

ON SELF-INTERPRETERS FOR SYSTEM T AND OTHER TYPED λ -CALCULI

ANDREJ BAUER

ABSTRACT. Matt Brown and Jens Palsberg constructed a self-interpreter for System F_ω . They require a self-interpreter to have an injective quoting function and that the source codes be β -normal. I show that under the same constraints there is a self-interpreter for Gödel's System T , and in fact for any strongly normalizing typed λ -calculus expressive enough to encode natural numbers and pairs. The proof is a swindle which makes one think that the constraints by Brown and Palsberg are not good enough, but I show that they are essentially the best we could hope for in terms of type complexity. What we need is a better definition of a self-interpreter that takes into account structural properties of self-interpreters.

1. INTRODUCTION

An interpreter is a program which reads source code and executes it. A self-interpreter is an interpreter implemented in the language that it interprets. The first one was Turing's universal machine, but it was Steve Russell's implementation of `eval` in Lisp that made self-interpreters a popular programming exercise.

In defiance of the received wisdom that a total programming language cannot have a self-interpreter, Matt Brown and Jens Palsberg [3] implemented a self-interpreter for System F_ω , which is a strongly normalizing typed λ -calculus and thus certainly a total language. In order to avoid the trivial self-interpreter they imposed certain constraints on what a self-interpreter is. I show that under the same conditions already Gödel's System T has a self-interpreter (Theorem 3.2). The construction is trivial, which makes one wonder whether something is at fault with the notion of self-interpreter used by Brown and Palsberg. However, I show that there cannot be a significant improvement (Corollary 3.7) because the type of a code must be at least as complex as the type of the program it encodes.

2. UNTYPED AND TYPED SELF-INTERPRETERS

2.1. Untyped self-interpreters. We work first in the untyped λ -calculus [2]. A self-interpreter is a program (closed term) u which takes as an argument the (*source*) *code* $\ulcorner e \urcorner$ of a program e and outputs a term which is equivalent to e ,

$$u \ulcorner e \urcorner \equiv_\beta e.$$

In order to avoid the trivial case $u = \lambda x. x$ and $\ulcorner e \urcorner = e$ we should specify what counts as source code. A programmer would certainly expect codes to be strings of characters or abstract syntax trees, and a logician would point out that such data can be coded by

This material is based upon work supported by the Air Force Office of Scientific Research, Air Force Materiel Command, USAF under Award No. FA9550-14-1-0096. Any opinions, findings, and conclusions or recommendations expressed in this publication are those of the author(s) and do not necessarily reflect the views of the Air Force Office of Scientific Research, Air Force Materiel Command, USAF..

numbers. Assuming that all variables are of the form x_0, x_1, x_2, \dots , we define the *Gödel encoding* g of terms as numbers by

$$\begin{aligned} g(x_i) &= 2^0 \cdot 3^i, \\ g(e_1 e_2) &= 2^1 \cdot 3^{g(e_1)} \cdot 5^{g(e_2)}, \\ g(\lambda x_i. e) &= 2^2 \cdot 3^i \cdot 5^{g(e)} \end{aligned}$$

and let the quoting function be $\ulcorner e \urcorner = g(e)$, where \underline{n} is the *Church numeral* representing number n . Implementing the corresponding interpreter is an edifying programming exercise which we leave to the enthusiasts. We note that u may be implemented using only *primitive* recursion because $\ulcorner e \urcorner$ calculates its result by structural recursion on e , and primitive recursion suffices to discern the syntax of e from its Gödel code $g(e)$.

2.2. Typed self-interpreters. Let us now move to the typed setting. As in the untyped case, we call closed terms *programs* and write $\text{Prg}(\tau)$ for the set of programs of type τ .

We consider extensions of the *simply typed λ -calculus* [2, §A.1], which we refer to simply as *calculi*. One is *Gödel's System T*, see [2, §A.2] and [1], which is the simply typed λ -calculus extended with binary products $\sigma \times \tau$, a ground type of natural numbers nat and *primitive* recursion at each type. It is expressive enough to manipulate syntax through Gödel encoding. System *T* is strongly normalizing [1, §4.3], and so every well-typed term e has a unique β -normal form $n(e)$.

Another extension is *PCF* [8], which is the simply typed λ -calculus extended with the natural numbers nat and *general* recursion: for every type τ it has a *fixed-point operator* $\text{fix}_\tau : (\tau \rightarrow \tau) \rightarrow \tau$ governed by the β -rule $\text{fix}_\tau f \rightarrow^\beta f (\text{fix}_\tau f)$. The fixed-point operators make the calculus non-normalizing because $\text{fix}_\tau(\lambda x : \tau. x)$ is a program of type τ with an infinite sequence of β -reductions.

Definition 2.1. A *typed self-interpreter* is given by a type ν of (*source*) *codes*, and for each type τ a *quoting function* $\ulcorner \cdot \urcorner_\tau : \text{Prg}(\tau) \rightarrow \text{Prg}(\nu)$ and an *interpreter* $u_\tau \in \text{Prg}(\nu \rightarrow \tau)$ such that $u_\tau \ulcorner e \urcorner_\tau \equiv_\beta e$ for all $e \in \text{Prg}(\tau)$.

Note that the quoting functions need not be λ -definable, i.e., there may be no program q_τ such that $\ulcorner e \urcorner_\tau \equiv_\beta q_\tau e$. The following theorem is the justification for the popular opinion that total languages do not have self-interpreters.

Theorem 2.2. *If a calculus has a self-interpreter then it has fixed-point operators at all types.*

Proof. The proof proceeds by diagonalization much like in Lawvere's fixed point theorem [4]. Let us first show that for every type τ every $f \in \text{Prg}(\tau \rightarrow \tau)$ has a fixed point. Define $g : \nu \rightarrow \tau$ by $g = \lambda x : \nu. f (u_{\nu \rightarrow \tau} x x)$, and let $n = \ulcorner g \urcorner_{\nu \rightarrow \tau}$ be its code. Then $u_{\nu \rightarrow \tau} n n$ is a fixed point of f because $u_{\nu \rightarrow \tau} n n \equiv_\beta g n \equiv_\beta f (u_{\nu \rightarrow \tau} n n)$. Now we obtain the fixed-point operator $\text{fix}_\tau : (\tau \rightarrow \tau) \rightarrow \tau$ as the fixed point of the program

$$\lambda y : (\tau \rightarrow \tau) \rightarrow \tau. \lambda f : \tau \rightarrow \tau. f (y f). \quad \square$$

Corollary 2.3. *System T does not have a self-interpreter.*

Proof. In System *T* successor $\text{succ} : \text{nat} \rightarrow \text{nat}$ has no fixed points. \square

The corollary holds for other kinds of calculi, as long as they possess endomaps without fixed points, which is typical of strongly normalizing calculi such as System F_ω .

We can ask whether Theorem 2.2 can be inverted: does a simply typed λ -calculus with natural numbers and fixed-point operators have a self-interpreter? As Alex Simpson

pointed out to me, the question was answered affirmatively by John Longley and Gordon Plotkin who described a rather intricate self-interpreter for PCF [6, Prop. 6].

2.3. Self-reducers. Our interpreters convert source code to actual code. An alternative is to have them operate exclusively with source code. In the following we shall use the type nat as the type of source codes, and the standard Gödel encoding for quoting, i.e., in this section we define the quoting function $\ulcorner \cdot \urcorner_\tau : \text{Prg}(\tau) \rightarrow \text{Prg}(\text{nat})$ by

$$\ulcorner e \urcorner_\tau = \underline{g}(e).$$

Here \underline{n} stands for the closed term

$$\underbrace{\text{succ} (\text{succ} \cdots (\text{succ } 0) \cdots)}_n$$

of type nat .

Definition 2.4. A *self-reducer* is given by a family of *reducers* $\mathbf{r}_\tau \in \text{Prg}(\text{nat} \rightarrow \text{nat})$, one for each type τ , such that $\mathbf{r}_\tau \ulcorner e \urcorner_\tau \equiv_\beta \ulcorner \mathbf{n}(e) \urcorner_\tau$ for all $e \in \text{Prg}(\tau)$.

Whereas a self interpreter turns the code $\ulcorner e \urcorner_\tau$ to the corresponding program e , a self-reducer turns the code $\ulcorner e \urcorner_\tau$ to the code $\ulcorner \mathbf{n}(e) \urcorner_\tau$ of the normal form of e .

Theorem 2.5. *If there is a self-reducer then every $f \in \text{Prg}(\text{nat} \rightarrow \text{nat})$ has a fixed point.*

Proof. Let us first observe that it only takes primitive recursion (which is available in System T) to define

- (1) a function $\text{enc} \in \text{Prg}(\text{nat} \rightarrow \text{nat})$ such that $\text{enc } \underline{n} \equiv_\beta \ulcorner \underline{n} \urcorner_{\text{nat}}$ for every $n \in \mathbb{N}$,
- (2) a function $\text{dec} \in \text{Prg}(\text{nat} \rightarrow \text{nat})$ such that $\text{dec } \ulcorner \underline{n} \urcorner_{\text{nat}} \equiv_\beta \underline{n}$ for every $n \in \mathbb{N}$,
- (3) a function $\text{app} \in \text{Prg}(\text{nat} \rightarrow \text{nat} \rightarrow \text{nat})$ such that $\text{app } \ulcorner e_1 \urcorner_{\text{nat} \rightarrow \text{nat}} \ulcorner e_2 \urcorner_{\text{nat}} \equiv_\beta \ulcorner e_1 e_2 \urcorner_{\text{nat}}$, for all $e_1 \in \text{Prg}(\text{nat} \rightarrow \text{nat})$ and $e_2 \in \text{Prg}(\text{nat})$.

Now consider any $f \in \text{Prg}(\text{nat} \rightarrow \text{nat})$ and define $g \in \text{Prg}(\text{nat} \rightarrow \text{nat})$ by

$$g = \lambda x : \text{nat} . f (\text{dec } (\mathbf{r}_{\text{nat}} (\text{app } x (\text{enc } x))))$$

Let $k = \ulcorner g \urcorner_{\text{nat} \rightarrow \text{nat}}$ and observe that $g k$ normalizes to a numeral, i.e., there is $\ell \in \mathbb{N}$ such that $\mathbf{n}(g k) = \underline{\ell}$. Now compute:

$$\begin{aligned} \underline{\ell} &\equiv_\beta g k \\ &\equiv_\beta f (\text{dec } (\mathbf{r}_{\text{nat}} (\text{app } k (\text{enc } k)))) \\ &\equiv_\beta f (\text{dec } (\mathbf{r}_{\text{nat}} (\text{app } (\ulcorner g \urcorner_{\text{nat} \rightarrow \text{nat}}) \ulcorner k \urcorner_{\text{nat}}))) \\ &\equiv_\beta f (\text{dec } (\mathbf{r}_{\text{nat}} \ulcorner g k \urcorner_{\text{nat}})) \\ &\equiv_\beta f (\text{dec } (\ulcorner \underline{\ell} \urcorner_{\text{nat}})) \\ &\equiv_\beta f \underline{\ell} \end{aligned} \quad \square$$

Corollary 2.6. *System T does not have a self-reducer.*

Proof. Once again, $\text{succ} : \text{nat} \rightarrow \text{nat}$ has no fixed points. \square

It would be interesting to generalize the non-existence of self-reducers to general quoting functions.

3. BROWN-PALSBERG SELF-INTERPRETERS

To obtain a self-interpreter for System T we need a notion of source code with varying type. The following definition is fashioned after Brown and Palsberg [3].

Definition 3.1. A *weak Brown-Palsberg self-interpreter* is given by, for each type τ , a type of codes $\Box\tau$, a *quoting function* $\ulcorner \cdot \urcorner_\tau : \text{Prg}(\tau) \rightarrow \text{Prg}(\Box\tau)$, and an *interpreter* $\mathsf{u}_\tau : \text{Prg}(\Box\tau \rightarrow \tau)$ such that $\mathsf{u}_\tau \ulcorner e \urcorner_\tau \equiv_\beta e$ for all $e \in \text{Prg}(\tau)$. Such an interpreter is *strong* when for every type τ , the quoting function $\ulcorner \cdot \urcorner_\tau$ is

- (1) *normal*: $\ulcorner e \urcorner_\tau$ is β -normal for all $e \in \text{Prg}(\tau)$, and
- (2) *acceptable*: there is $\mathsf{g}_\tau : \Box\tau \rightarrow \text{nat}$ such that $\mathsf{g}_\tau \ulcorner e \urcorner_\tau = \underline{\mathsf{g}}(e)$ for all $e \in \text{Prg}(\tau)$.

The normality condition expresses the idea that codes should be values, as opposed to programs that still need to be evaluated, while acceptability says that the syntax of an expression is discernible from its code. Brown and Palsberg also require that the quoting function be injective, which follows from our definition because $\underline{\mathsf{g}}$ is injective. They do not explicitly postulate acceptability, although they provide programs that extract the syntax of an expression from its code.

A Brown-Palsberg self-interpreter cannot have a trivial quoting function $\ulcorner a \urcorner_\tau = a$ because codes must be β -normal, while injectivity of the quoting function prevents coding by β -normal forms $\ulcorner a \urcorner_\tau = \mathsf{n}(a)$.

Theorem 3.2. *System T has a Brown-Palsberg self-interpreter.*

Proof. Define

$$\begin{aligned} \Box\tau &= \text{nat} \times \tau, & \ulcorner e \urcorner_\tau &= \langle \underline{\mathsf{g}}(e), \mathsf{n}(e) \rangle, \\ \mathsf{u}_\tau &= \lambda x : \text{nat} \times \tau . \text{snd } x, & \mathsf{g}_\tau &= \lambda x : \text{nat} \times \tau . \text{fst } x. \end{aligned}$$

Clearly we have $\mathsf{u}_\tau \ulcorner e \urcorner_\tau \equiv_\beta \mathsf{n}(e) \equiv_\beta e$, $\langle \underline{\mathsf{g}}(e), \mathsf{n}(e) \rangle$ is β -normal because numerals are β -normal and so are pairs of β -normal terms, and the quoting function is acceptable by fiat. \square

It is clear that the same proof applies to any calculus that has binary products, natural numbers, and any notion of normal form. System F_ω is an example, so we could replace the elegant self-interpreter by Brown and Palsberg with our technologically primitive solution.

The proof of Theorem 3.2 abuses the fact that Brown-Palsberg interpreters allow codes to be as complex as the programs they encode. Theorem 2.2 prevents us from using a fixed type of codes, but perhaps $\Box\tau$ can at least be less complex than τ ? We shall show that this is not possible, but first we need a precise measure of complexity.

Definition 3.3. The *level* of a type τ is defined inductively as follows:

$$\begin{aligned} \text{lev}(\text{nat}) &= 0 \\ \text{lev}(\sigma \times \tau) &= \max(\text{lev}(\sigma), \text{lev}(\tau)) \\ \text{lev}(\sigma \rightarrow \tau) &= \max(1 + \text{lev}(\sigma), \text{lev}(\tau)). \end{aligned}$$

The level of τ is the deepest nesting of \rightarrow to the left. The *pure types* ν_0, ν_1, \dots defined by

$$\nu_0 = \text{nat} \quad \text{and} \quad \nu_{k+1} = \nu_k \rightarrow \text{nat}$$

have arbitrarily high levels because $\text{lev}(\nu_k) = k$.

Definition 3.4. A type σ is a *retract* of τ , written $\sigma \triangleleft \tau$, if there are programs $s : \sigma \rightarrow \tau$ and $r : \tau \rightarrow \sigma$, respectively called a *section* and a *retraction*, such that

$$\lambda x : \sigma . r (s x) \equiv_{\beta\eta} \lambda x : \sigma . x.$$

If additionally $\lambda y : \tau . s (r y) \equiv_{\beta\eta} \lambda y : \tau . y$ then we say that σ and τ are *isomorphisms* and that σ and τ are *isomorphic*, written $\sigma \cong \tau$.

Some basic observations about retracts are:

- (1) every type is a retract of itself, $\tau \triangleleft \tau$,
- (2) being a retract is transitive: if $\rho \triangleleft \sigma$ and $\sigma \triangleleft \tau$ then $\rho \triangleleft \tau$,
- (3) if $\sigma \triangleleft \sigma'$ and $\tau \triangleleft \tau'$ then $\sigma \times \tau \triangleleft \sigma' \times \tau'$ and $\sigma' \rightarrow \tau \triangleleft \sigma' \rightarrow \tau'$,
- (4) $\sigma \triangleleft \sigma \times \tau$ and $\tau \triangleleft \sigma \times \tau$.

The last observation relies on the fact that every type in System T is inhabited. The following lemma has been known for a long time, e.g. [9, §1.8].

Lemma 3.5. *If $\text{lev}(\sigma) \leq \text{lev}(\tau)$ then σ is a retract of τ .*

Proof. We first verify that $\nu_k \triangleleft \nu_{k+1}$ for all $k \in \mathbb{N}$. Define programs $s_k : \nu_k \rightarrow \nu_{k+1}$ and $r_k : \nu_{k+1} \rightarrow \nu_k$ by

$$\begin{aligned} s_0 &= \lambda n : \text{nat} . \lambda m : \text{nat} . n, \\ r_0 &= \lambda f : \text{nat} \rightarrow \text{nat} . f \text{ zero}, \end{aligned}$$

and

$$\begin{aligned} s_{k+1} &= \lambda f : \nu_{k+1} . \lambda g : \nu_{k+1} . f (r_k g), \\ r_{k+1} &= \lambda f : \nu_{k+2} . \lambda g : \nu_k . f (s_k g). \end{aligned}$$

We check by induction on k that s_k and r_k form a section-retraction pair. The base case is confirmed by

$$\lambda x : \text{nat} . r_0 (s_0 x) \equiv_{\beta} \lambda x : \text{nat} . (\lambda m : \text{nat} . x) \text{ zero} \equiv_{\beta} \lambda x : \text{nat} . x.$$

and the induction step by

$$\begin{aligned} \lambda f : \nu_{k+1} . r_{k+1} (s_{k+1} f) &\equiv_{\eta} \lambda f : \nu_{k+1} . \lambda g : \nu_k . r_{k+1} (s_{k+1} f) g \\ &\equiv_{\beta} \lambda f : \nu_{k+1} . \lambda g : \nu_k . (s_{k+1} f) (s_k g) \\ &\equiv_{\beta} \lambda f : \nu_{k+1} . \lambda g : \nu_k . f (r_k (s_k g)) \\ &\equiv_{\beta\eta} \lambda f : \nu_{k+1} . \lambda g : \nu_k . f g \\ &\equiv_{\eta} \lambda f : \nu_{k+1} . f. \end{aligned}$$

We used the induction step to pass from the third to the fourth line.

Next, let $p : \text{nat} \rightarrow \text{nat} \times \text{nat}$ and $q : \text{nat} \rightarrow \text{nat} \times \text{nat}$ be isomorphisms witnessing $\text{nat} \cong \text{nat} \times \text{nat}$. Then $p_k : \nu_k \rightarrow \nu_k \times \nu_k$ and $q_k : \nu_k \times \nu_k \rightarrow \nu_k$ defined by $p_0 = p$, $q_0 = q$ and

$$\begin{aligned} p_{k+1} &= \lambda f : \nu_{k+1} . \langle (\lambda g : \nu_k . \text{fst } (q (f g))), (\lambda g : \nu_k . \text{snd } (q (f g))) \rangle \\ q_{k+1} &= \lambda f : \nu_{k+1} \times \nu_{k+1} . \lambda g : \nu_k . p \langle \text{fst } f g, \text{snd } f g \rangle \end{aligned}$$

are isomorphism between ν_k and $\nu_k \times \nu_k$.

To prove the lemma it suffices to show that the types τ and $\nu_{\text{lev}(\tau)}$ are retracts of each other for any type τ . We proceed by induction on the structure of τ .

The base case nat is trivial.

For a product type $\sigma \times \tau$ we have a chain of retractions

$$\sigma \times \tau \triangleleft \nu_{\text{lev}(\sigma)} \times \nu_{\text{lev}(\tau)} \triangleleft \nu_{\text{lev}(\sigma \times \tau)} \times \nu_{\text{lev}(\sigma \times \tau)} \cong \nu_{\text{lev}(\sigma \times \tau)}.$$

In the opposite direction, if $\text{lev}(\sigma) \leq \text{lev}(\tau)$ then

$$\nu_{\text{lev}(\sigma \times \tau)} = \nu_{\text{lev}(\tau)} \triangleleft \tau \triangleleft \sigma \times \tau.$$

The case $\text{lev}(\tau) \leq \text{lev}(\sigma)$ is symmetric.

Finally consider a function type $\sigma \rightarrow \tau$. If $\text{lev}(\tau) = 0$ then $\tau = \text{nat}$, hence

$$\sigma \rightarrow \tau \triangleleft \nu_{\text{lev}(\sigma)} \rightarrow \text{nat} = \nu_{\text{lev}(\sigma)+1} = \nu_{\text{lev}(\sigma \rightarrow \tau)}$$

and

$$\nu_{\text{lev}(\sigma \rightarrow \tau)} = \nu_{\text{lev}(\sigma)+1} = \nu_{\text{lev}(\sigma)} \rightarrow \text{nat} \triangleleft \sigma \rightarrow \text{nat} = \sigma \rightarrow \tau.$$

If $\text{lev}(\tau) = n + 1$ then

$$\begin{aligned} \sigma \rightarrow \tau \triangleleft \nu_{\text{lev}(\sigma)} \rightarrow \nu_{n+1} &\cong \nu_{\text{lev}(\sigma)} \times \nu_n \rightarrow \text{nat} \triangleleft \\ \nu_{\max(\text{lev}(\sigma), n)} \rightarrow \text{nat} &= \nu_{1+\max(\text{lev}(\sigma), n)} = \nu_{\text{lev}(\sigma \rightarrow \tau)}. \end{aligned}$$

and

$$\begin{aligned} \nu_{\text{lev}(\sigma \rightarrow \tau)} &= \nu_{1+\max(\text{lev}(\sigma), n)} = \nu_{\max(\text{lev}(\sigma), n)} \rightarrow \text{nat} \triangleleft \\ \nu_{\text{lev}(\sigma)} \times \nu_n \rightarrow \text{nat} &\cong \nu_{\text{lev}(\sigma)} \rightarrow \nu_{n+1} \triangleleft \sigma \rightarrow \tau. \quad \square \end{aligned}$$

With the lemma in hand we can prove an obstruction theorem for self-interpreters.

Theorem 3.6. *For any type τ , if a weak Brown-Palsberg self-interpreter satisfies $\text{lev}(\Box\tau) < \text{lev}(\tau)$ then every $f \in \text{Prg}(\tau \rightarrow \tau)$ has a fixed point with respect to $\beta\eta$ -equivalence.*

Proof. If $\text{lev}(\Box\tau) < \text{lev}(\tau)$ then $\text{lev}(\Box\tau \rightarrow \tau) \leq \text{lev}(\tau)$, hence Lemma 3.5 gives us a section and a retraction

$$s : (\Box\tau \rightarrow \tau) \rightarrow \tau \quad \text{and} \quad r : \tau \rightarrow (\Box\tau \rightarrow \tau).$$

Define $g : \Box\tau \rightarrow \tau$ by $g = \lambda x : \Box\tau. f(r(u_\tau x))$ and let $n = \ulcorner s g \urcorner_\tau$. Then

$$r(u_\tau n) n \equiv_\beta r(s g) n \equiv_{\beta\eta} g n \equiv_\beta f(r(u_\tau n) n). \quad \square$$

The following corollary squashes any hope of a significant improvement of the notion of Brown-Palsberg interpreters, as far as complexity of the types of codes is concerned.

Corollary 3.7. *A weak Brown-Palsberg self-interpreter for System T satisfies $\text{lev}(\Box\tau) \geq \text{lev}(\tau)$ for every type τ .*

Proof. It suffices to show that every type τ has a program $f : \tau \rightarrow \tau$ without a $\beta\eta$ -fixed point. By Lemma 3.5 there is a section $s : \text{nat} \rightarrow \tau$ and a corresponding retraction $r : \tau \rightarrow \text{nat}$. Define

$$f = \lambda x : \tau. s(\text{succ}(r x)).$$

To see that f does not have a fixed point, suppose $f e \equiv_{\beta\eta} e$ for some $e \in \text{Prg}(\tau)$. Then

$$r e \equiv_{\beta\eta} r(f e) \equiv_\beta r(s(\text{succ}(r e))) \equiv_{\beta\eta} \text{succ}(r e).$$

Let n be the unique number such that $n(r e) \equiv_\beta \underline{n}$. Then

$$\underline{n+1} = \text{succ } \underline{n} \equiv_\beta \text{succ}(r e) \equiv_{\beta\eta} (r e) \equiv_\beta \underline{n},$$

which is a contradiction because it implies $n+1 = n$. To see this, we could presumably find a reference to a suitable theorem about $\beta\eta$ -normal forms for System T , but here is a semantic proof. The set-theoretic model of System T validates both β - and η -equality, thus

$\underline{k} \equiv_{\beta\eta} \underline{m}$ implies that in the set-theoretic model the interpretations of \underline{k} and \underline{m} are equal, but they are of course k and m , respectively. \square

4. THE MORAL

The self-interpreter for F_ω given by Brown and Palsberg has important structural properties that Definition 3.1 fails to capture. For instance, their encoding of types commutes with substitution [3, Thm. 5.2] and is a congruence with respect to type equality [3, Thm. 5.3]. In the original work [7] on meta-circularity Frank Pfenning and Peter Lee called such phenomena *reflexivity*. Unfortunately they spoke of it at an informal level and did not provide a definition.

In order to shed further light on self-interpreters for total languages we need a definition of self-interpreters which takes into account structural properties of self-interpreters that distinguishes the interpreter by Brown and Palsberg from the one given in Theorem 3.2. However, it will not do to simply require that $\tau \mapsto \Box\tau$ be a congruence which commutes with substitution, as that is not a general enough idea. I would instead expect to see a definition of a structure-preserving *homomorphism* between calculi, like in algebra. A good starting point might be John Longley’s notion of computability structures and simulations [5].

REFERENCES

- [1] Jeremy Avigad and Solomon Feferman. Chapter V: Gödel’s Functional (“Dialectica”) Interpretation. In Samuel R. Buss, editor, *Handbook of Proof Theory*, volume 137 of *Studies in Logic and the Foundations of Mathematics*, pages 337–405. Elsevier, 1998.
- [2] H. P. Barendregt. *The Lambda Calculus: Its Syntax and Semantics*, volume 103 of *Studies in Logic and the Foundations of Mathematics*. College Publications, 1984.
- [3] Matt Brown and Jens Palsberg. Breaking through the normalization barrier: A self-interpreter for F-omega. In *Principles of Programming Languages (POPL)*, January 2016.
- [4] F. William Lawvere. Diagonal arguments and cartesian closed categories. *Lecture Notes in Mathematics*, 92:134–145, 1969. Republished in: *Reprints in Theory and Applications of Categories*, No. 15 (2006), 1–13.
- [5] John Longley. Computability structures, simulations and realizability. *Mathematical Structures in Computer Science*, 24(2), 2014.
- [6] John Longley and Gordon Plotkin. Logical full abstraction and pcf. In *Tbilisi Symposium on Language, Logic and Computation. SiLLI/CSLI*, pages 333–352. SiLLI/CSLI, 1996.
- [7] Frank Pfenning and Peter Lee. Metacircularity in the polymorphic λ -calculus. *Theoretical Computer Science*, 89(1):137 – 159, 1991.
- [8] G.D. Plotkin. Lcf considered as a programming language. *Theoretical Computer Science*, 5(3):223 – 255, 1977.
- [9] A. S. Troelstra. *Metamathematical Investigation of Intuitionistic Arithmetic and Analysis*. Number 344 in *Lecture Notes in Mathematics*. Springer-Verlag, 1973.

ANDREJ BAUER, UNIVERSITY OF LJUBLJANA, SLOVENIA
E-mail address: Andrej.Bauer@andrej.com