

## Week 1: Undecidability and Decision Problems

*Dylan Hendrickson**MIT Educational Studies Program*

## 1.1 The Halting Problem

We'll start by considering whether programs can solve certain problems at all. Most problems you encounter can be solved by a program, though it may take a huge amount of time, perhaps brute forcing its way through every potential solution. But there are some problems that it's not clear how to write a program to solve. For example, the halting problem: given a program and some input, does the program eventually halt on that input? Some programs will keep running in an infinite loop, and it's not always obvious when this will happen. (Often, it will be easy to tell that a particular program will halt, or that it will loop, but we want an algorithm that *always* tells us whether a program halts.)

An attempt at an algorithm is to run the program in question; if it halts, then you can answer 'yes.' But if it doesn't halt, then you'll never find out; maybe it will keep running forever or maybe it will suddenly stop if you simulate it for a little longer.

**Theorem 1.1.** *There is no program that solves the halting problem. In other words, the halting problem is undecidable.*

*Proof.* Suppose there were a program  $H$  which solved the halting problem; that is, for any program  $P$  and input  $x$ ,  $H(P, x)$  answers either 'yes' or 'no' depending on whether  $P(x)$  halts. We'll show that this assumption leads to a contradiction, so it must be false. Use  $H$  to construct a new program  $G$  which does the following when input a program  $P$ :

- Run  $H(P, P)$ ; this tells us whether  $P$  halts when given its own code as input.
- If  $H(P, P)$  is 'yes,' enter an infinite loop.
- If  $H(P, P)$  is 'no,' halt.

Now  $G(P)$  halts exactly when  $H(P, P)$  is 'no,' which is when  $P(P)$  loops. Let's consider what happens when we run  $G$  on its own code, by plugging in  $G$  for  $P$ . We find that  $G(G)$  halts exactly when  $G(G)$  loops, which is impossible. So our assumption that  $H$  exists must be wrong.  $\square$

## 1.2 Proving other problems undecidable

Next let's consider the following question: given a program  $P$ , is there any input on which  $P$  outputs 'yes'?

We'll show this problem is undecidable by converting instances of the halting problem to it. Suppose we had a program  $E$  which tells us whether an arbitrary program can output 'yes.' We solve the halting problem as follows:

- We want to know whether  $P(x)$  halts.

- Create a new program  $Q$  based on  $P$  and  $x$ .
- Use the assumed program  $E$  to see whether  $Q$  outputs ‘yes’ on any inputs.
- This answers the original question of whether  $P(x)$  halts.

Here’s how we construct  $Q$ . It’s the program which takes an input  $y$  and does the following:

- If  $y \neq 0$ , output ‘no.’
- If  $y = 0$ , run  $P(x)$ .
- When it halts, output ‘yes.’

Clearly  $Q(0)$  outputs ‘yes’ iff  $P(x)$  halts, and  $Q$  otherwise never outputs ‘yes.’ So we can use our program  $E$  to solve the halting problem; since the halting problem is undecidable,  $E$  can’t exist.

Notice the general strategy here: we know a problem  $A$  (here, the halting problem) is hard, and we use this to show another problem  $B$  hard. To do so, we describe a way to convert instances of  $A$  into instances of  $B$ ; then if we can solve  $B$  quickly, we’d be able to also solve  $A$  quickly. In general, this is called a *reduction* from  $A$  to  $B$ , and it means  $B$  is at least as hard to solve as  $A$ . We’ll talk more about this next week.

### 1.3 Decision Problems

When we talk about what programs are capable of, we usually mean the classes of yes/no questions they can answer. Such a class of yes/no questions is called a *decision problem*. Notice that it’s uninteresting to consider only a single yes/no question, since it has an answer of either ‘yes’ or ‘no,’ and either the program that outputs ‘yes’ or the program that outputs ‘no’ solves it (even if we don’t know which one). We instead consider infinite families of yes/no questions, and look for programs that can answer all of them. When a program answers ‘yes,’ we say it *accepts* the input, and when it answers ‘no,’ we say it *rejects* the input.

If there is a program that answers every yes/no question in a decision problem, we say that the program *decides* the decision problem, and that the decision problem is *decidable*. Most decision problems you care about are decidable, but we’ve just seen that the halting problem and the problem of whether a program accepts any input and both undecidable.

If there is a program that accepts whenever the answer is ‘yes,’ and either rejects or runs forever whenever the answer is ‘no,’ we say that the program *recognizes* the decision problem, and that the decision problem is *recognizable*. The algorithm which, given inputs  $P$  and  $x$ , runs  $P(x)$  until it halts and then accepts, recognizes the halting problem. If  $P(x)$  halts, this algorithm will eventually accept, and if  $P(x)$  doesn’t halt, this algorithm runs forever. So the halting problem is recognizable but not decidable.

Deciders have to always given an answer, but recognizers only have to answer when the answer is ‘yes,’ and are allowed to loop when the answer is ‘no.’ What if we allow the program to loop when the answer is ‘yes’ but not when it’s ‘no’? Then we say that the program *corecognizes* the decision problem, and the decision problem is *corecognizable*. For example, take the negation of the halting problem: given a program  $P$  and input  $x$ , does  $P(x)$  run forever? The program which runs  $P(x)$  and rejects when it halts corecognizes this decision problem; if the answer is ‘no,’ then  $P(x)$  halts and the program rejects, and if the answer is ‘yes’ then it runs forever.

**Exercise 1.** Is the problem shown undecidable in Section 1.2 (determining whether a program accepts any inputs) recognizable? Is it corecognizable?

We are now in a position to prove some results about decidability, recognizability, and corecognizability, but I'd like to introduce some notation first. If  $L$  is a decision problem, we write  $x \in L$  to mean  $x$  is in  $L$ , or the answer to the question corresponding to  $x$  is 'yes.' We write  $x \notin L$  if the answer to the question corresponding to  $x$  is 'no.' The *complement*  $\bar{L}$  of  $L$  is the decision problem with all answers opposite those of  $L$ , so  $x \in L$  exactly when  $x \notin \bar{L}$  and vice-versa. Note that  $\bar{\bar{L}} = L$ . I'll also use the shorthand 'iff' for 'if and only if.'

(I might sometimes say 'language' instead of 'decision problem'; the terms are essentially equivalent, but I think it's easier to think in terms of decision problems. A language is a set of inputs, which you think of as the inputs for which the answer is 'yes.')

**Lemma 1.2.** *A decision problem  $L$  is recognizable iff  $\bar{L}$  is corecognizable.*

**Exercise 2.** Prove Lemma 1.2. Once you understand the notation, you should be able to see that it's not saying anything particularly deep or surprising.

**Lemma 1.3.** *A decision problem  $L$  is decidable iff it is both recognizable and corecognizable.*

*Proof.* There are two directions to prove; one is easy. First suppose  $L$  is decidable. The program that decides  $L$  both recognizes and corecognizes  $L$ , so  $L$  is recognizable and corecognizable.

Now suppose  $L$  is recognizable and corecognizable. Then there are programs  $F$  and  $G$  which recognize and corecognize  $L$ , respectively. That means that whenever  $x \in L$ ,  $F$  accepts  $x$  (and  $G$  either accepts or loops on  $x$ ), and whenever  $x \notin L$ ,  $G$  rejects  $x$  (and  $F$  either rejects or loops on  $x$ ). We construct a new program  $P$  which simulates  $F$  and  $G$  in parallel. The exact nature of the simulation isn't important; maybe it alternates running single lines of code from  $F$  and  $G$ , keeping them separate from each other. If the simulation of  $F$  accepts,  $P$  accepts, and if the simulation of  $G$  rejects,  $P$  rejects. We know that exactly one of these will happen, depending on whether  $x \in L$ . If  $x \in L$ ,  $P$  will accept  $x$ , and if  $x \notin L$ ,  $P$  will reject  $x$ . Thus  $P$  decides  $L$ , so  $L$  is decidable.  $\square$

Since we know the halting problem is recognizable but not decidable, it must not be corecognizable.

There's another equivalent way to define recognizability (actually multiple, but one that is important to us). Let  $L$  be a decision problem. We will see that  $L$  is recognizable if there is a program  $P$  which always halts, such that for every  $x \in L$ , there is some positive integer  $n$  such that  $P(x, n)$  accepts, and if  $P(x, n)$  accepts, then  $x \in L$ . We say that  $n$  is a *certificate* for  $x$ ; we can use  $n$  convince the program that the answer to the question corresponding to  $x$  is 'yes.'

**Theorem 1.4.** *This is an equivalent definition of recognizability. That is, a decision problem  $L$  is recognizable iff there is a decidable decision problem  $D$  such that*

- If  $x \in L$ , there is a positive integer  $n$  such that  $(x, n) \in D$
- If  $(x, n) \in D$ , then  $x \in L$ .

Before reading the proof, make sure you understand what this theorem is saying. If you're feeling up to it, try to prove it yourself. The idea of the proof is that the certificate  $n$  that convinces you that  $x \in L$  is the number of steps the program that recognizes  $L$  takes to accept on input  $x$ .

*Proof.* Suppose  $L$  is recognizable. Let  $F$  be a program that recognizes  $L$ . We construct a decision problem  $D$  satisfying the conditions: given  $x$  and  $n$ , does  $F(x)$  accept within  $n$  steps? If  $x \in L$ , then  $F(x)$  must accept at some point, so there is some  $n$  such that  $(x, n) \in D$ . If  $(x, n) \in D$ , then  $F(x)$  accepts within  $n$

steps, so in particular  $F(x)$  accepts, and thus  $x \in L$ . Finally,  $D$  is decidable: the program which, on input  $x$  and  $n$ , simulates  $F(x)$  for  $n$  steps and accepts if the simulation accepted and rejects otherwise, decides  $D$ .

Now suppose there is a decision problem  $D$  satisfying the conditions decided by a program  $P$ . We write a program that recognizes  $L$ . On input  $x$ , for each positive integer  $n$  in order, run  $P(x, n)$ . If  $P(x, n)$  accepts, accept; otherwise move on to the next  $n$ .

If  $x \in L$ , then this program will eventually find an  $n$  such that  $(x, n) \in D$ , and accept. If  $x \notin L$ , then no such  $n$  will work, and the program will run forever.  $\square$

**Exercise 3.** Give an equivalent definition of corecognizability along the lines of Theorem 1.4, and prove that it's equivalent (you can, but don't have to, use Theorem 1.4 in your proof).

Decision problems are organized into *complexity classes*. A complexity class contains all the decision problems that can be 'easily solved' by a program; the meaning of 'easily solved' is what distinguishes different complexity classes. When 'easily solved' means 'decided,' we have the complexity class of decidable problems, called **R** (for 'recursive,' an older word for 'decidable'). When 'easily solved' means 'recognized' or 'corecognized,' we have the classes **RE** and **coRE** (for 'recursively enumerable,' an older term describing an alternate definition of 'recognizable'). There are hundreds of named complexity classes, a few of which we'll get to know in this class. A table of some of the important ones for our purposes is at the end of these notes; if you want to see more, visit [complexityzoo.uwaterloo.ca](https://www.math.ucdavis.edu/~greg/zoology/diagram.xml) for descriptions of complexity classes or <https://www.math.ucdavis.edu/~greg/zoology/diagram.xml> for an interactive visualization.

We define the complement  $\bar{\mathbf{S}}$  of a complexity class **S** as the collection of the complements of languages in **S**. Lemma 1.2 can then be expressed as **coRE** =  $\bar{\mathbf{RE}}$ . Lemma 1.3 can be expressed as **R** = **RE**  $\cap$  **coRE** (the symbol  $\cap$  is for intersection;  $A \cap B$  is the collection of things that are in both  $A$  and  $B$ ). The relation between **R** and **RE** is analogous to the relation between **P** and **NP** (complexity classes we'll see more of later), so it would be reasonable to use **NR** for **RE**, the N standing for 'nondeterministic,' which sort of means 'with guessing.'

Name	Stands for	Meaning of 'easily solved'	Complete problem
<b>L</b>	log space	answered in logarithmic space	undirected reachability
<b>NL</b>	nondeterministic <b>L</b>	verified in logarithmic space	directed reachability
<b>P</b>	polynomial	answered in polynomial time	formula evaluation
<b>NP</b>	nondeterministic <b>P</b>	verified in polynomial time	formula satisfiability
<b>coNP</b>	complement of <b>NP</b>	verified false in polynomial time	formula unsatisfiability
<b>PH</b>	polynomial hierarchy	generalization of <b>NP</b>	
<b>PSPACE</b>	polynomial space	answered in polynomial space	quantified boolean formula
<b>EXP</b>	exponential	answered in exponential time	halt in $k$ steps
<b>NEXP</b>	nondeterministic <b>EXP</b>	verified in exponential time	succinct Hamiltonian path
<b>EXSPACE</b>	exponential space	answered in exponential space	regex equivalence
<b>R</b>	recursive	answered at all	
<b>RE</b>	recursively enumerable	verified at all	program halts
<b>coRE</b>	complement of <b>RE</b>	verified false at all	program loops
<b>AH</b>	arithmetical hierarchy	generalization of <b>RE</b>	