

COS 487
Fall 2012
Tarjan

Notes on the Recursion Theorem

We are interested in Turing machines as algorithms to compute functions: if M is a Turing machine, $M(x)$ denotes the output of machine M on input x , which is the value at x of the function computed by M .

Question: Can we construct a machine that produces a description of itself? Such a machine can be turned into a self-reproducing machine. More generally, is there a way to give a Turing machine access to its own description?

Here is a simple construction of a machine that writes a description of itself on an empty tape. Let M be a one-input Turing machine. A diagonal machine $D.M$ for M has the following behavior: On an initially empty tape, $D.M$ runs M on $\langle M \rangle$. The following is an implementation of $D.M$. Write $\langle M \rangle$. Transfer control to M . A diagonal machine D has the following behavior: $D(\langle M \rangle) = \langle D.M \rangle$. That is, given a description $\langle M \rangle$ of a Turing machine M , D writes a description of a diagonal machine $D.M$ for M . The following is an implementation of D . On input $\langle M \rangle$, construct a machine $D.M$ that writes $\langle M \rangle$ on an empty tape and then runs M , and write a description $\langle D.M \rangle$.

Now let SELF be the machine such that $D(\langle D \rangle) = \langle \text{SELF} \rangle$. That is, let SELF be the machine whose description is written by D when run on input $\langle D \rangle$. What does SELF do? On empty tape, it writes $\langle D \rangle$ and then runs D on $\langle D \rangle$, producing $D(\langle D \rangle) = \langle \text{SELF} \rangle$. That is, SELF writes a description of itself!

Recursion Theorem: For any Turing machine T that computes a function of two inputs, the first of which is a Turing machine description, there is a Turing machine $R.T$ such that $R.T(x) = T(\langle R.T \rangle, x)$.

Furthermore, there is a Turing machine R such that $R(\langle T \rangle) = \langle R.T \rangle$. That is, the proof of the Recursion Theorem can be made constructive.

Proof:

We start by generalizing $D.M$ and D to machines M with two inputs. Let M be a two-input Turing machine. A diagonal machine $D.M$ for M has the following behavior: $D.M(x) = M(\langle M \rangle, x)$. The following is an implementation of $D.M$. On input x : write $\langle M \rangle$ next to x , then run M on $\langle M \rangle, x$. This implementation builds both a machine to write $\langle M \rangle$, and a copy of M itself, into its finite control. An alternative is to write two copies of $\langle M \rangle$ and then run U on $\langle M \rangle, \langle M \rangle, x$, where U is a universal two-input Turing machine. This implementation is more generic in that the only part that depends on $\langle M \rangle$ is the part that writes $\langle M \rangle$, but either implementation will do.

A diagonal machine D has the following behavior: $D(\langle M \rangle) = \langle D.M \rangle$. That is, on input $\langle M \rangle$, it produces a description of a diagonal machine for M . The following is an implementation of D . On input $\langle M \rangle$: construct $D.M$, where $D.M$ is implemented as in the preceding paragraph, and write $\langle D.M \rangle$.

A meta machine $E.T$ for T has the following behavior: $E.T(\langle M \rangle, x) = T(D(\langle M \rangle), x)$. The following is an implementation of $E.T$. On input $\langle M \rangle, x$: run D on $\langle M \rangle$ to write $\langle D.M \rangle$, then run T on $\langle D.M \rangle, x$.

Let $R.T$ be the machine such that $D(\langle E.T \rangle) = \langle R.T \rangle$; that is, the machine whose description D produces on input $\langle E.T \rangle$, where $E.T$ is some meta machine for T . Then $R.T(x) = E.T(\langle E.T \rangle, x) = T(D(\langle E.T \rangle), x) = T(\langle R.T \rangle, x)$, as desired.

[End of proof]

To make this proof constructive, we build a meta machine E such that $E(\langle T \rangle) = \langle E.T \rangle$. That is, on input $\langle T \rangle$, E produces a description of $E.T$. An implementation of E is as follows. On input $\langle T \rangle$: construct $E.T$, where $E.T$ is implemented as described above, and write $\langle E.T \rangle$. Given E , we build a machine R such that $R(\langle T \rangle) = D(E(\langle T \rangle)) = \langle R.T \rangle$.

Instead of directly constructing a machine that writes a description of itself, as we did at the beginning of these notes, we can obtain such a machine by applying the recursion theorem to the trivial machine HALT that halts immediately. Indeed, let SELF' be such that $\langle \text{SELF}' \rangle = R(\langle \text{HALT} \rangle)$. That is, SELF' is the machine whose description is written by R on input $\langle \text{HALT} \rangle$. Then SELF' on empty input writes $\langle \text{SELF}' \rangle$ and halts. (Why?)