. That is, given any standard Turing machine  $M$ , describe a tag-Turing machine  $M'$  that accepts and rejects exactly the same input strings as  $M$ .

9. \* (a) Prove that any standard Turing machine can be simulated by a standard Turing machine with only three states. [*Hint: Keep an encoding of the state of the simulated machine on the tape of the simulating machine.*]
- ★ (b) Prove that any standard Turing machine can be simulated by a standard Turing machine with only two states.

## Simulating “Stronger” Machines

10. A **two-dimensional** Turing machine uses an infinite two-dimensional grid of cells as the tape; at each transition, the head can move from its current cell to any of its four neighbors on the grid. The transition function of such a machine has the form  $\delta: Q \times \Gamma \rightarrow Q \times \Gamma \times \{\uparrow, \leftarrow, \downarrow, \rightarrow\}$ , where the arrows indicate which direction the head should move.
  - (a) Prove that any two-dimensional Turing machine can be simulated by a standard Turing machine.
  - (b) Suppose further that we endow our two-dimensional Turing machine with the following additional actions, in addition to moving the head:
    - Insert row: Move all symbols on or above the row containing the head up one row, leaving the head's row blank.
    - Insert column: Move all symbols on or to the right of the column containing the head one column to the right, leaving the head's column blank.
    - Delete row: Move all symbols above the row containing the head down one row, deleting the head's row of symbols.
    - Delete column: Move all symbols the right of the column containing the head one column to the right, deleting the head's column of symbols.

Show that any two-dimensional Turing machine that can add and delete rows can be simulated by a standard Turing machine.

11. A **binary-tree** Turing machine uses an infinite binary tree as its tape; that is, *every* cell in the tape has a left child and a right child. At each step, the head moves from its current

cell to its **P**arent, its **L**eft child, or to its **R**ight child. Thus, the transition function of such a machine has the form  $\delta : Q \times \Gamma \rightarrow Q \times \Gamma \times \{\mathbf{P}, \mathbf{L}, \mathbf{R}\}$ . The input string is initially given along the left spine of the tape.

Show that any binary-tree Turing machine can be simulated by a standard Turing machine.

12. A **stack-tape** Turing machine uses an semi-infinite tape, where every cell is actually the top of an independent stack. The behavior of the machine at each iteration is governed by its internal state and the symbol *at the top* of the current cell's stack. At each transition, the head can optionally push a new symbol onto the stack, or pop the top symbol off the stack. (If a stack is empty, its "top symbol" is a blank and popping has no effect.)

Show that any stack-tape Turing machine can be simulated by a standard Turing machine. (Compare with Problem 4!)

13. A **tape-stack** Turing machine has two actions that modify its work tape, in addition to simply writing individual cells: it can **save** the entire tape by pushing in onto a stack, and it can **restore** the entire tape by popping it off the stack. Restoring a tape returns the content of every cell to its content when the tape was saved. Saving and restoring the tape do not change the machine's state or the position of its head. If the machine attempts to "restore" the tape when the stack is empty, the machine crashes.

Show that any tape-stack Turing machine can be simulated by a standard Turing machine.

Rewrite in the language of algorithms instead of the language of Turing machines, using “source code” instead of “encoding” everywhere. Formulation in terms of TMs makes almost everything much more complicated than it needs to be. (The dovetail/product construction in the proof of Lemma 4 may be an exception.)

## 7 Undecidability

Perhaps the single most important result in Turing’s remarkable 1936 paper that introduces Turing machines is his observation that there are problems that cannot be solved by *any* algorithm. Turing’s canonical example of an undecidable problem was the **halting problem**, which asks whether a given Turing machine halts when given a particular input string. Among other consequences, Turing’s undecidability result provided an elegant negative solution to Hilbert’s *Entscheidungsproblem*, which asked for an algorithm to decide whether a given statement of first-order logic is true—no such algorithm exists.

### 7.1 Acceptable versus Decidable

Recall that there are three possible outcomes for a Turing machine  $M$  running on any particular input string  $w$ : acceptance, rejection, and divergence. Every Turing machine  $M$  immediately defines four different languages (over the input alphabet  $\Sigma$  of  $M$ ):

- The *accepting* language  $\text{ACCEPT}(M) := \{w \in \Sigma^* \mid M \text{ accepts } w\}$
- The *rejecting* language  $\text{REJECT}(M) := \{w \in \Sigma^* \mid M \text{ rejects } w\}$
- The *halting* language  $\text{HALT}(M) := \text{ACCEPT}(M) \cup \text{REJECT}(M)$
- The *diverging* language  $\text{DIVERGE}(M) := \Sigma^* \setminus \text{HALT}(M)$

For any language  $L$ , the sentence “ **$M$  accepts  $L$** ” means  $\text{ACCEPT}(M) = L$ , and the sentence “ **$M$  decides  $L$** ” means  $\text{ACCEPT}(M) = L$  and  $\text{DIVERGE}(M) = \emptyset$ .

Now let  $L$  be an arbitrary language. We say that  $L$  is **acceptable** (or *semi-computable*, or *semi-decidable*, or *recognizable*, or *listable*, or *recursively enumerable*) if some Turing machine accepts  $L$ , and **unacceptable** otherwise. Similarly,  $L$  is **decidable** (or *computable*, or *recursive*) if some Turing machine decides  $L$ , and **undecidable** otherwise.

## 7.2 Lo, I Have Become Death, Stealer of Pie

There is a subtlety in the definitions of “acceptable” and “decidable” that many beginners miss: A language can be decidable even if we can’t exhibit a specific Turing machine that decides it. As a canonical example, consider the language  $\Pi = \{w \mid 1^{|w|} \text{ appears in the binary expansion of } \pi\}$ . Despite appearances, this language is decidable! There are only two cases to consider:

- Suppose there is an integer  $N$  such that the binary expansion of  $\pi$  contains the substring  $1^N$  but does not contain the substring  $1^{N+1}$ . Let  $M_N$  be the Turing machine with  $N + 3$  states  $\{0, 1, \dots, N, \text{accept}, \text{reject}\}$ , start state 0, and the following transition function:

$$\delta(q, a) = \begin{cases} \text{accept} & \text{if } a = \square \\ \text{reject} & \text{if } a \neq \square \text{ and } q = n \\ (q + 1, a, +1) & \text{otherwise} \end{cases}$$

This machine correctly decides  $\Pi$ .

- Suppose the binary expansion of  $\pi$  contains arbitrarily long substrings of  $1$ s. Then any Turing machine that accepts all inputs correctly decides  $\Pi$ .

We have no idea which of these machines correctly decides  $\Pi$ , but one of them does, and that’s enough!

## 7.3 Useful Properties

This subsection lists several simple but useful properties of (un)decidable and (un)acceptable languages. Almost all of these properties follow from routine definition-chasing; readers are strongly encouraged to try to prove each lemma themselves before reading ahead.

One might reasonably ask why we don’t also define “rejectable” and “halttable” languages. The following lemma, whose proof is an easy exercise (hint, hint), implies that these sets are both identical to the acceptable languages.

**Lemma 1.** *Let  $M$  be an arbitrary Turing machine.*

- There is a Turing machine  $M^R$  such that  $\text{ACCEPT}(M^R) = \text{REJECT}(M)$  and  $\text{REJECT}(M^R) = \text{ACCEPT}(M)$ .*
- There is a Turing machine  $M^A$  such that  $\text{ACCEPT}(M^A) = \text{ACCEPT}(M)$  and  $\text{REJECT}(M^A) = \emptyset$ .*
- There is a Turing machine  $M^H$  such that  $\text{ACCEPT}(M^H) = \text{HALT}(M)$  and  $\text{REJECT}(M^H) = \emptyset$ .*

The decidable languages have several fairly obvious useful closure properties.

**Lemma 2.** *If  $L$  and  $L'$  are decidable, then  $L \cup L'$ ,  $L \cap L'$ ,  $L \setminus L'$ , and  $L' \setminus L$  are also decidable.*

**Proof:** Let  $M$  and  $M'$  be Turing machines that decide  $L$  and  $L'$ , respectively. We can build a Turing machine  $M_\cup$  that decides  $L \cup L'$  as follows. First,  $M_\cup$  copies its input string  $w$  onto a second tape. Then  $M_\cup$  runs  $M$  on input  $w$  (on the first tape), and then runs  $M'$  on input  $w$  (on the second tape). If either  $M$  or  $M'$  accepts, then  $M_\cup$  accepts; if both  $M$  and  $M'$  reject, then  $M_\cup$  rejects.

The other three languages are similar. □



**Corollary 3.** *The following hold for all languages  $L$  and  $L'$ .*

- (a) *If  $L \cap L'$  is undecidable and  $L'$  is decidable, then  $L$  is undecidable.*
- (b) *If  $L \cup L'$  is undecidable and  $L'$  is decidable, then  $L$  is undecidable.*
- (c) *If  $L \setminus L'$  is undecidable and  $L'$  is decidable, then  $L$  is undecidable.*
- (d) *If  $L' \setminus L$  is undecidable and  $L'$  is decidable, then  $L$  is undecidable.*

Unfortunately, acceptable languages are not quite as well-behaved as decidable languages, thanks to the subtle distinction between *rejecting* a string and *not accepting* a string.

**Lemma 4.** *For all acceptable languages  $L$  and  $L'$ , the languages  $L \cup L'$  and  $L \cap L'$  are also acceptable.*

**Proof:** Let  $M$  and  $M'$  be Turing machines that decide  $L$  and  $L'$ , respectively. We can build a Turing machine  $M_\cap$  that decides  $L \cap L'$  as follows. First,  $M_\cap$  copies its input string  $w$  onto a second tape. Then  $M_\cap$  runs  $M$  on input  $w$  using the first tape, and then runs  $M'$  on input  $w$  using the second tape. If both  $M$  and  $M'$  accept, then  $M_\cap$  accepts; if either  $M$  or  $M'$  reject, then  $M_\cap$  rejects; if either  $M$  or  $M'$  diverge, then  $M_\cap$  diverges (automatically).

The construction for  $L \cup L'$  is more subtle; instead of running  $M$  and  $M'$  in series, we must run them in parallel. Like  $M_\cap$ , the new machine  $M_\cup$  starts by copying its input string  $w$  onto a second tape. But then  $M_\cup$  runs  $M$  and  $M'$  simultaneously; with each step of  $M_\cup$  simulating both one step of  $M$  on the first tape and one step of  $M'$  on the second. Ignoring the states and transitions needed for initialization, the state set of  $M_\cup$  is the product of the state sets of  $M$  and  $M'$ , and the transition function is

$$\delta_\cup(q, a, q', a') = \begin{cases} \text{accept}_\cup & \text{if } q = \text{accept} \text{ or } q' = \text{accept}' \\ \text{reject}_\cup & \text{if } q = \text{reject} \text{ and } q' = \text{reject}' \\ (\delta(q, a), \delta'(q', a')) & \text{otherwise} \end{cases}$$

Thus,  $M_\cup$  accepts as soon as either  $M$  or  $M'$  accepts, and rejects only after both  $M$  or  $M'$  reject. □

**Lemma 5.** *An acceptable language  $L$  is decidable if and only if  $\Sigma^* \setminus L$  is also acceptable.*

**Proof:** Let  $M$  and  $\overline{M}$  be Turing machines that accept  $L$  and  $\Sigma^* \setminus L$ , respectively. Following the previous proof, we construct a new Turing machine  $M^*$  that copies its input onto a second tape, and then simulates  $M$  and  $\overline{M}$  in parallel on the two tapes. If  $M$  accepts, then  $M^*$  accepts; if  $\overline{M}$  accepts, then  $M^*$  rejects. Since every string is accepted by either  $M$  or  $\overline{M}$ , we conclude that  $M^*$  decides  $L$ .

The other direction follows immediately from Lemma 1. □

## 7.4 Code is Data; Data is Code

Perhaps the single most important observation in developing these undecidability results—and one of the most important observations in computer science more broadly—is that Turing machines can be encoded as strings. At one level, this observation is completely trivial: Any written description of a Turing machine is a string, and modern code is just a sequence of bytes, stored in a file like any other data. But this apparently trivial observation is actually incredibly powerful.

Most natural encodings of Turing machines have three important properties.

- **Unique:** Different Turing machines are encoded as different strings.
- **Modifiable:** We can algorithmically modify any Turing machine  $M$ , given the encoding of  $M$  as input. For example, there are algorithms to swap the **accept** and **reject** states of any Turing machine, or to add new states and transitions representing pre- and post-processing phases, or to build a new machine that calls  $M$  as a subroutine, or to build a new machine that runs several copies of  $M$  in parallel.
- **Executable:** There is a fixed **universal** Turing machine  $U$  that can simulate the behavior of an arbitrary Turing machine  $M$ , given the encodings of  $M$  and  $w$  as input. For example, if we decided to encode Turing machines as Python programs, then  $U$  would be a Python interpreter.

The precise details of the encoding are unimportant, but for the sake of concreteness, let me describe a natural encoding of Turing machines as strings over the six-character alphabet  $\{0, 1, \{, \cdot, \$, \square\}$ . Let  $M = (\Gamma, \square, \Sigma, Q, \text{start}, \text{accept}, \text{reject}, \delta)$  be an arbitrary Turing machine, with a single half-infinite tape and a single read-write head. (I will consistently indicate the states and tape symbols of  $M$  in *slanted green* to distinguish them from the *upright red* symbols in the encoding alphabet.)

- We encode each symbol  $a \in \Gamma$  as a unique string  $\langle a \rangle$  of  $\lceil \lg(|\Gamma|) \rceil$  bits. For example, if  $\Gamma = \{0, 1, \$, x, \square\}$ , we might use the following encoding:

$$\langle 0 \rangle = 001, \quad \langle 1 \rangle = 010, \quad \langle \$ \rangle = 011, \quad \langle x \rangle = 100, \quad \langle \square \rangle = 000.$$

- Similarly, we encode each state  $q \in Q$  as a distinct string  $\langle q \rangle$  of  $\lceil \lg|Q| \rceil$  bits. Without loss of generality, we encode the start state with all **1s** and the reject state with all **0s**. For example, if  $Q = \{\text{start}, \text{seek1}, \text{seek0}, \text{reset}, \text{verify}, \text{accept}, \text{reject}\}$ , we might use the following encoding:

$$\begin{array}{llll} \langle \text{start} \rangle = 111 & \langle \text{seek1} \rangle = 010 & \langle \text{seek0} \rangle = 011 & \langle \text{reset} \rangle = 100 \\ \langle \text{verify} \rangle = 101 & \langle \text{accept} \rangle = 110 & \langle \text{reject} \rangle = 000 & \end{array}$$

- Finally, we encode the machine  $M$  itself as the string  $\langle M \rangle = \{\langle \text{reject} \rangle \cdot \langle \square \rangle\} \langle \delta \rangle$ , where  $\langle \delta \rangle$  is the concatenation of substrings  $\{\langle p \rangle \cdot \langle a \rangle \cdot \langle q \rangle \cdot \langle b \rangle \cdot \langle \Delta \rangle\}$  encoding each transition  $\delta(p, a) = (q, b, \Delta)$  such that  $q \neq \text{reject}$ . We encode the actions  $\Delta = \pm 1$  by defining  $\langle -1 \rangle := 0$  and  $\langle +1 \rangle := 1$ . Conveniently, every transition string has exactly the same length. For example, with the symbol and state encodings described above, the transition  $\delta(\text{reset}, \$) = (\text{start}, \$, +1)$  would be encoded as the string

$$\{100 \cdot 011 \cdot 001 \cdot 011 \cdot 1\}.$$

Our first example Turing machine for recognizing  $\{0^n 1^n 0^n \mid n \geq 0\}$  would be represented by the following string (broken into multiple lines for readability):

```
{000•000}{001•001•010•011•1}{001•100•101•011•1}{010•001•010•001•1}
{010•100•010•100•1}{010•010•011•100•1}{011•010•011•010•1}
{011•100•011•100•1}{011•001•100•100•1}{100•001•100•001•0}
{100•010•100•010•0}{100•100•100•100•0}{100•011•001•011•1}
{101•100•101•011•1}{101•000•110•000•0}
```

Building a universal Turing machine  $U$  that uses this encoding is more a matter of careful bookkeeping than real insight. We can encode any *configuration* of  $M$  on  $U$ 's work tape by encoding each cell of  $M$ 's tape as a string  $\{\langle q \rangle \bullet \langle a \rangle\}$  indicating that (1) the cell contains symbol  $a$ ; (2) if  $q \neq \text{reject}$ , then  $M$ 's head is located at this cell, and  $M$  is in state  $q$ ; and (3) if  $q = \text{reject}$ , then  $M$ 's head is located somewhere else. We also surround the entire tape encoding with brackets  $\{$  and  $\}$ . For example, the initial configuration  $(\text{start}, \text{00110}, 0)$  for our example Turing machine would be encoded as follows.

$$\langle \text{start}, \text{00110}, 0 \rangle = \underbrace{\{\{111 \bullet 001\}\}}_{\text{start } 0} \underbrace{\{\{000 \bullet 001\}\}}_{\text{reject } 0} \underbrace{\{\{000 \bullet 010\}\}}_{\text{reject } 1} \underbrace{\{\{000 \bullet 010\}\}}_{\text{reject } 1} \underbrace{\{\{000 \bullet 001\}\}}_{\text{reject } 0}$$

Similarly, the intermediate configuration  $(\text{reset}, \$0x1x, 3)$  would be encoded as follows:

$$\langle \text{reset}, \$\$x1x, 3 \rangle = \underbrace{\{\{000 \bullet 011\}\}}_{\text{reject } \$} \underbrace{\{\{000 \bullet 011\}\}}_{\text{reject } 0} \underbrace{\{\{000 \bullet 100\}\}}_{\text{reject } x} \underbrace{\{\{010 \bullet 010\}\}}_{\text{reset } 1} \underbrace{\{\{000 \bullet 100\}\}}_{\text{reject } x}$$

To simulate one step of  $M$ 's execution, we (1) find the location of the head (or reject if the head has vanished), (2) look up the transition for the state-symbol pair at the head, and (3) update the current cell and one of its neighbors to reflect the transition. The remaining grungy details are left as an exercise.

## 7.5 Self-Haters Gonna Self-Hate

A Turing machine encoding  $\langle M \rangle$  is just a string, and any string (over the correct alphabet) can be used as the input to a Turing machine. Thus, a suitable encoding of *any* Turing machine can be used as the input to *any* Turing machine. In particular:

**The encoding  $\langle M \rangle$  of Turing machine  $M$   
can be used as input to *the same* Turing machine  $M$ .**

Turing used this observation about self-reference to derive his first undecidable language as follows. Let's say that a Turing machine  $M$  is **self-rejecting** if it rejects its own encoding  $\langle M \rangle$ . Let **SELFREJECT** be the set of all encodings of self-rejecting Turing machines:

$$\text{SELFREJECT} := \{ \langle M \rangle \mid M \text{ rejects } \langle M \rangle \}$$

**Theorem 6.** *SELFREJECT is undecidable.*

**Proof:** Suppose to the contrary that there is a Turing machine  $SR$  that decides **SELFREJECT**. Then by definition,  $\text{ACCEPT}(SR) = \text{SELFREJECT}$  and  $\text{DIVERGE}(SR) = \emptyset$ . More explicitly, for *any* Turing machine  $M$ ,

- $SR$  accepts  $\langle M \rangle \iff M$  rejects  $\langle M \rangle$ , and
- $SR$  rejects  $\langle M \rangle \iff M$  does not reject  $\langle M \rangle$ .

In particular, these equivalences must hold when  $M$  is the machine  $SR$ . Thus,

- $SR$  accepts  $\langle SR \rangle \iff SR$  rejects  $\langle SR \rangle$ , and
- $SR$  rejects  $\langle SR \rangle \iff SR$  does not reject  $\langle SR \rangle$ .

In short,  $SR$  accepts  $\langle SR \rangle$  if and only if  $SR$  rejects  $\langle SR \rangle$ , which is impossible! The only logical conclusion is that the Turing machine  $SR$  does not exist. □

## 7.6 Aside: Uncountable Barbers

Turing’s proof by contradiction is an avatar of the famous *diagonalization argument* that uncountable sets exist, published by Georg Cantor in 1891. Indeed, SELFREJECT is sometimes called “the diagonal language”. Recall that a function  $f : A \rightarrow B$  is a **surjection**<sup>1</sup> if  $f(A) = \{f(a) \mid a \in A\} = B$ .

**Cantor’s Theorem.** *Let  $f : X \rightarrow 2^X$  be an **arbitrary** function from an **arbitrary** set  $X$  to its power set. This function  $f$  is not a surjection.*

**Proof:** Fix an arbitrary function  $f : X \rightarrow 2^X$ . Call an element  $x \in X$  **happy** if  $x \in f(x)$  and **sad** if  $x \notin f(x)$ . Let  $Y$  be the set of all sad elements of  $X$ ; that is, for every element  $x \in X$ , we have

$$x \in Y \iff x \notin f(x).$$

For the sake of argument, suppose  $f$  is a surjection. Then (by definition of surjection) there must be an element  $y \in X$  such that  $f(y) = Y$ . Then for every element  $x \in X$ , we have

$$x \in f(y) \iff x \notin f(x).$$

In particular, the previous equivalence must hold when  $x = y$ :

$$y \in f(y) \iff y \notin f(y).$$

We have a contradiction! We conclude that  $f$  is not a surjection after all. □

Now let  $X = \Sigma^*$ , and define the function  $f : X \rightarrow 2^X$  as follows:

$$f(w) := \begin{cases} \text{ACCEPT}(M) & \text{if } w = \langle M \rangle \text{ for some Turing machine } M \\ \emptyset & \text{if } w \text{ is not the encoding of a Turing machine} \end{cases}$$

Cantor’s theorem immediately implies that not all languages are acceptable.

Alternatively, let  $X$  be the set of all Turing machines that halt on all inputs. For any Turing machine  $M \in X$ , let  $f(M)$  be the set of all Turing machines  $N \in X$  such that  $M$  accepts the encoding  $\langle N \rangle$ . Then a Turing machine  $M$  is **sad** if it rejects its own encoding  $\langle M \rangle$ ; thus,  $Y$  is essentially the set SELFREJECT. Cantor’s argument now immediately implies that no Turing machine decides the language SELFREJECT.

The core of Cantor’s diagonalization argument also appears in the “barber paradox” popularized by Bertrand Russell in the 1910s. In a certain small town, every resident has a haircut on Haircut Day. Some residents cut their own hair; others have their hair cut by another resident of the same town. To obtain an official barber’s license, a resident must cut the hair of all residents who don’t cut their own hair, and no one else. Given these assumptions, we can immediately conclude that there are no licensed barbers. After all, who would cut the barber’s hair?

To map Russell’s barber paradox back to Cantor’s theorem, let  $X$  be the set of residents, and let  $f(x)$  be the set of residents who have their hair cut by  $x$ ; then a resident is **sad** if they do not cut their own hair. To prove that SELFREJECT is undecidable, replace “resident” with “a Turing machine that halts on all inputs”, and replace “ $A$  cuts  $B$ ’s hair” with “ $A$  accepts  $\langle B \rangle$ ”.

---

<sup>1</sup>more commonly, flouting all reasonable standards of grammatical English, “an onto function”

## 7.7 Just Don't Know What to Do with Myself

Similar diagonal arguments imply that three other languages are also undecidable:

$$\begin{aligned}\text{SELFACCEPT} &:= \{ \langle M \rangle \mid M \text{ accepts } \langle M \rangle \} \\ \text{SELFHALT} &:= \{ \langle M \rangle \mid M \text{ halts on } \langle M \rangle \} \\ \text{SELF DIVERGE} &:= \{ \langle M \rangle \mid M \text{ diverges on } \langle M \rangle \}\end{aligned}$$

The proofs for these three languages are not quite as direct as the proof for **SELFREJECT**; each fictional deciding machine requires a small modification to create the contradiction.

**Theorem 7.** *SELFACCEPT is undecidable.*

**Proof:** For the sake of argument, suppose there is a Turing machine  $SA$  such that  $\text{ACCEPT}(SA) = \text{SELFACCEPT}$  and  $\text{DIVERGE}(M) = \emptyset$ . Let  $SA^R$  be the Turing machine obtained from  $SA$  by swapping its **accept** and **reject** states (as in the proof of Lemma 1). Then  $\text{REJECT}(SA^R) = \text{SELFACCEPT}$  and  $\text{DIVERGE}(SA^R) = \emptyset$ . It follows that  $SA^R$  rejects  $\langle SA^R \rangle$  if and only if  $SA^R$  accepts  $\langle SA^R \rangle$ , which is impossible.  $\square$

**Theorem 8.** *SELFHALT is undecidable.*

**Proof:** Suppose to the contrary that there is a Turing machine  $SH$  such that  $\text{ACCEPT}(SH) = \text{SELFHALT}$  and  $\text{DIVERGE}(SH) = \emptyset$ . Let  $SH^X$  be the Turing machine obtained from  $SH$  by redirecting every transition to **accept** to a new hanging state **hang**, and then redirecting every transition to **reject** to **accept**. Then  $\text{ACCEPT}(SH^X) = \Sigma^* \setminus \text{SELFHALT}$  and  $\text{REJECT}(SH^X) = \emptyset$ . It follows that  $SH^X$  accepts  $\langle SH^X \rangle$  if and only if  $SH^X$  does not halt on  $\langle SH^X \rangle$ , and we have a contradiction.  $\square$

**Theorem 9.** *SELF DIVERGE is unacceptable and therefore undecidable.*

**Proof:** Suppose to the contrary that there is a Turing machine  $SD$  such that  $\text{ACCEPT}(M) = \text{SELF DIVERGE}$ . Let  $SD^A$  be the Turing machine obtained from  $M$  by redirecting every transition to **reject** to a new hanging state **hang** such that  $\delta(\text{hang}, a) = (\text{hang}, a, +1)$  for every symbol  $a$ . Then  $\text{ACCEPT}(SD^A) = \text{SELF DIVERGE}$  and  $\text{REJECT}(SD^A) = \emptyset$ . It follows that  $SD^A$  accepts  $\langle SD^A \rangle$  if and only if  $SD^A$  does not halt on  $\langle SD^A \rangle$ , which is impossible.  $\square$

## \*7.8 Nevertheless, Acceptable

Our undecidability argument for **SELF DIVERGE** actually implies the stronger result that **SELF DIVERGE** is unacceptable; we never assumed that the hypothetical accepting machine  $SD$  halts on all inputs. However, we can use or modify our universal Turing machine  $U$  to *accept* the other three self-referential languages.

**Theorem 10.** *SELFACCEPT is acceptable.*

**Proof:** We describe a Turing machine  $SA$  that accepts the language **SELFACCEPT**. Given any string  $w$  as input,  $SA$  first verifies that  $w$  is the encoding of a Turing machine. If  $w$  is not the encoding of a Turing machine, then  $SA$  diverges. Otherwise,  $w = \langle M \rangle$  for some Turing machine  $M$ ; in this case,  $SA$  writes the string  $ww = \langle M \rangle \langle M \rangle$  onto its tape and passes control to the universal Turing machine  $U$ .  $U$  then simulates  $M$  (the machine encoded by the first half of

its input) on the string  $\langle M \rangle$  (the second half of its input).<sup>2</sup> In particular,  $U$  accepts  $\langle M, M \rangle$  if and only if  $M$  accepts  $\langle M \rangle$ . We conclude that  $SR$  accepts  $\langle M \rangle$  if and only if  $M$  accepts  $\langle M \rangle$ .  $\square$

**Theorem 11.** *SELFREJECT is acceptable.*

**Proof:** Let  $U^R$  be the Turing machine obtained from our universal machine  $U$  by swapping the **accept** and **reject** states. We describe a Turing machine  $SR$  that accepts the language SELFREJECT as follows.  $SR$  first verifies that its input string  $w$  is the encoding of a Turing machine and diverges if not. Otherwise,  $SR$  writes the string  $ww = \langle M, M \rangle$  onto its tape and passes control to the reversed universal Turing machine  $U^R$ . Then  $U^R$  accepts  $\langle M, M \rangle$  if and only if  $M$  rejects  $\langle M \rangle$ . We conclude that  $SR$  accepts  $\langle M \rangle$  if and only if  $M$  rejects  $\langle M \rangle$ .  $\square$

Finally, because SELFHALT is the union of two acceptable languages, SELFHALT is also acceptable.

## 7.9 The Halting Problem via Reduction

Now consider the following related languages:<sup>3</sup>

$$\begin{aligned}\text{ACCEPT} &:= \{ \langle M, w \rangle \mid M \text{ accepts } w \} \\ \text{REJECT} &:= \{ \langle M, w \rangle \mid M \text{ rejects } w \} \\ \text{HALT} &:= \{ \langle M, w \rangle \mid M \text{ halts on } w \} \\ \text{DIVERGE} &:= \{ \langle M, w \rangle \mid M \text{ diverges on } w \}\end{aligned}$$

Deciding the language HALT is usually called the **halting problem**: Given a program  $M$  and an input  $w$  to that program, does the program halt? This problem may seem trivial; why not just run the program and see? More formally, why not just pass the input string  $\langle M, x \rangle$  to our universal Turing machine  $U$ ? That strategy works perfectly if we just want to *accept* HALT, but we actually want to *decide* HALT; if  $M$  is not going to halt on  $w$ , we still want an answer in a finite amount of time. Sadly, we can't always get what we want.

**Theorem 12.** *HALT is undecidable.*

**Proof:** Suppose to the contrary that there is a Turing machine  $H$  that decides HALT. Then we can use  $H$  to build another Turing machine  $SH$  that decides the language SELFHALT. Given any string  $w$ , the machine  $SH$  first verifies that  $w = \langle M \rangle$  for some Turing machine  $M$  (rejecting if not), then writes the string  $ww = \langle M, M \rangle$  onto the tape, and finally passes control to  $H$ . But SELFHALT is undecidable, so no such machine  $SH$  exists. We conclude that  $H$  does not exist either.  $\square$

Nearly identical arguments imply that the languages ACCEPT, REJECT, and DIVERGE are undecidable.

---

<sup>2</sup>To simplify the presentation, I am implicitly assuming here that  $\langle M \rangle = \langle \langle M \rangle \rangle$ . Without this assumption, we need a Turing machine that transforms an arbitrary string  $w \in \Sigma_M^*$  into its encoding  $\langle w \rangle \in \Sigma_U^*$ ; building such a Turing machine is straightforward.

<sup>3</sup>Many sources including Sipser and Wikipedia uses the shorter name  $A_{TM}$  instead of ACCEPT, but uses  $HALT_{TM}$  instead of HALT. I have no idea why Sipser thought four-letter names are okay, but six-letter names are not. The subscript TM is just a reminder that these are languages of *Turing machine* encodings, as opposed to encodings of DFAs or some other machine model.

Here we have our first example of an undecidability proof by **reduction**. Specifically, we **reduced** the language SELFHALT to the language HALT. More generally, to reduce one language  $X$  to another language  $Y$ , we assume (for the sake of argument) that there is a program  $P_Y$  that decides  $Y$ , and we write another program that decides  $X$ , using  $P_Y$  as a black-box subroutine. If later we discover that  $Y$  is decidable, we can immediately conclude that  $X$  is decidable. Equivalently, if we later discover that  $X$  is undecidable, we can immediately conclude that  $Y$  is undecidable.

**To prove that a language  $L$  is undecidable,  
reduce a known undecidable language to  $L$ .**

Perhaps the most confusing aspect of reduction arguments is that the *languages* we want to prove undecidable nearly (but not quite) always involve encodings of Turing machines, while at the same time, the *programs* that we build to prove them undecidable are also Turing machines. Our proof that HALT is undecidable involved three different machines:

- The hypothetical Turing machine  $H$  that decides HALT.
- The new Turing machine  $SH$  that decides SELFHALT, using  $H$  as a subroutine.
- The Turing machine  $M$  whose encoding is the input to  $H$ .

It is *incredibly* easy to get confused about which machines are playing each in the proof. Therefore, it is absolutely *vital* that we give each machine in a reduction proof a unique and mnemonic name, and then **always** refer to each machine **by name**. Never write, say, or even *think* “the Turing machine” or “the state” or “the tape” or “the input” or (gods forbid) “it”. You also may find it useful to think of the working **programs** we are trying to construct ( $H$  and  $SH$  in this proof) as being written in a different language than the arbitrary **source code** that we want those programs to analyze ( $\langle M \rangle$  in this proof).

## 7.10 One Million Years Dungeon!

As a more complex set of examples, consider the following languages:

$$\begin{aligned}\text{NEVERACCEPT} &:= \{ \langle M \rangle \mid \text{ACCEPT}(M) = \emptyset \} \\ \text{NEVERREJECT} &:= \{ \langle M \rangle \mid \text{REJECT}(M) = \emptyset \} \\ \text{NEVERHALT} &:= \{ \langle M \rangle \mid \text{HALT}(M) = \emptyset \} \\ \text{NEVERDIVERGE} &:= \{ \langle M \rangle \mid \text{DIVERGE}(M) = \emptyset \}\end{aligned}$$

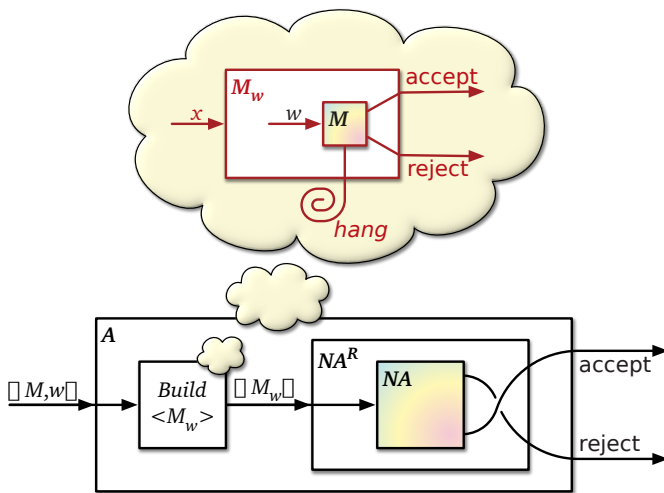
**Theorem 13.** *NEVERACCEPT is undecidable.*

**Proof:** Suppose to the contrary that there is a Turing machine  $NA$  that decides NEVERACCEPT. Then by swapping the **accept** and **reject** states, we obtain a Turing machine  $NA^R$  that decides the complementary language  $\Sigma^* \setminus \text{NEVERACCEPT}$ .

To reach a contradiction, we construct a Turing machine  $A$  that decides ACCEPT as follows. Given the encoding  $\langle M, w \rangle$  of an arbitrary machine  $M$  and an arbitrary string  $w$  as input,  $A$  writes the encoding  $\langle M_w \rangle$  of a new Turing machine  $M_w$  that ignores its input, writes  $w$  onto the tape, and then passes control to  $M$ . Finally,  $A$  passes the new encoding  $\langle M_w \rangle$  as input to  $NA^R$ . The following cartoon tries to illustrate the overall construction.

Before going any further, it may be helpful to list the various Turing machines that appear in this construction.





A reduction from from ACCEPT to NEVERACCEPT, which proves NEVERACCEPT undecidable.

- The hypothetical Turing machine  $NA$  that decides NEVERACCEPT.
- The Turing machine  $NA^R$  that decides  $\Sigma^* \setminus \text{NEVERACCEPT}$ , which we constructed by modifying  $NA$ .
- The Turing machine  $A$  that we are building, which decides ACCEPT using  $NA^R$  as a black-box subroutine.
- The Turing machine  $M$ , whose encoding is part of the input to  $A$ .
- The Turing machine  $M_w$  whose encoding  $A$  constructs from  $\langle M, w \rangle$  and then passes to  $NA^R$  as input.

Now let  $M$  be an arbitrary Turing machine and  $w$  be an arbitrary string, and suppose we run our new Turing machine  $A$  on the encoding  $\langle M, w \rangle$ . To complete the proof, we need to consider two cases: Either  $M$  accepts  $w$  or  $M$  does not accept  $w$ .

- First, suppose  $M$  accepts  $w$ .
  - Then for all strings  $x$ , the machine  $M_w$  accepts  $x$ .
  - So  $\text{ACCEPT}(M_w) = \Sigma^*$ , by the definition of  $\text{ACCEPT}(M_w)$ .
  - So  $\langle M_w \rangle \notin \text{NEVERACCEPT}$ , by definition of NEVERACCEPT.
  - So  $NA$  rejects  $\langle M_w \rangle$ , because  $NA$  decides NEVERACCEPT.
  - So  $NA^R$  accepts  $\langle M_w \rangle$ , by construction of  $NA^R$ .
  - We conclude that  $A$  accepts  $\langle M, w \rangle$ , by construction of  $A$ .
- On the other hand, suppose  $M$  does not accept  $w$ , either rejecting or diverging instead.
  - Then for all strings  $x$ , the machine  $M_w$  does not accept  $x$ .
  - So  $\text{ACCEPT}(M_w) = \emptyset$ , by the definition of  $\text{ACCEPT}(M_w)$ .
  - So  $\langle M_w \rangle \in \text{NEVERACCEPT}$ , by definition of NEVERACCEPT.
  - So  $NA$  accepts  $\langle M_w \rangle$ , because  $NA$  decides NEVERACCEPT.
  - So  $NA^R$  rejects  $\langle M_w \rangle$ , by construction of  $NA^R$ .
  - We conclude that  $A$  rejects  $\langle M, w \rangle$ , by construction of  $A$ .



In short,  $A$  decides the language **ACCEPT**, which is impossible. We conclude that  $NA$  does not exist.  $\square$

Again, similar arguments imply that the languages **NEVERREJECT**, **NEVERHALT**, and **NEVERDIVERGE** are undecidable. In each case, the core of the argument is describing how to transform the incoming machine-and-input encoding  $\langle M, w \rangle$  into the encoding of an appropriate new Turing machine  $\langle M_w \rangle$ .

Now that we know that **NEVERACCEPT** and its relatives are undecidable, we can use them as the basis of further reduction proofs. Here is a typical example:

**Theorem 14.** *The language  $DIVERGESAME := \{ \langle M_1 \rangle \langle M_2 \rangle \mid DIVERGE(M_1) = DIVERGE(M_2) \}$  is undecidable.*

**Proof:** Suppose for the sake of argument that there is a Turing machine  $DS$  that decides **DIVERGESAME**. Then we can build a Turing machine  $ND$  that decides **NEVERDIVERGE** as follows. Fix a Turing machine  $Y$  that accepts  $\Sigma^*$  (for example, by defining  $\delta(\text{start}, a) = (\text{accept}, \cdot, \cdot)$  for all  $a \in \Gamma$ ). Given an arbitrary Turing machine encoding  $\langle M \rangle$  as input,  $ND$  writes the string  $\langle M \rangle \langle Y \rangle$  onto the tape and then passes control to  $DS$ . There are two cases to consider:

- If  $DS$  accepts  $\langle M \rangle \langle Y \rangle$ , then  $DIVERGE(M) = DIVERGE(Y) = \emptyset$ , so  $\langle M \rangle \in \text{NEVERDIVERGE}$ .
- If  $DS$  rejects  $\langle M \rangle \langle Y \rangle$ , then  $DIVERGE(M) \neq DIVERGE(Y) = \emptyset$ , so  $\langle M \rangle \notin \text{NEVERDIVERGE}$ .

In short,  $ND$  accepts  $\langle M \rangle$  if and only if  $\langle M \rangle \in \text{NEVERDIVERGE}$ , which is impossible. We conclude that  $DS$  does not exist.  $\square$

## 7.11 Rice's Theorem

In 1953, Henry Rice proved the following extremely powerful theorem, which essentially states that **every** interesting question about the language accepted by a Turing machine is undecidable.

The following formulation is closer to the proof and may be (slightly) easier to use:

**Rice's Theorem.** *For any set  $\mathcal{L}$  of languages, if  $\emptyset \notin \mathcal{L}$  and there is a Turing machine  $M$  such that  $\text{ACCEPT}(M) \in \mathcal{L}$ , then the language  $\text{ACCEPTIN}(\mathcal{L}) := \{ \langle M \rangle \mid \text{ACCEPT}(M) \in \mathcal{L} \}$  is undecidable.*

The only downside of this formulation is that when  $\emptyset \in \mathcal{L}$ , we need to consider either the complementary property  $\bar{\mathcal{L}} = 2^{\Sigma^*} \setminus \mathcal{L}$  or the complementary language  $\{ \langle M \rangle \mid \text{ACCEPT}(M) \notin \mathcal{L} \}$ .

**Rice's Theorem.** *Let  $\mathcal{L}$  be any set of languages that satisfies the following conditions:*

- *There is a Turing machine  $Y$  such that  $\text{ACCEPT}(Y) \in \mathcal{L}$ .*
- *There is a Turing machine  $N$  such that  $\text{ACCEPT}(N) \notin \mathcal{L}$ .*

*The language  $\text{ACCEPTIN}(\mathcal{L}) := \{ \langle M \rangle \mid \text{ACCEPT}(M) \in \mathcal{L} \}$  is undecidable.*

**Proof:** Without loss of generality, suppose  $\emptyset \notin \mathcal{L}$ . (A symmetric argument establishes the theorem in the opposite case  $\emptyset \in \mathcal{L}$ .) Fix an arbitrary Turing machine  $Y$  such that  $\text{ACCEPT}(Y) \in \mathcal{L}$ .

Suppose to the contrary that there is a Turing machine  $A_{\mathcal{L}}$  that decides  $\text{ACCEPTIN}(\mathcal{L})$ . To derive a contradiction, we describe a Turing machine  $H$  that decides the halting language **HALT**, using  $A_{\mathcal{L}}$  as a black-box subroutine. Given the encoding  $\langle M, w \rangle$  of an arbitrary Turing machine  $M$  and an arbitrary string  $w$  as input,  $H$  writes the encoding  $\langle WTF \rangle$  of a new Turing machine  $WTF$  that executes the following algorithm:

$WTF(x)$ :

run  $M$  on input  $w$  (and discard the result)

run  $Y$  on input  $x$

$H$  then passes the new encoding  $\langle WTF \rangle$  to  $A_{\mathcal{L}}$ .

Now let  $M$  be an arbitrary Turing machine and  $w$  be an arbitrary string, and suppose we run our new Turing machine  $H$  on the encoding  $\langle M, w \rangle$ . There are two cases to consider.

- Suppose  $M$  halts on input  $w$ .
  - Then for all strings  $x$ , the machine  $WTF$  accepts  $x$  if and only if  $Y$  accepts  $x$ .
  - So  $\text{ACCEPT}(WTF) = \text{ACCEPT}(Y)$ , by definition of  $\text{ACCEPT}(\cdot)$ .
  - So  $\text{ACCEPT}(WTF) \in \mathcal{L}$ , by definition of  $Y$ .
  - So  $A_{\mathcal{L}}$  accepts  $\langle WTF \rangle$ , because  $A_{\mathcal{L}}$  decides  $\text{ACCEPTIN}(\mathcal{L})$ .
  - So  $H$  accepts  $\langle M, w \rangle$ , by definition of  $H$ .
- Suppose  $M$  does not halt on input  $w$ .
  - Then for all strings  $x$ , the machine  $WTF$  does not halt on input  $x$ , and therefore does not accept  $x$ .
  - So  $\text{ACCEPT}(WTF) = \emptyset$ , by definition of  $\text{ACCEPT}(WTF)$ .
  - So  $\text{ACCEPT}(WTF) \notin \mathcal{L}$ , by our assumption that  $\emptyset \notin \mathcal{L}$ .
  - So  $A_{\mathcal{L}}$  rejects  $\langle WTF \rangle$ , because  $A_{\mathcal{L}}$  decides  $\text{ACCEPTIN}(\mathcal{L})$ .
  - So  $H$  rejects  $\langle M, w \rangle$ , by definition of  $H$ .

In short,  $H$  decides the language  $\text{HALT}$ , which is impossible. We conclude that  $A_{\mathcal{L}}$  does not exist.  $\square$

The set  $\mathcal{L}$  in the statement of Rice's Theorem is often called a **property** of languages, rather than a *set*, to avoid the inevitable confusion about sets of sets of finite sequences of characters. We can also think of  $\mathcal{L}$  as a **decision problem** about languages, where the languages are represented by Turing machines that accept or decide them. Rice's theorem states that the **only** properties of languages that are decidable are the trivial properties "Does this Turing machine accept an acceptable language?" (Answer: Yes, by definition.) and "Does this Turing machine accept Discover?" (Answer: No, because Discover is a credit card, not a language.)

Rice's Theorem makes it incredibly easy to prove that language properties are undecidable; we only need to exhibit one acceptable language that has the property and another acceptable language that does not. In fact, **every** proof using Rice's theorem can use at least one of the following Turing machines:

- $M_{\text{ACCEPT}}$  accepts every string, by defining  $\delta(\text{start}, a) = \text{accept}$  for every tape symbol  $a$ .
- $M_{\text{REJECT}}$  rejects every string, by defining  $\delta(\text{start}, a) = \text{reject}$  for every tape symbol  $a$ .
- $M_{\text{DIVERGE}}$  diverges on every string, by defining  $\delta(\text{start}, a) = (\text{start}, a, +1)$  for every tape symbol  $a$ .

**Corollary 15.** *Each of the following languages is undecidable.*

- (a)  $\{\langle M \rangle \mid M \text{ accepts given an empty initial tape}\}$
- (b)  $\{\langle M \rangle \mid M \text{ accepts the string UIUC}\}$
- (c)  $\{\langle M \rangle \mid M \text{ accepts exactly three strings}\}$

- (d)  $\{\langle M \rangle \mid M \text{ accepts all palindromes}\}$
- (e)  $\{\langle M \rangle \mid \text{ACCEPT}(M) \text{ is regular}\}$
- (f)  $\{\langle M \rangle \mid \text{ACCEPT}(M) \text{ is not regular}\}$
- (g)  $\{\langle M \rangle \mid \text{ACCEPT}(M) \text{ is undecidable}\}$
- (h)  $\{\langle M \rangle \mid \text{ACCEPT}(M) = \text{ACCEPT}(N)\}$ , for some arbitrary fixed Turing machine  $N$ .

**Proof:** In all cases, undecidability follows from Rice's theorem.

- (a) Let  $\mathcal{L}$  be the set of all languages that contain the empty string. Then  $\text{ACCEPTIN}(\mathcal{L}) = \{\langle M \rangle \mid M \text{ accepts given an empty initial tape}\}$ .
  - Given an empty initial tape,  $M_{\text{ACCEPT}}$  accepts, so  $\text{ACCEPT}(M_{\text{ACCEPT}}) \in \mathcal{L}$ .
  - Given an empty initial tape,  $M_{\text{DIVERGE}}$  does not accept, so  $\text{ACCEPT}(M_{\text{DIVERGE}}) \notin \mathcal{L}$ .

Therefore, Rice's Theorem implies that  $\text{ACCEPTIN}(\mathcal{L})$  is undecidable.

- (b) Let  $\mathcal{L}$  be the set of all languages that contain the string **UIUC**.

- $M_{\text{ACCEPT}}$  accepts **UIUC**, so  $\text{ACCEPT}(M_{\text{ACCEPT}}) \in \mathcal{L}$ .
- $M_{\text{DIVERGE}}$  does not accept **UIUC**, so  $\text{ACCEPT}(M_{\text{DIVERGE}}) \notin \mathcal{L}$ .

Therefore,  $\text{ACCEPTIN}(\mathcal{L}) = \{\langle M \rangle \mid M \text{ accepts the string UIUC}\}$  is undecidable by Rice's Theorem.

- (c) There is a Turing machine that accepts the language **{larry, curly, moe}**. On the other hand,  $M_{\text{REJECT}}$  does not accept exactly three strings.
- (d)  $M_{\text{ACCEPT}}$  accepts all palindromes, and  $M_{\text{REJECT}}$  does not accept all palindromes.
- (e)  $M_{\text{REJECT}}$  accepts the regular language  $\emptyset$ , and there is a Turing machine  $M_{0^n 1^n}$  that accepts the non-regular language  $\{0^n 1^n \mid n \geq 0\}$ .
- (f)  $M_{\text{REJECT}}$  accepts the regular language  $\emptyset$ , and there is a Turing machine  $M_{0^n 1^n}$  that accepts the non-regular language  $\{0^n 1^n \mid n \geq 0\}$ .<sup>4</sup>
- (g)  $M_{\text{REJECT}}$  accepts the decidable language  $\emptyset$ , and there is a Turing machine that *accepts* the undecidable language **SELFREJECT**.
- (h) The Turing machine  $N$  accepts  $\text{ACCEPT}(N)$  by definition. For the negative Turing machine  $M_{\text{ACCEPT}}$  accepts  $\Sigma^*$  and the Turing machine  $M_{\text{REJECT}}$  accepts  $\emptyset$ , so at least one of those two machines does not accept  $\text{ACCEPT}(N)$ .  $\square$

We can also use Rice's theorem as a component in more complex undecidability proofs, where the target language consists of more than just a single Turing machine encoding.

**Theorem 16.** *The language  $L := \{\langle M, w \rangle \mid M \text{ accepts } w^k \text{ for every integer } k \geq 0\}$  is undecidable.*

**Proof:** Fix an arbitrary string  $w$ , and let  $\mathcal{L}$  be the set of all languages that contain  $w^k$  for all  $k$ . Then  $\text{ACCEPT}(M_{\text{ACCEPT}}) = \Sigma^* \in \mathcal{L}$  and  $\text{ACCEPT}(M_{\text{REJECT}}) = \emptyset \notin \mathcal{L}$ . Thus, *even if the string  $w$  is fixed in advance*, no Turing machine can decide  $L$ .  $\square$

Nearly identical reduction arguments imply the following variants of Rice's theorem. (The names of these theorems are not standard.)

<sup>4</sup>Yes, parts (e) and (f) have exactly the same proof.

**Rice's Rejection Theorem.** Let  $\mathcal{L}$  be any set of languages that satisfies the following conditions:

- There is a Turing machine  $Y$  such that  $\text{REJECT}(Y) \in \mathcal{L}$
- There is a Turing machine  $N$  such that  $\text{REJECT}(N) \notin \mathcal{L}$ .

The language  $\text{REJECTIN}(\mathcal{L}) := \{\langle M \rangle \mid \text{REJECT}(M) \in \mathcal{L}\}$  is undecidable.

**Rice's Halting Theorem.** Let  $\mathcal{L}$  be any set of languages that satisfies the following conditions:

- There is a Turing machine  $Y$  such that  $\text{HALT}(Y) \in \mathcal{L}$
- There is a Turing machine  $N$  such that  $\text{HALT}(N) \notin \mathcal{L}$ .

The language  $\text{HALTIN}(\mathcal{L}) := \{\langle M \rangle \mid \text{HALT}(M) \in \mathcal{L}\}$  is undecidable.

**Rice's Divergence Theorem.** Let  $\mathcal{L}$  be any set of languages that satisfies the following conditions:

- There is a Turing machine  $Y$  such that  $\text{DIVERGE}(Y) \in \mathcal{L}$
- There is a Turing machine  $N$  such that  $\text{DIVERGE}(N) \notin \mathcal{L}$ .

The language  $\text{DIVERGEIN}(\mathcal{L}) := \{\langle M \rangle \mid \text{DIVERGE}(M) \in \mathcal{L}\}$  is undecidable.

**Rice's Decision Theorem.** Let  $\mathcal{L}$  be any set of languages that satisfies the following conditions:

- There is a Turing machine  $Y$  that **decides** a language in  $\mathcal{L}$ .
- There is a Turing machine  $N$  that **decides** a language not in  $\mathcal{L}$ .

The language  $\text{DECIDEIN}(\mathcal{L}) := \{\langle M \rangle \mid M \text{ **decides** a language in } \mathcal{L}\}$  is undecidable.

As easy as it is to use Rice's theorem and its variants, they cannot be used for all undecidability proofs; these theorems only apply to properties of *languages*. For example, the language  $\text{THISISSPARTA} := \{\langle M \rangle \mid M \text{ accepts the string **SPARTA** after exactly 300 steps}\}$  is decidable, even though there are Turing machines that accept the string **SPARTA** after exactly 300 steps and there are other Turing machines that do not.

More subtly, Rice's theorem cannot be applied to self-referential languages like  $\text{REVACCEPT} := \{\langle M \rangle \mid M \text{ accepts } \langle M \rangle^R\}$ , because membership depends on details of the encoded machine and not just the language that the encoded machine accepts. To be clear: **REVACCEPT is undecidable**; you just can't use Rice's theorem to prove that fact.

## \*7.12 The Rice-McNaughton-Myhill-Shapiro Theorem

The following subtle generalization of Rice's theorem precisely characterizes which properties of acceptable languages are *acceptable*. This result was partially proved by Henry Rice in 1953, in the same paper that proved Rice's Theorem; Robert McNaughton, John Myhill, and Norman Shapiro completed the proof a few years later, each independently from the other two.<sup>5</sup>

**The Rice-McNaughton-Myhill-Shapiro Theorem.** Let  $\mathcal{L}$  be an arbitrary set of acceptable languages. The language  $\text{ACCEPTIN}(\mathcal{L}) := \{\langle M \rangle \mid \text{ACCEPT}(M) \in \mathcal{L}\}$  is **acceptable** if and only if  $\mathcal{L}$  satisfies the following conditions:

- (a)  $\mathcal{L}$  is **monotone**: For any language  $L \in \mathcal{L}$ , every superset of  $L$  is also in  $\mathcal{L}$ .
- (b)  $\mathcal{L}$  is **compact**: Every language in  $\mathcal{L}$  has a finite subset that is also in  $\mathcal{L}$ .

---

<sup>5</sup>McNaughton never published his proof (although he did announce the result); consequently, this theorem is sometimes called "The Rice-Myhill-Shapiro Theorem". Even more confusingly, Myhill published his proof twice, once in a paper with John Shepherdson and again in a later paper with Jacob Dekker. So maybe it should be called the Rice-Dekker-Myhill-(McNaughton-)Myhill-Shepherdson-Shapiro Theorem.

(c)  $\mathcal{L}$  is **finitely acceptable**: The language  $\{\langle L \rangle \mid L \in \mathcal{L} \text{ and } L \text{ is finite}\}$  is acceptable.<sup>6</sup>

I won't give a complete proof of this theorem (in part because it requires techniques I haven't introduced), but the following lemma is arguably the most interesting component:

**Lemma 17.** *Let  $\mathcal{L}$  be a set of acceptable languages. If  $\mathcal{L}$  is not monotone, then  $\text{ACCEPTIN}(\mathcal{L})$  is unacceptable.*

**Proof:** Suppose to the contrary that there is a Turing machine  $AI_{\mathcal{L}}$  that accepts  $\text{ACCEPTIN}(\mathcal{L})$ . Using this Turing machine as a black box, we describe a Turing machine  $SD$  that accepts the unacceptable language  $\text{SELF DIVERGE}$ . Fix two Turing machines  $Y$  and  $N$  such that

$$\begin{aligned} \text{ACCEPT}(Y) &\in \mathcal{L}, \\ \text{ACCEPT}(N) &\notin \mathcal{L}, \\ \text{and } \text{ACCEPT}(Y) &\subseteq \text{ACCEPT}(N). \end{aligned}$$

Let  $w$  be the input to  $SD$ . After verifying that  $w = \langle M \rangle$  for some Turing machine  $M$  (and rejecting otherwise),  $SD$  writes the encoding  $\langle WTF \rangle$  or a new Turing machine  $WTF$  that implements the following algorithm:

$WTF(x)$ :  
 write  $x$  to second tape  
 write  $\langle M \rangle$  to third tape  
 in parallel:  
     run  $Y$  on the first tape  
     run  $N$  on the second tape  
     run  $M$  on the third tape  
 if  $Y$  accepts  $x$   
     **accept**  
 if  $N$  accepts  $x$  and  $M$  halts on  $\langle M \rangle$   
     **accept**

Finally,  $SD$  passes the new encoding  $\langle WTF \rangle$  to  $AI_{\mathcal{L}}$ . There are two cases to consider:

- If  $M$  halts on  $\langle M \rangle$ , then  $\text{ACCEPT}(WTF) = \text{ACCEPT}(N) \notin \mathcal{L}$ , and therefore  $AI_{\mathcal{L}}$  does not accept  $\langle WTF \rangle$ .
- If  $M$  does not halt on  $\langle M \rangle$ , then  $\text{ACCEPT}(WTF) = \text{ACCEPT}(Y) \in \mathcal{L}$ , and therefore  $AI_{\mathcal{L}}$  accepts  $\langle WTF \rangle$ .

In short,  $SD$  accepts  $\text{SELF DIVERGE}$ , which is impossible. We conclude that  $SD$  does not exist.  $\square$

**Corollary 18.** *Each of the following languages is unacceptable.*

- (a)  $\{\langle M \rangle \mid \text{ACCEPT}(M) \text{ is finite}\}$
- (b)  $\{\langle M \rangle \mid \text{ACCEPT}(M) \text{ is infinite}\}$
- (c)  $\{\langle M \rangle \mid \text{ACCEPT}(M) \text{ is regular}\}$
- (d)  $\{\langle M \rangle \mid \text{ACCEPT}(M) \text{ is not regular}\}$
- (e)  $\{\langle M \rangle \mid \text{ACCEPT}(M) \text{ is decidable}\}$

<sup>6</sup>Here the encoding  $\langle L \rangle$  of a finite language  $L \subseteq \Sigma^*$  is exactly the string that you would write down to explicitly describe  $L$ . Formally,  $\langle L \rangle$  is the unique string over the alphabet  $\Sigma \cup \{\{, \bullet, \}, \mathbf{\epsilon}\}$  that contains the strings in  $L$  in lexicographic order, separated by commas  $\bullet$  and surrounded by braces  $\{\}$ , with  $\mathbf{\epsilon}$  representing the empty string. For example,  $\{\{\epsilon, 0, 01, 0110, 01101001\}\} = \{\mathbf{\epsilon} \bullet 0 \bullet 01 \bullet 0110 \bullet 01101001\}$ .

- (f)  $\{\langle M \rangle \mid \text{ACCEPT}(M) \text{ is undecidable}\}$
- (g)  $\{\langle M \rangle \mid M \text{ accepts at least one string in SELF DIVERGE}\}$
- (h)  $\{\langle M \rangle \mid \text{ACCEPT}(M) = \text{ACCEPT}(N)\}$ , for some arbitrary fixed Turing machine  $N$ .

**Proof:** (a) The set of finite languages is not monotone:  $\emptyset$  is finite;  $\Sigma^*$  is not finite; both  $\emptyset$  and  $\Sigma^*$  are acceptable (in fact decidable); and  $\emptyset \subset \Sigma^*$ .

(b) The set of infinite acceptable languages is not compact: No finite subset of the infinite acceptable language  $\Sigma^*$  is infinite!

(c) The set of regular languages is not monotone: Consider the languages  $\emptyset$  and  $\{0^n 1^n \mid n \geq 0\}$ .

(d) The set of non-regular acceptable languages is not monotone: Consider the languages  $\{0^n 1^n \mid n \geq 0\}$  and  $\Sigma^*$ .

(e) The set of decidable languages is not monotone: Consider the languages  $\emptyset$  and SELF REJECT.

(f) The set of undecidable acceptable languages is not monotone: Consider the languages SELF REJECT and  $\Sigma^*$ .

(g) The set  $\mathcal{L} = \{L \mid L \cap \text{SELF DIVERGE} \neq \emptyset\}$  is not finitely acceptable. For any string  $w$ , deciding whether  $\{w\} \in \mathcal{L}$  is equivalent to deciding whether  $w \in \text{SELF DIVERGE}$ , which is impossible.

(h) If  $\text{ACCEPT}(N) \neq \Sigma^*$ , then the set  $\{\text{ACCEPT}(N)\}$  is not monotone. On the other hand, if  $\text{ACCEPT}(N) = \Sigma^*$ , then the set  $\{\text{ACCEPT}(N)\}$  is not compact: No finite subset of  $\Sigma^*$  is equal to  $\Sigma^*$ !

□

## 7.13 Turing Machine Behavior: It's Complicated

Rice's theorems imply that every interesting question about the language that a Turing machine accepts—or more generally, the function that a program computes—is undecidable. A more subtle question is whether we can recognize Turing machines that exhibit certain *internal behavior*. Some behaviors we can recognize; others we can't.

**Theorem 19.** *The language  $\text{NEVERLEFT} := \{\langle M, w \rangle \mid \text{Given } w \text{ as input, } M \text{ never moves left}\}$  is decidable.*

**Proof:** Given the encoding  $\langle M, w \rangle$ , we simulate  $M$  with input  $w$  using our universal Turing machine  $U$ , but with the following termination conditions. If  $M$  ever moves its head to the left, then we **reject**. If  $M$  halts without moving its head to the left, then we **accept**. Finally, if  $M$  reads more than  $|Q|$  blanks, where  $Q$  is the state set of  $M$ , then we **accept**. If the first two cases do not apply,  $M$  only moves to the right; moreover, after reading the entire input string,  $M$  only reads blanks. Thus, after reading  $|Q|$  blanks, it must repeat some state, and therefore loop forever without moving to the left. The three cases are exhaustive. □

**Theorem 20.** *The language  $\text{LEFTTHREE} := \{\langle M, w \rangle \mid \text{Given } w \text{ as input, } M \text{ eventually moves left three times in a row}\}$  is undecidable.*

**Proof:** Given  $\langle M \rangle$ , we build a new Turing machine  $M'$  that accepts the same language as  $M$  and moves left three times in a row if and only if it accepts, as follows. For each non-accepting state  $p$

of  $M$ , the new machine  $M'$  has three states  $p_1, p_2, p_3$ , with the following transitions:

$$\begin{aligned}\delta'(p_1, a) &= (q_2, b, \Delta), & \text{where } (q, b, \Delta) &= \delta(p, a) \text{ and } q \neq \text{accept} \\ \delta'(p_2, a) &= (p_3, a, +1) \\ \delta'(p_3, a) &= (p_1, a, -1)\end{aligned}$$

In other words, after each non-accepting transition,  $M'$  moves once to the right and then once to the left. For each transition to **accept**,  $M'$  has a sequence of seven transitions: three steps to the right, then three steps to the left, and then finally **accept'**, all without modifying the tape. (The three steps to the right ensure that  $M'$  does not fall off the left end of the tape.)

Finally,  $M'$  moves left three times in a row if and only if  $M$  accepts  $w$ . Thus, if we could decide LEFTTHREE, we could also decide ACCEPT, which is impossible.  $\square$

There is no hard and fast rule like Rice's theorem to distinguish decidable behaviors from undecidable behaviors, but I can offer two rules of thumb.

- If it is possible to simulate an arbitrary Turing machine while avoiding the target behavior, then the behavior is not decidable. For example, there is no algorithm to determine whether a given Turing machine reenters its **start** state, or revisits the left end of the tape, or writes a blank.
- If a Turing machine with the target behavior is limited to a finite number of configurations, or is guaranteed to force an infinite loop after a finite number of transitions, then the behavior is likely to be decidable. For example, there *are* algorithms to determine whether a given Turing machine ever leaves its **start** state, or reads its entire input string, or writes a non-blank symbol over a blank.

## Exercises

1. Let  $M$  be an arbitrary Turing machine.
  - (a) Describe a Turing machine  $M^R$  such that

$$\text{ACCEPT}(M^R) = \text{REJECT}(M) \quad \text{and} \quad \text{REJECT}(M^R) = \text{ACCEPT}(M).$$

- (b) Describe a Turing machine  $M^A$  such that

$$\text{ACCEPT}(M^A) = \text{ACCEPT}(M) \quad \text{and} \quad \text{REJECT}(M^A) = \emptyset.$$

- (c) Describe a Turing machine  $M^H$  such that

$$\text{ACCEPT}(M^H) = \text{HALT}(M) \quad \text{and} \quad \text{REJECT}(M^H) = \emptyset.$$

2.
  - (a) Prove that ACCEPT is undecidable.
  - (b) Prove that REJECT is undecidable.
  - (c) Prove that DIVERGE is undecidable.
3.
  - (a) Prove that NEVERREJECT is undecidable.



- (b) Prove that NEVERHALT is undecidable.
  - (c) Prove that NEVERDIVERGE is undecidable.
4. Prove that each of the following languages is undecidable.
- (a)  $\text{ALWAYSACCEPT} := \{\langle M \rangle \mid \text{ACCEPT}(M) = \Sigma^*\}$
  - (b)  $\text{ALWAYSREJECT} := \{\langle M \rangle \mid \text{REJECT}(M) = \Sigma^*\}$
  - (c)  $\text{ALWAYSHALT} := \{\langle M \rangle \mid \text{HALT}(M) = \Sigma^*\}$
  - (d)  $\text{ALWAYSDIVERGE} := \{\langle M \rangle \mid \text{DIVERGE}(M) = \Sigma^*\}$
5. Let  $\mathcal{L}$  be a non-empty proper subset of the set of acceptable languages. Prove that the following languages are undecidable:
- (a)  $\text{REJECTIN}(\mathcal{L}) := \{\langle M \rangle \mid \text{REJECT}(M) \in \mathcal{L}\}$
  - (b)  $\text{HALTIN}(\mathcal{L}) := \{\langle M \rangle \mid \text{HALT}(M) \in \mathcal{L}\}$
  - (c)  $\text{DIVERGEIN}(\mathcal{L}) := \{\langle M \rangle \mid \text{DIVERGE}(M) \in \mathcal{L}\}$
6. For each of the following decision problems, either *sketch* an algorithm or prove that the problem is undecidable. Recall that  $w^R$  denotes the reversal of string  $w$ . For each problem, the input is the encoding  $\langle M \rangle$  of a Turing machine  $M$ .
- (a) Does  $M$  reject the empty string?
  - (b) Does  $M$  accept  $\langle M \rangle^R$ ?
  - (c) Does  $M$  accept  $\langle M \rangle \langle M \rangle$ ?
  - (d) Does  $M$  accept  $\langle M \rangle^k$  for any integer  $k$ ?
  - (e) Does  $M$  accept the encoding of any Turing machine?
  - (f) Is there a Turing machine that accepts  $\langle M \rangle$ ?
  - (g) Is  $\langle M \rangle$  a palindrome?
  - (h) Does  $M$  reject any palindrome?
  - (i) Does  $M$  accept all palindromes?
  - (j) Does  $M$  diverge only on palindromes?
  - (k) Is there an input string that forces  $M$  to move left?
  - (l) Is there an input string that forces  $M$  to move left three times in a row?
  - (m) Does  $M$  accept the encoding of any Turing machine  $N$  such that  $\text{ACCEPT}(N) = \text{SELFDIVERGE}$ ?
7. For each of the following decision problems, either *sketch* an algorithm or prove that the problem is undecidable. Recall that  $w^R$  denotes the reversal of string  $w$ . For each problem, the input is an encoding  $\langle M, w \rangle$  of a Turing machine  $M$  and its input string  $w$ .
- (a) Does  $M$  accept the string  $ww^R$ ?



- (b) Does  $M$  accept either  $w$  or  $w^R$ ?
- (c) Does  $M$  either accept  $w$  or reject  $w^R$ ?
- (d) Does  $M$  accept the string  $w^k$  for some integer  $k$ ?
- (e) Does  $M$  accept  $w$  in at most  $2^{|w|}$  steps?
- (f) If we run  $M$  on input  $w$ , does  $M$  ever change a symbol on its tape?
- (g) If we run  $M$  on input  $w$ , does  $M$  ever move to the right?
- (h) If we run  $M$  on input  $w$ , does  $M$  ever move to the right twice in a row?
- (i) If we run  $M$  on input  $w$ , does  $M$  move its head to the right more than  $2^{|w|}$  times (not necessarily consecutively)?
- (j) If we run  $M$  with input  $w$ , does  $M$  ever change a  $\square$  on the tape to any other symbol?
- (k) If we run  $M$  with input  $w$ , does  $M$  ever change a  $\square$  on the tape to  $1$ ?
- (l) If we run  $M$  with input  $w$ , does  $M$  ever write a  $\square$ ?
- (m) If we run  $M$  with input  $w$ , does  $M$  ever leave its **start** state?
- (n) If we run  $M$  with input  $w$ , does  $M$  ever reenter its **start** state?
- (o) If we run  $M$  with input  $w$ , does  $M$  ever reenter a state that it previously left? That is, are there states  $p \neq q$  such that  $M$  moves from state  $p$  to state  $q$  and then later moves back to state  $p$ ?

8. Let  $M$  be a Turing machine, let  $w$  be an arbitrary input string, and let  $s$  and  $t$  be positive integers. We say that  $M$  accepts  $w$  **in space**  $s$  if  $M$  accepts  $w$  after accessing at most the first  $s$  cells on the tape, and  $M$  accepts  $w$  **in time**  $t$  if  $M$  accepts  $w$  after at most  $t$  transitions.

- (a) Prove that the following languages are decidable:
  - i.  $\{\langle M, w \rangle \mid M \text{ accepts } w \text{ in time } |w|^2\}$
  - ii.  $\{\langle M, w \rangle \mid M \text{ accepts } w \text{ in space } |w|^2\}$
- (b) Prove that the following languages are undecidable:
  - i.  $\{\langle M \rangle \mid M \text{ accepts at least one string } w \text{ in time } |w|^2\}$
  - ii.  $\{\langle M \rangle \mid M \text{ accepts at least one string } w \text{ in space } |w|^2\}$

9. Let  $L_0$  be an arbitrary language. For any integer  $i > 0$ , define the language

$$L_i := \{\langle M \rangle \mid M \text{ decides } L_{i-1}\}.$$

For which integers  $i > 0$  is  $L_i$  decidable? Obviously the answer depends on the initial language  $L_0$ ; give a complete characterization of all possible cases. Prove your answer is correct. *[Hint: This question is a lot easier than it looks!]*

10. Argue that each of the following decision problems about programs in your favorite programming language are undecidable.

- (a) Does this program correctly compute Fibonacci numbers?

- (b) Can this program fall into an infinite loop?
- (c) Will the value of this variable ever change?
- (d) Will this program every attempt to deference a null pointer?
- (e) Does this program free every block of memory that it dynamically allocates?
- (f) Is any statement in this program unreachable?
- (g) Do these two programs compute the same function?

★11. Call a Turing machine **conservative** if it never writes over its input string. More formally, a Turing machine is conservative if for every transition  $\delta(p, a) = (q, b, \Delta)$  where  $a \in \Sigma$ , we have  $b = a$ ; and for every transition  $\delta(p, a) = (q, b, \Delta)$  where  $a \notin \Sigma$ , we have  $b \neq \Sigma$ .

- (a) Prove that if  $M$  is a conservative Turing machine, then  $\text{ACCEPT}(M)$  is a regular language.
- (b) Prove that the language  $\{\langle M \rangle \mid M \text{ is conservative and } M \text{ accepts } \varepsilon\}$  is undecidable.

Together, these two results imply that every conservative Turing machine accepts the same language as some DFA, but it is impossible to determine *which* DFA.

- ★12. (a) Prove that it is undecidable whether a given C++ program is syntactically correct.  
[Hint: Use templates!]
- (b) Prove that it is undecidable whether a given ANSI C program is syntactically correct.  
[Hint: Use the preprocessor!]
- (c) Prove that it is undecidable whether a given Perl program is syntactically correct.  
[Hint: Does that slash character / delimit a regular expression or represent division?]

## 8 Universal Models of Computation

Remind about the Church-Turing thesis.

There is some confusion here between **universal models of computation** and the somewhat wider class of **undecidable problems/languages**.

### 8.1 Universal Turing Machines

The pinnacle of Turing machine constructions is the *universal* Turing machine. For modern computer scientists, it's useful to think of a universal Turing machine as a "Turing machine *interpreter* written in Turing machine". Just as the input to a Python interpreter is a string of Python source code, the input to our universal Turing machine  $U$  is a string  $\langle M, w \rangle$  that encodes both an arbitrary Turing machine  $M$  and a string  $w$  in the input alphabet of  $M$ . Given these encodings,  $U$  simulates the execution of  $M$  on input  $w$ ; in particular,

- $U$  accepts  $\langle M, w \rangle$  if and only if  $M$  accepts  $w$ .
- $U$  rejects  $\langle M, w \rangle$  if and only if  $M$  rejects  $w$ .

In the next few pages, I will sketch a universal Turing machine  $U$  that uses the input alphabet  $\{0, 1, [, ], \bullet, | \}$  and a somewhat larger tape alphabet. However, I do *not* require that the Turing machines that  $U$  simulates have similarly small alphabets, so we first need a method to encode *arbitrary* input and tape alphabets.

### Encodings

Let  $M = (\Gamma, \square, \Sigma, Q, \textit{start}, \textit{accept}, \textit{reject}, \delta)$  be an arbitrary Turing machine, with a single half-infinite tape and a single read-write head. (I will consistently indicate the states and tape symbols of  $M$  in *slanted green* to distinguish them from the *upright red* states and tape symbols of  $U$ .)

We encode each symbol  $a \in \Gamma$  as a unique string  $|a|$  of  $\lceil \lg(|\Gamma|) \rceil$  bits. Thus, if  $\Gamma = \{0, 1, \$, x, \square\}$ , we might use the following encoding:

$$\langle 0 \rangle = 001, \quad \langle 1 \rangle = 010, \quad \langle \$ \rangle = 011, \quad \langle x \rangle = 100, \quad \langle \square \rangle = 000.$$

The input string  $w$  is encoded by its sequence of symbol encodings, with separators  $\bullet$  between every pair of symbols and with brackets  $[$  and  $]$  around the whole string. For example, with this encoding, the input string *001100* would be encoded on the input tape as

$$\langle 001100 \rangle = [001 \bullet 001 \bullet 010 \bullet 010 \bullet 001 \bullet 001]$$

Similarly, we encode each state  $q \in Q$  as a distinct string  $\langle q \rangle$  of  $\lceil \lg|Q| \rceil$  bits. Without loss of generality, we encode the start state with all 1s and the reject state with all 0s. For example, if  $Q = \{\text{start}, \text{seek1}, \text{seek0}, \text{reset}, \text{verify}, \text{accept}, \text{reject}\}$ , we might use the following encoding:

$$\begin{array}{llll} \langle \text{start} \rangle = 111 & \langle \text{seek1} \rangle = 010 & \langle \text{seek0} \rangle = 011 & \langle \text{reset} \rangle = 100 \\ \langle \text{verify} \rangle = 101 & \langle \text{accept} \rangle = 110 & \langle \text{reject} \rangle = 000 & \end{array}$$

We encode the machine  $M$  itself as the string  $\langle M \rangle = [\langle \text{reject} \rangle \bullet \langle \square \rangle] \langle \delta \rangle$ , where  $\langle \delta \rangle$  is the concatenation of substrings  $[\langle p \rangle \bullet \langle a \rangle \mid \langle q \rangle \bullet \langle b \rangle \bullet \langle \Delta \rangle]$  encoding each transition  $\delta(p, a) = (q, b, \Delta)$  such that  $q \neq \text{reject}$ . We encode the actions  $\Delta = \pm 1$  by defining  $\langle -1 \rangle := 0$  and  $\langle +1 \rangle := 1$ . Conveniently, every transition string has exactly the same length. For example, with the symbol and state encodings described above, the transition  $\delta(\text{reset}, \$) = (\text{start}, \$, +1)$  would be encoded as

$$[100 \bullet 011 \mid 001 \bullet 011 \bullet 1].$$

Our first example Turing machine for recognizing  $\{0^n 1^n 0^n \mid n \geq 0\}$  would be represented by the following string (here broken into multiple lines for readability):

$$\begin{aligned} [000 \bullet 000] & [001 \bullet 001 \mid 010 \bullet 011 \bullet 1] [001 \bullet 100 \mid 101 \bullet 011 \bullet 1] \\ & [010 \bullet 001 \mid 010 \bullet 001 \bullet 1] [010 \bullet 100 \mid 010 \bullet 100 \bullet 1] \\ & [010 \bullet 010 \mid 011 \bullet 100 \bullet 1] [011 \bullet 010 \mid 011 \bullet 010 \bullet 1] \\ & [011 \bullet 100 \mid 011 \bullet 100 \bullet 1] [011 \bullet 001 \mid 100 \bullet 100 \bullet 1] \\ & [100 \bullet 001 \mid 100 \bullet 001 \bullet 0] [100 \bullet 010 \mid 100 \bullet 010 \bullet 0] \\ & [100 \bullet 100 \mid 100 \bullet 100 \bullet 0] [100 \bullet 011 \mid 001 \bullet 011 \bullet 1] \\ & [101 \bullet 100 \mid 101 \bullet 011 \bullet 1] [101 \bullet 000 \mid 110 \bullet 000 \bullet 0] \end{aligned}$$

Finally, we encode any *configuration* of  $M$  on  $U$ 's work tape by alternating between encodings of states and encodings of tape symbols. Thus, each tape cell is represented by the string  $[\langle q \rangle \bullet \langle a \rangle]$  indicating that (1) the cell contains symbol  $a$ ; (2) if  $q \neq \text{reject}$ , then  $M$ 's head is located at this cell, and  $M$  is in state  $q$ ; and (3) if  $q = \text{reject}$ , then  $M$ 's head is located somewhere else. Conveniently, each cell encoding uses exactly the same number of bits. We also surround the entire tape encoding with brackets  $[$  and  $]$ .

For example, with the encodings described above, the initial configuration  $(\text{start}, \uparrow 001100, 0)$  for our first example Turing machine would be encoded on  $U$ 's tape as follows.

$$\langle \text{start}, \uparrow 001100, 0 \rangle = \underbrace{[111 \bullet 001]}_{\text{start } 0} \underbrace{[000 \bullet 001]}_{\text{reject } 0} \underbrace{[000 \bullet 010]}_{\text{reject } 1} \underbrace{[000 \bullet 010]}_{\text{reject } 1} \underbrace{[000 \bullet 001]}_{\text{reject } 0} \underbrace{[000 \bullet 000]}_{\text{reject } 0}$$

Similarly, the intermediate configuration  $(\text{reset}, \$0x\uparrow x0, 3)$  would be encoded as follows:

$$\langle \text{reset}, \$\$x\uparrow x0, 3 \rangle = \underbrace{[000 \bullet 011]}_{\text{reject } \$} \underbrace{[000 \bullet 011]}_{\text{reject } 0} \underbrace{[000 \bullet 100]}_{\text{reject } x} \underbrace{[010 \bullet 010]}_{\text{reset } 1} \underbrace{[000 \bullet 100]}_{\text{reject } x} \underbrace{[000 \bullet 000]}_{\text{reject } 0}$$

## Input and Execution

Without loss of generality, we assume that the input to our universal Turing machine  $U$  is given on a separate read-only *input tape*, as the encoding of an arbitrary Turing machine  $M$  followed by an encoding of its input string  $x$ . Notice the substrings  $[[$  and  $]]$  each appear only once on the input tape, immediately before and after the encoded transition table, respectively.  $U$  also has a read-write *work tape*, which is initially blank.

We start by initializing the work tape with the encoding  $\langle \text{start}, x, 0 \rangle$  of the initial configuration of  $M$  with input  $x$ . First, we write  $[[\langle \text{start} \rangle \bullet]$ . Then we copy the encoded input string  $\langle x \rangle$  onto the work tape, but we change the punctuation as follows:

- Instead of copying the left bracket  $[$ , write  $[[\langle \text{start} \rangle \bullet]$ .
- Instead of copying each separator  $\bullet$ , write  $][\langle \text{reject} \rangle \bullet]$ .
- Instead of copying the right bracket  $]$ , write two right brackets  $]]$ .

The state encodings  $\langle \text{start} \rangle$  and  $\langle \text{reject} \rangle$  can be copied directly from the beginning of  $\langle M \rangle$  (replacing 0s for 1s for  $\langle \text{start} \rangle$ ). Finally, we move the head back to the start of  $U$ 's tape.

At the start of each step of the simulation,  $U$ 's head is located at the start of the work tape. We scan through the work tape to the unique encoded cell  $[[\langle p \rangle \bullet \langle a \rangle]]$  such that  $p \neq \text{reject}$ . Then we scan through the encoded transition function  $\langle \delta \rangle$  to find the unique encoded tuple  $[[\langle p \rangle \bullet \langle a \rangle | \langle q \rangle \bullet \langle b \rangle \bullet \langle \Delta \rangle]]$  whose left half matches our the encoded tape cell. If there is no such tuple, then  $U$  immediately halts and rejects. Otherwise, we copy the right half  $\langle q \rangle \bullet \langle b \rangle$  of the tuple to the work tape. Now if  $q = \text{accept}$ , then  $U$  immediately halts and accepts. (We don't bother to encode  $\text{reject}$  transformations, so we know that  $q \neq \text{reject}$ .) Otherwise, we transfer the state encoding to either the next or previous encoded cell, as indicated by  $M$ 's transition function, and then continue with the next step of the simulation.

During the final state-copying phase, we ever read two right brackets  $]]$ , indicating that we have reached the right end of the tape encoding, we replace the second right bracket with  $][\langle \text{reject} \rangle \bullet \langle \square \rangle]]$  (mostly copied from the beginning of the machine encoding  $\langle M \rangle$ ) and then scan back to the left bracket we just wrote. This trick allows our universal machine to *pretend* that its tape contains an infinite sequence of *encoded* blanks  $][\langle \text{reject} \rangle \bullet \langle \square \rangle]]$  instead of *actual* blanks  $\square$ .

## Example

As an illustrative example, suppose  $U$  is simulating our first example Turing machine  $M$  on the input string  $001100$ . The execution of  $M$  on input  $w$  eventually reaches the configuration  $(\text{seek1}, \$\$x1x0, 3)$ . At the start of the corresponding step in  $U$ 's simulation,  $U$  is in the following configuration:

$\uparrow [[000 \bullet 011][000 \bullet 011][000 \bullet 100][010 \bullet 010][000 \bullet 100][000 \bullet 001]]$

First  $U$  scans for the first encoded tape cell whose state is not  $\text{reject}$ . That is,  $U$  repeatedly compares the first half of each encoded state cell on the work tape with the prefix  $][\langle \text{reject} \rangle \bullet$  of the machine encoding  $\langle M \rangle$  on the input tape.  $U$  finds a match in the fourth encoded cell.

$[[000 \bullet 011][000 \bullet 011][000 \bullet 100][010 \bullet 010][000 \bullet 100][000 \bullet 001]]$   
 $\uparrow$

Next,  $U$  scans the machine encoding  $\langle M \rangle$  for the substring  $[010 \bullet 010]$  matching the current encoded cell.  $U$  eventually finds a match in the left side of the the encoded transition  $[010 \bullet 010 | 011 \bullet 100 \bullet 1]$ .  $U$  copies the state-symbol pair  $011 \bullet 100$  from the right half of this encoded transition into the current encoded cell. (The underline indicates which symbols are changed.)

$[[000 \bullet 011][000 \bullet 011][000 \bullet 100][\underline{011 \bullet 100}][000 \bullet 100][000 \bullet 001]]$   
 $\uparrow$

The encoded transition instructs  $U$  to move the current state encoding one cell to the right. (The underline indicates which symbols are changed.)

$[[000 \bullet 011] [000 \bullet 011] [000 \bullet 100] [\underline{000} \bullet 100] [\underline{011} \bullet 100] [000 \bullet 001]]$

Finally,  $U$  scans left until it reads two left brackets  $[[$ ; this returns the head to the left end of the work tape to start the next step in the simulation.  $U$ 's tape now holds the encoding of  $M$ 's configuration ( $\text{seek}0, \$\$xx\text{seek}0, 4$ ), as required.

$[[000 \bullet 011] [000 \bullet 011] [000 \bullet 100] [000 \bullet 100] [011 \bullet 100] [000 \bullet 001]]$

## 8.2 Two-Stack Machines

A **two-stack machine** is a Turing machine with two tapes with the following restricted behavior. At all times, on each tape, every cell to the right of the head is blank, and every cell at or to the left of the head is non-blank. Thus, a head can only move right by writing a non-blank symbol into a blank cell; symmetrically, a head can only move left by erasing the rightmost non-blank cell. Thus, each tape behaves like a stack. To avoid underflow, there is a special symbol at the start of each tape that cannot be overwritten. Initially, one tape contains the input string, with the head at its *last* symbol, and the other tape is empty (except for the start-of-tape symbol).

Simulate a doubly-infinite tape with two stacks, one holding the tape contents to the left of the head, the other holding the tape contents to the right of the head. For each transition of a standard Turing machine  $M$ , the stack machine pops the top symbol off the (say) left stack, changes its internal state according to the transition  $\delta$ , and then either pushes a new symbol onto the right stack, or pushes a new symbol onto the left stack and then moves the top symbol from the right stack to the left stack.

## 8.3 Counter Machines

A configuration of a  $k$ -counter machine consists of  $k$  non-negative integers and an internal state from some finite set  $Q$ . The transition function  $\delta: Q \times \{0, +1\}^k \rightarrow Q \times \{-1, 0, +1\}^k$  takes an internal state and the signs of the counters as input, and produces a new internal state and changes to counters as output.

- Prove that any Turing machine can be simulated by a three-counter machine. One counter holds the binary representation of the tape after the head; another counter holds the reversed binary representation of the tape before the head. Implement transitions via halving, doubling, and parity, using the third counter for scratch work.
- Prove that two counters can simulate three. Store  $2^a 3^b 5^c$  in one counter, use the other for scratch work.
- Prove that a three-counter machine can compute any computable function: Given input  $(n, 0, 0)$ , we can compute  $(f(n), 0, 0)$  for *any* computable function  $f$ . First transform  $(n, 0, 0)$  to  $(2^n, 0, 0)$  using all three counters; then run two- (or three-) counter TM simulation to obtain  $(2^{f(n)}, 0, 0)$ ; and finally transform  $(2^{f(n)}, 0, 0)$  to  $(f(n), 0, 0)$  using all three counters.
- **HARD:** Prove that a two-counter machine cannot transform  $(n, 0)$  to  $(2^n, 0)$ . [Barzhdin 1963, Yao 1971, Schröpel 1972]

## 8.4 FRACTRAN

FRACTRAN [Conway 1987]: A one-counter machine whose “program” is a sequence of rational numbers. The counter is initially 1. At each iteration, multiply the counter by the first rational number that yields an integer; if there is no such number, halt.

- Prove that for any computable function  $f : \mathbb{N} \rightarrow \mathbb{N}$ , there is a FRACTRAN program that transforms  $2^{n+1}$  into  $3^{f(n)+1}$ , for all natural numbers  $n$ .
- Prove that every FRACTRAN program, given the integer 1 as input, either outputs 1 or loops forever. It follows that there is no FRACTRAN program for the increment function  $n \mapsto n + 1$ .

## 8.5 Post Correspondence Problem

Given  $n$  of pairs of strings  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ , is there a finite sequence of integers  $(i_1, i_2, \dots, i_k)$  such that  $x_{i_1}x_{i_2} \cdots x_{i_k} = y_{i_1}y_{i_2} \cdots y_{i_k}$ ? For notation convenience, we write each pair vertically as  $\begin{bmatrix} x \\ y \end{bmatrix}$  instead of horizontally as  $(x, y)$ . For example, given the string pairs

$$a = \begin{bmatrix} 0 \\ 100 \end{bmatrix}, b = \begin{bmatrix} 01 \\ 00 \end{bmatrix}, c = \begin{bmatrix} 110 \\ 11 \end{bmatrix},$$

we should answer TRUE, because

$$cbca = \begin{bmatrix} 110 \\ 11 \end{bmatrix} \begin{bmatrix} 01 \\ 00 \end{bmatrix} \begin{bmatrix} 110 \\ 11 \end{bmatrix} \begin{bmatrix} 0 \\ 100 \end{bmatrix}$$

gives us **110110100** for both concatenations. As more extreme examples, the shortest solutions for the input

$$a = \begin{bmatrix} 0 \\ 001 \end{bmatrix}, b = \begin{bmatrix} 001 \\ 1 \end{bmatrix}, c = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

have length 75; one such solution is *aacaacabbabccaaccaaacbaabbaacbacbbccbbacbacbcbacbbacbacbbacbbacccbabbbccbaacaacaacbabbaacacbccbbabacbcaaccbacabbbbabcccbcaababaaccbcbbbacccbabbbcc*. The shortest solution for the instance

$$a = \begin{bmatrix} 0 \\ 000 \end{bmatrix}, b = \begin{bmatrix} 0 \\ 0101 \end{bmatrix}, c = \begin{bmatrix} 01 \\ 1 \end{bmatrix}, d = \begin{bmatrix} 1111 \\ 10 \end{bmatrix}$$

is the unbelievable  $a^2b^8a^4c^{16}ab^4a^2b^4ad^4b^3c^8a^6c^8b^2c^4bc^6d^2a^{18}d^2c^4dcad^2cb^{54}c^3dca^2c^{111}dc^6d^{28}cb^{17}c^{63}d^{16}c^{16}d^4c^4dc$ , which has total length 451. Finally, the shortest solution for the instance

$$a = \begin{bmatrix} 0 \\ 00010 \end{bmatrix}, b = \begin{bmatrix} 010 \\ 01 \end{bmatrix}, c = \begin{bmatrix} 100 \\ 0 \end{bmatrix},$$

has length 528.

The simplest universality proof simulates a tag-Turing machine.

## 8.6 Matrix Mortality

Given a set of integer matrices  $A_1, \dots, A_k$ , is the product of any sequence of these matrices (with repetition) equal to 0? Undecidable by reduction from PCP, even for two  $15 \times 15$  matrices or six  $3 \times 3$  matrices [Cassaigne, Halava, Harju, Nicolas 2014]

## 8.7 Dynamical Systems

Ray Tracing [Reif, Tygar, and Yoshida 1994] The configuration of a Turing machine is encoded as the  $(x, y)$  coordinates of a light path crossing the unit square  $[0, 1] \times [0, 1]$ , where the  $x$ - (resp.  $y$ -)coordinate encodes the tape contents to the left (resp. right) of the head. Need either quadratic-surface mirrors or refraction to simulate transitions.

N-body problem [Smith 2006]: Similar idea

Skolem-Pisot reachability: Given an integer vector  $x$  and an integer matrix  $A$ , does  $A^n x = (0, \dots)$  for any integer  $n$ ? [Halava, Harju, Hirvensalo, Karhumäki 2005] It's surprising that this problem is undecidable; the similar mortality problem for one matrix is not.

## 8.8 Wang Tiles

Turing machine simulation is straightforward. *Small* Turing-complete tile sets via affine maps (via two-stack machines) are a little harder.

## 8.9 Combinator Calculus

In the 1920s, Moses Schönfinkel developed what can now be interpreted as a model of computation now called *combinator calculus* or *combinatory logic*. Combinator calculus operates on **terms**, where every term is either one of a finite number of **combinators** (represented here by upper case letters) or an ordered pair of terms. For notational convenience, we omit commas between components of every pair and parentheses around the *left* term in every pair. Thus, **SKK(IS)** is shorthand for the term  $((S, K), K), (I, S)$ .

We can “evaluate” any term by a sequence of rewriting rules that depend on its first primitive combinator. Schönfinkel defined three primitive combinators with the following evaluation rules:

- Identity: **I** $x \mapsto x$
- Constant: **K** $xy \mapsto x$
- Substitution: **S** $xyz \mapsto xz(yz)$

Here,  $x$ ,  $y$ , and  $z$  are variables representing unknown but arbitrary terms. “Computation” in the combinator calculus is performed by repeatedly evaluating arbitrary (sub)terms with one of these three structures, until all such (sub)terms are gone.

For example, the term **S(K(SI))Kxy** (for any terms  $x$  and  $y$ ) evaluates as follows:

$\underline{S(K(SI))Kxy} \mapsto \underline{K(SI)} \underline{x(Kx)y}$	Substitution
$\mapsto \underline{SI(Kx)} y$	Constant
$\mapsto \underline{Iy(Kxy)}$	Substitution
$\mapsto y(\underline{Kxy})$	Identity
$\mapsto yx$	Constant

Thus, we can define a new combinator **R**  $:= S(K(SI))K$  that upon evaluation reverses the next two terms: **R** $xy \mapsto yx$ .



On the other hand, evaluating  $SII(S(KI)(SII))$  leads to an infinite loop:

$\underline{SII(S(KI)(SII))} \mapsto \underline{I(S(KI)(SII))}(I(S(KI)(SII)))$	Substitution
$\mapsto S(KI)(SII)(\underline{I(S(KI)(SII))})$	Identity
$\mapsto \underline{S(KI)(SII)}(S(KI)(SII))$	Identity
$\mapsto \underline{KI(S(KI)(SII))}(SII(S(KI)(SII)))$	Substitution
$\mapsto \underline{I(SII(S(KI)(SII)))}$	Constant
$\mapsto SII(S(KI)(SII))$	Identity

Wikipedia sketches a direct *undecidability* proof. Is there a Turing-completeness proof that avoids  $\lambda$ -calculus?

## Exercises

1. A **tag**-Turing machine has two heads: one can only read, the other can only write. Initially, the read head is located at the left end of the tape, and the write head is located at the first blank after the input string. At each transition, the read head can either move one cell to the right or stay put, but the write head *must* write a symbol to its current cell and move one cell to the right. Neither head can ever move to the left.

Prove that any standard Turing machine can be simulated by a tag-Turing machine. That is, given any standard Turing machine  $M$ , describe a tag-Turing machine  $M'$  that accepts and rejects exactly the same input strings as  $M$ .

2.  $\star$ (a) Prove that any standard Turing machine can be simulated by a Turing machine with only three states. [Hint: Use the tape to store an encoding of the state of the machine yours is simulating.]
- $\star$ (b) Prove that any standard Turing machine can be simulated by a Turing machine with only two states.
3. A **two-dimensional** Turing machine uses an infinite two-dimensional grid of cells as the tape; at each transition, the head can move from its current cell to any of its four neighbors on the grid. The transition function of such a machine has the form  $\delta: Q \times \Gamma \rightarrow Q \times \Gamma \times \{\uparrow, \leftarrow, \downarrow, \rightarrow\}$ , where the arrows indicate which direction the head should move.

- (a) Prove that any two-dimensional Turing machine can be simulated by a standard Turing machine.
- (b) Suppose further that we endow our two-dimensional Turing machine with the following additional actions, in addition to moving the head:
  - Insert row: Move all symbols on or above the row containing the head up one row, leaving the head's row blank.
  - Insert column: Move all symbols on or to the right of the column containing the head one column to the right, leaving the head's column blank.
  - Delete row: Move all symbols above the row containing the head down one row, deleting the head's row of symbols.
  - Delete column: Move all symbols the right of the column containing the head one column to the right, deleting the head's column of symbols.

Show that any two-dimensional Turing machine that can add and delete rows can be simulated by a standard Turing machine.

4. A **binary-tree** Turing machine uses an infinite binary tree as its tape; that is, *every* cell in the tape has a left child and a right child. At each step, the head moves from its current cell to its **P**arent, its **L**eft child, or to its **R**ight child. Thus, the transition function of such a machine has the form  $\delta : Q \times \Gamma \rightarrow Q \times \Gamma \times \{P, L, R\}$ . The input string is initially given along the left spine of the tape.

Show that any binary-tree Turing machine can be simulated by a standard Turing machine.

5. A **stack-tape** Turing machine uses an semi-infinite tape, where every cell is actually the top of an independent stack. The behavior of the machine at each iteration is governed by its internal state and the symbol *at the top* of the current cell's stack. At each transition, the head can optionally push a new symbol onto the stack, or pop the top symbol off the stack. (If a stack is empty, its “top symbol” is a blank and popping has no effect.)

Show that any stack-tape Turing machine can be simulated by a standard Turing machine. (Compare with Problem ??!)

6. A **tape-stack** Turing machine has two actions that modify its work tape, in addition to simply writing individual cells: it can **save** the entire tape by pushing in onto a stack, and it can **restore** the entire tape by popping it off the stack. Restoring a tape returns the content of every cell to its content when the tape was saved. Saving and restoring the tape do not change the machine's state or the position of its head. If the machine attempts to “restore” the tape when the stack is empty, the machine crashes.

Show that any tape-stack Turing machine can be simulated by a standard Turing machine.

- **Tape alphabet** =  $\mathbb{N}$ .
  - Read: zero or positive. Write: +1, −1
  - Read: even or odd. Write: +1, −1, ×2, ÷2
  - Read: positive, negative, or zero. Write:  $x + y$  (merge),  $x - y$  (merge), 1, 0
- Never three times in a row in the same direction
- Hole-punch TM: tape alphabet  $\{\square, \blacksquare\}$ , and only  $\square \mapsto \blacksquare$  transitions allowed.

*If first you don't succeed, then try and try again.*

*And if you don't succeed again, just try and try and try.*

— Marc Blitzstein, “Useless Song”, *The Three Penny Opera* (1954)

Adaptation of Bertold Brecht, “Das Lied von der Unzulänglichkeit menschlichen Strebens” *Die Dreigroschenoper* (1928)

*Children need encouragement.*

*If a kid gets an answer right, tell him it was a lucky guess.*

*That way he develops a good, lucky feeling.*

— Jack Handey, “Deep Thoughts”, *Saturday Night Live* (March 21, 1992)

## 9 Nondeterministic Turing Machines

### 9.1 Definitions

In his seminal 1936 paper, Turing also defined an extension of his “automatic machines” that he called **choice machines**, which are now more commonly known as **nondeterministic Turing machines**. The execution of a nondeterministic Turing machine is *not determined* entirely by its input and its transition function; rather, at each step of its execution, the machine can *choose* from a set of possible transitions. The distinction between deterministic and nondeterministic Turing machines exactly parallels the distinction between deterministic and nondeterministic finite-state automata.

Formally, a nondeterministic Turing machine has all the components of a standard deterministic Turing machine—a finite tape alphabet  $\Gamma$  that contains the input alphabet  $\Sigma$  and a blank symbol  $\square$ ; a finite set  $Q$  of internal states with special **start**, **accept**, and **reject** states; and a transition function  $\delta$ . However, the transition function now has the signature

$$\delta: Q \times \Gamma \rightarrow 2^{Q \times \Gamma \times \{-1, +1\}}.$$

That is, for each state  $p$  and tape symbol  $a$ , the output  $\delta(p, a)$  of the transition function is a *set* of triples of the form  $(q, b, \Delta) \in Q \times \Gamma \times \{-1, +1\}$ . Whenever the machine finds itself in state  $p$  reading symbol  $a$ , the machine *chooses* an arbitrary triple  $(q, b, \Delta) \in \delta(p, a)$ , and then changes its state to  $q$ , writes  $b$  to the tape, and moves the head by  $\Delta$ . If the set  $\delta(p, a)$  is empty, the machine moves to the **reject** state and halts.

The set of all possible transition sequences of a nondeterministic Turing machine  $N$  on a given input string  $w$  define a rooted tree, called a **computation tree**. The initial configuration (**start**,  $w, 0$ ) is the root of the computation tree, and the children of any configuration  $(q, x, i)$  are the configurations that can be reached from  $(q, x, i)$  in one transition. In particular, any configuration whose state is **accept** or **reject** is a leaf. For deterministic Turing machines, this computation tree is just a single path, since there is at most one valid transition from every configuration.

## 9.2 Acceptance and Rejection

Unlike deterministic Turing machines, there is a fundamental asymmetry between the acceptance and rejection criteria for nondeterministic Turing machines. Let  $N$  be any nondeterministic Turing machine, and let  $w$  be any string.

- $N$  **accepts**  $w$  if and only if there is *at least one* sequence of valid transitions from the initial configuration (**start**,  $w$ , 0) that leads to the **accept** state. Equivalently,  $N$  accepts  $w$  if the computation tree contains at least one **accept** leaf.
- $N$  **rejects**  $w$  if and only if *every* sequence of valid transitions from the initial configuration (**start**,  $w$ , 0) leads to the **reject** state. Equivalently,  $N$  rejects  $w$  if every path through the computation tree ends with a **reject** leaf.

In particular,  $N$  can accept  $w$  even when there are choices that allow the machine to run forever, but rejection requires  $N$  to halt after only a finite number of transitions, no matter what choices it makes along the way. Just as for deterministic Turing machines, it is possible that  $N$  neither accepts nor rejects  $w$ .

Acceptance and rejection of *languages* are defined exactly as they are for deterministic machines. A non-deterministic Turing machine  $N$  **accepts** a language  $L \subseteq \Sigma^*$  if  $N$  accepts all strings in  $L$  and nothing else;  $N$  **rejects**  $L$  if  $N$  rejects every string in  $L$  and nothing else; and finally,  $N$  **decides**  $L$  if  $N$  accepts  $L$  and rejects  $\Sigma^* \setminus L$ .

## 9.3 Time and Space Complexity

- Define “time” and “space”.
- $\text{TIME}(f(n))$  is the class of languages that can be decided by a deterministic multi-tape Turing machine in  $O(f(n))$  time.
- $\text{NTIME}(f(n))$  is the class of languages that can be decided by a nondeterministic multi-tape Turing machine in  $O(f(n))$  time.
- $\text{SPACE}(f(n))$  is the class of languages that can be decided by deterministic multi-tape Turing machine in  $O(f(n))$  space.
- $\text{NSPACE}(f(n))$  is the class of languages that can be decided by a nondeterministic multi-tape Turing machine in  $O(f(n))$  space.
- Why multi-tape TMs? Because  $t$  steps on any  $k$ -tape Turing machine can be simulated in  $O(t \log t)$  steps on a two-tape machine [Hennie and Stearns 1966, essentially using lazy counters and amortization], and in  $O(t^2)$  steps on a single-tape machine [Hartmanis and Stearns 1965; realign multiple tracks at every simulation step]. Moreover, the latter quadratic bound is tight [Hennie 1965 (palindromes, via communication complexity)].

## 9.4 Deterministic Simulation via Dovetailing

**Theorem 1.** For any nondeterministic Turing machine  $N$ , there is a deterministic Turing machine  $M$  that accepts exactly the same strings and  $N$  and rejects exactly the same strings as  $N$ . Moreover, if every computation path of  $N$  on input  $x$  halts after at most  $t$  steps, then  $M$  halts on input  $x$  after at most  $O(t^2 r^{2t})$  steps, where  $r$  is the maximum size of any transition set in  $N$ .

**Proof:** I’ll describe a deterministic machine  $M$  that performs a breadth-first search of the computation tree of  $N$ . (The depth-first search performed by a standard recursive backtracking algorithm won’t work here. If  $N$ ’s computation tree contains an infinite path, a depth-first search would get stuck in that path without exploring the rest of the tree.)

At the beginning of each simulation round,  $M$ 's tape contains a string of the form

$$\square\square\cdots\square\bullet\bullet y_1q_1z_1\bullet y_2q_2z_2\bullet\cdots\bullet y_kq_kz_k\bullet\bullet$$

where each substring  $y_iq_iz_i$  encodes a configuration  $(q_i, y_iz_i, |y_i|)$  of some computation path of  $N$ , and  $\bullet$  is a new symbol not in the tape alphabet of  $N$ . The machine  $M$  interprets this sequence of encoded configurations as a queue, with new configurations inserted on the right and old configurations removed from the left. The double-separators  $\bullet\bullet$  uniquely identify the start and end of this queue; outside this queue, the tape is entirely blank.

Specifically, in each round, first  $M$  appends the encodings of all configurations than  $N$  can reach in one transition from the first encoded configuration  $(q_1, y_1z_1, |y_1|)$ ; then  $M$  erases the first encoded configuration.

$$\begin{array}{c} \cdots\square\square\bullet\bullet y_1q_1z_1\bullet y_2q_2z_2\bullet\cdots\bullet y_rq_rz_r\bullet\bullet\square\square\cdots \\ \Downarrow \qquad\qquad\qquad \Downarrow \\ \cdots\square\square\square\square\bullet\bullet y_2q_2z_2\bullet\cdots\bullet y_kq_kz_k\bullet\tilde{y}_1\tilde{q}_1\tilde{z}_1\bullet\tilde{y}_2\tilde{q}_2\tilde{z}_2\bullet\cdots\bullet\tilde{y}_r\tilde{q}_r\tilde{z}_r\bullet\bullet\square\square\cdots \end{array}$$

Suppose each transition set  $\delta_N(q, a)$  has size at most  $r$ . Then after simulating  $t$  steps of  $N$ , the tape string of  $M$  encoding  $O(r^t)$  different configurations of  $N$  and therefore has length  $L = O(tr^t)$  (not counting the initial blanks). If  $M$  begins each simulation phase by moving the initial configuration from the beginning to the end of the tape string, which takes  $O(t^2r^t)$  time, the time for the rest of the the simulation phase is negligible. Altogether, simulating all  $r^t$  possibilities for the the  $t$ th step of  $N$  requires  $O(t^2r^{2t})$  time. We conclude that  $M$  can simulate the first  $t$  steps of every computation path of  $N$  in  $O(t^2r^{2t})$  time, as claimed.  $\square$

The running time of this simulation is dominated by the time spent reading from one end of the tape string and writing to the other. It is fairly easy to reduce the running time to  $O(tr^t)$  by using either two tapes (a “read tape” containing  $N$ -configurations at time  $t$  and a “write tape” containing  $N$ -configurations at time  $t + 1$ ) or two independent heads on the same tape (one at each end of the queue).

## 9.5 Nondeterminism as Advice

Any nondeterministic Turing machine  $N$  can also be simulated by a *deterministic* machine  $M$  with *two* inputs: the user input string  $w \in \Sigma^*$ , and a so-called **advice** string  $x \in \Omega^*$ , where  $\Omega$  is another finite alphabet. Only the first input string  $w$  is actually given by the user. At least for now, we assume that the advice string  $x$  is given on a separate read-only tape.

The deterministic machine  $M$  simulates  $N$  step-by-step, but whenever  $N$  has a choice of how to transition,  $M$  reads a new symbol from the advice string, and that symbol determines the choice. In fact, without loss of generality, we can assume that  $M$  reads a new symbol from the advice string and moves the advice-tape's head to the right on *every* transition. Thus,  $M$ 's transition function has the form  $\delta_M: Q \times \Gamma \times \Omega \rightarrow Q \times \Gamma \times \{-1, +1\}$ , and we require that

$$\delta_N(q, a) = \{\delta_M(q, a, \omega) \mid \omega \in \Omega\}$$

For example, if  $N$  has a binary choice

$$\delta_N(\text{branch}, ?) = \{(\text{left}, L, -1), (\text{right}, R, +1)\},$$

then  $M$  might determine this choice by defining

$$\delta_M(\text{branch}, ?, 0) = (\text{left}, L, -1) \quad \text{and} \quad \delta_M(\text{branch}, ?, 1) = (\text{right}, R, +1)$$

More generally, if every set  $\delta_N(p, a)$  has size  $r$ , then we let  $\Omega = \{1, 2, \dots, r\}$  and define  $\delta_M(q, a, i)$  to be the  $i$ th element of  $\delta_N(q, a)$  in some canonical order.

Now observe that  $N$  accepts a string  $w$  if and only if  $M$  accepts the pair  $(w, x)$  for *some* string  $x \in \Omega^*$ , and  $N$  rejects  $w$  if and only if  $M$  rejects the pair  $(w, x)$  for *all* strings  $x \in \Omega^*$ .

The “advice” formulation of nondeterminism allows a different strategy for simulation by a standard deterministic Turing machine, which is often called **dovetailing**. Consider all possible advice strings  $x$ , in increasing order of length; listing these advice strings is equivalent to repeatedly incrementing a base- $r$  counter. For each advice string  $x$ , simulate  $M$  on input  $(w, x)$  for exactly  $|x|$  transitions.

```
DOVETAILM(w):
  for t ← 1 to ∞
    done ← TRUE
    for all strings x ∈ Ωt
      if M accepts (w, x) in at most t steps
        accept
      if M(w, x) does not halt in at most t steps
        done ← FALSE
    if done
      reject
```

The most straightforward Turing-machine implementation of this algorithm requires three tapes: A read-only input tape containing  $w$ , an advice tape containing  $x$  (which is also used as a timer for the simulation), and the work tape. This simulation requires  **$O(tr^t)$  time** to simulate all possibilities for  $t$  steps of the original non-deterministic machine  $N$ .

If we insist on using a standard Turing machine with a single tape and a single head, the simulation becomes slightly more complex, but (unlike our earlier queue-based strategy) not significantly slower. This standard machine  $S$  maintains a string of the form  $\bullet w \bullet x \bullet z$ , where  $z$  is the current work-tape string of  $M$  (or equivalently, of  $N$ ), with marks (on a second track) indicating the current positions of the heads on  $M$ ’s work tape and  $M$ ’s advice tape. Simulating a single transition of  $M$  now requires  $O(|x|)$  steps, because  $S$  needs to shuttle its single head between these two marks. Thus,  $S$  requires  **$O(t^2 r^t)$  time** to simulate all possibilities for  $t$  steps of the original non-deterministic machine  $N$ . This is significantly faster than the queue-based simulation, because we don’t record (and therefore don’t have to repeatedly scan over) intermediate configurations; recomputing everything from scratch is actually cheaper!

## 9.6 The Cook-Levin Theorem

Define SAT and CIRCUITSAT. Non-determinism is fundamentally different from other Turing machine extensions, in that it seems to provide an exponential speedup for some problems, just like NFAs can use exponentially fewer states than DFAs for the same language.

**The Cook-Levin Theorem.** *If SAT ∈ P, then P=NP.*

**Proof:** Let  $L \subseteq \Sigma^*$  be an arbitrary language in NP, over some fixed alphabet  $\Sigma$ . There must be an integer  $k$  and Turing machine  $M$  that satisfies the following conditions:

- For all strings  $w \in L$ , there is at least one string  $x \in \Sigma^*$  such that  $M$  accepts the string  $w \square x$ .
- For all strings  $w \notin L$  and  $x \in \Sigma^*$ ,  $M$  rejects the string  $w \square x$ .
- For all strings  $w, x \in \Sigma^*$ ,  $M$  halts on input  $w \square x$  after at most  $\max\{1, |w|^k\}$  steps.

Now suppose we are given a string  $w \in \Sigma^*$ . Let  $n = |w|$  and let  $N = \max\{1, |w|^k\}$ . We construct a boolean formula  $\Phi_w$  that is satisfiable if and only if  $w \in L$ , by following the execution of  $M$  on input  $w \square x$  for some unknown advice string  $x$ . Without loss of generality, we can assume that  $|x| = N - n - 1$  (since we can extend any shorter string  $x$  with blanks.) Our formula  $\Phi_w$  uses the following boolean variables for all symbols  $a \in \Gamma$ , all states  $q \in Q$ , and all integers  $0 \leq t \leq N$  and  $0 \leq i \leq N + 1$ .

- $Q_{t,i,q}$  —  $M$  is in state  $q$  with its head at position  $i$  after  $t$  transitions.
- $T_{t,i,a}$  — The  $k$ th cell of  $M$ 's work tape contains  $a$  after  $t$  transitions.

The formula  $\Phi_w$  is the conjunction of the following constraints:

- **Boundaries:** To simplify later constraints, we include artificial boundary variables just past both ends of the tape:

$$\begin{aligned} Q_{t,i,q} = Q_{t,N+1,q} &= \text{FALSE} && \text{for all } 0 \leq t \leq N \text{ and } q \in Q \\ T_{t,0,a} = T_{t,N+1,a} &= \text{FALSE} && \text{for all } 0 \leq t \leq N \text{ and } a \in \Gamma \end{aligned}$$

- **Initialization:** We have the following values for variables with  $t = 0$ :

$$\begin{aligned} Q_{0,1,\text{start}} &= \text{TRUE} \\ Q_{0,1,q} &= \text{FALSE} && \text{for all } q \neq \text{start} \\ H_{0,i,q} &= \text{FALSE} && \text{for all } i \neq 1 \text{ and } q \in Q \\ T_{0,i,w_i} &= \text{TRUE} && \text{for all } 1 \leq i \leq n \\ T_{0,i,a} &= \text{FALSE} && \text{for all } 1 \leq i \leq n \text{ and } a \neq w_i \\ T_{0,n+1,\square} &= \text{TRUE} \\ T_{0,n+1,a} &= \text{FALSE} && \text{for all } a \neq \square \end{aligned}$$

- **Uniqueness:** The variables  $T_{0,i,a}$  with  $n+2 \leq i \leq N$  represent the unknown advice string  $x$ ; these are the “inputs” to  $\Phi_w$ . We need some additional constraints ensure that for each  $i$ , *exactly one* of these variables is TRUE:

$$\left( \bigvee_{a \in \Gamma} T_{0,j,a} \right) \wedge \bigwedge_{a \neq b} (\overline{T_{0,j,a}} \vee \overline{T_{0,j,b}})$$

- **Transitions:** For all  $1 \leq t \leq N$  and  $1 \leq i \leq N$ , the following constraints simulate the transition from time  $t-1$  to time  $t$ .

$$\begin{aligned} Q_{t,i,q} &= \bigvee_{\delta(p,a)=(q, \cdot, +1)} (Q_{t-1,i-1,p} \wedge T_{t-1,i,a}) \vee \bigvee_{\delta(p,a)=(q, \cdot, -1)} (Q_{t-1,i+1,p} \wedge T_{t-1,i,a}) \\ T_{t,i,b} &= \bigvee_{\delta(p,a)=(\cdot, b, \cdot)} (Q_{t-1,i,p} \wedge T_{t-1,i,a}) \vee \left( \bigwedge_{q \in Q} \overline{Q_{t-1,i,q}} \wedge T_{t-1,i,b} \right) \end{aligned}$$

- **Output:** We have one final constraint that indicates acceptance.

$$z = \bigvee_{t=0}^N \bigvee_{i=1}^N Q_{t,i,\text{accept}}$$

By definition,  $\Phi_w$  is satisfiable if and only if some input values  $T_{0,i,a}$ , all constraints are satisfied, including acceptance. A straightforward induction argument implies that *even without the acceptance constraint*, any assignment of values to the unknown variables  $T_{0,i,a}$  that satisfies the uniqueness constraints determines *unique* values for the other variables in  $\Phi_w$ , which consistently describe the execution of  $M$ . To satisfy the acceptance constraint, this execution of  $M$  must lead to the **accept** state. Thus,  $\Phi_w$  is satisfiable if and only if there is a string  $x \in \Gamma^*$  such that  $M$  accepts the input  $w \sqcup x$ . We conclude that  $\Phi_w$  is satisfiable if and only if  $w \in L$ .

It remains only to argue that the reduction requires only polynomial time. For any input string  $w$  of length  $n$ , the formula  $\Phi_w$  has  $O(N^2)$  variables and  $O(N^2)$  constraints (where the hidden constants depend on the machine  $M$ ). Every constraint except acceptance has constant length, so altogether  $\Phi_w$  has length  $O(N^2)$ . Moreover, we can construct  $\Phi_w$  in  $O(N^2) = O(n^{2k})$  time.

In conclusion, if we could decide SAT for formulas of size  $M$  in  $O(M^c)$  time, then we could decide membership in  $L$  in  $O(n^{2kc})$  time, which implies that  $L \in P$ .  $\square$

## Exercises

1. Prove that the following problem is NP-hard, *without* using the Cook-Levin Theorem. Given a string  $\langle M, w \rangle$  that encodes a non-deterministic Turing machine  $M$  and a string  $w$ , does  $M$  accept  $w$  in at most  $|w|$  transitions?

More exercises!