

THEORETICAL PEARLS

Self-interpretation in lambda calculus

HENK BARENDREGT

Faculty of Mathematics and Computer Science, Catholic University of Nijmegen, The Netherlands

This editorial emphasizes results from the theory behind functional programming (including lambda calculus, type theory and term-rewriting systems) that are particularly beautiful, and which have short and elegant proofs. Readers are encouraged to send comments or contributions to:

Henk Barendregt
Faculty of Mathematics and Computer Science,
Catholic University, Nijmegen,
Toernooiveld 1, 6525 ED Nijmegen,
The Netherlands
E-mail: henk@cs.kun.nl

Programming languages which are capable of interpreting themselves have been fascinating computer scientists. Indeed, if this is possible then a ‘strange loop’ (in the sense of Hofstadter, 1979) is involved. Nevertheless, the phenomenon is a direct consequence of the existence of universal languages. Indeed, if all computable functions can be captured by a language, then so can the particular job of interpreting the code of a program of that language. Self-interpretation will be shown here to be possible in lambda calculus.

The set of λ -terms, notation Λ , is defined by the following abstract syntax

$$\Lambda = V | \Lambda \Lambda | \lambda V. \Lambda$$

where $V = v | V'$

is the set $\{v, v', v'', v''', \dots\}$ of *variables*. Arbitrary variables are usually denoted by x, y, z, \dots and λ -terms by M, N, L, \dots . A *redex* is a λ -term of the form

$$(\lambda x. M) N$$

and has as *contractum* $M[x := N]$,

that is, the result of substituting N for (the free occurrences of) x in M . Stylistically, it can be said that λ -terms represent functional programs including their input. A *reduction machine* executes such terms by trying to reduce them to normal form; that is, redexes are continuously replaced by their contracta until hopefully no more redexes are present. If such a normal form can be reached, then this is the output of the functional program; otherwise, the program diverges.

From the point of view of a reduction machine, a λ -term M can be considered as an executable. It ‘itches’ at many places: all redexes want to be reduced.

1 Definition

(i) Each λ -term M has a unique natural number $\#M$ as code. One way of coding is

$$\begin{aligned}\#v^{(i)} &= \langle 0, i \rangle, \\ \#(MN) &= \langle 1, \langle \#M, \#N \rangle \rangle, \\ \#(\lambda x. M) &= \langle 2, \langle \#x, \#M \rangle \rangle,\end{aligned}$$

where $\langle -, - \rangle$ is some effective coding of pairs of numbers as a single number, for example,

$$\langle n, m \rangle = \frac{1}{2}(n+m)(n+m+1) + m.$$

(ii) Let $\ulcorner 0 \urcorner, \ulcorner 1 \urcorner, \ulcorner 2 \urcorner, \dots$ be some set of numerals (λ -terms representing the natural numbers). We take the *Church numerals* $\ulcorner n \urcorner \equiv \lambda f x. f^n(x)$.

Write $\ulcorner M \urcorner \equiv \ulcorner \#M \urcorner$, the *internal λ -code* of M . Now $\ulcorner M \urcorner$ does not itch: being a Church numeral it is in normal form.

Write $FV(M)$ for the set of free variables of M . A λ -term M is *closed* if $FV(M) = \emptyset$. the set of closed λ -terms is denoted by Λ^0 .

2 Definition

(i) An *interpreter* (or *evaluator*) is an (external) function $E: \Lambda \rightarrow \Lambda$ such that

$$E(\ulcorner M \urcorner) \equiv M.$$

(ii) A *self-interpreter* is a λ -term E such that for $M \in \Lambda^0$ one has

$$E \ulcorner M \urcorner =_{\beta} M. \quad (1)$$

Here $=_{\beta}$ (or simply $=$) denotes convertibility between elements of Λ .

3 Remarks

(i) Equation (1) cannot hold for open terms containing free variables. Indeed, E has at most a finite number of free variables and $\ulcorner M \urcorner$ being a numeral has none, but on the right-hand side M may have arbitrarily many free variables.

(ii) Define the *quote* to be the function $Q: \Lambda \rightarrow \Lambda$ such that

$$Q(M) \equiv \ulcorner M \urcorner.$$

A *self-quote* is a λ -term Q such that (say for closed terms M)

$$QM =_{\beta} \ulcorner M \urcorner.$$

Such a self-quote does not exist, however. Indeed, the existence of Q implies

$$\ulcorner \ulcorner I \urcorner \urcorner =_{\beta} Q(\ulcorner I \urcorner) =_{\beta} QI =_{\beta} \ulcorner I \urcorner.$$

Since numerals are in β -normal form it follows that $\ulcorner \ulcorner I \urcorner \urcorner \equiv \ulcorner I \urcorner$, so $\#(\ulcorner I \urcorner) = \#I$. However, $\ulcorner I \urcorner$ and I are different terms, and so have different codes, which is a contradiction.

Kleene already in (1936) showed that there is a self-interpreter E for the lambda calculus. One would think that E is defined by recursion on the structure of its argument. There is, however, a difficulty: closed terms are not built up inductively,

but formed as a subset of the wider class of open terms. Kleene avoided this problem by building up closed terms from combinators S, K and I (actually, he worked with the λ I-calculus and used combinators like S, B, C and I). The construction was as follows

$$\begin{array}{c} \text{CL} \quad \quad \text{E}_{\text{CL}} \\ \lceil M \rceil \rightarrow \lceil M_{\text{CL}} \rceil \rightarrow M_{\text{CL}} =_{\beta} M \end{array}$$

where CL is a compiler from λ -terms to combinatory terms, and E_{CL} is an interpreter for combinatory terms. The translation CL gives, for example,

$$(\lambda z. zz)_{\text{CL}} \equiv \text{SII} (=_{\beta} \lambda z. \text{Iz}(\text{Iz}) =_{\beta} \lambda z. zz).$$

P. de Bruin (my former student) gave an essentially simpler construction of a self-interpreter for the λ -calculus. He used an idea from denotational semantics. In the following construction, F plays the role of an environment in the sense that it determines the values of the free variables.

4 Theorem (Kleene, 1936)

There exists a self-interpreter E for the lambda calculus.

Proof (de Bruin)

By the representability of computable functions there is a term E_0 such that

$$\begin{aligned} \text{E}_0 \lceil x \rceil F &=_{\beta} F \lceil x \rceil, \\ \text{E}_0 \lceil MN \rceil F &=_{\beta} F(\text{E}_0 \lceil M \rceil F)(\text{E}_0 \lceil N \rceil F), \\ \text{E}_0 \lceil \lambda x. M \rceil F &=_{\beta} \lambda x. (\text{E}_0 \lceil M \rceil F_{\lceil \lceil x \rceil \mapsto x \rceil}), \end{aligned}$$

where $F_{\lceil \lceil x \rceil \mapsto x \rceil} = F'x$, with

$$\begin{aligned} F'x \lceil x \rceil &=_{\beta} x, \\ F'x \lceil y \rceil &=_{\beta} F \lceil y \rceil, \quad \text{if } y \neq x. \end{aligned}$$

By induction on the structure of $M \in \Lambda$ it can be shown that

$$\text{E}_0 \lceil M \rceil F =_{\beta} M[x_1 := F \lceil x_1 \rceil, \dots, x_n := F \lceil x_n \rceil] \quad (2)$$

(simultaneous substitution), where $\{x_1, \dots, x_n\} = \text{FV}(M)$. Now we can take

$$\text{E} \equiv \lambda m. \text{E}_0 m \text{I}.$$

Indeed, for closed M it follows by equation (2) that

$$\text{E} \lceil M \rceil =_{\beta} \text{E}_0 \lceil M \rceil \text{I} =_{\beta} M. \quad \square$$

Using the self-interpreter E it can be shown that certain λ -terms exist without giving details. We first introduce some λ -terms inspired by the language LISP.

5 Definition

$$\begin{aligned} \text{cons} &\equiv \lambda xyz. zxy; \\ \text{nil} &\equiv \lambda xyz. y; \\ \text{null} &\equiv \lambda x. x(\lambda \text{abcd}. d). \end{aligned}$$

6 Proposition

- (i) $\text{null nil} = \lambda xy. x \equiv \text{true};$
- (ii) $\text{null (cons a b)} = \lambda xy. y \equiv \text{false}.$
- (iii) Moreover, there exist terms car and cdr such that

$$\begin{aligned}\text{car (cons a b)} &= a; \\ \text{cdr (cons a b)} &= b.\end{aligned}$$

Proof

(i), (ii). Easy. (iii) Take $\text{car} \equiv \lambda x. x(\lambda ab. a)$ and $\text{cdr} \equiv \lambda x. x(\lambda ab. b)$. \square

7 Notation

Write

$$\begin{aligned}a : b &\equiv \text{cons } a \ b; \\ \langle \rangle &\equiv \text{nil}; \\ \langle x_1, \dots, x_{n+1} \rangle &\equiv x_1 : \langle x_2, \dots, x_{n+1} \rangle.\end{aligned}$$

For example, $\langle a, b \rangle \equiv a : b : \text{nil} \equiv (\text{cons } a (\text{cons } b \text{ nil}))$.

The following problem was raised by Dr Wim Vree of the University of Amsterdam.

8 Problem

Does there exist a λ -term F such that for all $n \in \mathbb{N}$ one has

$$F \ulcorner n \urcorner = \lambda x_1 \dots x_n. \langle x_1, \dots, x_n \rangle? \quad (3)$$

Solution

Write $M_n \equiv \lambda x_1 \dots x_n. \langle x_1, \dots, x_n \rangle$. Clearly, $\#M_n$ is computable from n , say $\#M_n = g(n)$ with g recursive. Let g be λ -defined by G , say. Then

$$G \ulcorner n \urcorner = \ulcorner g(n) \urcorner = \ulcorner M_n \urcorner.$$

Then $F = \lambda n. E(Gn)$ satisfies equation (3)

$$F \ulcorner n \urcorner = E(G \ulcorner n \urcorner) = E \ulcorner M_n \urcorner = M_n. \quad \square$$

At first, Vree thought the answer to Problem 8 was negative. After seeing the positive answer, he came up with a more constructive solution.

Constructive solution

One can find a λ -term rev such that for all n

$$\text{rev } \langle x_1, \dots, x_n \rangle = \langle x_n, \dots, x_1 \rangle.$$

(For example,

$$\text{rev} = \lambda L_1. \text{rev}' L_1 \langle \rangle,$$

with

$$\begin{aligned}\text{rev}' (a : b) L_2 &= \text{rev}' b (a : L_2) \\ \text{rev}' \text{nil } L_2 &= L_2.\end{aligned}$$

So take $\text{rev}' \equiv Y(\lambda r L_1 L_2. \text{if}[\text{null } L_1] \text{ then } [L_2] \text{ else } [r (\text{cdr } L_1) ((\text{car } L_1) : L_2)])$, where Y is the fixed point combinator and if X then Y else Z is simply XYZ .)

Construct a λ -term V such that

$$\begin{aligned} V^{\ulcorner n+1 \urcorner} &= \lambda Lx. (V^{\ulcorner n \urcorner}(x; L)), \\ V^{\ulcorner 0 \urcorner} &= \text{rev}. \end{aligned}$$

(For example, $V = Y(\lambda vn. \text{if}[\text{zero? } n] \text{ then rev else } [\lambda Lx. v(\text{pred } n)(x; L)])$, where pred represents the predecessor function.)

Then $F = \lambda n. V n \text{ nil}$ satisfies equation (3). Indeed

$$\begin{aligned} F^{\ulcorner n \urcorner} x_1 \dots x_n &= V^{\ulcorner n \urcorner} \text{nil } x_1 \dots x_n \\ &= (\lambda Lx. V^{\ulcorner n-1 \urcorner}(x; L)) \text{nil } x_1 \dots x_n \\ &= V^{\ulcorner n-1 \urcorner}(x_1; \text{nil}) x_2 \dots x_n \\ &= (\lambda Lx. V^{\ulcorner n-2 \urcorner}(x; L)) (x_1; \text{nil}) x_2 \dots x_n \\ &= V^{\ulcorner n-2 \urcorner}(x_2; x_1; \text{nil}) x_3 \dots x_n \\ &\dots \\ &= V^{\ulcorner 0 \urcorner}(x_n; \dots; x_1; \text{nil}) \\ &= \text{rev} \langle x_n, \dots, x_1 \rangle \\ &= \langle x_1, \dots, x_n \rangle. \end{aligned}$$

□

A concrete λ -term satisfying equation (3) is the following

$$\begin{aligned} F &\equiv \lambda n. (\lambda ab. b(aab)) (\lambda ab. b(aab)) \\ &\quad (\lambda vn. [n(\lambda xyz. y) (\lambda xy. x)] \\ &\quad \quad [\lambda L. (\lambda ab. b(aab)) (\lambda ab. b(aab)) \\ &\quad \quad \quad (\lambda rL_1 L_2. [L_1(\lambda abcd. d)] [L_2] [r(L_1(\lambda ab. b)) (\lambda z. z(L_1(\lambda ab. a)) L_2)]) \\ &\quad \quad \quad L(\lambda xyz. y)]) \\ &\quad \quad [\lambda Lx. v(\lambda yz. n(\lambda pq. q(py)) (\lambda w. z) (\lambda t. t)) (\lambda z. zxL)]) \\ &\quad n(\lambda xyz. y)). \end{aligned}$$

9 Exercises

(i) Show that there is no λ -term G such that for all $n \in \mathbb{N}$ one has

$$Gx_1 \dots x_n = \langle x_1, \dots, x_n \rangle. \quad (4)$$

(ii) Construct a λ -term H such that for all $n \in \mathbb{N}$ one has

$$H^{\ulcorner n \urcorner} x_1 \dots x_n = \lambda z. zx_1 \dots x_n. \quad (5)$$

References

- Hofstadter, D. R. 1979. *Gödel, Escher, Bach: an Eternal Golden Braid*, Basic Books.
 Kleene, S. C. 1936. λ -definability and recursiveness. *Duke Math. J.* 2, 340–353.