
Expressions That Talk About Themselves

MAARTEN FOKKINGA

*University of Twente, Department INF, PO Box 217, 7500 AE Enschede, The Netherlands,
Email: fokkinga@cs.utwente.nl*

In this paper we explain, independent of any formalism, how to construct an expression in a formal system that ‘does something with itself’. As an application of this procedure we show how to construct a program that writes its own text, lay-out included. The emphasis is on the systematic development in concrete programming languages, rather than on the abstract theory or the end result (which, indeed, are not new at all).

Received January 26, 1995; revised June 17, 1996

INTRODUCTION

Self-reference occurs frequently in theoretical investigations of formal systems. In predicate logic self-reference is used in Gödel’s Incompleteness Theorem: a *formula* is constructed that asserts, roughly said, ‘I am not formally provable’; in a programming language self-reference manifests itself in the existence of a *program* that writes its own text, a so-called self-reproducing program. Self-reference also plays a role in Nature, in the ability of self-reproduction of the formal system of DNA molecules.

In this paper we explain, independent of any formalism, how to construct an expression in a formal system that does something with itself. As an application of this procedure we show how to construct a program that writes its own text, lay-out included. (The self-reproducing Miranda¹ [1] script is, after some optimization, exactly 36 characters long.)

To appreciate the problems of self-reference, the reader may try to write a self-reproducing program from scratch. It will soon be clear that the program tends to be infinitely long, and therefore is not a program at all. This holds so even when recursion is used in an attempt to avoid the blow-up of the text:

```
program Self;
begin write ("program Self;
            begin write ( ...
                        ... )
            end Self.")
end Self.
```

Moreover, some trick seems to be needed in order to write special symbols (like quotes). Our method needs no tricks and works in any programming language, for any way to encode special characters within strings.

The technique of self-reference and self-reproduction is by no means new, and very many solutions have been reported throughout the years. In a technical setting Odifreddi [2, pp. 165–174] dealt with self-reference within logic, automata and DNA molecules and he

gives many good references to the literature. For a popular account about self-reference and Gödel’s Theorem read Hofstadter’s famous book [3] or the classic ‘Gödel’s Proof’ by Nagel and Newman [4]. Applications of the well-known theoretical constructions are almost always impractical, in the sense that the self-referring construction (device, formula, program) is too large and takes too much time to perform its task. This has been recognized quite early, and attempts for improvement have been made, for instance by Thatcher [5]. Later on Kozen [6] investigates in an abstract (axiomatic) setting self-reference theorems whose proofs yield, theoretically, time-efficient constructions. More recently Jones and his colleagues [7] provided, in a rather concrete setting, alternative constructive proofs that do make applications practical: they seek to exploit the self-reference theory for compiler generation and the like.

1. HOW TO OBTAIN SELF-REFERENCE

1.1. Notation

Identifiers in italic font and math symbols like $(,)$, $[,]$, $/$, \ll , \gg are used to *describe* expressions from the formal system. Expressions and symbols from the formal system itself are written in type writer font. We shall assume that X is a symbol of the formal system (a variable, maybe) that may occur in expressions; we will use it to indicate a place where a replacement has to take place.

1.2. Coding and substitution

In order that an expression can do something sensibly with itself, the formal system must provide a way to code expressions. For example, expressions might be coded as numerals, or as character sequences. The code of an expression E is denoted $\ll E \gg$.

We assume that various operations on codes can be expressed in the formal system. For a programming language this comes as no surprise; for Predicate Logic this is not at all obvious (and we shall not attempt to

¹ Miranda is a trademark of Research Software Ltd.

prove that here: it is the heart of Gödel's theorem). Specifically, we shall want a way to express a kind of substitution. So we assume that there exists a term $subst[,]$ (with two open places where codes can be put) such that for each code c and expression E the expression $subst[c, \ll E \gg]$ (that is, $subst[,]$ with the two places filled by c and $\ll E \gg$) denotes the code of the expression that results from replacing the symbol X by c in E :

$$subst[c, \ll E \gg] \approx \ll [c/X]E \gg .$$

The sign \approx is semantic equivalence, meaning that the left-hand side and the right-hand side may be interchanged in all contexts, if only the 'standard meaning' of expressions and terms is considered. So, in predicate logic relation \approx is formally provable equivalence and in a programming language relation \approx is value-equality. Actually, in all uses of $subst$ in the sequel the two arguments are equal. Therefore we will assume that $subst$ has only one open place; the property that we assume to hold of $subst$ thus reads:

$$subst[\ll E \gg] \approx \ll [\ll E \gg /X]E \gg .$$

Without loss of generality, we may assume that symbol X does not occur in $subst$.

1.3. Self-reference

The expressibility of substitution is a key to obtaining self-reference. It allows an expression to use a code several times whereas it occurs only once in the expression. (In a programming language, there exist alternative means to obtain this effect, namely repetitive constructs like for-loops, list comprehension and replication of strings [8, 9].) In this way the infinite blow-up as suggested for the self-reproducing program is avoided. Indeed, take $self$ and $code$ as follows:

$$self = subst[\ll \langle \text{the entire right-hand side but with } X \text{ instead of this box} \rangle \gg],$$

or, written more formally:

$$self = subst[code]$$

$$code = \ll subst[X] \gg .$$

Then we have that $self$ is semantically equivalent to its own code:

$$self = \{ \text{definition } self \}$$

$$subst[code]$$

$$\approx \{ \text{property } subst, \text{definition } code \}$$

$$\ll [code/X](subst[X]) \gg$$

$$= \{ \text{substitution, } subst \text{ doesn't contain } X \}$$

$$\ll subst[code] \gg$$

$$= \{ \text{definition } self \}$$

$$\ll self \gg .$$

It only remains to adapt this device to other purposes.

1.4. A general purpose self-reference construction

Let $Expr[]$ be an expression with an empty place ... where a code may be put; we write $Expr[c]$ for $Expr$ with the empty place ... filled by c . Here are the two main examples:

$$Expr[] \approx \langle \text{a formula asserting that the formula with code ... is not provable} \rangle$$

$$Expr[] \approx \langle \text{a program that writes the text of the program with code ...} \rangle$$

We set out to construct an expression $Self$ that is semantically equivalent to statement $Expr$ about itself:

$$Self \approx Expr[\ll Self \gg].$$

For the example expressions above we thus get:

$$Self \approx \langle \text{a formula asserting that the formula with code } \ll Self \gg \text{ is not provable} \rangle$$

$$Self \approx \langle \text{a program that writes the text of the program with code } \ll Self \gg \rangle$$

Following the construction above of $self$, we take:

$$Self = Expr[subst[\ll \langle \text{the entire rhs but } X \text{ instead of this box} \rangle \gg]],$$

or, written more formally:

$$Self = Expr[subst[code]]$$

$$code = \ll Expr[subst[X]] \gg .$$

Indeed, assuming that $Expr$ does not contain X , we find in exactly the same way as in the calculation above (Yes, try it!) that $Self \approx Expr[\ll Self \gg]$, as desired.

The attentive reader may have noticed that up to now the coding $\ll \gg$ may be trivial. Both the identity and each constant function satisfy all of the assumptions and the coding need not be invertible. However, if $Expr[c]$ is to do something non-trivial with c , it has to be able to retrieve from c the constituent pieces of the expression encoded by c . It is only at this point that further properties of codes are relevant. (Thanks to Peter Asveld for this observation.)

(In Lisp, programs and data are the same, namely S-expressions, and it is possible to retrieve the constituent text parts from a program E even if E is given as a program. So, in Lisp we can take the identity as coding.)

2. APPLICATION: A SELF-REPRODUCING PROGRAM

Using the above technique for constructing self-referring expressions, we shall construct a self-reproducing program *Self*:

Self \approx

{ a program that writes the text of
the program with code $\ll Self \gg$ }.

Thus the formal system is a programming language.

Since in all programming languages a text can be written if it is first encoded as a string (a character sequence) and then subject to `write`, say, we choose the string encoding as the coding $\ll \gg$. Thus, writing is an inverse of coding, that is, for any text *x*:

`write ($\ll x \gg$)` yields *x*.

For the definiteness, let us assume that a text is coded as a string in the following way:

$\ll x \gg =$
text *x* with each quote, new line, and backslash
replaced by `\q`, `\n`, and `\b`, respectively,
and the whole enclosed by a quote " at both sides.

For example:

x = `aaa"bbb(newline)ccc`
 $\ll x \gg$ = `"aaa\bqbbb\bnccc"`
 $\ll\ll x \gg\gg$ = `"\qaaa\bqbbb\bnccc\q"`.

We call the quote, backslash and new line the special symbols, since these, and only these, undergo a change when subject to coding. In the following it will be obvious how to adapt our construction to alternative ways of encoding special symbols; for example, encoding a quote by doubling it, or encoding a quote by `\"` and a backslash by `\\`. Later on we shall use the concatenation $\#$ of strings; it is specified by the following property:

$\ll x \gg \# \ll y \gg = \ll xy \gg$.

In the right-hand side the texts *x* and *y* are juxtapositioned to form one long text.

2.1. The program

Following the general purpose construction of a self-referring *Self*, we take:

```
Expr[...] = program p;
          . . . . .
          begin write (...)
          end p.

Self      = Expr[subst[code]]
code      =  $\ll Expr[subst[X]] \gg$ .
```

It remains to express in the programming language the expression `subst[...]`. For this we choose to use a

procedure call `S (args, ...)`, where procedure `S` is declared within *Expr* in the declaration part of program *p*, and *args* is a collection of auxiliary arguments that we might need in the future:

```
procedure S (parms; c: string): string;
(* pre: c =  $\ll E \gg$  result:  $\ll [c/X]E \gg$  *)
begin . . . end S;
```

In order to synthesize a definition for procedure `S`, we try to express in programming language terms its result for an argument *code* $= \ll E \gg$ where $E = E'XE''$ and *X* does not occur in E', E'' :

```
result of S (args, code)
=
subst[code]
=
subst[ $\ll E \gg$ ]
=
 $\ll [code/X](E'XE'') \gg$ 
=
 $\ll E' code E'' \gg$ 
=
 $\ll E' \gg \# \ll code \gg \# \ll E'' \gg$ 
(*) = {defining firstpart and lastpart as below}
      firstpart code  $\# \ll code \gg \#$  lastpart code.
```

Notice that procedure `S` has to do an extra coding in the middle part of its result. For step (*) operations *firstpart* and *lastpart* should satisfy:

firstpart $\ll E'XE'' \gg = \ll E' \gg$
lastpart $\ll E'XE'' \gg = \ll E'' \gg$.

If it is possible to define these operations and $\ll \gg$ and $\#$ (as is the case in most languages), their definitions can be placed before procedure `S`, in the declaration part of program *p*. However, we shall show later how to eliminate these operations; so they are not needed in the end. Having the operations available, it is easy to complete the definition of procedure `S`, and the complete program thus reads:

```
program p;
auxiliary procedures
procedure S (parms; c: string): string;
begin return firstpart c  $\# \ll c \gg \#$  lastpart c
end S;
begin write (S (args,
 $\ll$  (the surrounding text with X for this box)  $\gg$ ))
end p.
```

As said before, argument *args* may contain auxiliary data to ease the task of procedure `S`.

(In Miranda, it so happens that the pseudo-function `show` is equal to the coding function $\ll \gg$. Moreover, the program can be written in such a way that the box occurs last, so *firstpart* *c* is `init c` and *lastpart* *c* is empty. The optimization below eliminates the need for *firstpart* altogether.)

2.2. Optimization

It is possible to eliminate the need for operations $\#$, *firstpart*, *lastpart* and procedures that implement the coding $\ll \gg$. We achieve this by an extreme form of *eager evaluation*; for several operations the time of evaluation is shifted from run-time through compile-time to construction-time of the program. (The extreme opposite is *lazy evaluation*, where the evaluation of operations is delayed as long as possible.)

First, instead of programming the disassembling of *code* within procedure *S*, we ourselves disassemble *code* and write the parts as separate arguments to *S*. This eliminates the need for disassembling procedures like *firstpart*. It also eliminates the need for symbol *X*; its sole purpose was to indicate one point where a replacement had to occur, and we know at which point.

Second, the need for the assembling operation $\#$ is eliminated by exploiting the following equivalence:

$$\text{write}(\ll x \gg \# \ll y \gg) \approx \text{write}(\ll x \gg); \text{write}(\ll y \gg)$$

So, instead of passing a concatenation of several codes to a write procedure, procedure *S* itself can perform write actions on the codes separately.

Together these simplifications also enable an easy way to implement the coding $\ll \gg$ that *S* has to perform: we ourselves split *code* at each special symbol and for each splitting point we ourselves know what special code has to be written and where, and so we let procedure *S* do just that.

To work out this idea in an illustrative way, we shall first assume that the above program does not contain special symbols outside the box. The optimization then gives the following text:

```
program p;
auxiliary procedures
procedure S (parms; c0,c1: string);
begin write (c0);
      wq; write (c0); wc; write (c1); wq;
      write (c1)
end S;
begin S (args,
```

‘ ‘ (the text surrounding this box
with ‘ ‘, ‘ ‘ instead of this box) ’ ’

```
end p.
```

Here, *wq* and *wc* have to be defined in such a way that their effect is the writing of a quote " and a quoted comma ", respectively. Notice that the box stands for two arguments, *code*₀ and *code*₁ say; the former *code* is equal to *code*₀ $\#$ $\ll X \gg$ $\#$ *code*₁.

There are several ways to define *wq* and *wc*. The simplest way is to write the quote and quoted comma using the ASCII code. The ASCII codes do not contain special symbols, and if program *p* is written as one line, our simplifying assumption is true. Then we are readily

done, and the *auxiliary procedures* parameters *parms* and arguments *args* are not needed at all. However, we refrain from using this specific coding property. Instead, we strive for a generally applicable solution.

Another way is to *name* the quote and quoted comma in some way or another, and to use the names instead of the symbols themselves to write them. One option is to name them as constants, another option is to name *and write* them in a procedure, and yet another option is to name them via the argument-parameter binding of *args* and *parms*. Quite arbitrarily we choose the latter option. So we take:

```
args = ", ", "\q", ...
parms = c,q: string, ...
```

Then *wc* is: write (q); write (c); write (q) and *wq* is: write (q). However, this *args* does contain special symbols, and our simplifying assumption (no special symbols outside the box) is definitely invalid. So, now we really have to split the text surrounding the box at all special symbols, taking care that *S* writes the special symbols at the splitting points. We also have to include a naming for the backslash and newline. For brevity we say write (x,y,...) instead of write (x); write (y); ... Thus we obtain our final program:

```
program p;
procedure S (c,q,bs,nl: string;
            c0,c1,...,cn+1: string);
begin
  write (c0,sym1,c1,...,symk,ck);
  write (q,c0,q,c,q,c1,q,c,...,c,q,cn+1,q);
  write (ck+1, symk+1,...,symn,cn+1)
end S;
begin S (" ", "\q", "\b", "\n",
```

‘ ‘ (the text surrounding this box with ‘ ‘, ‘ ‘
for this box and for every special symbol) ’ ’

```
end p.
```

Here each *sym*_{*i*} stands for q, bs or nl according to the special symbol at which the *i*th splitting occurs; *k* is the number of special symbols preceding the box, and *n* is the total number of special symbols outside the box. Note that all *sym*_{*i*} are texts without special symbols; were this not the case, we would have a hard time!

The box now stands for a series of codes, *code*₀, *code*₁, ..., *code*_{*n+1*} say. The reader may wish to check that, assumed the program is written as one long line, we have:

```
code0 = "(the text before the first quote)"
k,n    = 11,11
code12 = "(the text after the last quote)",
```

and further:

```
sym1 = q      code1 = ", "
sym2 = q      code2 = ", "
sym3 = q      code3 = ""
sym4 = bs     code4 = "q"
```

```

sym5  = q   code5 = ", "
sym6  = q   code6 = ""
sym7  = bs  code7 = "b"
sym8  = q   code8 = ", "
sym9  = q   code9 = ""
sym10 = bs  code10 = "n"
sym11 = q   code11 = ", "

```

Further optimization is possible: several codes are equal, some are empty, and since there are no new lines, parameter `n1` is nowhere used.

2.3. Exercise

Write a program that writes "Yes, that was me." if it reads its own text on a given file and that writes "Rubbish!" otherwise.

REFERENCES

- [1] Turner, D. A. (1986) An overview of Miranda. In *Research Topics in Functional Programming*, pp. 1–16. Addison-Wesley, Reading, MA. Also: SIGPLAN Notices, Volume 21, number 3, pp. 158–166.
- [2] Odifreddi, P. (1989) *Classical Recursion Theory*, vol. 125 *Studies in Logic*. North-Holland, Amsterdam.
- [3] Hofstadter, D. R. (1979) *Gödel, Escher, Bach: an Eternal Golden Braid*. The Harvester Press, Stanford Terrace, Hassocks, UK.
- [4] Nagel, E. and Newman, J. R. (1958) *Gödel's Proof*. New York University Press, New York.
- [5] Thatcher, J. W. (1963) The construction of a self-describing Turing machine. In Fox, J. (ed.) *Proc. of the Symp. on Mathematical Theory of Automata*, volume XII of *Micro-wave Research Institute Symposia Series*, pp. 164–171, Brooklyn, NY. Polytechnic Press, Brooklyn, NY.
- [6] Kozen, D. (1980) Indexings of subrecursive classes. *Theoretical Computer Science*, **11**, 277–301.
- [7] Jones, N. D. (1992) Computer implementation and applications of Kleene's S-m-n and Recursion Theorems. In Moschovakis Y. N. (ed.), *Logic From Computer Science*, volume 21 of *Mathematical Sciences Research Institute publications*, pp. 243–263. Springer-Verlag, New York.
- [8] Fokkinga, M. M. (1973) A self-reproducing Algol 60 program. *Algol Bulletin*, **35**, 24–26.
- [9] Fokkinga, M. M. (1979) *Some self-reproducing Algol-like programs and Kleene's recursion theorem formulated in concrete programming languages*. Technical Report TW-memo 281, University of Twente, Enschede, Netherlands.