## Theory of Computation: The Recursion Theorem

Arjun Chandrasekhar



#### This is my favorite lecture

- I love this topic
- This is the first lecture I ever made from scratch
- For once, Turing machines will be more convenient than modern programming languages
- This lecture got me my first full-time teaching job
- This lecture lead to one of my favorite teaching stories

#### Pop quiz!

For extra credit: write a non-empty program that prints out its own source code (without using any I/O operations). You have 20 minutes.

#### Recall

- Certain problems are not decidable, or even recognizable
- Initially we proved this with Diagonalization
- Usually we prove this using reduction
  - ► If A is undecidable/unrecognizable, and A ≤<sub>M</sub> B, then B is undecidable/unrecognizable
- Is there an alternative (perhaps easier) way to prove these results?

### Today's goals

- Understand how to construct a self-reproducing Turing machine
- Describe Turing machines that analyze (and contradict) their own behavior
- Use the recursion theorem for concise impossiblity proofs

#### Self-reproducing robots

- Can we build a robot that builds another robot?
- Can we build a robot that builds another robot-building robot?
- Can we build a robot-building robot that builds an <u>identical</u> robot-building robot?

- Can we build write a robot program that builds another robot program?
- Can we build write a robot program that builds another robot-building robot program-building program?
- Can we build a robot-building robot program-building program that builds an identical robot-building robot program-building program?

#### Statement of the recursion theorem

#### The recursion theorem:

- 1. There exists a Turing machine that prints out its own description
- 2. A Turing machine can be programmed to obtain and then analyze its own description.

We will prove both parts by construction.

In doing so, we will see why the simplicity of the TM model can be convenient!

#### Idea behind a self-reproducing programs

Print the following sentence twice, the second time in quotes "Print the following sentence twice, the second time in quotes"

#### Review of notation

- If M is a Turing machine, then (M) refers to the string description of M
- Can think of (M) as the source code, and M is the executable file
- Think back to the first two programming assignments
  - On the first assignment, you read a string description (D) of a DFA and parsed it into an actual DFA object D
  - On the second assignment, you read a string description (N) of an NFA, converted it into an equivalent DFA object D, and output a string description (D)

## Building blocks of our self-reproducing program

To construct a Turing machine that prints out its own description, we will split the program up into two sub-machines Q and  $P_{\langle Q \rangle}$ 

- Machine Q takes an input string w and creates a machine P<sub>w</sub> that always prints out the same string w
- Machine P<sub>(Q)</sub> is a Turing machine that ignores its input and always prints out (Q), i.e. a description of the machine Q

#### Machines that only print one string



#### Machines that only print one string

Let  $P_w$  be a Turing machine that erases whatever is on the tape and prints out string w.

- w is a constant that we decide on ahead of time
- Can a human construct this machine?
- Can we write a computer program to construct this machine?

Let  $q(w) = \langle P_w \rangle$ . That is, the function q takes as input a string w. It outputs the description of a Turing machine  $P_w$  that always prints out the string w.

- Is q a Turing-computable function?
- If you understand what q is doing and convince yourself that q is computable, then the proof of the recursion theorem will be straightforward

- Consider the Turing machine P<sub>computability</sub> ignores its input and always prints out the same string "computability"
- Can we construct a machine Q<sub>computability</sub> that prints out (P<sub>computability</sub>)?
  - Yes!
  - We can treat (P<sub>computability</sub>) like any other pre-determined string w

15

- Consider the Turing machine P<sub>recursion</sub> ignores its input and always prints out the same string "recursion"
- Can we construct a machine  $Q_{\text{recursion}}$  that prints out  $\langle P_{\text{recursion}} \rangle$ ?
  - Yes!
  - It's just like the last slide. We use the string "recursion" instead of the string "computability"

- Can we construct a machine Q<sub>w</sub> that prints out (P<sub>w</sub>) for any arbitrary string w?
  - Yes!
  - It's just like the last two slides. Instead of the string "recursion" or "computability", we substitute our string of choice w

- Can we construct a machine Q that takes as input w and outputs (P<sub>w</sub>), the description of the machine P<sub>w</sub>?
- Yes!
- We know what 99% of the program P<sub>w</sub> will look like
- If we are given w, we know how to finish writing the program
  - Substitute in w for the part where P prints out its output string
  - We basically convert w from a constant to a parameter of the method Q

Now we are ready to describe the machine SELF, a machine which prints out its own description. SELF consists of two sub-machines,  $P_{\langle Q \rangle}$  and Q.

- ▶ Q takes as input ⟨M⟩, a Turing machine description, and does the following:
  - 1. Compute  $\langle P_{\langle M \rangle} \rangle$ , i.e. a description of a machine that ignores its input and prints out the string  $\langle M \rangle$  (where  $\langle M \rangle$  is also a Turing machine description)
  - 2. Combine  $\langle P_{\langle M \rangle} 
    angle$  and  $\langle M 
    angle$  into a single machine  $M^*$

19

- 3. Print out  $\langle M^* \rangle$  and halt
- ► P<sub>(Q)</sub> ignores whatever input is on the tape and prints out the description of machine Q.
  SELF first runs P<sub>(Q)</sub>, and then runs Q. What happens?

SELF first runs  $P_{\langle Q \rangle}$ , and then runs Q. What happens?

- 1. Machine  $P_{\langle Q \rangle}$  prints  $\langle Q \rangle$  on the tape, and passes control to Q.
- 2. Machine Q reads  $\langle Q \rangle$  on the tape. That is, it encounters its own description, which was produced by  $P_{\langle Q \rangle}$
- 3. Machine Q constructs the machine  $P_w$ , which in this case is  $P_{\langle Q \rangle}$ .
- 4. Machine Q then combines the machines  $\langle P_{\langle Q \rangle} \rangle$ and  $\langle Q \rangle$  into one machine  $M^*$

Does that last step loop familiar?

SELF first runs  $P_{\langle Q \rangle}$ , and then runs Q. What happens?

- SELF produces a machine (M\*) which is a combination of P(Q) and Q
- But SELF also of a combination of  $P_{\langle Q \rangle}$  and Q!
- SELF has reproduced its own description!

21

#### SELF = $\langle P_{(Q)} Q \rangle$

Initially  $P_{(Q)}$  starts with control of the tape, with an input string w

- 1.  $P_{\langle Q \rangle}$  erases w and prints  $\langle Q \rangle$  on the tape, then passes control to Q
- Q reads (Q) on the tape. That is, it encounters its own description which was produced by P<sub>(Q)</sub>
- Q reads ⟨Q⟩ and uses this to construct P<sub>(Q)</sub>
- Q combines (P<sub>(Q)</sub>) and (Q) into a single machine (M<sup>\*</sup>)



SELF =  $\langle \mathsf{P}_{(\mathsf{Q})} \mathsf{Q} \rangle$ 

Initially  $P_{(Q)}$  starts with control of the tape, with an input string w

- P<sub>⟨Q⟩</sub> erases w and prints ⟨Q⟩ on the tape, then passes control to Q
- Q reads (Q) on the tape. That is, it encounters its own description which was produced by P<sub>(Q)</sub>
- Q reads ⟨Q⟩ and uses this to construct P<sub>(Q)</sub>
- Q combines (P<sub>(Q</sub>) and (Q) into a single machine (M<sup>\*</sup>)





SELF =  $\langle \mathsf{P}_{(\mathsf{Q})} \mathsf{Q} \rangle$ 

Initially  $P_{(Q)}$  starts with control of the tape, with an input string w

- 1.  $P_{\langle Q \rangle}$  erases w and prints  $\langle Q \rangle$  on the tape, then passes control to Q
- Q reads (Q) on the tape. That is, it encounters its own description which was produced by P<sub>(Q)</sub>
- Q reads ⟨Q⟩ and uses this to construct P<sub>(Q)</sub>
- Q combines (P<sub>(Q</sub>) and (Q) into a single machine (M<sup>\*</sup>)



SELF =  $\langle \mathsf{P}_{(\mathsf{Q})} \mathsf{Q} \rangle$ 

Initially  $P_{(Q)}$  starts with control of the tape, with an input string w

- 1.  $P_{\langle Q \rangle}$  erases w and prints  $\langle Q \rangle$  on the tape, then passes control to Q
- Q reads (Q) on the tape. That is, it encounters its own description which was produced by P<sub>(Q)</sub>
- Q reads ⟨Q⟩ and uses this to construct P<sub>(Q)</sub>
- Q combines (P<sub>(Q</sub>) and (Q) into a single machine (M<sup>\*</sup>)





SELF =  $\langle \mathsf{P}_{(\mathsf{Q})} \mathsf{Q} \rangle$ 

Initially  $P_{(Q)}$  starts with control of the tape, with an input string w

- 1.  $P_{\langle Q \rangle}$  erases w and prints  $\langle Q \rangle$  on the tape, then passes control to Q
- Q reads (Q) on the tape. That is, it encounters its own description which was produced by P<sub>(Q)</sub>
- Q reads ⟨Q⟩ and uses this to construct P<sub>(Q)</sub>
- Q combines (P<sub>⟨Q⟩</sub>) and (Q) into a single machine (M<sup>\*</sup>)

$$\begin{array}{c} \mathsf{Q} \\ \mathsf{$$

- Let T be a Turing machine that computes a function  $t : \Sigma^* \times \Sigma^* \to \Sigma^*$ .
  - T represents a machine that receives a machine description as one of its inputs and analyzes that machine
- There is a Turing machine R that computes a function  $r : \Sigma^* \to \Sigma^*$ , where for every w $r(w) = t(\langle R \rangle, w)$ 
  - R performs the same analysis as T does, but it analyzes its own description rather than taking a machine description as one of its inputs

23

We construct R in three parts:  $P_{\langle QT \rangle}, Q, T$ 

- 1.  $P_{\langle QT \rangle}$  prints out the description  $\langle QT \rangle$ , a combination of machines Q and T.
  - Instead of erasing whatever input it gets on the tape, it writes (QT) following its input w
- 2. *Q* reads  $\langle QT \rangle$  on the tape and uses this to construct  $P_{\langle QT \rangle}$ .
- 3. *Q* then combines  $\langle P_{\langle QT \rangle} \rangle$ ,  $\langle Q \rangle$ , and  $\langle T \rangle$  into one machine  $\langle M^* \rangle = \langle P_{\langle QT \rangle} QT \rangle = \langle R \rangle$
- 4. Q writes  $\langle R \rangle$  onto the tape (again, following w) and passes control to T
- 5. T reads w and  $\langle R \rangle$  on. the tape and computes  $t(\langle R \rangle, w)$

#### $\mathsf{R} = \langle \mathsf{P}_{\langle \mathsf{QT} \rangle} \mathsf{QT} \rangle$

### Initially ${\rm P}_{_{\!\!\!(QT)}}$ starts with control of the tape, with an input string w

- P<sub>(QT)</sub> prints (QT) on the tape (following the original input w), and then passes control to Q
- 2. Q reads (QT) and uses it to construct P<sub>(QT)</sub>
- 3. Q combines  $P_{(QT)}$ ,  $\langle Q \rangle$ , and  $\langle T \rangle$  into a single machine  $\langle M^* \rangle$
- Q writes ⟨R⟩ onto the tape (following w) and passes control to ⊤
- T reads w and ⟨R⟩ on the tape and computes t(w, ⟨R⟩)



 $\mathsf{R} = \langle \mathsf{P}_{\langle \mathsf{QT} \rangle} \mathsf{QT} \rangle$ 

Initially  $\mathsf{P}_{\text{(QT)}}$  starts with control of the tape, with an input string w

- P<sub>(QT)</sub> prints (QT) on the tape (following the original input w), and then passes control to Q
- Q reads ⟨QT⟩ and uses it to construct P<sub>(QT)</sub>
- 3. Q combines  $P_{(QT)}$ ,  $\langle Q \rangle$ , and  $\langle T \rangle$  into a single machine  $\langle M^* \rangle$
- Q writes ⟨R⟩ onto the tape (following w) and passes control to ⊤
- T reads w and ⟨R⟩ on the tape and computes t(w, ⟨R⟩)





 $\mathsf{R} = \langle \mathsf{P}_{\langle \mathsf{QT} \rangle} \mathsf{QT} \rangle$ 

Initially  $\mathsf{P}_{_{\!(\!Q\!T\!)}}$  starts with control of the tape, with an input string w

- P<sub>(QT)</sub> prints (QT) on the tape (following the original input w), and then passes control to Q
- 2. Q reads (QT) and uses it to construct P<sub>(QT)</sub>
- 3. Q combines  $P_{\langle QT \rangle}$ ,  $\langle Q \rangle$ , and  $\langle T \rangle$  into a single machine  $\langle M^* \rangle$
- Q writes ⟨R⟩ onto the tape (following w) and passes control to ⊤
- T reads w and ⟨R⟩ on the tape and computes t(w, ⟨R⟩)





 $\mathsf{R} = \langle \mathsf{P}_{\langle \mathsf{QT} \rangle} \mathsf{QT} \rangle$ 

Initially  $\mathsf{P}_{_{\!(\!Q\!T\!)}}$  starts with control of the tape, with an input string w

- P<sub>(QT)</sub> prints (QT) on the tape (following the original input w), and then passes control to Q
- 2. Q reads (QT) and uses it to construct P<sub>(QT)</sub>
- Q combines P<sub>(QT)</sub>, (Q), and (T) into a single machine (M<sup>\*</sup>)
- Q writes ⟨R⟩ onto the tape (following w) and passes control to ⊤
- T reads w and ⟨R⟩ on the tape and computes t(w, ⟨R⟩)

$$\langle M^* \rangle = \langle P_{\langle Q \rangle} Q T \rangle = \langle R \rangle$$

 $\mathsf{R} = \langle \mathsf{P}_{\langle \mathsf{QT} \rangle} \mathsf{QT} \rangle$ 

Initially  $\mathsf{P}_{_{\!(\!Q\!T\!)}}$  starts with control of the tape, with an input string w

- P<sub>(QT)</sub> prints (QT) on the tape (following the original input w), and then passes control to Q
- 2. Q reads (QT) and uses it to construct P<sub>(QT)</sub>
- 3. Q combines  $P_{(QT)}$ ,  $\langle Q \rangle$ , and  $\langle T \rangle$  into a single machine  $\langle M^* \rangle$
- Q writes ⟨R⟩ onto the tape (following w) and passes control to T
- T reads w and ⟨R⟩ on the tape and computes t(w, ⟨R⟩)



 $\mathsf{R} = \langle \mathsf{P}_{\langle \mathsf{QT} \rangle} \mathsf{QT} \rangle$ 

Initially  $\mathsf{P}_{_{\!(\!Q\!T\!)}}$  starts with control of the tape, with an input string w

- P<sub>(QT)</sub> prints (QT) on the tape (following the original input w), and then passes control to Q
- 2. Q reads (QT) and uses it to construct P<sub>(QT)</sub>
- 3. Q combines  $P_{(QT)}$ ,  $\langle Q \rangle$ , and  $\langle T \rangle$  into a single machine  $\langle M^* \rangle$
- Q writes ⟨R⟩ onto the tape (following w) and passes control to ⊤
- T reads w and ⟨R⟩ on the tape and computes t(w, ⟨R⟩)





- When constructing a Turing machine *M*, you can include the command "Obtain M's description ⟨*M*⟩" as part of the pseudocode
- Whatever you were going to do with (M), the recursion theorem tells us that there exists a machine that finds a way to put its own description on the tape before processing that description

HALT is undecidable: another proof Theorem: HALT = { $\langle M, w \rangle | M$  halts on w} is undecidable

**Proof:** AFSOC machine *H* decides HALT. Create the following machine *M*:

- 1. M takes a string w as input
- 2. Obtain its own description  $\langle M \rangle$
- 3. Run *H* on  $\langle M, w \rangle$
- 4. "Do the opposite"
  - If H accepts  $\langle M, w \rangle$ , go into a loop
  - If *H* rejects  $\langle M, w \rangle$ , immediately halt

M halts on input w if H says it should loop, and loops if H says it should halt.

#### HALT is undecidable: another proof

#### HALT = {<M, w> | M halts on w} A paradoxical machine M



If we can decide HALT, we create a paradoxical machine

Undecidability proofs using the recursion theorem

We can use the following recipe to prove that a language L is undecidable

- 1. AFSOC a machine D can decide L
- 2. Create a machine *M* that obtains its own description, uses *D* to analyze itself, and "does the opposite" of what it should
- 3. Conclude that D is not deciding L correctly

#### Recursion theorem recipe

#### L = an undecidable language

A paradoxical machine M



If we can decide L, we create a paradoxical machine

A<sub>TM</sub> is undecidable: an alternate proof **Theorem:** A<sub>TM</sub> = { $\langle M, w \rangle | w \in L(M)$ } is undecidable Let's try proving this using the recursion theorem

A<sub>TM</sub> is undecidable: an alternate proof Theorem:  $A_{TM} = \{\langle M, w \rangle | w \in L(M)\}$  is undecidable

**Proof:** AFSOC machine *A* decides ATM. Create a machine *M* that does the following:

- 1. M takes a string w as input
- 2. Obtain its own description  $\langle M \rangle$
- 3. Run *A* on  $\langle M, w \rangle$
- 4. "Do the opposite"
  - If A accepts  $\langle M, w \rangle$ , reject
  - If A rejects  $\langle M, w \rangle$ , accept

M accepts w if A says it should reject, and rejects if A says it should accept.

#### $A_{\rm TM}$ is undecidable: an alternate proof

A<sub>TM</sub> = {<M, w> | M accepts w} A paradoxical machine M



If we can decide  ${\rm A}_{_{\rm TM}}\!,$  we create a paradoxical machine

### The language $\mathrm{REG}_\mathrm{TM}$

Consider the following language:

$$\operatorname{REG}_{\operatorname{TM}} = \{ \langle {\it M} 
angle | {\it L}({\it M}) ext{ is regular} \}$$

- We receive a TM description  $\langle M \rangle$  as input
- We seek to determine if *M* recognizes a regular language
  - "Can this TM be converted to a DFA?"

**Theorem**: REG<sub>TM</sub> = { $\langle M \rangle | L(M)$  is regular} is undecidable

Hint 1: Use the recursion theorem

Hint 2: Make a machine that is regular when it shouldn't be, and vice-versa

# $\begin{array}{l} \operatorname{REG}_{\mathrm{TM}} \text{ is undecidable} \\ \text{Theorem: } \operatorname{REG}_{\mathrm{TM}} = \{ \langle M \rangle | L(M) \text{ is regular} \} \text{ is} \\ \text{undecidable} \end{array}$

**Proof:** AFSOC machine R decides  $REG_{TM}$ . Construct a machine M that does the following:

- 1. M takes a string w as input
- 2. Obtain its own description  $\langle M \rangle$
- 3. Run *R* on  $\langle M \rangle$
- 4. "Do the opposite"
  - If R accepts (M), simulate a machine that recognizes 0<sup>n</sup>1<sup>n</sup>
  - If R accepts (M), simulate a machine that recognizes 0\*1\*

#### $\mathrm{REG}_{\mathrm{TM}}$ is undecidable

- If R says M is regular, M recognizes 0<sup>n</sup>1<sup>n</sup> a non-regular language.
- If *M* says *A* is not regular, it recognizes 0\*1\* a regular language.
- Thus, R is not deciding  $REG_{TM}$  correctly

#### $\mathrm{REG}_{\mathrm{TM}}$ is undecidable

#### REG<sub>TM</sub> = {<M> | L(M) is regular}

A paradoxical machine M



If we can decide  $\text{REG}_{\text{TM}}$ , we create a paradoxical machine

#### Minimal Turing machines

- ► Let *M* be a Turing machine, with string description ⟨*M*⟩.
- We say *M* is **minimal** if there is no Turing machine *M*<sub>2</sub> such that:

1.  $L(M) = L(M_2)$ 2.  $|\langle M_2 \rangle| < |\langle M \rangle|$ 

i.e., there is no machine with a shorter description than M that recognizes the same language.

### The language $\mathrm{MIN}_\mathrm{TM}$

#### $MIN_{TM} = \{ \langle M \rangle | M \text{ is a minimal TM} \}$

- We receive a TM description/source code (M)
  We want to know if M is minimal
  - "Can this source code be rewritten more concisely?"

**Proof:** AFSOC some enumerator E enumerates  $MIN_{TM}$ . Create a machine M that does the following:

- 1. M takes a string w as input
- 2. Obtain its own description  $\langle M \rangle$
- 3. Run *E* until it lists a machine  $\langle M_2 \rangle$  such that  $|\langle M \rangle| < |\langle M_2 \rangle|$  (i.e. a machine with a longer description)
- 4. Simulate  $M_2$  on w

#### $\rm MIN_{\rm TM}$ is not recursively enumerable

- ▶ Note that  $L(M) = L(M_2)$ , and  $|\langle M \rangle| < |\langle M_2 \rangle|$
- But this is not supposed to be possible since M<sub>2</sub> is supposed to be a minimal TM.
- Thus, our enumerator did not enumerate MIN<sub>TM</sub> correctly.