

**Caveat lector:** This is the zeroth (draft) edition of this lecture note. Please send bug reports and suggestions to jeffe@illinois.edu.

*I said in my haste, All men are liars.*

— Psalms 116:11 (King James Version)

*yields falsehood when preceded by its quotation.*

— William V. Quine, “Paradox”, *Scientific American* (1962)

*Some problems are so complex that you have to be highly intelligent and well informed just to be undecided about them.*

— Laurence Johnston Peter, *Peter’s Almanac* (September 24, 1982)

*“Proving or disproving a formula—once you’ve encrypted the formula into numbers, that is—is just a calculation on that number. So it means that the answer to the question is, no! Some formulas cannot be proved or disproved by any mechanical process! So I guess there’s some point in being human after all!”*

*Alan looked pleased until Lawrence said this last thing, and then his face collapsed. “Now there you go making unwarranted assumptions.”*

— Neal Stephenson, *Cryptonomicon* (1999)

*No matter how P might perform, Q will scoop it:*

*Q uses P’s output to make P look stupid.*

*Whatever P says, it cannot predict Q:*

*P is right when it’s wrong, and is false when it’s true!*

— Geoffrey S. Pullum, “[Scooping the Loop Sniffer](#)” (2000)

*This castle is in unacceptable condition! **UNACCEPTABLE!!***

— Earl of Lemongrab [Justin Poiland], “Too Young”  
*Adventure Time* (August 8, 2011)

## 37 Undecidability

Perhaps the single most important result in Turing’s remarkable 1936 paper is his solution to Hilbert’s *Entscheidungsproblem*, which asked for a general automatic procedure to determine whether a given statement of first-order logic is *provable*. Turing proved that no such procedure exists; there is no systematic way to distinguish between statements that cannot be proved even in principle and statements whose proofs we just haven’t found yet.

### 37.1 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.

### 37.2 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 1s. 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!

### 37.3 Useful Lemmas

This subsection contains several lemmas that are useful for proving that languages are (un)decidable or (un)acceptable. For almost all of these lemmas, the proofs are straightforward; 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 “haltable” languages. The following lemma, whose proof is an easy exercise (hint, hint), implies that these 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 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.*

The asymmetry between acceptance and rejection implies that merely acceptable languages are not quite as well-behaved as decidable languages.

**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.  $\square$

**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.  $\square$

### 37.4 Self-Haters Gonna Self-Hate

Let  $U$  be an arbitrary fixed universal Turing machine. Any Turing machine  $M$  can be encoded as a string  $\langle M \rangle$  of symbols from  $U$ 's input alphabet, so that  $U$  can simulate the execution of  $M$  on any suitably encoded input string. Different universal Turing machines require different encodings.<sup>1</sup>

A Turing machine encoding is just a string, and any string (over the correct alphabet) can be used as the input to a Turing machine. Thus, we can use the encoding  $\langle M \rangle$  of any Turing machine  $M$  as the input to another Turing machine. We've already seen an example of this ability in our universal Turing

<sup>1</sup>In fact, these undecidability proofs never actually use the universal Turing machine; all we really need is an encoding function that associates a unique string  $\langle M \rangle$  with every Turing machine  $M$ . However, we *do* need the encoding to be compatible with a universal Turing machine for the results in Section 37.7.

machine  $U$ , but more significantly, we can use  $\langle M \rangle$  as the input to *the same Turing machine*  $M$ . Thus, each of the following languages is well-defined:

$$\begin{aligned}\text{SELFACCEPT} &:= \{ \langle M \rangle \mid M \text{ accepts } \langle M \rangle \} \\ \text{SELFREJECT} &:= \{ \langle M \rangle \mid M \text{ rejects } \langle M \rangle \} \\ \text{SELFHALT} &:= \{ \langle M \rangle \mid M \text{ halts on } \langle M \rangle \} \\ \text{SELFDIVERGE} &:= \{ \langle M \rangle \mid M \text{ diverges on } \langle M \rangle \}\end{aligned}$$

One of Turing's key observations is that **SELFREJECT is undecidable**; Turing proved this theorem by contradiction as follows:

Suppose to the contrary that there is a Turing machine  $SR$  such that  $\text{ACCEPT}(SR) = \text{SELFREJECT}$  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 equal to  $SR$ . 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!

### 37.5 Aside: Uncountable Barbers

Turing's proof by contradiction is nearly identical to 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>2</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. □

<sup>2</sup>more commonly, flouting all reasonable standards of grammatical English, "an onto function"

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  $\langle B \rangle$ ".

### 37.6 Just Don't Know What to Do with Myself

Similar diagonal arguments imply that the other three languages `SELFACCEPT`, `SELFHALT`, and `SELF-DIVERGE` are also undecidable. The proofs are not quite as direct for these three languages as the proof for `SELFREJECT`; each fictional deciding machine requires a small modification to create the contradiction.

**Theorem 6.** *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 7.** *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 8.** *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.  $\square$

### \*37.7 Nevertheless, Acceptable

Our undecidability argument for SELFDIVERGE actually implies the stronger result that SELFDIVERGE 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 to accept the other three languages.

**Theorem 9.** *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>3</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 10.** *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.

### 37.8 The Halting Problem via Reduction

Consider the following related languages:<sup>4</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 what is usually meant by 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.

<sup>3</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  $\langle w \rangle$  for  $U$ ; building such a Turing machine is straightforward.

<sup>4</sup>Sipser uses the shorter name  $A_{TM}$  instead of ACCEPT, but uses  $HALT_{TM}$  instead of HALT. I have no idea why he thought four-letter names are okay, but six-letter names are not. His 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.

**Theorem 11.** *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.

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 machine” or “that machine” 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).

### 37.9 One Million Years Dungeon!

As a more complex set of examples, consider the following languages:

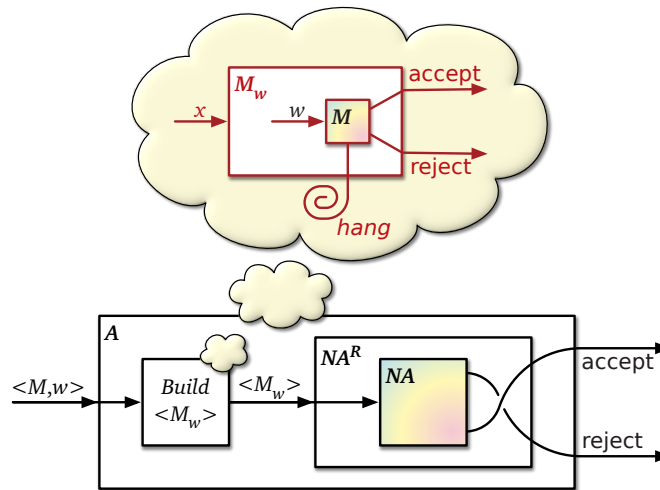
$$\begin{aligned} \text{NEVERACCEPT} &:= \{ \langle M \rangle \mid \text{ACCEPT}(M) = \emptyset \} \\ \text{NEVERREJECT} &:= \{ \langle M \rangle \mid \text{REJECT}(M) = \emptyset \} \\ \text{NEVERHALT} &:= \{ \langle M \rangle \mid \text{HALT}(M) = \emptyset \} \\ \text{NEVERDIVERGE} &:= \{ \langle M \rangle \mid \text{DIVERGE}(M) = \emptyset \} \end{aligned}$$

**Theorem 12.** *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  $\Sigma^* \setminus NEVERACCEPT$ .

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.



A reduction from from  $ACCEPT$  to  $NEVERACCEPT$ , which proves  $NEVERACCEPT$  undecidable.

Before going any further, it may be helpful to list the various Turing machines that appear in this construction.

- The hypothetical Turing machine  $NA$  that decides  $NEVERACCEPT$ .
- The Turing machine  $NA^R$  that decides  $\Sigma^* \setminus NEVERACCEPT$ , 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  $ACCEPT(M_w) = \Sigma^*$ , by the definition of  $ACCEPT(M_w)$ .
  - So  $\langle M_w \rangle \notin NEVERACCEPT$ , 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  $\text{NEVERACCEPT}$ .
  - So  $NA$  accepts  $\langle M_w \rangle$ , because  $NA$  decides  $\text{NEVERACCEPT}$ .
  - 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  $\text{ACCEPT}$ , which is impossible. We conclude that  $NA$  does not exist.  $\square$

Again, similar arguments imply that the languages  $\text{NEVERREJECT}$ ,  $\text{NEVERHALT}$ , and  $\text{NEVERDIVERGE}$  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  $\text{NEVERACCEPT}$  and its relatives are undecidable, we can use them as the basis of further reduction proofs. Here is a typical example:

**Theorem 13.** *The language  $\text{DIVERGESAME} := \{ \langle M_1 \rangle \langle M_2 \rangle \mid \text{DIVERGE}(M_1) = \text{DIVERGE}(M_2) \}$  is undecidable.*

**Proof:** Suppose for the sake of argument that there is a Turing machine  $DS$  that decides  $\text{DIVERGESAME}$ . Then we can build a Turing machine  $ND$  that decides  $\text{NEVERDIVERGE}$  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  $\text{DIVERGE}(M) = \text{DIVERGE}(Y) = \emptyset$ , so  $\langle M \rangle \in \text{NEVERDIVERGE}$ .
- If  $DS$  rejects  $\langle M \rangle \langle Y \rangle$ , then  $\text{DIVERGE}(M) \neq \text{DIVERGE}(Y) = \emptyset$ , 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.  $\square$

### 37.10 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.

**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  $\text{HALT}$ , 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. 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, most proofs 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 14.** *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  $\text{ACCEPTIN}(\mathcal{L}) = \{\langle M \rangle \mid M \text{ accept given an empty initial tape}\}$ .

- Given an empty initial tape,  $M_{\text{ACCEPT}}$  accepts, so  $\text{HALT}(M_{\text{ACCEPT}}) \in \mathcal{L}$ .
- Given an empty initial tape,  $M_{\text{DIVERGE}}$  does not accept, so  $\text{HALT}(M_{\text{DIVERGE}}) \notin \mathcal{L}$ .

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  $\text{HALT}(M_{\text{ACCEPT}}) \in \mathcal{L}$ .
- $M_{\text{DIVERGE}}$  does not accept **UIUC**, so  $\text{HALT}(M_{\text{DIVERGE}}) \notin \mathcal{L}$ .

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  $\{\text{larry, curly, moe}\}$ . 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>5</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. The Turing machine  $N^R$ , obtained by swapping the **accept** and **reject** states of  $N$ , accepts the language  $\text{HALT}(L) \setminus \text{ACCEPT}(N) \neq \text{ACCEPT}(N)$ .  $\square$

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 15.** *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.)

**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:*

<sup>5</sup>Yes, parts (e) and (f) have exactly the same proof.

- There is a Turing machine  $Y$  such that  $DIVERGE(Y) \in \mathcal{L}$
- There is a Turing machine  $N$  such that  $DIVERGE(N) \notin \mathcal{L}$ .

The language  $DIVERGEIN(\mathcal{L}) := \{\langle M \rangle \mid DIVERGE(M) \in \mathcal{L}\}$  is undecidable.

**Rice's Decision Theorem.** Let  $\mathcal{L}$  be any set of languages that satisfies the following conditions:

- There is a Turing machine  $Y$  such that **decides** an language in  $\mathcal{L}$ .
- There is a Turing machine  $N$  such that **decides** an language not in  $\mathcal{L}$ .

The language  $DECIDEIN(\mathcal{L}) := \{\langle M \rangle \mid M \text{ decides a language in } \mathcal{L}\}$  is undecidable.

As a final sanity check, always be careful to distinguish the following objects:

- The string  $\varepsilon$
- The language  $\emptyset$
- The language  $\{\varepsilon\}$
- The language property  $\emptyset$
- The language property  $\{\emptyset\}$
- The language property  $\{\{\varepsilon\}\}$
- The Turing machine  $M_{\text{REJECT}}$  that rejects every string and therefore **decides** the language  $\emptyset$ .
- The Turing machine  $M_{\text{DIVERGE}}$  that diverges on every string and therefore **accepts** the language  $\emptyset$ .

### \*37.11 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>6</sup>

**The Rice-McNaughton-Myhill-Shapiro Theorem.** Let  $\mathcal{L}$  be an arbitrary set of acceptable languages. The language  $ACCEPTIN(\mathcal{L}) := \{\langle M \rangle \mid ACCEPT(M) \in \mathcal{L}\}$  is **acceptable** if and only if  $\mathcal{L}$  satisfies the following conditions:

- $\mathcal{L}$  is **monotone**: For any language  $L \in \mathcal{L}$ , every superset of  $L$  is also in  $\mathcal{L}$ .
- $\mathcal{L}$  is **compact**: Every language in  $\mathcal{L}$  has a finite subset that is also in  $\mathcal{L}$ .
- $\mathcal{L}$  is **finitely acceptable**: The language  $\{\langle L \rangle \mid L \in \mathcal{L} \text{ and } L \text{ is finite}\}$  is acceptable.<sup>7</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 16.** Let  $\mathcal{L}$  be a set of acceptable languages. If  $\mathcal{L}$  is not monotone, then  $ACCEPTIN(\mathcal{L})$  is unacceptable.

<sup>6</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.

<sup>7</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  $\Sigma \cup \{\{, \cdot, \}, \varepsilon\}$  that contains the strings in  $L$  in lexicographic order, separated by commas  $\cdot$ , and surrounded by braces  $\{ \}$ , with  $\varepsilon$  representing the empty string. For example,  $\langle \{\varepsilon, 0, 01, 0110, 01101001\} \rangle = \{\varepsilon, 0, 01, 0110, 01101001\}$ .

**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. □

**Corollary 17.** *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)  $\{\langle M \rangle \mid \text{ACCEPT}(M) \text{ is decidable}\}$
- (f)  $\{\langle M \rangle \mid \text{ACCEPT}(M) \text{ is undecidable}\}$
- (g)  $\{\langle M \rangle \mid M \text{ accepts at least one string in } \text{SELF DIVERGE}\}$
- (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 `SELFREJECT`.
- (f) The set of undecidable acceptable languages is not monotone: Consider the languages `SELFREJECT` 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  $\Sigma^*$ !  $\square$

### 37.12 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 18.** *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.  $\square$

**Theorem 19.** *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

- Let  $M$  be an arbitrary Turing machine.
  - Describe a Turing machine  $M^R$  such that  $\text{ACCEPT}(M^R) = \text{REJECT}(M)$  and  $\text{REJECT}(M^R) = \text{ACCEPT}(M)$ .
  - Describe is a Turing machine  $M^A$  such that  $\text{ACCEPT}(M^A) = \text{ACCEPT}(M)$  and  $\text{REJECT}(M^A) = \emptyset$ .
  - Describe is a Turing machine  $M^H$  such that  $\text{ACCEPT}(M^H) = \text{HALT}(M)$  and  $\text{REJECT}(M^H) = \emptyset$ .
- Prove that **ACCEPT** is undecidable.
  - Prove that **REJECT** is undecidable.
  - Prove that **DIVERGE** is undecidable.
- Prove that **NEVERREJECT** is undecidable.
  - Prove that **NEVERHALT** is undecidable.
  - Prove that **NEVERDIVERGE** is undecidable.
- Prove that each of the following languages is undecidable.
  - $\text{ALWAYSACCEPT} := \{\langle M \rangle \mid \text{ACCEPT}(M) = \Sigma^*\}$
  - $\text{ALWAYSREJECT} := \{\langle M \rangle \mid \text{REJECT}(M) = \Sigma^*\}$
  - $\text{ALWAYSHALT} := \{\langle M \rangle \mid \text{HALT}(M) = \Sigma^*\}$
  - $\text{ALWAYSDIVERGE} := \{\langle M \rangle \mid \text{DIVERGE}(M) = \Sigma^*\}$
- Let  $\mathcal{L}$  be a non-empty proper subset of the set of acceptable languages. Prove that the following languages are undecidable:
  - $\text{REJECTIN}(\mathcal{L}) := \{\langle M \rangle \mid \text{REJECT}(M) \in \mathcal{L}\}$
  - $\text{HALTIN}(\mathcal{L}) := \{\langle M \rangle \mid \text{HALT}(M) \in \mathcal{L}\}$
  - $\text{DIVERGEIN}(\mathcal{L}) := \{\langle M \rangle \mid \text{DIVERGE}(M) \in \mathcal{L}\}$
- 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  $M$ .
  - Does  $M$  accept  $\langle M \rangle^R$ ?
  - Does  $M$  reject any palindrome?
  - Does  $M$  accept all palindromes?



- (d) Does  $M$  diverge only on palindromes?
- (e) Is there an input string that forces  $M$  to move left?
- (f) Is there an input string that forces  $M$  to move left three times in a row?
- (g) Does  $M$  accept the encoding of any Turing machine  $N$  such that  $\text{ACCEPT}(N) = \text{SELF DIVERGE}$ ?
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  $\mathbf{1}$ ?
- (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:
- $\{ \langle M, w \rangle \mid M \text{ accepts } w \text{ in time } |w|^2 \}$
  - $\{ \langle M, w \rangle \mid M \text{ accepts } w \text{ in space } |w|^2 \}$
- (b) Prove that the following languages are undecidable:
- $\{ \langle M \rangle \mid M \text{ accepts at least one string } w \text{ in time } |w|^2 \}$
  - $\{ \langle M \rangle \mid M \text{ accepts at least one string } w \text{ in space } |w|^2 \}$

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.
- Does this program correctly compute Fibonacci numbers?
  - Can this program fall into an infinite loop?
  - Will the value of this variable ever change?
  - Will this program every attempt to dereference a null pointer?
  - Does this program free every block of memory that it dynamically allocates?
  - Is any statement in this program unreachable?
  - 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  $b \neq \Sigma$ .

- Prove that if  $M$  is a conservative Turing machine, then  $\text{ACCEPT}(M)$  is a regular language.
- 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?]*