

These notes are adapted from [10].

1 Turing Machines and Effective Computability

In these notes we will introduce Turing machines (TMs), named after Alan Turing, who invented them in 1936. Turing machines can compute any function normally considered computable; in fact, it is quite reasonable to define *computable* to mean computable by a TM.

TMs were invented in the 1930s, long before real computers appeared. They came at a time when mathematicians were trying to come to grips with the notion of *effective computation*. They knew various algorithms for computing things effectively, but they weren't quite sure how to define "effectively computable" in a general way that would allow them to distinguish between the computable and the noncomputable. Several alternative formalisms evolved, each with its own peculiarities, in the attempt to nail down this notion:

- Turing machines (Alan Turing [20]);
- Post systems (Emil Post [12, 13]);
- μ -recursive functions (Kurt Gödel [7], Jacques Herbrand);
- λ -calculus (Alonzo Church [2], Stephen C. Kleene [8]); and
- combinatory logic (Moses Schönfinkel [18], Haskell B. Curry [4]).

All of these systems embody the idea of *effective computation* in one form or another. They work on various types of data; for example, Turing machines manipulate strings over a finite alphabet, μ -recursive functions manipulate the natural numbers, the λ -calculus manipulates λ -terms, and combinatory logic manipulates terms built from combinator symbols.

However, there are natural translations between all these different types of data. For example, let $\{0, 1\}^*$ denote the set of all finite-length strings over the alphabet $\{0, 1\}$. There is a simple one-to-one correspondence between $\{0, 1\}^*$ and the natural numbers $\mathbb{N} = \{0, 1, 2, \dots\}$ defined by

$$x \mapsto \#(1x) - 1, \tag{1}$$

where $\#y$ is the natural number represented by the binary string y . Conversely, it is easy to encode just about anything (natural numbers, λ -terms, strings in $\{0, 1, 2, \dots, 9\}^*$, trees, graphs, ...) as strings in $\{0, 1\}^*$. Under these natural encodings of the data, it turns out that all the formalisms above can simulate one another, so despite their superficial differences they are all computationally equivalent.

Nowadays we can take unabashed advantage of our more modern perspective and add programming languages such as Java or C (or idealized versions of them) to this list—a true luxury compared to what Church and Gödel had to struggle with.

Of the classical systems listed above, the one that most closely resembles a modern computer is the Turing machine. Besides the off-the-shelf model we will define below, there are also many custom variations (non-deterministic, multitape, multidimensional tape, two-way infinite tapes, and so on) that all turn out to be computationally equivalent in the sense that they can all simulate one another.

Church's Thesis

Because these vastly dissimilar formalisms are all computationally equivalent, the common notion of computability that they embody is extremely robust, which is to say that it is invariant under fairly radical perturbations in the model. All these mathematicians with their pet systems turned out to be looking at the same thing from different angles. This was too striking to be mere coincidence. They soon came to the realization that the commonality among all these systems must be the elusive notion of *effective computability* that they had sought for so long. Computability is not just Turing machines, nor the λ -calculus, nor the μ -recursive functions, nor the Pascal programming language, but the common spirit embodied by them all.

Alonzo Church [3] gave voice to this thought, and it has since become known as *Church's thesis* (or the *Church-Turing thesis*). It is not a theorem, but rather a declaration that all these formalisms capture precisely our intuition about what it means to be effectively computable in principle, no more and no less. Church's thesis may not seem like such a big deal in retrospect, since by now we are thoroughly familiar with the capabilities of modern computers; but keep in mind that at the time it was first formulated, computers and programming languages had yet to be invented. Coming to this realization was an enormous intellectual leap.

Probably the most compelling development leading to the acceptance of Church's thesis was the Turing machine. It was the first model that could be considered readily programmable. If someone laid one of the other systems out in front of you and declared, "*This* system captures exactly what we mean by effectively computable," you might harbor some skepticism. But it is hard to argue with Turing machines. One can rightly challenge Church's thesis on the grounds that there are aspects of computation that are not addressed by Turing machines (for example, randomized or interactive computation), but no one could dispute that the notion of effective computability as captured by Turing machines is robust and important.

Universality and Self-Reference

One of the most intriguing aspects of these systems, and a pervasive theme in our study of them, is the idea of *programs as data*. Each of these programming systems is powerful enough that programs can be written that understand and manipulate other programs that are encoded as data in some reasonable way. For example, in the λ -calculus, λ -terms act as both programs and data; combinator symbols in combinatory logic manipulate other combinator symbols; there is a so-called Gödel numbering of the μ -recursive functions in which each function has a number that can be used as input to other μ -recursive functions; and Turing machines can interpret their input strings as descriptions of other Turing machines. It is not a far step from this idea to the notion of *universal simulation*, in which a universal program or machine U is constructed to take an encoded description of another program or machine M and a string x as input and perform a step-by-step simulation of M on input x . A modern-day example of this phenomenon would be a Scheme interpreter written in Scheme.

One far-reaching corollary of universality is the notion of *self-reference*. It is exactly this capability that led to the discovery of natural uncomputable problems. If you know some set theory, you can convince yourself that uncomputable problems must exist by a cardinality argument: there are uncountably many decision problems but only countably many Turing machines. However, self-reference allows us to construct very simple and natural examples of uncomputable problems. For example, there do not exist general procedures that can determine whether a given block of code in a given C program is ever going to be executed, or whether a given C program will ever halt. These are important problems that compiler builders would like to solve; unfortunately, one can give a formal proof that they are unsolvable.

Perhaps the most striking example of the power of self-reference is the incompleteness theorem of Kurt Gödel. Starting near the beginning of the twentieth century with Whitehead and Russell's *Principia Mathematica* [21], there was a movement to reduce all of mathematics to pure symbol manipulation, independent of semantics. This was in part to understand and avoid the set-theoretic paradoxes discovered by Russell and others. This movement was advocated by the mathematician David Hilbert and became known as the *formalist program*. It attracted a lot of followers and fed the development of formal logic as we know it today. Its proponents believed that all mathematical truths could be derived in some fixed formal system,

just by starting with a few axioms and applying rules of inference in a purely mechanical way. This view of mathematical proof is highly computational. The formal deductive system most popular at the time for reasoning about the natural numbers, called *Peano arithmetic* (PA), was believed to be adequate for expressing and deriving mechanically all true statements of number theory. The incompleteness theorem showed that this was wrong: there exist even fairly simple statements of number theory that are true but not provable in PA. This holds not only for PA but for *any* reasonable extension of it. This revelation was a significant setback for the formalist program and sent shock waves throughout the mathematical world.

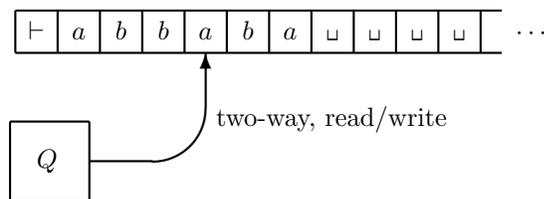
Gödel proved the incompleteness theorem using self-reference. The basic observation needed here is that the language of number theory is expressive enough to talk about itself and about proofs in PA. For example, one can write down a number-theoretic statement that says that a certain other number-theoretic statement has a proof in PA, and one can reason about such statements using PA itself. Now, by a tricky argument involving substitutions, one can actually construct statements that talk about whether they themselves are provable. Gödel actually constructed a sentence that said, “I am not provable.”

The consequences of universality are not only philosophical but also practical. Universality was in a sense the germ of the idea that led to the development of computers as we know them today: the notion of a *stored program*, a piece of software that can be read and executed by hardware. This programmability is what makes computers so versatile. Although it was only realized in physical form several years later, the notion was definitely present in Turing’s theoretical work in the 1930s.

Informal Description of Turing Machines

We describe here a deterministic, one-tape Turing machine. This is the standard off-the-shelf model. There are many variations, apparently more powerful or less powerful but in reality not. We will consider some of these in §3.

A TM has a finite set of states Q , a semi-infinite tape that is delimited on the left end by an endmarker \vdash and is infinite to the right, and a head that can move left and right over the tape, reading and writing symbols.



The input string is of finite length and is initially written on the tape in contiguous tape cells snug up against the left endmarker. The infinitely many cells to the right of the input all contain a special blank symbol \sqcup .

The machine starts in its start state s with its head scanning the left endmarker. In each step it reads the symbol on the tape under its head. Depending on that symbol and the current state, it writes a new symbol on that tape cell, moves its head either left or right one cell, and enters a new state. The action it takes in each situation is determined by a transition function δ . It *accepts* its input by entering a special accept state t and *rejects* by entering a special reject state r . On some inputs it may run infinitely without ever accepting or rejecting, in which case it is said to *loop* on that input.

Formal Definition of Turing Machines

Formally, a *deterministic one-tape Turing machine* is a 9-tuple

$$M = (Q, \Sigma, \Gamma, \vdash, \sqcup, \delta, s, t, r),$$

where

- Q is a finite set (the *states*);

- Σ is a finite set (the *input alphabet*);
- Γ is a finite set (the *tape alphabet*) containing Σ as a subset;
- $\sqcup \in \Gamma - \Sigma$, the *blank symbol*;
- $\vdash \in \Gamma - \Sigma$, the *left endmarker*;
- $\delta : Q \times \Gamma \rightarrow Q \times \Gamma \times \{L, R\}$, the *transition function*;
- $s \in Q$, the *start state*;
- $t \in Q$, the *accept state*; and
- $r \in Q$, the *reject state*, $r \neq t$.

Intuitively, $\delta(p, a) = (q, b, d)$ means, “When in state p scanning symbol a , write b on that tape cell, move the head in direction d , and enter state q .” The symbols L and R stand for left and right, respectively.

We restrict TMs so that the left endmarker is never overwritten with another symbol and the machine never moves off the tape to the left of the endmarker; that is, we require that for all $p \in Q$ there exists $q \in Q$ such that

$$\delta(p, \vdash) = (q, \vdash, R). \quad (2)$$

We also require that once the machine enters its accept state, it never leaves it, and similarly for its reject state; that is, for all $b \in \Gamma$ there exist $c, c' \in \Gamma$ and $d, d' \in \{L, R\}$ such that

$$\delta(t, b) = (t, c, d), \quad \delta(r, b) = (r, c', d'). \quad (3)$$

We sometimes refer to the state set and transition function collectively as the *finite control*.

Example 1. Here is a TM that accepts the non-context-free set $\{a^n b^n c^n \mid n \geq 0\}$. Informally, the machine starts in its start state s , then scans to the right over the input string, checking that it is of the form $a^* b^* c^*$. It doesn't write anything on the way across (formally, it writes the same symbol it reads). When it sees the first blank symbol \sqcup , it overwrites it with a right endmarker \dashv . Now it scans left, erasing the first c it sees, then the first b it sees, then the first a it sees, until it comes to the \vdash . It then scans right, erasing one a , one b , and one c . It continues to sweep left and right over the input, erasing one occurrence of each letter in each pass. If on some pass it sees at least one occurrence of one of the letters and no occurrences of another, it rejects. Otherwise, it eventually erases all the letters and makes one pass between \vdash and \dashv seeing only blanks, at which point it accepts.

Formally, this machine has

$$Q = \{s, q_1, \dots, q_{10}, t, r\}, \quad \Sigma = \{a, b, c\}, \quad \Gamma = \Sigma \cup \{\vdash, \sqcup, \dashv\}.$$

There is nothing special about \dashv ; it is just an extra useful symbol in the tape alphabet. The transition function δ is specified by the following table:

	\vdash	a	b	c	\sqcup	\dashv
s	(s, \vdash, R)	(s, a, R)	(q_1, b, R)	(q_2, c, R)	(q_3, \dashv, L)	—
q_1	—	$(r, -, -)$	(q_1, b, R)	(q_2, c, R)	(q_3, \dashv, L)	—
q_2	—	$(r, -, -)$	$(r, -, -)$	(q_2, c, R)	(q_3, \dashv, L)	—
q_3	$(t, -, -)$	$(r, -, -)$	$(r, -, -)$	(q_4, \sqcup, L)	(q_3, \sqcup, L)	—
q_4	$(r, -, -)$	$(r, -, -)$	(q_5, \sqcup, L)	(q_4, c, L)	(q_4, \sqcup, L)	—
q_5	$(r, -, -)$	(q_6, \sqcup, L)	(q_5, b, L)	—	(q_5, \sqcup, L)	—
q_6	(q_7, \vdash, R)	(q_6, a, L)	—	—	(q_6, \sqcup, L)	—
q_7	—	(q_8, \sqcup, R)	$(r, -, -)$	$(r, -, -)$	(q_7, \sqcup, R)	$(t, -, -)$
q_8	—	(q_8, a, R)	(q_9, \sqcup, R)	$(r, -, -)$	(q_8, \sqcup, R)	$(r, -, -)$
q_9	—	—	(q_9, b, R)	(q_{10}, \sqcup, R)	(q_9, \sqcup, R)	$(r, -, -)$
q_{10}	—	—	—	(q_{10}, c, R)	(q_{10}, \sqcup, R)	(q_3, \dashv, L)

The symbol $-$ in the table above means “don’t care.” The transitions for t and r are not included in the table—just define them to be anything satisfying the restrictions (2) and (3).

Configurations and Acceptance

At any point in time, the read/write tape of the Turing machine M contains a semi-infinite string of the form $y\sqcup^\omega$, where $y \in \Gamma^*$ (y is a finite-length string) and \sqcup^ω denotes the semi-infinite string

$$\sqcup \sqcup \sqcup \sqcup \sqcup \sqcup \sqcup \sqcup \dots$$

(Here ω denotes the smallest infinite ordinal.) Although the string is infinite, it always has a finite representation, since all but finitely many of the symbols are the blank symbol \sqcup .

We define a *configuration* to be an element of $Q \times \{y\sqcup^\omega \mid y \in \Gamma^*\} \times \mathbb{N}$, where $\mathbb{N} = \{0, 1, 2, \dots\}$. A configuration is a global state giving a snapshot of all relevant information about a TM computation at some instant in time. The configuration (p, z, n) specifies a current state p of the finite control, current tape contents z , and current position of the read/write head $n \geq 0$. We usually denote configurations by α, β, γ .

The *start configuration* on input $x \in \Sigma^*$ is the configuration

$$(s, \vdash x\sqcup^\omega, 0).$$

The last component 0 means that the machine is initially scanning the left endmarker \vdash .

One can define a *next configuration relation* $\xrightarrow[M]{1}$. For a string $z \in \Gamma^\omega$, let z_n be the n th symbol of z (the leftmost symbol is z_0), and let $z[n/b]$ denote the string obtained from z by substituting b for z_n at position n . For example,

$$\vdash b a a c a b c a \dots [4/b] = \vdash b a a b c a b c a \dots$$

The relation $\xrightarrow[M]{1}$ is defined by

$$(p, z, n) \xrightarrow[M]{1} \begin{cases} (q, z[n/b], n - 1) & \text{if } \delta(p, z_n) = (q, b, L), \\ (q, z[n/b], n + 1) & \text{if } \delta(p, z_n) = (q, b, R). \end{cases}$$

Intuitively, if the tape contains z and if M is in state p scanning the n th tape cell, and δ says to print b , go left, and enter state q , then after that step the tape will contain $z[n/b]$, the head will be scanning the $(n - 1)$ st tape cell, and the new state will be q .

We define the reflexive transitive closure $\xrightarrow[M]{*}$ of $\xrightarrow[M]{1}$ inductively:

- $\alpha \xrightarrow[M]{0} \alpha$,
- $\alpha \xrightarrow[M]{n+1} \beta$ if $\alpha \xrightarrow[M]{n} \gamma \xrightarrow[M]{1} \beta$ for some γ , and
- $\alpha \xrightarrow[M]{*} \beta$ if $\alpha \xrightarrow[M]{n} \beta$ for some $n \geq 0$.

The machine M is said to *accept* input $x \in \Sigma^*$ if

$$(s, \vdash x\sqcup^\omega, 0) \xrightarrow[M]{*} (t, y, n)$$

for some y and n and is said to *reject* x if

$$(s, \vdash x\sqcup^\omega, 0) \xrightarrow[M]{*} (r, y, n)$$

for some y and n . It is said to *halt* on input x if it either accepts x or rejects x . This is just a mathematical definition; the machine doesn't really grind to a halt in the literal sense. It is possible that it neither accepts nor rejects, in which case it is said to *loop* on input x . A Turing machine is said to be *total* if it halts on all inputs; that is, if for all inputs it either accepts or rejects. The set $L(M)$ denotes the set of strings accepted by M .

We call a set of strings

- *recursively enumerable* (r.e.) if it is $L(M)$ for some Turing machine M ,
- *co-r.e.* if its complement is r.e., and
- *recursive* if it is $L(M)$ for some *total* Turing machine M .

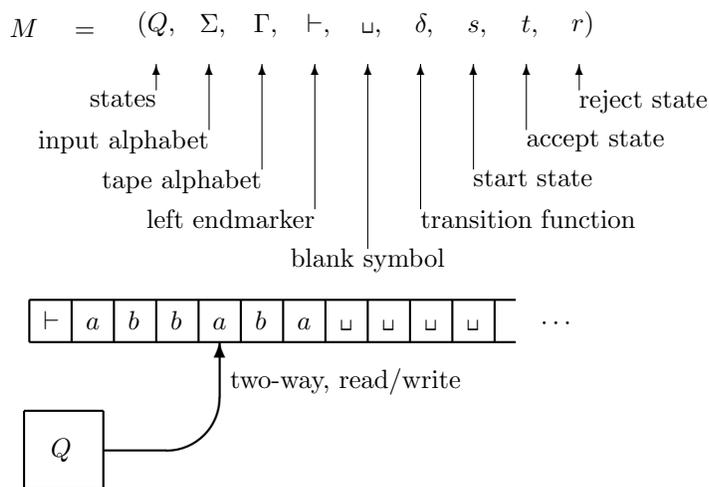
In common parlance, the term “recursive” usually refers to an algorithm that calls itself. The definition above has nothing to do with this usage. As used here, it is just a name for a set accepted by a Turing machine that always halts.

Historical Notes

Church's thesis is often referred to as the Church–Turing thesis, although Alonzo Church was the first to formulate it explicitly [3]. The thesis was based on Church and Kleene's observation that the λ -calculus and the μ -recursive functions of Gödel and Herbrand were computationally equivalent [3]. Church was apparently unaware of Turing's work at the time of the writing of [3], or if he was, he failed to mention it. Turing, on the other hand, cited Church's paper [3] explicitly in [20], and apparently considered his machines to be a much more compelling definition of computability. In an appendix to [20], Turing outlined a proof of the computational equivalence of Turing machines and the λ -calculus.

2 Examples of Turing Machines

In the last section we defined deterministic one-tape Turing machines:



In each step, based on the current tape symbol it is reading and its current state, it prints a new symbol on the tape, moves its head either left or right, and enters a new state. This action is specified formally by the transition function

$$\delta : Q \times \Gamma \rightarrow Q \times \Gamma \times \{L, R\}.$$

Intuitively, $\delta(p, a) = (q, b, d)$ means, “When in state p scanning symbol a , write b on that tape cell, move the head in direction d , and enter state q .”

We defined a *configuration* to be a triple (p, z, n) where p is a state, z is a semi-infinite string of the form $y\sqcup^\omega$, $y \in \Sigma^*$, describing the contents of the tape, and n is a natural number denoting a tape head position.

The transition function δ was used to define the *next configuration relation* $\xrightarrow[M]{1}$ on configurations and its reflexive transitive closure $\xrightarrow[M]{*}$. The machine M *accepts* input $x \in \Sigma^*$ if

$$(s, \vdash x\sqcup^\omega, 0) \xrightarrow[M]{*} (t, y, n)$$

for some y and n , and *rejects* input x if

$$(s, \vdash x\sqcup^\omega, 0) \xrightarrow[M]{*} (r, y, n)$$

for some y and n . The left configuration above is the *start configuration* on input x . Recall that we restricted TMs so that once a TM enters its accept state, it may never leave it, and similarly for its reject state. If M never enters its accept or reject state on input x , it is said to *loop* on input x . It is said to *halt* on input x if it either accepts or rejects. A TM that halts on all inputs is called *total*.

Define the set

$$L(M) = \{x \in \Sigma^* \mid M \text{ accepts } x\}.$$

This is called the set *accepted by M*. A subset of Σ^* is called *recursively enumerable* (r.e.) if it is $L(M)$ for some M . A set is called *recursive* if it is $L(M)$ for some total M .

For now, the terms *r.e.* and *recursive* are just technical terms describing the sets accepted by TMs and total TMs, respectively; they have no other significance. We will discuss the origin of this terminology in §3.

Example 2. Consider the set $\{ww \mid w \in \{a, b\}^*\}$, the set of all strings over the alphabet $\{a, b\}$ consisting of two concatenated copies of some string w . It is a recursive set, because we can give a total TM M for it. The machine M works as follows. On input x , it scans out to the first blank symbol \sqcup , counting the number of symbols mod 2 to make sure x is of even length and rejecting immediately if not. It lays down a right endmarker \dashv , then repeatedly scans back and forth over the input. In each pass from right to left, it marks the first unmarked a or b it sees with $\acute{}$. In each pass from left to right, it marks the first unmarked a or b it sees with $\grave{}$. It continues this until all symbols are marked. For example, on input

$$a a b b a a a b b a$$

the initial tape contents are

$$\vdash a a b b a a a b b a \sqcup \sqcup \sqcup \dots$$

and the following are the tape contents after the first few passes.

$$\begin{aligned} &\vdash a a b b a a a b b \acute{a} \dashv \sqcup \sqcup \dots \\ &\vdash \grave{a} a b b a a a b b \acute{a} \dashv \sqcup \sqcup \dots \\ &\vdash \grave{a} a b b a a a b \acute{b} \acute{a} \dashv \sqcup \sqcup \dots \\ &\vdash \grave{a} \grave{a} b b a a a b \acute{b} \acute{a} \dashv \sqcup \sqcup \dots \\ &\vdash \grave{a} \grave{a} b b a a a \acute{b} \acute{b} \acute{a} \dashv \sqcup \sqcup \dots \end{aligned}$$

Marking a with $\grave{}$ formally means writing the symbol $\grave{a} \in \Gamma$; thus

$$\Gamma = \{a, b, \vdash, \sqcup, \dashv, \grave{a}, \grave{b}, \acute{a}, \acute{b}\}.$$

When all symbols are marked, we have the first half of the input string marked with $\grave{}$ and the second half marked with $\acute{}$.

$$\vdash \grave{a} \grave{a} \grave{b} \grave{b} \grave{a} \acute{a} \acute{a} \acute{b} \acute{b} \acute{a} \dashv \sqcup \sqcup \dots$$

The reason we did this was to find the center of the input string.

Decidability and Semidecidability

A property P of strings is said to be *decidable* if the set of all strings having property P is a recursive set; that is, if there is a total Turing machine that accepts input strings that have property P and rejects those that do not. A property P is said to be *semidecidable* if the set of strings having property P is an r.e. set; that is, if there is a Turing machine that on input x accepts if x has property P and rejects or loops if not. For example, it is decidable whether a given string x is of the form ww , because we can construct a Turing machine that halts on all inputs and accepts exactly the strings of this form.

Although you often hear them switched, the adjectives *recursive* and *r.e.* are best applied to sets and *decidable* and *semidecidable* to properties. The two notions are equivalent, since

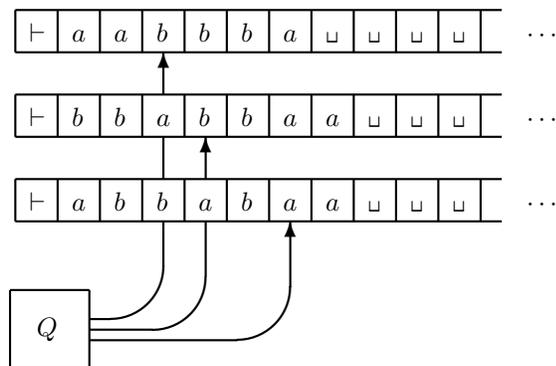
$$\begin{aligned}
 P \text{ is decidable} &\Leftrightarrow \{x \mid P(x)\} \text{ is recursive.} \\
 A \text{ is recursive} &\Leftrightarrow "x \in A" \text{ is decidable,} \\
 P \text{ is semidecidable} &\Leftrightarrow \{x \mid P(x)\} \text{ is r.e.,} \\
 A \text{ is r.e.} &\Leftrightarrow "x \in A" \text{ is semidecidable.}
 \end{aligned}$$

3 Equivalent Models

As mentioned, the concept of computability is remarkably robust. As evidence of this, we will present several different flavors of Turing machines that at first glance appear to be significantly more or less powerful than the basic model defined in §2 but in fact are computationally equivalent.

Multiple Tapes

First, we show how to simulate multitape Turing machines with single-tape Turing machines. Thus extra tapes don't add any power. A three-tape machine is similar to a one-tape TM except that it has three semi-infinite tapes and three independent read/write heads. Initially, the input occupies the first tape and the other two are blank. In each step, the machine reads the three symbols under its heads, and based on this information and the current state, it prints a symbol on each tape, moves the heads (they don't all have to move in the same direction), and enters a new state.



Its transition function is of type

$$\delta : Q \times \Gamma^3 \rightarrow Q \times \Gamma^3 \times \{L, R\}^3.$$

Say we have such a machine M . We build a single-tape machine N simulating M as follows. The machine N will have an expanded tape alphabet allowing us to think of its tape as divided into three tracks. Each track will contain the contents of one of M 's tapes. We also mark exactly one symbol on each track to indicate that this is the symbol currently being scanned on the corresponding tape of M . The configuration of M illustrated above might be simulated by the following configuration of N .

⊢	⊢	a	a	\widehat{b}	b	b	a	␣	␣	␣	␣	...
	⊢	b	b	a	\widehat{b}	b	a	a	␣	␣	␣	
	⊢	a	b	b	a	b	\widehat{a}	a	␣	␣	␣	

A tape symbol of N is either \vdash , an element of Σ , or a triple

c
d
e

where c, d, e are tape symbols of M , each either marked or unmarked. Formally, we might take the tape alphabet of N to be

$$\Sigma \cup \{\vdash\} \cup (\Gamma \cup \Gamma')^3,$$

where

$$\Gamma' \stackrel{\text{def}}{=} \{\widehat{c} \mid c \in \Gamma\}.$$

The three elements of $\Gamma \cup \Gamma'$ stand for the symbols in corresponding positions on the three tapes of M , either marked or unmarked, and

␣
␣
␣

is the blank symbol of N .

On input $x = a_1 a_2 \cdots a_n$, N starts with tape contents

$$\vdash a_1 a_2 a_3 \cdots a_n \text{␣␣␣} \cdots.$$

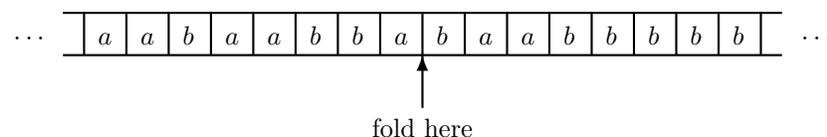
It first copies the input to its top track and fills in the bottom two tracks with blanks. It also shifts everything right one cell so that it can fill in the leftmost cells on the three tracks with the simulated left endmarker of M , which it marks with $\widehat{}$ to indicate the position of the heads in the starting configuration of M .

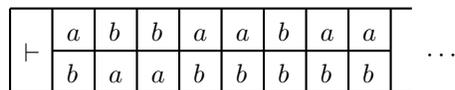
⊢	$\widehat{\vdash}$	a_1	a_2	a_3	a_4	...	a_n	␣	␣	...
	$\widehat{\vdash}$	␣	␣	␣	␣		␣	␣	␣	
	$\widehat{\vdash}$	␣	␣	␣	␣		␣	␣	␣	

Each step of M is simulated by several steps of N . To simulate one step of M , N starts at the left of the tape, then scans out until it sees all three marks, remembering the marked symbols in its finite control. When it has seen all three, it determines what to do according to M 's transition function δ , which it has encoded in its finite control. Based on this information, it goes back to all three marks, rewriting the symbols on each track and moving the marks appropriately. It then returns to the left end of the tape to simulate the next step of M .

Two-Way Infinite Tapes

Two-way infinite tapes do not add any power. Just fold the tape someplace and simulate it on two tracks of a one-way infinite tape:

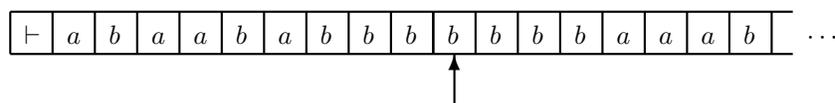




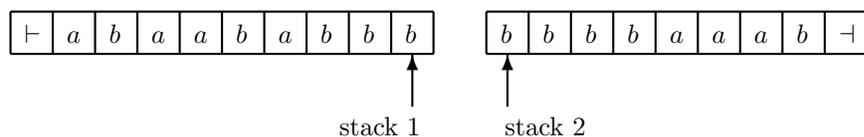
The bottom track is used to simulate the original machine when its head is to the right of the fold, and the top track is used to simulate the machine when its head is to the left of the fold, moving in the opposite direction.

Two Stacks

A machine with a two-way, *read-only* input head and two stacks is as powerful as a Turing machine. Intuitively, the computation of a one-tape TM can be simulated with two stacks by storing the tape contents to the left of the head on one stack and the tape contents to the right of the head on the other stack. The motion of the head is simulated by popping a symbol off one stack and pushing it onto the other. For example,



is simulated by



Counter Automata

A *k-counter automaton* is a machine equipped with a two-way read-only input head and *k* integer counters. Each counter can store an arbitrary nonnegative integer. In each step, the automaton can independently increment or decrement its counters and test them for 0 and can move its input head one cell in either direction. It cannot write on the tape.

A stack can be simulated with two counters as follows. We can assume without loss of generality that the stack alphabet of the stack to be simulated contains only two symbols, say 0 and 1. This is because we can encode finitely many stack symbols as binary numbers of fixed length, say *m*; then pushing or popping one stack symbol is simulated by pushing or popping *m* binary digits. Then the contents of the stack can be regarded as a binary number whose least significant bit is on top of the stack. The simulation maintains this number in the first of the two counters and uses the second to effect the stack operations. To simulate pushing a 0 onto the stack, we need to double the value in the first counter. This is done by entering a loop that repeatedly subtracts one from the first counter and adds two to the second until the first counter is 0. The value in the second counter is then twice the original value in the first counter. We can then transfer that value back to the first counter, or just switch the roles of the two counters. To push 1, the operation is the same, except the value of the second counter is incremented once at the end. To simulate popping, we need to divide the counter value by two; this is done by decrementing one counter while incrementing the other counter every second step. Testing the parity of the original counter contents tells whether a simulated 1 or 0 was popped.

Since a two-stack machine can simulate an arbitrary TM, and since two counters can simulate a stack, it follows that a four-counter automaton can simulate an arbitrary TM.

However, we can do even better: a two-counter automaton can simulate a four-counter automaton. When the four-counter automaton has the values *i, j, k, l* in its counters, the two-counter automaton will have the value $2^i 3^j 5^k 7^l$ in its first counter. It uses its second counter to effect the counter operations of the four-counter automaton. For example, if the four-counter automaton wanted to add one to *k* (the value of the

third counter), then the two-counter automaton would have to multiply the value in its first counter by 5. This is done in the same way as above, adding 5 to the second counter for every 1 we subtract from the first counter. To simulate a test for zero, the two-counter automaton has to determine whether the value in its first counter is divisible by 2, 3, 5, or 7, respectively, depending on which counter of the four-counter automaton is being tested.

Combining these simulations, we see that two-counter automata are as powerful as arbitrary Turing machines. However, as you can imagine, it takes an enormous number of steps of the two-counter automaton to simulate one step of the Turing machine.

One-counter automata are not as powerful as arbitrary TMs.

Enumeration Machines

We defined the recursively enumerable (r.e.) sets to be those sets accepted by Turing machines. The term *recursively enumerable* comes from a different but equivalent formalism embodying the idea that the elements of an r.e. set can be enumerated one at a time in a mechanical fashion.

Define an *enumeration machine* as follows. It has a finite control and two tapes, a read/write *work tape* and a write-only *output tape*. The work tape head can move in either direction and can read and write any element of Γ . The output tape head moves right one cell when it writes a symbol, and it can only write symbols in Σ . There is no input and no accept or reject state. The machine starts in its start state with both tapes blank. It moves according to its transition function like a TM, occasionally writing symbols on the output tape as determined by the transition function. At some point it may enter a special *enumeration state*, which is just a distinguished state of its finite control. When that happens, the string currently written on the output tape is said to be *enumerated*. The output tape is then automatically erased and the output head moved back to the beginning of the tape (the work tape is left intact), and the machine continues from that point. The machine runs forever. The set $L(E)$ is defined to be the set of all strings in Σ^* that are ever enumerated by the enumeration machine E . The machine might never enter its enumeration state, in which case $L(E) = \emptyset$, or it might enumerate infinitely many strings. The same string may be enumerated more than once.

Enumeration machines and Turing machines are equivalent in computational power:

Theorem 1. *The family of sets enumerated by enumeration machines is exactly the family of r.e. sets. In other words, a set is $L(E)$ for some enumeration machine E if and only if it is $L(M)$ for some Turing machine M .*

Proof. We show first that given an enumeration machine E , we can construct a Turing machine M such that $L(M) = L(E)$. Let M on input x copy x to one of three tracks on its tape, then simulate E , using the other two tracks to record the contents of E 's work tape and output tape. For every string enumerated by E , M compares this string to x and accepts if they match. Then M accepts its input x iff x is ever enumerated by E , so the set of strings accepted by M is exactly the set of strings enumerated by E .

Conversely, given a TM M , we can construct an enumeration machine E such that $L(E) = L(M)$. We would like E somehow to simulate M on all possible strings in Σ^* and enumerate those that are accepted.

Here is an approach that doesn't quite work. The enumeration machine E writes down the strings in Σ^* one by one on the bottom track of its work tape in some order. For every input string x , it simulates M on input x , using the top track of its work tape to do the simulation. If M accepts x , E copies x to its output tape and enters its enumeration state. It then goes on to the next string.

The problem with this procedure is that M might not halt on some input x , and then E would be stuck simulating M on x forever and would never move on to strings later in the list (and it is impossible to determine in general whether M will ever halt on x , as we will see in §4). Thus E should not just list the strings in Σ^* in some order and simulate M on those inputs one at a time, waiting for each simulation to halt before going on to the next, because the simulation might never halt.

The solution to this problem is *timesharing*. Instead of simulating M on the input strings one at a time,

the enumeration machine E should run several simulations at once, working a few steps on each simulation and then moving on to the next. The work tape of E can be divided into segments separated by a special marker $\# \in \Gamma$, with a simulation of M on a different input string running in each segment. Between passes, E can move way out to the right, create a new segment, and start up a new simulation in that segment on the next input string. For example, we might have E simulate M on the first input for one step, then the first and second inputs for one step each, then the first, second, and third inputs for one step each, and so on. If any simulation needs more space than initially allocated in its segment, the entire contents of the tape to its right can be shifted to the right one cell. In this way M is eventually simulated on all input strings, even if some of the simulations never halt. \square

Historical Notes

Turing machines were invented by Alan Turing [20]. Originally they were presented in the form of enumeration machines, since Turing was interested in enumerating the decimal expansions of computable real numbers and values of real-valued functions. Turing also introduced the concept of nondeterminism in his original paper, although he did not develop the idea.

The basic properties of the r.e. sets were developed by Kleene [9] and Post [13, 14].

Counter automata were studied by Fischer [5], Fischer et al. [6], and Minsky [11].

4 Universal Machines and Diagonalization

A Universal Turing Machine

Now we come to a crucial observation about the power of Turing machines: there exist Turing machines that can simulate other Turing machines whose descriptions are presented as part of the input. There is nothing mysterious about this; it is the same as writing an OCaml interpreter in OCaml.

First we need to fix a reasonable encoding scheme for Turing machines over the alphabet $\{0, 1\}$. This encoding scheme should be simple enough that all the data associated with a machine M —the set of states, the transition function, the input and tape alphabets, the endmarker, the blank symbol, and the start, accept, and reject states—can be determined easily by another machine reading the encoded description of M . For example, if the string begins with the prefix

$$0^n 10^m 10^k 10^s 10^t 10^r 10^u 10^v 1,$$

this might indicate that the machine has n states represented by the numbers 0 to $n - 1$; it has m tape symbols represented by the numbers 0 to $m - 1$, of which the first k represent input symbols; the start, accept, and reject states are s , t , and r , respectively; and the endmarker and blank symbol are u and v , respectively. The remainder of the string can consist of a sequence of substrings specifying the transitions in δ . For example, the substring

$$0^p 10^a 10^q 10^b 10$$

might indicate that δ contains the transition

$$(p, a) \rightarrow (q, b, L),$$

the direction to move the head encoded by the final digit. The exact details of the encoding scheme are not important. The only requirements are that it should be easy to interpret and able to encode all Turing machines up to isomorphism.

Once we have a suitable encoding of Turing machines, we can construct a *universal Turing machine* U such that

$$L(U) \stackrel{\text{def}}{=} \{M\#x \mid x \in L(M)\}.$$

In other words, presented with (an encoding over $\{0,1\}$ of) a Turing machine M and (an encoding over $\{0,1\}$ of) a string x over M 's input alphabet, the machine U accepts $M\#x$ iff M accepts x .¹ The symbol $\#$ is just a symbol in U 's input alphabet other than 0 or 1 used to delimit M and x .

The machine U first checks its input $M\#x$ to make sure that M is a valid encoding of a Turing machine and x is a valid encoding of a string over M 's input alphabet. If not, it immediately rejects.

If the encodings of M and x are valid, the machine U does a step-by-step simulation of M . This might work as follows. The tape of U is partitioned into three tracks. The description of M is copied to the top track and the string x to the middle track. The middle track will be used to hold the simulated contents of M 's tape. The bottom track will be used to remember the current state of M and the current position of M 's read/write head. The machine U then simulates M on input x one step at a time, shuttling back and forth between the description of M on its top track and the simulated contents of M 's tape on the middle track. In each step, it updates M 's state and simulated tape contents as dictated by M 's transition function, which U can read from the description of M . If ever M halts and accepts or halts and rejects, then U does the same.

As we have observed, the string x over the input alphabet of M and its encoding over the input alphabet of U are two different things, since the two machines may have different input alphabets. If the input alphabet of M is bigger than that of U , then each symbol of x must be encoded as a string of symbols over U 's input alphabet. Also, the tape alphabet of M may be bigger than that of U , in which case each symbol of M 's tape alphabet must be encoded as a string of symbols over U 's tape alphabet. In general, each step of M may require many steps of U to simulate.

Diagonalization

We now show how to use a universal Turing machine in conjunction with a technique called *diagonalization* to prove that the halting and membership problems for Turing machines are undecidable. In other words, the sets

$$\begin{aligned} \text{HP} &\stackrel{\text{def}}{=} \{M\#x \mid M \text{ halts on } x\}, \\ \text{MP} &\stackrel{\text{def}}{=} \{M\#x \mid x \in L(M)\} \end{aligned}$$

are not recursive.

The technique of diagonalization was first used by Cantor at the end of the nineteenth century to show that there does not exist a one-to-one correspondence between the natural numbers \mathbb{N} and its *power set*

$$2^{\mathbb{N}} = \{A \mid A \subseteq \mathbb{N}\},$$

the set of all subsets of \mathbb{N} . In fact, there does not even exist a function

$$f : \mathbb{N} \rightarrow 2^{\mathbb{N}}$$

that is onto. Here is how Cantor's argument went.

Suppose (for a contradiction) that such an onto function f did exist. Consider an infinite two-dimensional matrix indexed along the top by the natural numbers $0, 1, 2, \dots$ and down the left by the sets $f(0), f(1), f(2), \dots$.

¹Note that we are using the metasyMBOL M for both a Turing machine and its encoding over $\{0,1\}$ and the metasyMBOL x for both a string over M 's input alphabet and its encoding over $\{0,1\}$. This is for notational convenience.

Fill in the matrix by placing a 1 in position i, j if j is in the set $f(i)$ and 0 if $j \notin f(i)$.

	0	1	2	3	4	5	6	7	8	9	...
$f(0)$	1	0	0	1	1	0	1	0	1	1	
$f(1)$	0	0	1	1	0	1	1	0	0	1	
$f(2)$	0	1	1	0	0	0	1	1	0	1	
$f(3)$	0	1	0	1	1	0	1	1	0	0	
$f(4)$	1	0	1	0	0	1	0	0	1	1	...
$f(5)$	1	0	1	1	0	1	1	1	0	1	
$f(6)$	0	0	1	0	1	1	0	0	1	1	
$f(7)$	1	1	1	0	1	1	1	0	1	0	
$f(8)$	0	0	1	0	0	0	0	1	1	0	
$f(9)$	1	1	0	0	1	0	0	1	0	0	
\vdots					\vdots						\ddots

The i th row of the matrix is a bit string describing the set $f(i)$. For example, in the above picture, $f(0) = \{0, 3, 4, 6, 8, 9, \dots\}$ and $f(1) = \{2, 3, 5, 6, 9, \dots\}$. By our (soon to be proved fallacious) assumption that f is onto, every subset of \mathbb{N} appears as a row of this matrix.

But we can construct a new set that does not appear in the list by complementing the main diagonal of the matrix (hence the term *diagonalization*). Look at the infinite bit string down the main diagonal (in this example, 1011010010 \dots) and take its Boolean complement (in this example, 0100101101 \dots). This new bit string represents a set B (in this example, $B = \{1, 4, 6, 7, 9, \dots\}$). But the set B does not appear anywhere in the list down the left side of the matrix, since it differs from every $f(i)$ on at least one element, namely i . This is a contradiction, since every subset of \mathbb{N} was supposed to occur as a row of the matrix, by our assumption that f was onto.

This argument works not only for the natural numbers \mathbb{N} , but for any set A whatsoever. Suppose (for a contradiction) there existed an onto function from A to its power set:

$$f : A \rightarrow 2^A.$$

Let

$$B = \{x \in A \mid x \notin f(x)\}$$

(this is the formal way of *complementing the diagonal*). Then $B \subseteq A$. Since f is onto, there must exist $y \in A$ such that $f(y) = B$. Now we ask whether $y \in f(y)$ and discover a contradiction:

$$\begin{aligned} y \in f(y) &\Leftrightarrow y \in B && \text{since } B = f(y) \\ &\Leftrightarrow y \notin f(y) && \text{definition of } B. \end{aligned}$$

Thus no such f can exist.

Undecidability of the Halting Problem

We have discussed how to encode descriptions of Turing machines as strings in $\{0,1\}^*$ so that these descriptions can be read and simulated by a universal Turing machine U . The machine U takes as input an encoding of a Turing machine M and a string x and simulates M on input x , and

- halts and accepts if M halts and accepts x ,
- halts and rejects if M halts and rejects x , and
- loops if M loops on x .

The machine U doesn't do any fancy analysis on the machine M to try to determine whether or not it will halt. It just blindly simulates M step by step. If M doesn't halt on x , then U will just go on happily simulating M forever.

It is natural to ask whether we can do better than just a blind simulation. Might there be some way to analyze M to determine in advance, before doing the simulation, whether M would eventually halt on x ? If U could say for sure in advance that M would not halt on x , then it could skip the simulation and save itself a lot of useless work. On the other hand, if U could ascertain that M *would* eventually halt on x , then it could go ahead with the simulation to determine whether M accepts or rejects. We could then build a machine U' that takes as input an encoding of a Turing machine M and a string x , and

- halts and accepts if M halts and accepts x ,
- halts and rejects if M halts and rejects x , and
- halts and rejects if M loops on x .

This would say that $L(U') = L(U) = \text{MP}$ is a recursive set.

Unfortunately, this is not possible in general. There are certainly machines for which it is possible to determine halting by some heuristic or other: machines for which the start state is the accept state, for example. However, there is no general method that gives the right answer for all machines.

We can prove this using Cantor's diagonalization technique. For $x \in \{0, 1\}^*$, let M_x be the Turing machine with input alphabet $\{0, 1\}$ whose encoding over $\{0, 1\}^*$ is x . (If x is not a legal description of a TM with input alphabet $\{0, 1\}^*$ according to our encoding scheme, we take M_x to be some arbitrary but fixed TM with input alphabet $\{0, 1\}$, say a trivial TM with one state that immediately halts.) In this way we get a list

$$M_\varepsilon, M_0, M_1, M_{00}, M_{01}, M_{10}, M_{11}, M_{100}, M_{101}, \dots \tag{4}$$

containing all possible Turing machines with input alphabet $\{0, 1\}$ indexed by strings in $\{0, 1\}^*$. We make sure that the encoding scheme is simple enough that a universal machine can determine M_x from x for the purpose of simulation.

Now consider an infinite two-dimensional matrix indexed along the top by strings in $\{0, 1\}^*$ and down the left by TMs in the list (4).

The matrix contains an H in position x, y if M_x halts on input y and an L if M_x loops on input y .

	ε	0	1	00	01	10	11	000	001	010	\dots
M_ε	H	L	L	H	H	L	H	L	H	H	
M_0	L	L	H	H	L	H	H	L	L	H	
M_1	L	H	H	L	L	L	H	H	L	H	
M_{00}	L	H	L	H	H	L	H	H	L	L	
M_{01}	H	L	H	L	L	H	L	L	H	H	\dots
M_{10}	H	L	H	H	L	H	H	H	L	H	
M_{11}	L	L	H	L	H	H	L	L	H	H	
M_{000}	H	H	H	L	H	H	H	L	H	L	
M_{001}	L	L	H	L	L	L	L	H	H	L	
M_{010}	H	H	L	L	H	L	L	H	L	L	
\vdots	\vdots				\vdots						\ddots

The x th row of the matrix describes for each input string y whether or not M_x halts on y . For example, in the above picture, M_ε halts on inputs $\varepsilon, 00, 01, 11, 001, 010, \dots$ and does not halt on inputs $0, 1, 10, 000, \dots$

Suppose (for a contradiction) that there existed a *total* machine K accepting the set HP; that is, a machine that for any given x and y could determine the x, y th entry of the above table in finite time. Thus on input $M\#x$,

- K halts and accepts if M halts on x , and
- K halts and rejects if M loops on x .

Consider a machine N that on input $x \in \{0, 1\}^*$

- constructs M_x from x and writes $M_x\#x$ on its tape;
- runs K on input $M_x\#x$, accepting if K rejects and going into a trivial loop if K accepts.

Note that N is essentially complementing the diagonal of the above matrix. Then for any $x \in \{0, 1\}^*$,

$$\begin{aligned} N \text{ halts on } x &\Leftrightarrow K \text{ rejects } M_x\#x && \text{definition of } N \\ &\Leftrightarrow M_x \text{ loops on } x && \text{assumption about } K. \end{aligned}$$

This says that N 's behavior is different from every M_x on at least one string, namely x . But the list (4) was supposed to contain all Turing machines over the input alphabet $\{0, 1\}$, including N . This is a contradiction. \square

The fallacious assumption that led to the contradiction was that it was possible to determine the entries of the matrix effectively; in other words, that there existed a Turing machine K that given M and x could determine in a finite time whether or not M halts on x .

One can always simulate a given machine on a given input. If the machine ever halts, then we will know this eventually, and we can stop the simulation and say that it halted; but if not, there is no way in general to stop after a finite time and say for certain that it will never halt.

Undecidability of the Membership Problem

The membership problem is also undecidable. We can show this by *reducing* the halting problem to it. In other words, we show that if there were a way to decide membership in general, we could use this as a subroutine to decide halting in general. But we just showed above that halting is undecidable, so membership must be undecidable too.

Here is how we would use a total TM that decides membership as a subroutine to decide halting. Given a machine M and input x , suppose we wanted to find out whether M halts on x . Build a new machine N that is exactly like M , except that it accepts whenever M would either accept or reject. The machine N can be constructed from M simply by adding a new accept state and making the old accept and reject states transfer to this new accept state. Then for all x , N accepts x iff M halts on x . The membership problem for N and x (asking whether $x \in L(N)$) is therefore the same as the halting problem for M and x (asking whether M halts on x). If the membership problem were decidable, then we could decide whether M halts on x by constructing N and asking whether $x \in L(N)$. But we have shown above that the halting problem is undecidable, therefore the membership problem must also be undecidable.

5 Decidable and Undecidable Problems

Here are some examples of decision problems involving Turing machines. Is it decidable whether a given Turing machine

- has at least 481 states?
- takes more than 481 steps on input ε ?
- takes more than 481 steps on *some* input?
- takes more than 481 steps on *all* inputs?

- (e) ever moves its head more than 481 tape cells away from the left endmarker on input ε ?
- (f) accepts the null string ε ?
- (g) accepts any string at all?
- (h) accepts every string?
- (i) accepts a finite set?
- (j) accepts a recursive set?
- (k) is equivalent to a Turing machine with a shorter description?

Problems (a) through (e) are decidable and problems (f) through (k) are undecidable (proofs below). We will show that problems (f) through (j) are undecidable by showing that a decision procedure for one of these problems could be used to construct a decision procedure for the halting problem, which we know is impossible. Problem (k) is a little more difficult, and we will leave that as an exercise. Translated into modern terms, problem (k) is the same as determining whether there exists a shorter Java program equivalent to a given one.

The best way to show that a problem is decidable is to give a total Turing machine that accepts exactly the “yes” instances. Because it must be total, it must also reject the “no” instances; in other words, it must not loop on any input.

Problem (a) is easily decidable, since the number of states of M can be read off from the encoding of M . We can build a Turing machine that, given the encoding of M written on its input tape, counts the number of states of M and accepts or rejects depending on whether the number is at least 481.

Problem (b) is decidable, since we can simulate M on input ε with a universal machine for 481 steps (counting up to 481 on a separate track) and accept or reject depending on whether M has halted by that time.

Problem (c) is decidable: we can just simulate M on all inputs of length at most 481 for 481 steps. If M takes more than 481 steps on some input, then it will take more than 481 steps on some input of length at most 481, since in 481 steps it can read at most the first 481 symbols of the input.

The argument for problem (d) is similar. If M takes more than 481 steps on all inputs of length at most 481, then it will take more than 481 steps on all inputs.

For problem (e), if M never moves more than 481 tape cells away from the left endmarker, then it will either halt or loop in such a way that we can detect the looping after a finite time. This is because if M has k states and m tape symbols, and never moves more than 481 tape cells away from the left endmarker, then there are only $482km^{481}$ configurations it could possibly ever be in, one for each choice of head position, state, and tape contents that fit within 481 tape cells. If it runs for any longer than that without moving more than 481 tape cells away from the left endmarker, then it must be in a loop, because it must have repeated a configuration. This can be detected by a machine that simulates M , counting the number of steps M takes on a separate track and declaring M to be in a loop if the bound of $482km^{481}$ steps is ever exceeded.

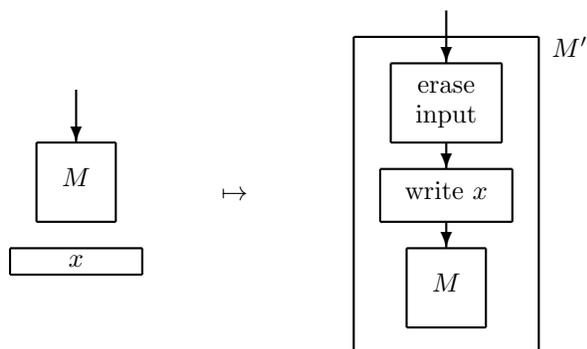
Problems (f) through (j) are undecidable. To show this, we show that the ability to decide any one of these problems could be used to decide the halting problem. Since we know that the halting problem is undecidable, these problems must be undecidable too. This is called a *reduction*.

Let's consider (f) first (although the same construction will take care of (g) through (i) as well). We will show that it is undecidable whether a given machine accepts ε , because the ability to decide this question would give the ability to decide the halting problem, which we know is impossible.

Suppose we could decide whether a given machine accepts ε . We could then decide the halting problem as follows. Say we are given a Turing machine M and string x , and we wish to determine whether M halts on x . Construct from M and x a new machine M' that does the following on input y :

- (i) erases its input y ;
- (ii) writes x on its tape (M' has x hard-wired in its finite control);

- (iii) runs M on input x (M' also has a description of M hard-wired in its finite control);
- (iv) accepts if M halts on x .



Note that M' does the same thing on all inputs y : if M halts on x , then M' accepts its input y ; and if M does not halt on x , then M' does not halt on y , therefore does not accept y . Moreover, this is true for every y . Thus

$$L(M') = \begin{cases} \Sigma^* & \text{if } M \text{ halts on } x, \\ \emptyset & \text{if } M \text{ does not halt on } x. \end{cases}$$

Now if we could decide whether a given machine accepts the null string ε , we could apply this decision procedure to the M' just constructed, and this would tell whether M halts on x . In other words, we could obtain a decision procedure for halting as follows: given M and x , construct M' , then ask whether M' accepts ε . The answer to the latter question is “yes” iff M halts on x . Since we know the halting problem is undecidable, it must also be undecidable whether a given machine accepts ε .

Similarly, if we could decide whether a given machine accepts any string at all, or whether it accepts every string, or whether the set of strings it accepts is finite, we could apply any of these decision procedures to M' and this would tell whether M halts on x . Since we know that the halting problem is undecidable, all of these problems must be undecidable too.

To show that (j) is undecidable, pick your favorite r.e. but nonrecursive set A (HP or MP will do) and modify the above construction as follows. Given M and x , build a new machine M'' that does the following on input y :

- (i) saves y on a separate track of its tape;
- (ii) writes x on a different track (x is hard-wired in the finite control of M'');
- (iii) runs M on input x (M is also hard-wired in the finite control of M'');
- (iv) if M halts on x , then M'' runs a machine accepting A on its original input y , and accepts if that machine accepts.

Either M does not halt on x , in which case the simulation in step (iii) never halts and M'' never accepts any string; or M does halt on x , in which case M'' accepts its input y iff $y \in A$. Thus

$$L(M'') = \begin{cases} A & \text{if } M \text{ halts on } x, \\ \emptyset & \text{if } M \text{ does not halt on } x. \end{cases}$$

Since A is not recursive and \emptyset is, if one could decide whether a given TM accepts a recursive set, then one could apply this decision procedure to M'' and this would tell whether M halts on x .

6 Reduction

There are two main techniques for showing that problems are undecidable: *diagonalization* and *reduction*. We saw examples of diagonalization in §4 and reduction in §5.

Once we have established that a problem such as HP is undecidable, we can show that another problem B is undecidable by *reducing* HP to B . Intuitively, this means we can manipulate instances of HP to make them look like instances of the problem B in such a way that “yes” instances of HP become “yes” instances of B and “no” instances of HP become “no” instances of B . Although we cannot tell effectively whether a given instance of HP is a “yes” instance, the manipulation preserves “yes”-ness and “no”-ness. If there existed a decision procedure for B , then we could apply it to the disguised instances of HP to decide membership in HP. In other words, combining a decision procedure for B with the manipulation procedure would give a decision procedure for HP. Since we have already shown that no such decision procedure for HP can exist, we can conclude that no decision procedure for B can exist.

We can give an abstract definition of reduction and prove a general theorem that will save us a lot of work in undecidability proofs from now on.

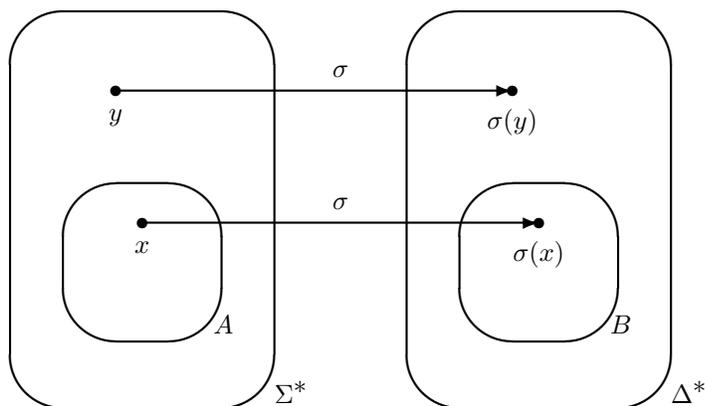
Given sets $A \subseteq \Sigma^*$ and $B \subseteq \Delta^*$, a (many-one) *reduction* of A to B is a computable function

$$\sigma : \Sigma^* \rightarrow \Delta^*$$

such that for all $x \in \Sigma^*$,

$$x \in A \Leftrightarrow \sigma(x) \in B. \quad (5)$$

In other words, strings in A must go to strings in B under σ , and strings not in A must go to strings not in B under σ .



The function σ need not be one-to-one or onto. It must, however, be *total* and *effectively computable*. This means σ must be computable by a total Turing machine that on any input x halts with $\sigma(x)$ written on its tape. When such a reduction exists, we say that A is *reducible* to B via the map σ , and we write $A \leq_m B$. The subscript m , which stands for “many-one,” is used to distinguish this relation from other types of reducibility relations.

The relation \leq_m of reducibility between languages is transitive: if $A \leq_m B$ and $B \leq_m C$, then $A \leq_m C$. This is because if σ reduces A to B and τ reduces B to C , then $\tau \circ \sigma$, the composition of σ and τ , is computable and reduces A to C .

Although we have not mentioned it explicitly, we have used reductions in the last few sections to show that various problems are undecidable.

Example 4. In showing that it is undecidable whether a given TM accepts the null string, we constructed from a given TM M and string x a TM M' that accepted the null string iff M halts on x . In this example,

$$\begin{aligned} A &= \{M\#x \mid M \text{ halts on } x\} = \text{HP}, \\ B &= \{M \mid \varepsilon \in L(M)\}, \end{aligned}$$

and σ is the computable map $M\#x \mapsto M'$.

Example 5. In showing that it is undecidable whether a given TM accepts a regular set, we constructed from a given TM M and string x a TM M'' such that $L(M'')$ is a nonregular set if M halts on x and \emptyset otherwise. In this example,

$$\begin{aligned} A &= \{M\#x \mid M \text{ halts on } x\} = \text{HP}, \\ B &= \{M \mid L(M) \text{ is regular}\}, \end{aligned}$$

and σ is the computable map $M\#x \mapsto M''$.

Here is a general theorem that will save us some work.

Theorem 2.

- (i) If $A \leq_m B$ and B is r.e., then so is A . Equivalently, if $A \leq_m B$ and A is not r.e., then neither is B .
- (ii) If $A \leq_m B$ and B is recursive, then so is A . Equivalently, if $A \leq_m B$ and A is not recursive, then neither is B .

Proof. (i) Suppose $A \leq_m B$ via the map σ and B is r.e. Let M be a TM such that $B = L(M)$. Build a machine N for A as follows: on input x , first compute $\sigma(x)$, then run M on input $\sigma(x)$, accepting if M accepts. Then

$$\begin{aligned} N \text{ accepts } x &\Leftrightarrow M \text{ accepts } \sigma(x) && \text{definition of } N \\ &\Leftrightarrow \sigma(x) \in B && \text{definition of } M \\ &\Leftrightarrow x \in A && \text{by (5)}. \end{aligned}$$

(ii) Recall from Lecture 2 that a set is recursive iff both it and its complement are r.e. Suppose $A \leq_m B$ via the map σ and B is recursive. Note that $\sim A \leq_m \sim B$ via the same σ (Check the definition!). If B is recursive, then both B and $\sim B$ are r.e. By (i), both A and $\sim A$ are r.e., thus A is recursive. \square

We can use Theorem 2(i) to show that certain sets are not r.e. and Theorem 2(ii) to show that certain sets are not recursive. To show that a set B is not r.e., we need only give a reduction from a set A we already know is not r.e. (such as $\sim\text{HP}$) to B . By Theorem 2(i), B cannot be r.e.

Example 6. Let's illustrate by showing that neither the set

$$\text{FIN} = \{M \mid L(M) \text{ is finite}\}$$

nor its complement is r.e. We show that neither of these sets is r.e. by reducing $\sim\text{HP}$ to each of them, where

$$\sim\text{HP} = \{M\#x \mid M \text{ does not halt on } x\} :$$

- (a) $\sim\text{HP} \leq_m \text{FIN}$,
- (b) $\sim\text{HP} \leq_m \sim\text{FIN}$.

Since we already know that $\sim\text{HP}$ is not r.e., it follows from Theorem 2(i) that neither FIN nor $\sim\text{FIN}$ is r.e. For (a), we want to give a computable map σ such that

$$M\#x \in \sim\text{HP} \Leftrightarrow \sigma(M\#x) \in \text{FIN}.$$

In other words, from $M\#x$ we want to construct a Turing machine $M' = \sigma(M\#x)$ such that

$$M \text{ does not halt on } x \Leftrightarrow L(M') \text{ is finite.} \tag{6}$$

Note that the description of M' can depend on M and x . In particular, M' can have a description of M and the string x hard-wired in its finite control if desired.

We have actually already given a construction satisfying (6). Given $M\#x$, construct M' such that on all inputs y , M' takes the following actions:

- (i) erases its input y ;
- (ii) writes x on its tape (M' has x hard-wired in its finite control);
- (iii) runs M on input x (M' also has a description of M hard-wired in its finite control);
- (iv) accepts if M halts on x .

If M does not halt on input x , then the simulation in step (iii) never halts, and M' never reaches step (iv). In this case M' does not accept its input y . This happens the same way for all inputs y , therefore in this case, $L(M) = \emptyset$. On the other hand, if M does halt on x , then the simulation in step (iii) halts, and y is accepted in step (iv). Moreover, this is true for all y . In this case, $L(M) = \Sigma^*$. Thus

$$\begin{aligned} M \text{ halts on } x &\Rightarrow L(M') = \Sigma^* \Rightarrow L(M') \text{ is infinite,} \\ M \text{ does not halt on } x &\Rightarrow L(M') = \emptyset \Rightarrow L(M') \text{ is finite.} \end{aligned}$$

Thus (6) is satisfied. Note that this is all we have to do to show that FIN is not r.e.: we have given the reduction (a), so by Theorem 2(i) we are done.

There is a common pitfall here that we should be careful to avoid. It is important to observe that the computable map σ that produces a description of M' from M and x does not need to execute the program (i) through (iv). It only produces the description of a machine M' that does so. The computation of σ is quite simple—it does not involve the simulation of any other machines or anything complicated at all. It merely takes a description of a Turing machine M and string x and plugs them into a general description of a machine that executes (i) through (iv). This can be done quite easily by a total TM, so σ is total and effectively computable.

Now (b). By definition of reduction, a map reducing $\sim\text{HP}$ to $\sim\text{FIN}$ also reduces HP to FIN, so it suffices to give a computable map τ such that

$$M\#x \in \text{HP} \Leftrightarrow \tau(M\#x) \in \text{FIN}.$$

In other words, from M and x we want to construct a Turing machine $M'' = \tau(M\#x)$ such that

$$M \text{ halts on } x \Leftrightarrow L(M'') \text{ is finite.} \tag{7}$$

Given $M\#x$, construct a machine M'' that on input y

- (i) saves y on a separate track;
- (ii) writes x on the tape;
- (iii) simulates M on x for $|y|$ steps (it erases one symbol of y for each step of M on x that it simulates);
- (iv) accepts if M has *not* halted within that time, otherwise rejects.

Now if M never halts on x , then M'' halts and accepts y in step (iv) after $|y|$ steps of the simulation, and this is true for all y . In this case $L(M'') = \Sigma^*$. On the other hand, if M does halt on x , then it does so after some finite number of steps, say n . Then M'' accepts y in (iv) if $|y| < n$ (since the simulation in (iii) has not finished by $|y|$ steps) and rejects y in (iv) if $|y| \geq n$ (since the simulation in (iii) does have time to complete). In this case M'' accepts all strings of length less than n and rejects all strings of length n or greater, so $L(M'')$ is a finite set. Thus

$$\begin{aligned} M \text{ halts on } x &\Rightarrow L(M'') = \{y \mid |y| < \text{running time of } M \text{ on } x\} \\ &\Rightarrow L(M'') \text{ is finite,} \end{aligned}$$

$$\begin{aligned} M \text{ does not halt on } x &\Rightarrow L(M'') = \Sigma^* \\ &\Rightarrow L(M'') \text{ is infinite.} \end{aligned}$$

Then (7) is satisfied.

It is important that the functions σ and τ in these two reductions can be computed by Turing machines that always halt.

Historical Notes

The technique of diagonalization was first used by Cantor [1] to show that there were fewer real algebraic numbers than real numbers.

Universal Turing machines and the application of Cantor's diagonalization technique to prove the undecidability of the halting problem appear in Turing's original paper [20].

Reducibility relations are discussed by Post [14]; see [17, 19].

7 Rice's Theorem

Rice's theorem says that undecidability is the rule, not the exception. It is a very powerful theorem, subsuming many undecidability results that we have seen as special cases.

Theorem 3 (Rice's theorem). *Every nontrivial property of the r.e. sets is undecidable.*

Yes, you heard right: that's *every* nontrivial property of the r.e. sets. So as not to misinterpret this, let us clarify a few things.

First, fix a finite alphabet Σ . A *property of the r.e. sets* is a map

$$P : \{\text{r.e. subsets of } \Sigma^*\} \rightarrow \{\mathbf{1}, \mathbf{0}\},$$

where $\mathbf{1}$ and $\mathbf{0}$ represent truth and falsity, respectively. For example, the property of emptiness is represented by the map

$$P(A) = \begin{cases} \mathbf{1} & \text{if } A = \emptyset, \\ \mathbf{0} & \text{if } A \neq \emptyset. \end{cases}$$

To ask whether such a property P is decidable, the set has to be presented in a finite form suitable for input to a TM. We assume that r.e. sets are presented by TMs that accept them. But keep in mind that the property is a property of *sets*, not of Turing machines; thus it must be true or false independent of the particular TM chosen to represent the set.

Here are some other examples of properties of r.e. sets: $L(M)$ is finite; $L(M)$ is recursive; M accepts 101001 (*i.e.*, $101001 \in L(M)$); $L(M) = \Sigma^*$. Each of these properties is a property of the set accepted by the Turing machine.

Here are some examples of properties of Turing machines that are *not* properties of r.e. sets: M has at least 481 states; M halts on all inputs; M rejects 101001; there exists a smaller machine equivalent to M . These are not properties of sets, because in each case one can give two TMs that accept the same set, one of which satisfies the property and the other of which doesn't.

For Rice's theorem to apply, the property also has to be *nontrivial*. This just means that the property is neither universally true nor universally false; that is, there must be at least one r.e. set that satisfies the property and at least one that does not. There are only two trivial properties, and they are both trivially decidable.

Proof of Rice's theorem. Let P be a nontrivial property of the r.e. sets. Assume without loss of generality that $P(\emptyset) = \mathbf{0}$ (the argument is symmetric if $P(\emptyset) = \mathbf{1}$). Since P is nontrivial, there must exist an r.e. set A such that $P(A) = \mathbf{1}$. Let K be a TM accepting A .

We reduce HP to the set $\{M \mid P(L(M)) = \mathbf{1}\}$, thereby showing that the latter is undecidable (Theorem 2(ii)). Given $M\#x$, construct a machine $M' = \sigma(M\#x)$ that on input y

- (i) saves y on a separate track someplace;
- (ii) writes x on its tape (x is hard-wired in the finite control of M');

- (iii) runs M on input x (a description of M is also hard-wired in the finite control of M');
- (iv) if M halts on x , M' runs K on y and accepts if K accepts.

Now either M halts on x or not. If M does not halt on x , then the simulation in (iii) will never halt, and the input y of M' will not be accepted. This is true for every y , so in this case $L(M') = \emptyset$. On the other hand, if M does halt on x , then M' always reaches step (iv), and the original input y of M' is accepted iff y is accepted by K ; that is, if $y \in A$. Thus

$$\begin{aligned} M \text{ halts on } x &\Rightarrow L(M') = A \Rightarrow P(L(M')) = P(A) = \mathbf{1}, \\ M \text{ does not halt on } x &\Rightarrow L(M') = \emptyset \Rightarrow P(L(M')) = P(\emptyset) = \mathbf{0}. \end{aligned}$$

This constitutes a reduction from HP to the set $\{M \mid P(L(M)) = \mathbf{1}\}$. Since HP is not recursive, by Theorem 2, neither is the latter set; that is, it is undecidable whether $L(M)$ satisfies P . \square

Rice's Theorem, Part II

A property $P : \{\text{r.e. sets}\} \rightarrow \{\mathbf{1}, \mathbf{0}\}$ of the r.e. sets is called *monotone* if for all r.e. sets A and B , if $A \subseteq B$, then $P(A) \leq P(B)$. Here \leq means less than or equal to in the order $\mathbf{0} \leq \mathbf{1}$. In other words, P is *monotone* if whenever a set has the property, then all supersets of that set have it as well. For example, the properties “ $L(M)$ is infinite” and “ $L(M) = \Sigma^*$ ” are monotone but “ $L(M)$ is finite” and “ $L(M) = \emptyset$ ” are not.

Theorem 4 (Rice's theorem, part II). *No nonmonotone property of the r.e. sets is semidecidable. In other words, if P is a nonmonotone property of the r.e. sets, then the set $T_P = \{M \mid P(L(M)) = \mathbf{1}\}$ is not r.e.*

Proof. Since P is nonmonotone, there exist TMs M_0 and M_1 such that $L(M_0) \subseteq L(M_1)$, $P(M_0) = \mathbf{1}$, and $P(M_1) = \mathbf{0}$.

We want to reduce \sim HP to T_P , or equivalently, HP to $\sim T_P = \{M \mid P(L(M)) = \mathbf{0}\}$. Since \sim HP is not r.e., neither will be T_P . Given $M \# x$, we want to show how to construct a machine M' such that $P(M') = \mathbf{0}$ iff M halts on x . Let M' be a machine that does the following on input y :

- (i) writes its input y on the top and middle tracks of its tape;
- (ii) writes x on the bottom track (it has x hard-wired in its finite control);
- (iii) simulates M_0 on input y on the top track, M_1 on input y on the middle track, and M on input x on the bottom track in a round-robin fashion; that is, it simulates one step of each of the three machines, then another step, and so on (descriptions of M_0 , M_1 , and M are all hard-wired in the finite control of M');
- (iv) accepts its input y if either of the following two events occurs:
 - (a) M_0 accepts y , or
 - (b) M_1 accepts y and M halts on x .

Either M halts on x or not, independent of the input y to M' . If M does not halt on x , then event (b) in step (iv) will never occur, so M' will accept y iff event (a) occurs, thus in this case $L(M') = L(M_0)$. On the other hand, if M does halt on x , then y will be accepted iff it is accepted by either M_0 or M_1 ; that is, if $y \in L(M_0) \cup L(M_1)$. Since $L(M_0) \subseteq L(M_1)$, this is equivalent to saying that $y \in L(M_1)$, thus in this case $L(M') = L(M_1)$. We have shown

$$\begin{aligned} M \text{ halts on } x &\Rightarrow L(M') = L(M_1) \\ &\Rightarrow P(L(M')) = P(L(M_1)) = \mathbf{0}, \end{aligned}$$

$$\begin{aligned} M \text{ does not halt on } x &\Rightarrow L(M') = L(M_0) \\ &\Rightarrow P(L(M')) = P(L(M_0)) = \mathbf{1}. \end{aligned}$$

The construction of M' from M and x constitutes a reduction from \sim HP to the set $T_P = \{M \mid P(L(M)) = \mathbf{1}\}$. By Theorem 2(i), the latter set is not r.e. \square

Historical Notes

Rice's theorem was proved by H. G. Rice [15, 16].

References

- [1] G. CANTOR, Über eine Eigenschaft des Inbegriffes aller reellen algebraischen Zahlen, *J. für die reine und angewandte Mathematik*, 77 (1874), pp. 258–262. Reprinted in *Georg Cantor Gesammelte Abhandlungen*, Berlin, Springer-Verlag, 1932, pp. 115–118.
- [2] A. CHURCH, A set of postulates for the foundation of logic, *Ann. Math.*, 33–34 (1933), pp. 346–366, 839–864.
- [3] ———, An unsolvable problem of elementary number theory, *Amer. J. Math.*, 58 (1936), pp. 345–363.
- [4] H.B. CURRY, An analysis of logical substitution, *Amer. J. Math.*, 51 (1929), pp. 363–384.
- [5] P.C. FISCHER, Turing machines with restricted memory access, *Information and Control*, 9 (1966), pp. 364–379.
- [6] P.C. FISCHER, A.R. MEYER, AND A.L. ROSENBERG, Counter machines and counter languages, *Math. Systems Theory*, 2 (1968), pp. 265–283.
- [7] K. GÖDEL, On undecidable propositions of formal mathematical systems, in *The Undecidable*, M. Davis, ed., Raven Press, Hewlett, NY, 1965, pp. 5–38.
- [8] S.C. KLEENE, A theory of positive integers in formal logic, *Amer. J. Math.*, 57 (1935), pp. 153–173, 219–244.
- [9] ———, Recursive predicates and quantifiers, *Trans. Amer. Math. Soc.*, 53 (1943), pp. 41–74.
- [10] D. KOZEN, *Automata and Computability*, Springer-Verlag, 1997.
- [11] M.L. MINSKY, Recursive unsolvability of Post's problem of 'tag' and other topics in the theory of Turing machines, *Ann. Math.*, 74 (1961), pp. 437–455.
- [12] E. POST, Finite combinatorial processes-formulation, I, *J. Symbolic Logic*, 1 (1936), pp. 103–105.
- [13] ———, Formal reductions of the general combinatorial decision problem, *Amer. J. Math.*, 65 (1943), pp. 197–215.
- [14] ———, Recursively enumerable sets of positive natural numbers and their decision problems, *Bull. Amer. Math. Soc.*, 50 (1944), pp. 284–316.
- [15] H.G. RICE, Classes of recursively enumerable sets and their decision problems, *Trans. Amer. Math. Soc.*, 89 (1953), pp. 25–59.
- [16] ———, On completely recursively enumerable classes and their key arrays, *J. Symbolic Logic*, 21 (1956), pp. 304–341.
- [17] H. ROGERS JR., *Theory of Recursive Functions and Effective Computability*, MIT Press, 1967.
- [18] M. SCHÖNFINKEL, Über die Bausteine der mathematischen Logik, *Math. Annalen*, 92 (1924), pp. 305–316.
- [19] R.I. SOARE, *Recursively Enumerable Sets and Degrees*, Springer-Verlag, Berlin, 1987.
- [20] A.M. TURING, On computable numbers with an application to the Entscheidungsproblem, *Proc. London Math. Soc.*, 42 (1936), pp. 230–265. Erratum: *Ibid.*, 43 (1937), pp. 544–546.
- [21] A.N. WHITEHEAD AND B. RUSSELL, *Principia Mathematica*, Cambridge University Press, Cambridge, 1910–1913. Three volumes.