# Writing a Self-Reproducing Program

Amisi Kiarie

December 18, 2021

## 1  The Problem Defined

The problem given by Ken Thompson[1] is as follows: Write a source program that, when compiled and executed, will produce as output an exact copy of its source.

To understand the difficulty involved in the problem, let us try to specify it somewhat rigorously. Let $\mathscr{T}$ be the set of all strings and $\mathscr{U} \subset \mathscr{T}$ be the set of source programmes in our language, and $\mathscr{C}$ be the set of compiled programs. Then we may denote the *compilation* and *execution maps*

$$\textbf{CMPL} : \mathscr{U} \to \mathscr{C} \text{ and } \textbf{EXEC} : \mathscr{C} \to \mathscr{T}$$

respectively, provided we restrict $\mathscr{U}$ to programs that, when compiled and executed, only produce text output.

A *self-reproducing program* $S \in \mathscr{U}$ is one with the property

$$S = \textbf{EXEC}(\textbf{CMPL}(S)). \tag{1}$$

This may motivate the definition of the *interpretation map* $\textbf{EVAL} : \mathscr{U} \to \mathscr{T}$ as

$$\textbf{EVAL}(P) \coloneqq \textbf{EXEC}(\textbf{CMPL}(P)),\ P \in \mathscr{U}.$$

This allows us to rewrite (1) as
$$S = \textbf{EVAL}(S).$$

In a word, $S$ is an identity element of $\textbf{EVAL}$. By the same token, we have

$$
\begin{aligned}
S &= \textbf{EVAL}(S) \\
&= \textbf{EVAL}(\textbf{EVAL}(S)) \\
&= \ldots \\
&= \textbf{EVAL}^{N}(S)
\end{aligned}
$$

for any positive integer $N$.

---

[1] *Reflections on Trusting Trust*, Turing Award Lecture, 1984.

## 2 Constructing $S$

### 2.1 A naive first attempt

Since we are focused on **EVAL** mainly, we may use the language python to arrive at a solution. Let $P_0$ be the one-line application

```python
print("print()")
```

Then the output **EVAL**$(P_0)$ is

```
print()
```

If we attempt to modify $P_0$ so that the portion interior to the **PRINT** command[2] is included in the output, we obtain $P_1$

```python
print("print(\"print()\")")
```

But this fails to solve our problem because **EVAL**$(P_1)$ is

```
print("print()")
```

However, we see that **EVAL**$(P_1) = P_0$, and it is not hard to see that we may easily construct the sequence $\{P_0, P_1, P_2, \dots\}$ with

$$\textbf{EVAL}(P_n) = P_{n-1}, \text{ for } n \geq 0.$$

Observe, then, that each step in writing a new 'upgrade' from $P_n \to P_{n+1}$ is something akin to applying **EVAL**$^{-1}$ to $P_n$. Now consider the following function:[3]

```python
def upgrade(s):
    escaped = s.replace("\"", "\\\"")
    return f"print(\"{escaped}\")"
```

Note that if we write

```python
upgrade("print()")
```

we obtain a string equivalent to $P_0$, i.e.

---

[2]The **PRINT** command maps strings from $\mathscr{T} \to \mathscr{T}$, and is defined in coordination with the representation map given below, so that $A = \textbf{PRINT}(\textbf{R}(A))$ for any string $A \in \mathscr{T}$.

[3]We are well aware that python offers both single- and double-quotes, but for the sake of illustrating our technique in a manner that is relatively language-agnostic, we have kept only to the latter.

```
print(upgrade("print()"))
```

outputs $P_0$. Moreover,

```
print(upgrade(upgrade("print()")))
```

outputs

```
print("print(\"print()\")")
```

which is simply $P_1$. The *upgrade* function thus maps $P_n \rightarrow P_{n+1}$.

Our aim is to obtain the program $S$ such that **EVAL**$(S) = S$, and we have the following:

- **EVAL**$(P_n) = P_{n-1}$

- *upgrade*$(P_n) = P_{n+1}$.

If we can somehow write a program that connects these two properties, we will obtain our self-reproducing $S$.

**Reflection**   This attempt fails because it grows into an infinite loop of expansions. Perhaps we may learn something from $\{P_n\}$, other than humility! Our exercise with $\{P_n\}$ teaches us it is futile to use a literal string alone to construct a self-reproducing program because this string must contain itself *and* other things: if nothing else, it must contain the code that contains the **PRINT** command itself. A self-reproducing program $S$, therefore, cannot be constructed by printing a literal string only.

## 2.2   Self-recursive construction of $S$

Let us consider another form a self-reproducing program $S$ may take. Suppose we have a block $\rho \in S$ which can be used *in some manner* to generate a literal representation of all of $S$ except itself. We could combine $\rho$ with another block $\varsigma \in S$ which can produce $\rho$ and make use of $\rho$ to produce itself.

More precisely, we can think of our program as having the form

$$S = \rho \bullet \varsigma, \tag{2}$$

where '$\bullet$' is the *concatenation operator*, and $\rho, \varsigma \in \mathscr{T}$:

- $\rho$ is a block which can be used to produce $\varsigma = S - \rho$

- $\varsigma$ is a block which uses $\rho$ to produce both $\rho$ and itself, thereby producing $S$.

Thus $\varsigma$ sees itself through $\rho$ and is a *mirror* or *self-reference* with respect to $\varsigma$. This is what is meant by *self-recursive construction*.[4]

To make this more concrete, define $\mathbf{R} : \mathscr{T} \to \mathscr{T}$ to be the *representation map* so that for any $A \in \mathscr{T}$, we have

$$A = \mathbf{PRINT}(\mathbf{R}(A)).$$

We can now understand $\rho$ as being the representation of $\varsigma$,[5]

$$\rho = \mathbf{R}(\varsigma) \tag{3}$$

and then

$$\varsigma = \mathbf{PRINT}(\rho) \cdot \mathbf{PRINT}(\mathbf{R}^{-1}(\rho)). \tag{4}$$

(3) and (4), together with (2), give us a recursive specification for $S$. We may be assured it terminates because the representation map does not evaluate the expressions it is given, i.e.

$$\rho = \mathbf{R}\big(\mathbf{PRINT}(\rho) \cdot \mathbf{PRINT}(\mathbf{R}^{-1}(\rho))\big)$$

will simply be a string, the $\rho$'s on the right will not be expanded further.

We now in a position to give this program in python. We begin by 'pretending' that we have a value for $\rho$, so we can construct $\varsigma$ on its basis. We must do this because according to (3), $\rho$ is defined in terms of $\varsigma$.

```python
rho = b'' # TODO: fill in value

# sigma:
print(f"rho = {rho}", end="") # print(rho)
print(rho.decode("utf-8"))    # print(R^-1(rho))
```

Ignoring line 1, we may identify in lines 2–4 the components of (4) above:

- On lines 2–4 we have the equivalent of the command $\mathbf{PRINT}(\rho)$

- Line 5 is $\mathbf{PRINT}(\mathbf{R}^{-1}(\rho))$.

Execution of the program as it is simply produces

---

[4]The code which is offered in this document is based upon Thompson's program; but an example is available in the repository which I constructed before looking at his solution. However, it is messier than his construction so I chose to make use of his ideas in this article.

[5]The definitions given for $\rho$ and $\varsigma$ are not meant to be exhaustive: the actual code may contain factors that are 'constant' with respect to these two, such as comments. So in effect we might have given a more full definition pair

$$\rho = C_0 \cdot \mathbf{R}(\varsigma) \cdot C_1$$

and

$$\varsigma = D_0 \cdot \mathbf{PRINT}(\rho) \cdot D_1 \cdot \mathbf{PRINT}(\mathbf{R}^{-1}(\rho)) \cdot D_2,$$

with constant factors $C_i, D_i \in \mathscr{T}$, but we felt that this would clutter the expressions.

```
rho = b''
```

Now we add the value for $\rho$. To obtain it we will take everything after the closing quote and encode it as a UTF-8 byte string. This is equivalent to applying the **R** in (3) above. Our final program is as follows:

```
1   rho = """
2
3   # sigma:
4   print(f"rho = {rho}", end="") # print(rho)
5   print(rho.decode('utf-8'))    # print(R^-1(rho))""".encode('utf-8')
6
7   # sigma:
8   print(f"rho = {rho}", end="") # print(rho)
9   print(rho.decode('utf-8'))    # print(R^-1(rho))
```

However, Thompson notes, purists will observe that this is not exactly a self-reproducing program: rather, it is the producer of one. Execution yields

```
1   rho = b'\n\n# sigma:\nprint(f"rho = {rho}", end="")...' # trimmed this line
2
3   # sigma:
4   print(f"rho = {rho}", end="") # print(rho)
5   print(rho.decode("utf-8"))    # print(R^-1(rho))
```

which is truly self-reproducing.