

# The Functional Machine Calculus

Willem Heijltjes

*Department of Computer Science  
University of Bath  
Bath, United Kingdom*

---

## Abstract

This paper presents the Functional Machine Calculus (FMC) as a simple model of higher-order computation with “reader/writer” effects: higher-order mutable store, input/output, and probabilistic and non-deterministic computation.

The FMC derives from the lambda-calculus by taking the standard operational perspective of a call-by-name stack machine as primary, and introducing two natural generalizations. One, “locations”, introduces multiple stacks, which each may represent an effect and so enable effect operators to be encoded into the abstraction and application constructs of the calculus. The second, “sequencing”, is known from kappa-calculus and concatenative programming languages, and introduces the imperative notions of “skip” and “sequence”. This enables the encoding of reduction strategies, including call-by-value lambda-calculus and monadic constructs.

The encoding of effects into generalized abstraction and application means that standard results from the lambda-calculus may carry over to effects. The main result is confluence, which is possible because encoded effects reduce algebraically rather than operationally. Reduction generates the familiar algebraic laws for state, and unlike in the monadic setting, reader/writer effects combine seamlessly. A system of simple types confers termination of the machine.

*Keywords:* lambda-calculus, computational effects, confluence, concatenative programming

---

## 1 Introduction

Higher-order programming and computational effects are ubiquitous in modern programs. Understanding them, and in particular their potent combination, is therefore an important challenge to computer science. Higher-order functional programming enjoys an elegant foundational theory in the  $\lambda$ -calculus, where  $\beta$ -reduction gives rise not only to operational semantics—by imposing an evaluation strategy—but also to an equational theory which may be regarded as definitive for higher-order functions. For computational effects, however, there are many approaches and, as yet, no single definitive theory. Such a theory would ideally include a convenient syntax, expressing a natural and convincing semantics, and supporting reasoning tools and methods such as type-systems and compile-time optimizations, while remaining amenable to refinement, extension and variation.

The rich history of approaches to the problem of computational effects in a higher-order setting includes Landin’s pioneering work [17], which cemented the central position of  $\lambda$ -calculus, and highlighted the difficulty of reconciling the drive for a complete theory with the practice of programming: Landin focussed on a call-by-value strategy and used *thunks* to delay evaluation where necessary. A more modular and flexible account was provided by Moggi’s use of *monads* [25], which has influenced not only theoretical work but also the design of the Haskell programming language. However, a fundamental problem with monads is that they don’t compose, and in practice multiple effects are combined by building a stack of monad transformers, which can become unwieldy for the programmer. Many alternatives and refinements have been proposed, including uniqueness types [32], continuations [10], encodings in (intuitionistic) linear logic [3,22] and in process calculi [23,13], premonoidal categories [31,30,20], Call-By-Push-Value [18] and its variants in linear logic [8,9], Arrows [14], algebraic effects [27,1], and effect handlers [29].

This paper offers a new solution to this challenge: the *Functional Machine Calculus* (FMC). It includes the effects of state, input/output, and probabilistic and non-deterministic computation—here referred to collectively as *reader/writer effects*. The FMC consists of two independent generalizations of the  $\lambda$ -calculus, *locations* and *sequencing*, that individually give two fragments, the *poly- $\lambda$ -calculus* and the *sequential  $\lambda$ -calculus*. It enjoys a clean equational theory supported by a confluent reduction semantics, and expresses both effects and higher-order features in the same terms, retaining the simplicity of the  $\lambda$ -calculus while being powerful enough to capture the reality of higher-order programming with multiple effects. We provide operational semantics in terms of an abstract machine, and a type system that confers termination of the machine. The remainder of this section will introduce both generalizations. Throughout the paper, proofs are omitted when they are straightforward.

### 1.1 Locations

The main objective of this work has been to preserve confluence, following the recent presentation of a confluent probabilistic  $\lambda$ -calculus by Dal Lago, Guerrieri, and Heijltjes [4]. This is perhaps surprising, as  $\lambda$ -calculi with effects are known to be non-confluent. The apparent contradiction disappears by disentangling the *operational* and the *algebraic* aspects of evaluation. In  $\lambda$ -calculus,  $\beta$ -reduction is algebraic (or more precisely,  *$\beta$ -equivalence* is), while stack machines such as Krivine’s [15] give an operational semantics. For effects, looking up a global variable or generating a random value is operational, while effect operators may interact algebraically via the laws of Plotkin and Power [27].

	global, operational	local, algebraic
$\lambda$ -calculus	stack machines	$\beta$ -reduction
effects	update, lookup, read, write, random	algebraic effect equations

Our starting point is the observation that for both the  $\lambda$ -calculus and reader/writer effects, the operational side can be given by push and pop actions on global stacks or streams. In a simple stack machine for the  $\lambda$ -calculus, *application*  $M N$  pushes its argument  $N$  to the stack and continues as  $M$ , and *abstraction*  $\lambda x. M$  pops a term  $N$  from the stack and binds it to  $x$ , to continue as  $\{N/x\}M$  (the capture-avoiding substitution of  $N$  for  $x$  in  $M$ ). Then, the following effects are also modelled via stacks or streams:

- reading from *input* is a pop from an input stream;
- writing to *output* is a push to an output stream;
- a memory cell  $c$  is modelled by a stack of depth one, where
  - *update*  $c := N$  pops from  $c$ , discarding the value, then pushes the new value  $N$ ;
  - *lookup*  $!c$  pops the value  $N$  from  $c$ , pushes  $N$  to reinstate  $c$ , and then returns  $N$ ;
- *probabilistic* and *non-deterministic* generators can be modelled as separate input streams.

This idea is captured in the *poly- $\lambda$ -calculus*: we introduce a set of *locations* to represent independent stacks or streams on the machine, and parameterize abstraction and application in this set to act as pop and push actions on the corresponding stack. Effect operators are then encoded in these constructs according to the above scheme. Beta-reduction, generalized to multiple locations, remains confluent, and for encoded effect operators it gives rise to the expected algebraic laws [27]. However, this encoding of effects forces their call-by-name semantics, while programming with effects requires control over when they are called. This is the purpose of the second generalization, *sequencing*.

### 1.2 Sequencing

The literature offers several ways to control reduction behaviour in higher-order languages, including continuation encodings between *cbv* and *cbn* [28], and call-by-push-value (*cbpv*) [18] which encodes both. To complement *locations*, our approach takes the stack machine as primary. Viewing the  $\lambda$ -calculus as a language of machine instruction sequences, the *sequencing* generalization extends it with composition and the empty sequence, analogous to imperative “sequence” and “skip”; not by introducing these as primitives, but again by generalizing the calculus in a subtle way so that they arise naturally. Such designs have arisen several times before: in the first-order  $\kappa$ -calculus of Hasegawa [11], generalized to higher-order in the context of premonoidal categories [30]; as the  $\Lambda_s$ -calculus in an analysis of compilers [7]; and in

higher-order stack programming languages, also called *concatenative* languages [12], such as *Joy* [36], *λ-FORTH* [21], *Cat* [6], and closest to our design, *Factor* [26].

Douence and Fradet demonstrate how their  $\Lambda_s$  encodes Plotkin’s *cbv*  $\lambda$ -calculus [28] as well as Moggi’s monadic constructs [24,25], illustrating how this design gives control over reduction. We will recall these encodings in Sections 3.2 and 3.3, and demonstrate the encoding of *cbpv* and *Arrows* [14].

### 1.3 The Functional Machine Calculus

The FMC combines both generalisations, *locations* and *sequencing*, in a simple model of higher-order computation with multiple effects. Its design and solid foundations in semantics mean that several important properties of the  $\lambda$ -calculus are preserved. The aim of this paper is introductory: it presents the syntax, operational semantics, and fundamental ideas and results with an emphasis on explanation and examples. In a forthcoming paper we will deepen these results with a strong normalization theorem, a domain-theoretic semantics of the untyped calculus, and semantics of the typed calculus in premonoidal and in Cartesian closed categories. The main results for the FMC as presented here are the following:

**Confluence** Beta-reduction is confluent, with evaluation behaviour expressed in syntax.

**Algebraic effects** The algebraic laws for reader/writer effects arise from reduction.

**Compositionality** Reader/writer effects combine seamlessly, due to the use of independent locations.

**Types** Simple types cover (higher-order) effect operations and confer termination of the machine.

## 2 The poly-lambda-calculus

We introduce a set of *locations*  $A$ , ranged over by  $a, b, c, \dots$ , to indicate the different stacks or streams for each effect. Abstraction and application are parameterized in  $A$  to give the corresponding pop and push actions. We write application  $M N$  as  $[N].M$  to emphasize the operational reading (push  $N$  and continue as  $M$ ), to easily attach a location  $a$ , and to give unique parsing—cf. De Bruijn [5]. Abstraction  $\lambda x. N$  is written  $\langle x \rangle. N$  to emphasize the duality with application.

**Definition 2.1** The *poly- $\lambda$ -calculus* is given by the grammar

$$M, N ::= x \mid [N]a.M \mid a\langle x \rangle.M$$

with from left to right a *variable*, an *application* or *push action* on location  $a$  with function  $M$  and argument  $N$ , and an *abstraction* or *pop action* on location  $a$  that binds  $x$  in  $M$ . Terms are considered modulo  $\alpha$ -equivalence. The regular  $\lambda$ -calculus embeds via a dedicated *main* location  $\lambda \in A$ , omitted from terms for brevity; so we may write  $\lambda x. M$  or  $\langle x \rangle. M$  for  $\lambda\langle x \rangle. M$ , and  $M N$  or  $[N].M$  for  $[N]\lambda. M$ .

The *poly-stack machine* is given by the following data. A *stack* of terms  $S$  is written with the top element to the right; we define them inductively below left, but they should better be considered as coinductive, to include streams. A *memory*  $S_A$  is a family of stacks or streams in  $A$ , defined below left. We write  $S_A; S_a$  to identify the stack for  $a$  in  $S_A$ . A *state* is a pair  $(S_A, M)$ , and the *transitions* or *steps* are given as top-to-bottom rules below centre. A *run* of the machine is a sequence of steps, written as  $(S_A, M) \Downarrow (T_A, N)$  or with a double line as below right.

$$\begin{array}{l} S ::= \varepsilon \mid S \cdot M \\ S_A ::= \{ S_a \mid a \in A \} \end{array} \quad \begin{array}{c} \frac{}{(S_A; S_a, [N]a.M)} \\ \frac{}{(S_A; S_a \cdot N, M)} \end{array} \quad \begin{array}{c} \frac{}{(S_A; S_a \cdot N, a\langle x \rangle.M)} \\ \frac{}{(S_A; S_a, \{N/x\}M)} \end{array} \quad \begin{array}{c} \frac{}{(S_A, M)} \\ \frac{}{(T_A, N)} \end{array}$$

### 2.1 Encoding effects

Consider the following  $\lambda$ -calculus with effects. We will encode it in the poly- $\lambda$ -calculus, with its *lazy* or *cbn* semantics. At the end of this section we will consider what would be needed to encode its *eager* or *cbv* semantics. (We assume familiarity with the operational semantics of  $\lambda$ -calculus and of effects; for an

introduction see e.g. Winskel [37].)

$M, N, P ::= x \mid MN \mid \lambda x.M$	$\lambda$ -calculus
read   write $N; M$	input/output
$c := N; M$   $!c$	state update and lookup
$N \oplus M$   $N + M$	probabilistic and non-deterministic sum

The *cbn*-encoding will follow the description in the introduction. The constructs of the above language are introduced as defined constructs (“sugar”) into the poly- $\lambda$ -calculus.

**Input/output:** Input uses a dedicated *input* location  $\mathbf{in} \in A$  and is encoded by  $\text{read} \triangleq \mathbf{in}\langle x \rangle.x$ . The machine is initialized with a stream  $S_{\text{in}} = \cdots N_3 \cdot N_2 \cdot N_1$  (infinite to the left), and the pop transition gives the expected operational semantics, below left. Writing to output uses a dedicated *output* location  $\mathbf{out} \in A$  and is encoded by  $\text{write } N; M \triangleq [N]\mathbf{out}.M$ . Evaluation then generates an output stream  $N_1, N_2, \dots$  (finite at any step) by the *push* machine transition, below right.

$$\frac{(S_A ; S_{\text{in}} \cdot N, \mathbf{in}\langle x \rangle.x)}{(S_A ; S_{\text{in}}, N)} \quad \frac{(S_A ; S_{\text{out}}, [N]\mathbf{out}.M)}{(S_A ; S_{\text{out}} \cdot N, M)}$$

**State:** A memory cell is modelled by a location  $c \in A$ . The associated stack is expected to hold at most one value, which is preserved by the encoding of the operators, and not enforced externally. Update and lookup are encoded by  $c := N; M \triangleq c(\_).[N]c.M$  and  $!c \triangleq c\langle x \rangle.[x]c.x$  where the underscore ( $\_$ ) represents a variable that does not occur in  $M$  or  $N$ . In the machine, the stack for each cell is initialized with a (dummy) value, and the transitions then give the expected operational semantics.

$$\frac{(S_A ; \varepsilon_c \cdot P, c(\_).[N]c.M)}{(S_A ; \varepsilon_c, [N]c.M)} \quad \frac{(S_A ; \varepsilon_c \cdot M, c\langle x \rangle.[x]c.x)}{(S_A ; \varepsilon_c, [M]c.M)}$$

$$\frac{(S_A ; \varepsilon_c \cdot N, M)}{(S_A ; \varepsilon_c \cdot M, M)}$$

**Probabilistic and non-deterministic sums:** Following the probabilistic case [4], probabilistic and non-deterministic sums are included via dedicated locations  $\mathbf{rnd}, \mathbf{nd} \in A$  by  $N \oplus M \triangleq \mathbf{rnd}\langle x \rangle.x M N$  and  $N + M \triangleq \mathbf{nd}\langle x \rangle.x M N$ . The machine is initialized with the corresponding streams of Church-encoded Booleans  $\lambda x. \lambda y. x$  and  $\lambda x. \lambda y. y$ , generated probabilistically for  $\mathbf{rnd}$  and non-deterministically for  $\mathbf{nd}$ . The machine steps are as expected.

**Example 2.2** Consider the following example term and its *cbn* encoding in the poly- $\lambda$ -calculus. (Numbers can be seen informally as constants, or as Church numerals.)

$$a := 2; (\lambda x. !a)(a := 3; 0) = a(\_).[2]a.[a(\_).[3]a.0].\langle x \rangle.a\langle y \rangle.[y]a.y$$

Its *cbn* reduction gives 2. It evaluates in the machine as follows (where the cell  $a$  is initialized with zero).

$$\frac{(\varepsilon_a \cdot 0; \varepsilon_\lambda, a(\_).[2]a.[a(\_).[3]a.0].\langle x \rangle.a\langle y \rangle.[y]a.y)}{(\varepsilon_a; \varepsilon_\lambda, [2]a.[a(\_).[3]a.0].\langle x \rangle.a\langle y \rangle.[y]a.y)}$$

$$\frac{(\varepsilon_a \cdot 2; \varepsilon_\lambda, [a(\_).[3]a.0].\langle x \rangle.a\langle y \rangle.[y]a.y)}{(\varepsilon_a \cdot 2; \varepsilon_\lambda \cdot a(\_).[3]a.0, \langle x \rangle.a\langle y \rangle.[y]a.y)}$$

$$\frac{(\varepsilon_a \cdot 2; \varepsilon_\lambda, a\langle y \rangle.[y]a.y)}{(\varepsilon_a \cdot 2; \varepsilon_\lambda, a\langle y \rangle.[y]a.y)}$$

$$\frac{(\varepsilon_a; \varepsilon_\lambda, [2]a.2)}{(\varepsilon_a \cdot 2; \varepsilon_\lambda, [2]a.2)}$$

$$(\varepsilon_a \cdot 2; \varepsilon_\lambda, 2)$$

## 2.2 Beta-reduction

In the  $\lambda$ -calculus,  $\beta$ -reduction lets successive push- and pop-actions interact. Generalizing to multiple locations, these must be actions on the same stack, while other stacks may be accessed in-between. The  $\beta$ -rule is then as below, where each  $X_i$  is an abstraction or application not on location  $a$ , and (if the former) not capturing in  $N$ . Reduction is closed under any context. (A formal definition is given in Section 4.)

$$[N]a. X_1 \dots X_n. a\langle x \rangle. M \rightarrow X_1 \dots X_n. \{N/x\}M$$

**Example 2.3** The term from example 2.2 reduces as follows, with reduced redexes underlined.

$$a\langle \_ \rangle. \underline{[2]a. [a\langle \_ \rangle. [3]a. 0]}. \langle x \rangle. \underline{a\langle y \rangle}. [y]a. y} \rightarrow a\langle \_ \rangle. \underline{[a\langle \_ \rangle. [3]a. 0]}. \langle x \rangle. [2]a. 2} \rightarrow a\langle \_ \rangle. [2]a. 2 = a := 2 ; 2$$

Analogous to  $\beta$ -reduction,  $\eta$ -reduction is as below, where each  $X_i$  is an abstraction or application not on location  $a$ , and  $x$  does not occur free in any of the  $X_i$  nor in  $M$ .

$$a\langle x \rangle. X_1 \dots X_n. [x]a. M \rightarrow_{\eta} X_1 \dots X_n. M$$

As an alternative to these rule schemes, terms may be taken modulo an equivalence  $\sim$  generated by the *permutations* below left, and  $\beta$ - and  $\eta$ -reduction defined only on adjacent operators, as below right. Observe that the machine semantics immediately validates these equivalences. We will use the below formulation to consider the relation with algebraic effects, but otherwise use the above formulation.

$$\begin{array}{l} [M]a. [N]b. P \sim [N]b. [M]a. P \\ \langle x \rangle a. [N]b. P \sim [N]b. \langle x \rangle a. P \quad \text{if } x \notin \text{fv}(N) \\ \langle x \rangle a. \langle y \rangle b. P \sim \langle y \rangle b. \langle x \rangle a. P \end{array} \quad \begin{array}{l} [N]a. a\langle x \rangle. M \rightarrow_{\beta} \{N/x\}M \\ a\langle x \rangle. [x]a. M \rightarrow_{\eta} M \quad \text{if } x \notin \text{fv}(M) \end{array}$$

Beta-reduction is confluent, as will be shown more generally for the FMC in Section 4. This is possible because it follows the algebraic laws for effects [27] instead of their operational semantics. For instance, laws for the interaction of *lookup* and *update* correspond to the following reductions.

$$\begin{array}{l} c := M ; c := N ; P = c\langle \_ \rangle. [M]c. c\langle \_ \rangle. [N]c. P \rightarrow c\langle \_ \rangle. [N]c. P = c := N ; P \\ c := M ; !c = c\langle \_ \rangle. [M]c. c\langle x \rangle. [x]c. x \rightarrow c\langle \_ \rangle. [M]c. M = c := M ; M \end{array}$$

The seven algebraic laws for global state of Plotkin and Power [27, p. 348] arise in our setting from  $\beta/\eta$ -reduction and  $\sim$ . Their notation has *update*  $u_{loc,v}(M)$  of location  $loc$  with value  $v$  in  $M$ , and *lookup*  $l_{loc}(M)_v$  of the value  $v$  at location  $loc$  using the value  $v$  as a parameter in  $M$ . These are encoded in the poly- $\lambda$ -calculus as below, using abstraction with  $x$  instead of parametrization in  $v$  in the lookup case. Values  $v$  may be taken as arbitrary poly- $\lambda$ -terms.

$$u_{a,v}(M) \triangleq a\langle \_ \rangle. [v]a. M \quad l_a(M)_x \triangleq a\langle x \rangle. [x]a. M$$

**Proposition 2.4** *The poly- $\lambda$ -calculus with  $\rightarrow_{\beta\eta} \cup \sim$  generates the algebraic laws for state.*

**Proof.** By the following equations, where  $a \neq b$  and in equation 7,  $x \notin \text{fv}(v)$ .

$$\begin{array}{l} 1. l_a(u_{a,y}(x))_y = a\langle y \rangle. [y]a. a\langle \_ \rangle. [y]a. x \rightarrow_{\beta} a\langle y \rangle. [y]a. x \rightarrow_{\eta} x \\ 2. l_a(l_a(M_{x,y})_x)_y = a\langle y \rangle. [y]a. a\langle x \rangle. [x]a. M_{x,y} \rightarrow_{\beta} a\langle y \rangle. [y]a. \{y/x\}M_{x,y} = l_a(M_{y,y})_y \\ 3. u_{a,v}(u_{a,v'}(x)) = a\langle \_ \rangle. [v]a. a\langle \_ \rangle. [v']a. x \rightarrow_{\beta} a\langle \_ \rangle. [v']a. x = u_{a,v'}(x) \\ 4. u_{a,v}(l_a(M_x)_x) = a\langle \_ \rangle. [v]a. a\langle x \rangle. [x]a. M_x \rightarrow_{\beta} a\langle \_ \rangle. [v]a. \{v/x\}M_x = u_{a,v}(M_v) \\ 5. l_a(l_b(M_{x,y})_y)_x = a\langle x \rangle. [x]a. b\langle y \rangle. [y]b. M_{x,y} \sim b\langle y \rangle. [y]b. a\langle x \rangle. [x]a. M_{x,y} = l_b(l_a(M_{x,y})_x)_y \\ 6. u_{a,v}(u_{b,v'}(x)) = a\langle \_ \rangle. [v]a. b\langle \_ \rangle. [v']b. x \sim b\langle \_ \rangle. [v']b. a\langle \_ \rangle. [v]a. x = u_{b,v'}(u_{a,v}(x)) \\ 7. u_{a,v}(l_b(M_x)_x) = a\langle \_ \rangle. [v]a. b\langle x \rangle. [x]b. M_x \sim b\langle x \rangle. [x]b. a\langle \_ \rangle. [v]a. M_x = l_b(u_{a,v}(M_x))_x \end{array}$$

□

### 2.3 Poly-types

A simple type for a poly-term represents its expected inputs, taken from multiple independent locations. Correspondingly, the antecedent of an implication is parameterized in a location, and implications on distinct locations may permute.

**Definition 2.5** *Simple poly-types* are given by the grammar below left, where  $o$  (omicron) is a *base* type and  $a(\sigma) \rightarrow \tau$  an *arrow* type, and considered modulo the congruence  $\sim$  given below right.

$$\rho, \sigma, \tau ::= o \mid a(\sigma) \rightarrow \tau \quad a(\rho) \rightarrow b(\sigma) \rightarrow \tau \sim b(\sigma) \rightarrow a(\rho) \rightarrow \tau \quad (a \neq b).$$

The typing rules are as follows, where a *context*  $\Gamma$  is a finite function from variables to types.

$$\frac{}{\Gamma, x : \tau \vdash x : \tau} \quad \frac{\Gamma, x : \sigma \vdash M : \tau}{\Gamma \vdash a(x). M : a(\sigma) \rightarrow \tau} \quad \frac{\Gamma \vdash N : \sigma \quad \Gamma \vdash M : a(\sigma) \rightarrow \tau}{\Gamma \vdash [N]a. M : \tau}$$

Observe that the congruence  $\sim$  means that a term  $M : a(\rho) \rightarrow b(\sigma) \rightarrow \tau$  may be prefixed by a push action  $[N]a$  where  $N : \rho$  or by one  $[P]b$  where  $P : \sigma$  (or both, in either order).

### 2.4 Towards encoding call-by-value semantics

The poly- $\lambda$ -calculus gives control over when effects are called, as we demonstrate by the following example.

**Example 2.6** Consider the following example term, which is a normal form with cbn semantics.

$$f (c := 2; 0) (!c) (c := 3; 1)$$

With cbv, the arguments may be evaluated left to right, reducing to  $f \ 0 \ 2 \ 1$ , or right to left, which gives  $f \ 0 \ 3 \ 1$ . The two readings are encoded as follows (using regular applications to  $f$  for readability).

$$\begin{aligned} c(\cdot). [2]c. c(x). [x]c. c(\cdot). [3]c. f \ 0 \ x \ 1 &\rightarrow c(\cdot). [3]c. f \ 0 \ 2 \ 1 \\ c(\cdot). [3]c. c(x). [x]c. c(\cdot). [2]c. f \ 0 \ x \ 1 &\rightarrow c(\cdot). [2]c. f \ 0 \ 3 \ 1 \end{aligned}$$

The encodings rely on repositioning an update  $c := 1$  as a prefix, and for a lookup  $!c$ , on separating the global actions  $c(x). [x]c$  from the variable  $x$  where the value is used. (The latter idea gives the cbv semantics in the probabilistic case [4].) However, it is unlikely that an encoding that only repositions effect operations can encode the cbv semantics of  $\lambda$ -calculus with effects. Consider the following example.

**Example 2.7** With a cbv semantics, the term

$$a := (\lambda x. b := 1; x) 0; !b$$

first reduces the redex, to give  $a := (b := 1; 0); !b$ , which then evaluates by updating  $b := 1$ , then  $a := 0$ , and reading  $!b$  as  $1$ . To obtain this semantics in the poly- $\lambda$ -calculus by manipulation of effect operations would require lifting  $b := 1$  out of a redex—a process which is likely undecidable in general.

The poly- $\lambda$ -calculus thus gives control over effects, but not evaluation behaviour in general. It is an open question whether this is sufficient for practical purposes—one we cannot answer here. Instead, we will consider *sequencing* as a natural way to include cbv semantics.

## 3 The sequential lambda-calculus

As an instruction sequence for a stack machine, a  $\lambda$ -term is a string of push and pop actions that must end in a variable. But the machine would naturally accept *any* sequence of actions and variables. Relaxing the variable restriction would further enable composition of sequences. This design of  $\lambda$ -calculus with sequential composition appears several times in the literature and in practice: as the calculus  $\Lambda_s$  [7], as the higher-order  $\kappa$ -calculus [30], and in concatenative programming languages such as Factor [26]. We call

this generalization of the  $\lambda$ -calculus *sequencing*, and implement it by introducing a *skip* (or *nil*) construct and making the variable a prefix.

**Definition 3.1** The *sequential  $\lambda$ -calculus* is given by the following grammar.

$$M, N, P ::= \star \mid x.M \mid [N].M \mid \langle x \rangle.M$$

We may omit the trailing  $\star$  from terms for readability. *Capture-avoiding composition*  $N;M$  is given by

$$\star;M = M \quad x.N;M = x.(N;M) \quad [P].N;M = [P].(N;M) \quad \langle y \rangle.N;M = \langle y \rangle.(N;M)$$

where in the last case  $y$  is not free in  $M$ . *Capture-avoiding substitution*  $\{M/x\}N$  is as follows.

$$\begin{aligned} \{M/x\}\star &= \star & \{M/x\}[P]a.N &= [\{M/x\}P]a.\{M/x\}N \\ \{M/x\}x.N &= M;\{M/x\}N & \{M/x\}a\langle x \rangle.N &= a\langle x \rangle.N \\ \{M/x\}y.N &= y.\{M/x\}N \quad (x \neq y) & \{M/x\}a\langle y \rangle.N &= a\langle y \rangle.\{M/x\}N \quad (x \neq y, y \notin \text{fv}(M)) \end{aligned}$$

*Beta-reduction* is otherwise standard, by closing the rule below left under all contexts. The abstract machine has *states*  $(S, M)$  of a stack and a term, and the transitions below right.

$$[N].\langle x \rangle.M \rightarrow \{N/x\}M \quad \frac{(S, [N].M)}{(S.N, M)} \quad \frac{(S.N, \langle x \rangle.M)}{(S, \{N/x\}M)}$$

**Example 3.2** Consider the following example terms.

$$\langle x \rangle.[x].[x] \quad \langle x \rangle.\langle y \rangle \quad [\langle x \rangle.[x]].\langle f \rangle.f.f.f$$

The first duplicates the top item on the stack; the second removes two items; the third pushes the term  $\langle x \rangle.[x]$  (which picks up and returns an item), pops it as  $f$ , and runs it three times.

Observe that the changes to evaluation are absorbed by substitution and composition, while  $\beta$ -reduction and machine evaluation remain unchanged. There is nevertheless a change in perspective from the  $\lambda$ -calculus, in that the return values or outputs of a computation are pushed to the stack, rather than left as the remainder of the term. Machine evaluation for a term  $M$  with input stack  $S$  is expected to terminate in  $\star$  (just as imperative computation successfully terminates in a *skip* command) with an output stack  $T$ , i.e.  $(S, M) \Downarrow (T, \star)$ . Then  $\star$  gives the identity run (of zero steps), and  $M;N$  gives composition of runs.

**Proposition 3.3** *If  $(R, M) \Downarrow (S, \star)$  and  $(S, N) \Downarrow (T, \star)$  then  $(R, M;N) \Downarrow (T, \star)$ .*

### 3.1 Sequential types

The type system for the sequential  $\lambda$ -calculus follows that of the  $\kappa$ -calculus [30]; similar type systems have also been studied for stack languages [33]. The type of a term describes the input/output behaviour of its machine evaluation: it consists of an implication between a vector of input types, one for each element consumed from the stack, and a vector of output types, one for each item returned to the stack.

**Definition 3.4** *Sequential types* are defined by:

$$\rho, \sigma, \tau, v ::= \vec{\sigma} \Rightarrow \vec{\tau} \quad \vec{\tau} ::= \tau_1 \dots \tau_n$$

Vector concatenation is by juxtaposition,  $\vec{\sigma}\vec{\tau}$ , and the reverse of a vector  $\vec{\tau} = \tau_1 \dots \tau_n$  is  $\overleftarrow{\tau} = \tau_n \dots \tau_1$ . Typing rules for the sequential  $\lambda$ -calculus are given below. A stack is typed by a type vector, where  $\Gamma \vdash \varepsilon.M_1 \dots M_n : \tau_1 \dots \tau_n$  if  $\Gamma \vdash M_i : \tau_i$  for each  $i \leq n$ .

$$\frac{}{\Gamma \vdash \star : \vec{\tau} \Rightarrow \vec{\tau}} \star \quad \frac{\Gamma, x : \vec{\rho} \Rightarrow \vec{\sigma} \vdash M : \vec{\sigma} \vec{\tau} \Rightarrow \vec{v}}{\Gamma, x : \vec{\rho} \Rightarrow \vec{\sigma} \vdash x.M : \vec{\rho} \vec{\tau} \Rightarrow \vec{v}} \text{var} \quad \frac{\Gamma, x : \rho \vdash M : \vec{\sigma} \Rightarrow \vec{\tau}}{\Gamma \vdash \langle x \rangle.M : \rho \vec{\sigma} \Rightarrow \vec{\tau}} \text{abs} \quad \frac{\Gamma \vdash N : \rho \quad \Gamma \vdash M : \rho \vec{\sigma} \Rightarrow \vec{\tau}}{\Gamma \vdash [N].M : \vec{\sigma} \Rightarrow \vec{\tau}} \text{app}$$

**Example 3.5** The terms in Example 3.2 can be typed as follows.

$$\langle x \rangle. [x]. [x] : \tau \Rightarrow \tau \tau \quad \langle x \rangle. \langle y \rangle : \tau \sigma \Rightarrow \quad \langle \langle x \rangle. [x] \rangle. \langle f \rangle. f. f. f : \tau \Rightarrow \tau$$

Observe that because stacks are last-in first-out, the identity function on the top two stack items is the term  $\langle x \rangle. \langle y \rangle. [y]. [x] : \tau \sigma \Rightarrow \sigma \tau$ , whereas the function that swaps them is  $\langle x \rangle. \langle y \rangle. [x]. [y] : \tau \sigma \Rightarrow \tau \sigma$ .

**Example 3.6** The term  $\lambda x. x x = \langle x \rangle. [x]. x$  can be typed by assigning  $x$  a type that does not consume input, i.e. of the form  $(\Rightarrow \bar{\tau})$ . The self-application  $x x = [x]. x$  then has the type  $\Rightarrow (\Rightarrow \bar{\tau}) \bar{\tau}$ , which reflects that the return values accumulate: if evaluating  $x$  returns the stack  $T : \bar{\tau}$ , then  $[x]. x$  returns the stack consisting of  $x$  prepended to  $T$ . The type derivation is below. Note that the term  $(\lambda x. x x)(\lambda y. y y)$  is not typeable: the argument  $\lambda y. y y$  needs a type that takes input, and then so should  $x$ .

$$\frac{\frac{x : \Rightarrow \bar{\tau} \vdash \star : \bar{\tau} \Rightarrow \bar{\tau}^*}{x : \Rightarrow \bar{\tau} \vdash x. \star : \Rightarrow \bar{\tau}} \text{var} \quad \frac{x : \Rightarrow \bar{\tau} \vdash \star : \bar{\tau} (\Rightarrow \bar{\tau}) \Rightarrow (\Rightarrow \bar{\tau}) \bar{\tau}^*}{x : \Rightarrow \bar{\tau} \vdash x. \star : (\Rightarrow \bar{\tau}) \Rightarrow (\Rightarrow \bar{\tau}) \bar{\tau}} \text{var}}{\frac{x : \Rightarrow \bar{\tau} \vdash [x. \star]. x. \star : \Rightarrow (\Rightarrow \bar{\tau}) \bar{\tau}}{\vdash \langle x \rangle. [x. \star]. x. \star : (\Rightarrow \bar{\tau}) \Rightarrow (\Rightarrow \bar{\tau}) \bar{\tau}} \text{abs}} \text{app}$$

**Remark 3.7** [Due to Chris Barrett] Observe that all sequential types  $\tau$  are inhabited by at least the element  $\perp_\tau : \tau$ , defined below (note that the base case  $n = m = 0$  gives  $\star : (\Rightarrow)$ ). This is in contrast with the simply-typed  $\lambda$ -calculus, where not all types are inhabited due to the presence of the uninhabited base type  $o$ .

$$\perp_\tau = \langle x_n \rangle \dots \langle x_1 \rangle. [\perp_{\tau_1}] \dots [\perp_{\tau_m}] \quad \text{where} \quad \tau = \sigma_n \dots \sigma_1 \Rightarrow \tau_1 \dots \tau_m$$

### 3.2 Encodings of call-by-name calculi

The (regular, call-by-name)  $\lambda$ -calculus is included as a fragment of the sequential  $\lambda$ -calculus. We extend this embedding with types and with products, and to Moggi's *computational metalanguage* [25] following Douce and Fradet [7]. The main observation is that implications embed as input-only sequential types, and products as output-only types, below left. The formal, inductive encoding is then below right.

$$\begin{array}{ll} \tau_1 \rightarrow \dots \rightarrow \tau_n \rightarrow o = \tau_1 \dots \tau_n \Rightarrow & o \triangleq (\Rightarrow) \quad \rho \rightarrow (\bar{\sigma} \Rightarrow \bar{\tau}) \triangleq \rho \bar{\sigma} \Rightarrow \bar{\tau} \\ \tau_1 \times \dots \times \tau_n = \Rightarrow \tau_n \dots \tau_1 & 1 \triangleq (\Rightarrow) \quad \sigma \times \tau \triangleq \Rightarrow \tau \sigma \end{array}$$

Following the types, product terms encode as follows.

$$(M, N) \triangleq [N]. [M] \quad \pi_i(P) \triangleq P ; \langle x_1 \rangle. \langle x_2 \rangle. x_i \quad () \triangleq \star$$

The computational metalanguage extends the  $\lambda$ -calculus with monadic type formers  $T(\sigma)$ , and a *return* construct  $[M]_T$  and a *let* construct  $\text{let}_T$  parameterized in  $T$ . In the interpretation in the sequential  $\lambda$ -calculus, the return value of a monadic function is pushed to the stack. A monadic function type is then interpreted as one with a single output, as below left. The language constructs are encoded as below right. It is easily verified that this extends correctly to type derivations and reductions.

$$\tau_1 \rightarrow \dots \rightarrow \tau_n \rightarrow T(\sigma) = \tau_1 \dots \tau_n \Rightarrow \sigma \quad [M]_T \triangleq [M] \quad \text{let}_T x \leftarrow N \text{ in } M \triangleq N ; \langle x \rangle. M$$

### 3.3 Encodings of call-by-value calculi

The cbv  $\lambda$ -calculus [28] and the computational  $\lambda$ -calculus  $\lambda_c$  [24], which extends the former with a monadic type constructor  $T$  and with the term constructs *return* and *let*, have an encoding in the  $\kappa$ -calculus [7,30]. We recall this for the sequential  $\lambda$ -calculus, and observe that it naturally extends to types. The cbv-interpretation of types is naturally viewed in two stages. First, types in isolation are translated as follows.

$$o_v = (\Rightarrow) \quad (\sigma \rightarrow \tau)_v = \sigma_v \Rightarrow \tau_v \quad T(\sigma)_v = \Rightarrow \sigma_v$$



Evaluation of a  $\lambda_c$ -term returns a value, which in the encoding is pushed to the stack. A typed term  $M : \tau$  will then translate as  $M_v : \Rightarrow \tau_v$ , with a single output type. Terms of the computational  $\lambda$ -calculus are then interpreted as below. Again it is easily verified that this extends correctly to type derivations and reduction. The machine behaviour of encoded terms can be seen to follow the SECD-machine [16].

$$\begin{aligned} x_v &= [x] & (\lambda x.M)_v &= [\langle x \rangle.M_v] & ([M]_T)_v &= [M_v] \\ (MN)_v &= N_v; M_v; \langle x \rangle.x & (\text{let}_T x \leftarrow N \text{ in } M)_v &= N_v; \langle x \rangle.M_v \end{aligned}$$

### 3.4 Arrows, call-by-push-value, and kappa-calculus

The sequential  $\lambda$ -calculus may encode the related formalisms of Arrows, **cbpv**, and  $\kappa$ -calculus. For reasons of space, we will not recall these calculi in detail and only provide an outline to the interested reader.

Hughes's *Arrows* [14] take the  $\lambda$ -calculus with products and extend it with a second implication  $\sigma \rightsquigarrow \tau$ , which we interpret directly as that of the sequential  $\lambda$ -calculus,  $\sigma \Rightarrow \tau$ . Arrow terms have three constructors, encoded as below. The first lifts a regular term  $M : \rho \rightarrow \sigma$  to an arrow term; the second composes two arrow terms  $P : \rho \rightsquigarrow \sigma$  and  $Q : \sigma \rightsquigarrow \tau$ ; and the third applies the arrow term  $P$  to the first element of a pair.

$$\begin{aligned} \text{arr } M : \rho \rightsquigarrow \sigma &\triangleq \langle x \rangle. [[x]. M] : \rho \Rightarrow \sigma \\ P \gg Q : \rho \rightsquigarrow \tau &\triangleq P; Q : \rho \Rightarrow \tau \\ \text{first } P : (\rho \times \tau) \rightsquigarrow (\sigma \times \tau) &\triangleq \langle x \rangle. [x.P] : (\Rightarrow \tau \rho) \Rightarrow (\Rightarrow \tau \sigma) \end{aligned}$$

The perspective that emerges from this encoding is that the arrow calculus corresponds to a version of the sequential  $\lambda$ -calculus with binary products instead of stacks (which may be considered  $n$ -ary products).

Characteristic of **cbpv** [18,19], and also featured in  $\kappa$ -calculus, is the separation of *computations* and *values*. In the sequential  $\lambda$ -calculus this distinction is present, too, if implicitly: values live on the stack, and computations run the machine. To make it explicit, we may extend the calculus with *think* and *force* constructs  $!M$  and  $?V$ , and their reduction rule, as below left. Term constructs of **cbpv** (without products or sums) then embed as below right.

$$\begin{array}{ll} V, W ::= x \mid !M & \text{think } M \triangleq !M \quad \text{force } V \triangleq ?V \\ M, N ::= \star \mid ?V.M \mid [V].M \mid \langle x \rangle.M & \text{return } V \triangleq [V] \quad N \text{ to } x.M \triangleq N; \langle x \rangle.M \\ !?N.M \rightarrow N; M & \lambda x.M \triangleq \langle x \rangle.M \quad V \text{ ' } M \triangleq [V].M \end{array}$$

Types for **cbpv** feature a monadic functor  $F$  and a *value* functor  $U$ . The latter could be introduced into sequential types as below left, though the structure of the arrow type  $\Rightarrow$  makes it redundant. Types then further encode as call-by-name types, below right.

$$\sigma, \tau ::= \vec{\sigma} \Rightarrow \vec{\tau} \quad \vec{\tau} ::= U\tau_1 \dots U\tau_n \quad \rho \rightarrow (\vec{\sigma} \Rightarrow \vec{\tau}) \triangleq \rho \vec{\sigma} \Rightarrow \vec{\tau} \quad F\sigma \triangleq \Rightarrow \sigma$$

The higher-order  $\kappa$ -calculus [30] is closely related to the sequential  $\lambda$ -calculus. Types are the same as sequential types: an implication between type vectors. Terms omit the unit  $\star$  and have composition  $M; N$  as a primitive (rather than prefixing). The remaining constructs encode as follows.

$$\text{push } V \triangleq [V] \quad \kappa x.M \triangleq \langle x \rangle.M \quad \text{mkthink } M \triangleq !M \quad \text{apply} \triangleq \langle x \rangle.?x$$

### 3.5 String diagrams

We may view typed sequential  $\lambda$ -terms as string diagrams. A term  $M : \rho_1 \dots \rho_m \Rightarrow \sigma_n \dots \sigma_1$  is rendered as below. The wires represent the input and output stacks, with the first element at the top.

$$\begin{array}{c} \rho_1 \text{---} \\ \vdots \\ \rho_m \text{---} \end{array} \boxed{M} \begin{array}{c} \text{---} \sigma_1 \\ \vdots \\ \text{---} \sigma_n \end{array}$$

We will use these diagrams to illustrate how types compose. First, *strict composition* is the composition of terms  $M : \tilde{\rho} \Rightarrow \tilde{\sigma}$  and  $N : \tilde{\sigma} \Rightarrow \tilde{\tau}$  into  $M ; N : \tilde{\rho} \Rightarrow \tilde{\tau}$ , below left. This does not give the most general form of composition. For that, we combine it with the following notion of *expansion*. If a term takes an input stack  $R$  to an output stack  $S$ , then when given a larger stack  $UR$  it returns  $US$ , with  $U$  untouched. Then if  $M : \tilde{\rho} \Rightarrow \tilde{\sigma}$  also  $M : \tilde{\rho}\tilde{v} \Rightarrow \tilde{v}\tilde{\sigma}$ , illustrated below right.

$$\begin{array}{c} \tilde{\rho} \text{---} \\ \vdots \\ \tilde{\rho} \text{---} \end{array} \boxed{M} \begin{array}{c} \vdots \\ \vdots \\ \vdots \end{array} \boxed{N} \begin{array}{c} \text{---} \tilde{\tau} \\ \vdots \\ \text{---} \tilde{\tau} \end{array} \qquad \begin{array}{c} \tilde{\rho} \text{---} \\ \vdots \\ \tilde{\rho} \text{---} \end{array} \boxed{M} \begin{array}{c} \text{---} \tilde{\sigma} \\ \vdots \\ \text{---} \tilde{\sigma} \end{array} \\ \tilde{v} \text{---} \text{---} \tilde{v} \end{array}$$

These constructions combine to give the general case, where in  $M ; N$  the type of  $M$  or  $N$  may be expanded. Note that in the regular  $\lambda$ -calculus only the first case arises, as the second case requires multiple outputs.

$$\begin{array}{c} \tilde{\rho} \text{---} \\ \vdots \\ \tilde{\rho} \text{---} \end{array} \boxed{M} \begin{array}{c} \vdots \\ \vdots \\ \vdots \end{array} \boxed{N} \begin{array}{c} \text{---} \tilde{\tau} \\ \vdots \\ \text{---} \tilde{\tau} \end{array} \qquad \begin{array}{c} \tilde{\rho} \text{---} \\ \vdots \\ \tilde{\rho} \text{---} \end{array} \boxed{M} \begin{array}{c} \text{---} \tilde{\tau} \\ \vdots \\ \text{---} \tilde{\tau} \end{array} \\ \tilde{v} \text{---} \text{---} \tilde{v} \end{array}$$

**Definition 3.8** *Type composition*  $\sigma \cdot \tau$  is the partial operation given below, and is undefined otherwise.

$$\begin{aligned} (\tilde{\rho} \Rightarrow \tilde{\sigma}) \cdot (\tilde{\sigma} \tilde{v} \Rightarrow \tilde{\tau}) &= (\tilde{\rho} \tilde{v} \Rightarrow \tilde{\tau}) \\ (\tilde{\rho} \Rightarrow \tilde{v} \tilde{\sigma}) \cdot (\tilde{\sigma} \Rightarrow \tilde{\tau}) &= (\tilde{\rho} \Rightarrow \tilde{v} \tilde{\tau}) \end{aligned}$$

Observe that the first equation in Definition 3.8 corresponds to the left composition diagram above, and the second equation to the right diagram, where in both cases  $\tilde{\sigma}$  gives the types of the connecting wires between  $M$  and  $N$ . The following proposition establishes these basic properties, as well as the familiar *subject reduction*.

**Proposition 3.9** *Typed terms satisfy the following properties:*

- *Strict composition:* if  $\Gamma \vdash M : \tilde{\rho} \Rightarrow \tilde{\sigma}$  and  $\Gamma \vdash N : \tilde{\sigma} \Rightarrow \tilde{\tau}$  then  $\Gamma \vdash M ; N : \tilde{\rho} \Rightarrow \tilde{\tau}$ .
- *Expansion:* if  $\Gamma \vdash M : \tilde{\rho} \Rightarrow \tilde{\sigma}$  then  $\Gamma \vdash M : \tilde{\rho}\tilde{v} \Rightarrow \tilde{v}\tilde{\sigma}$ .
- *Composition:* if  $\Gamma \vdash M : \sigma$  and  $\Gamma \vdash N : \tau$  and  $\sigma \cdot \tau$  is defined, then  $\Gamma \vdash M ; N : \sigma \cdot \tau$ .
- *Substitution:* if  $\Gamma \vdash M : \sigma$  and  $\Gamma, x : \sigma \vdash N : \tau$  then  $\Gamma \vdash \{M/x\}N : \tau$ .
- *Subject reduction:* if  $\Gamma \vdash N : \tau$  and  $N \rightarrow M$  then  $\Gamma \vdash M : \tau$ .

### 3.6 Machine termination

A remarkable aspect of the type system is how it gives a direct connection with termination of the machine. To expose this, we formalize the intuitive meaning of types as describing the initial and final stack of a run of the machine.

**Definition 3.10** The set  $\text{RUN}(\tilde{\sigma} \Rightarrow \tilde{\tau})$  is the set of terms  $M$  such that for any stack  $S \in \text{RUN}(\tilde{\sigma})$  there is a stack  $T \in \text{RUN}(\tilde{\tau})$  and a run of the machine  $(S, M) \Downarrow (T, \star)$ , where  $\text{RUN}(\tau_1 \dots \tau_n)$  is the set of stacks  $\varepsilon \cdot N_1 \dots N_n$  such that  $N_i \in \text{RUN}(\tau_i)$ .

Note that a successful run requires the term to be closed, so a set  $\text{RUN}(\tau)$  contains only closed terms. The following lemma shows that  $\text{RUN}(\tau)$  is always inhabited by the term  $\perp_\tau$  (see Remark 3.7).

**Lemma 3.11** *For any type  $\tau$  the set  $\text{RUN}(\tau)$  is inhabited:  $\perp_\tau \in \text{RUN}(\tau)$ .*

If  $M : \tau$  implies  $M \in \text{RUN}(\tau)$ , then a type derivation *is* a termination proof of the machine. This is Theorem 3.13, and proving it gives a concrete *Tait-style* reducibility proof [34], where  $\text{RUN}(\tau)$  takes the rôle of the reducibility set for  $\tau$ . By using the properties of machine runs the proof is then a simple, direct induction on type derivations.

Vector notation is extended to variables,  $\vec{x} = x_1 \dots x_n$ , to contexts as  $\vec{x} : \vec{\tau} = x_1 : \tau_1, \dots, x_n : \tau_n = \Gamma$ , and to simultaneous substitutions: if  $S = \varepsilon \cdot M_1 \dots M_n$  then  $\{S/\vec{x}\} = \{M_1/x_1, \dots, M_n/x_n\}$ . We write concatenation of stacks  $S$  and  $T$  by juxtaposition,  $ST$ .

**Lemma 3.12** *If  $\vec{w} : \vec{\omega} \vdash M : \tau$  then for any  $W \in \text{RUN}(\vec{\omega})$ ,  $\{W/\vec{w}\}M \in \text{RUN}(\tau)$ .*

**Proof.** By induction on the type derivation. In each case, let  $\Gamma = \vec{w} : \vec{\omega}$ , let  $W$  be a stack in  $\text{RUN}(\vec{\omega})$ , and let  $M' = \{W/\vec{w}\}M$ .

- If the derivation is a  $\star$ -rule for  $\Gamma \vdash \star : \vec{\tau} \Rightarrow \vec{\tau}$  then there is a trivial zero-step run  $(T, \star) \Downarrow (T, \star)$ .
- If the derivation ends in the variable rule below left, then for any  $N \in \text{RUN}(\vec{\rho} \Rightarrow \vec{\sigma})$ , the inductive hypothesis gives a run for  $\{N/x\}M'$  from any  $TS \in \text{RUN}(\vec{\tau} \Rightarrow \vec{\sigma})$  to some  $U \in \text{RUN}(\vec{v})$  (second item below). For  $N$  there is a run from any  $R \in \text{RUN}(\vec{\rho})$  to some  $S \in \text{RUN}(\vec{\sigma})$  (third item below). These runs compose into one for  $\{W/\vec{w}, N/x\}x.M = N; \{N/x\}M'$  as below right, expanding the stack on the run for  $N$  by  $T$ . Note that we may assume  $x \notin \text{fv}(W)$  (otherwise we rename  $x$ ).

$$\frac{\Gamma, x : \vec{\rho} \Rightarrow \vec{\sigma} \vdash M : \vec{\sigma} \Rightarrow \vec{\tau} \Rightarrow \vec{v} \quad \text{var}}{\Gamma, x : \vec{\rho} \Rightarrow \vec{\sigma} \vdash x.M : \vec{\rho} \Rightarrow \vec{\tau} \Rightarrow \vec{v}} \quad \frac{(TS, \{N/x\}M')}{(U, \star)} \quad \frac{(R, N)}{(S, \star)} \quad \frac{((TR, N; \{N/x\}M'))}{((TS, \{N/x\}M'))} \quad \frac{((TS, \{N/x\}M'))}{(U, \star)}$$

- If the derivation ends in the application rule below left, then by the inductive hypothesis for  $N$  we have  $N' = \{W/\vec{w}\}N \in \text{RUN}(\rho)$ , and for  $M$  we have a run from  $M'$  and any stack  $S \in \text{RUN}(\vec{\sigma})$  with  $N'$  added on top, to some  $T \in \text{RUN}(\vec{\tau})$ . This gives the run for  $\{W/\vec{w}\}[N].M = [N'].M'$  below right.

$$\frac{\Gamma \vdash N : \rho \quad \Gamma \vdash M : \rho \Rightarrow \vec{\sigma} \Rightarrow \vec{\tau}}{\Gamma \vdash [N].M : \vec{\sigma} \Rightarrow \vec{\tau}} \text{ app} \quad \frac{(S, [N'].M')}{(S \cdot N', M')} \quad \frac{(S \cdot N', M')}{(T, \star)}$$

- If the derivation ends in the abstraction rule below left, then for any  $N \in \text{RUN}(\rho)$  and  $S \in \text{RUN}(\vec{\sigma})$  the inductive hypothesis gives a run for  $\{N/x\}M'$  to some  $T \in \text{RUN}(\vec{\tau})$ . This gives the run for  $\{W/\vec{w}\}\langle x \rangle.M = \langle x \rangle.M'$  below right.

$$\frac{\Gamma, x : \rho \vdash M : \vec{\sigma} \Rightarrow \vec{\tau}}{\Gamma \vdash \langle x \rangle.M : \rho \Rightarrow \vec{\sigma} \Rightarrow \vec{\tau}} \text{ abs} \quad \frac{(S \cdot N, \langle x \rangle.M')}{(S, \{N/x\}M')} \quad \frac{(S, \{N/x\}M')}{(T, \star)}$$

□

The following theorem is then immediate.

**Theorem 3.13** *For a typed, closed term  $M : \vec{\sigma} \Rightarrow \vec{\tau}$  and stack  $S : \vec{\sigma}$  the machine terminates.*

## 4 The functional machine calculus

The combination of both generalizations, *locations* and *sequencing*, gives the Functional Machine Calculus.

**Definition 4.1** The *Functional Machine Calculus* (FMC) is given by the below grammar, with from left to right the constructors *nil*, a (sequential) *variable*, an *application* or *push action* on the location  $a$ , and an *abstraction* or *pop action* on  $a$  which binds  $x$  in  $M$ . Terms are considered modulo  $\alpha$ -equivalence.

$$M, N, P ::= \star \mid x.M \mid [N]a.M \mid a\langle x \rangle.M$$

Composition  $N; M$  and substitution  $\{N/x\}M$  are as for the sequential  $\lambda$ -calculus. The machine is as for the poly- $\lambda$ -calculus: a *state* is a pair  $(S_A, M)$  of a memory and a term, and the *transitions* are:

$$\frac{(\ S_A ; S_a \ , \ [N]a.M \ )}{(\ S_A ; S_a \cdot N \ , \ M \ )} \quad \frac{(\ S_A ; S_a \cdot N \ , \ a\langle x \rangle.M \ )}{(\ S_A ; S_a \ , \ \{N/x\}M \ )}$$

Beta-reduction is as for the poly- $\lambda$ -calculus: a redex consists of a successive application and abstraction *on the same location*, separated by any number of actions on other locations. We will now make this formal.

*Head contexts*  $H$  are defined as below left. The term obtained by replacing the hole  $\{\}$  in  $H$  with  $M$  is denoted  $H.M$ , where a binder  $a\langle x \rangle$  in  $H$  captures in  $M$ . The *binding variables*  $\text{bv}(H)$  of  $H$  are those variables  $x$  where  $H$  is constructed over  $a\langle x \rangle$ . The set of *locations* used in a term or context is denoted  $\text{loc}(M)$  respectively  $\text{loc}(H)$ . Then *Beta-reduction* is defined by the rewrite rule schema below right, where  $a \notin \text{loc}(H)$  and  $\text{bv}(H) \cap \text{fv}(N) = \emptyset$ , and is closed under all contexts.

$$H ::= \{\} \mid [M]a.H \mid a\langle x \rangle.H \quad [N]a.H.a\langle x \rangle.M \rightarrow H.\{N/x\}M$$

We will first consider the untyped calculus. We give the  $\text{cbv}$  encoding of effects and provide an intuition for programming in the FMC, then establish confluence and connect machine evaluation to  $\beta$ -reduction. We then consider simple types. Constants will be used informally, in examples.

#### 4.1 Call-by-value with effects

We extend the encoding  $(-)_v$  of the computational  $\lambda$ -calculus of Section 3.3 to effects as follows. (The case for  $N \oplus M$  is the same as for  $\text{cbn}$ , as it expects a Church boolean for  $x$ .)

$$\begin{aligned} \text{read}_v &= \text{in}\langle x \rangle.[x] \\ (\text{write } N; M)_v &= N_v; \langle x \rangle.[x]\text{out}.M_v \\ (c := N; M)_v &= N_v; \langle x \rangle.c\langle \_ \rangle.[x]c.M_v \\ !c_v &= c\langle x \rangle.[x]c.[x] \\ (N \oplus M)_v &= \text{rnd}\langle x \rangle.[M_v].[N_v].x \end{aligned}$$

**Example 4.2** The term from Example 2.2, below, is given a  $\text{cbv}$  interpretation as follows.

$$a := 2; (\lambda x. !a) (a := 3; 0) \rightarrow_{\text{cbv}} 3$$

Integers are values, and the translation will use  $i_v = [i]$ . An update with an integer then simplifies by:

$$(a := i; M)_v = [i].\langle x \rangle.a\langle \_ \rangle.[x]a; M \rightarrow a\langle \_ \rangle.[i]a; M$$

The  $\text{cbv}$ -translated term, after applying this reduction to the two updates, further reduces as follows.

$$a\langle \_ \rangle.[2]a.a\langle \_ \rangle.[3]a.[0].[\langle x \rangle.a\langle y \rangle.[y]a.[y]].\langle z \rangle.z \rightarrow a\langle \_ \rangle.[3]a.[3]$$

**Example 4.3** The term from Example 2.7,

$$a := (\lambda x. b := 1; x) 0; !b \rightarrow_{\text{cbv}} 1$$

translates and reduces as follows (with the same shortcut for update as above).

$$\begin{aligned}
& [0]. [\langle x \rangle. b \langle - \rangle. [1]b. [x]. \langle y \rangle. y. \langle z \rangle. a \langle - \rangle. [z]a. b \langle u \rangle. [u]b. [u]] \\
& \rightarrow [0]. \langle x \rangle. b \langle - \rangle. [1]b. [x]. \langle z \rangle. a \langle - \rangle. [z]a. b \langle u \rangle. [u]b. [u] \\
& \rightarrow b \langle - \rangle. [1]b. a \langle - \rangle. [0]a. b \langle u \rangle. [u]b. [u] \\
& \rightarrow b \langle - \rangle. a \langle - \rangle. [0]a. [1]b. [1] \\
& \sim a := 0; b := 1; 1
\end{aligned}$$

## 4.2 Programming in the FMC

As in the sequential  $\lambda$ -calculus, programming in the FMC naturally follows the concatenative paradigm. A term  $M$  is viewed as a function taking an input memory  $S_A$  to an output memory  $T_A$  by a run of the machine  $(S_A, M) \Downarrow (T_A, \star)$ . Functions standardly operate on the main stack  $\lambda$ , and it is then natural to consider effect operators that transfer values between the main stack and other locations, as the cbv translations of effect operators do. We introduce the following operations for input and output, a random generator, and a memory cell  $c$ . We further add *definitions* or *let*, as a redex.

$$\begin{aligned}
\text{print} & \triangleq \langle x \rangle. [x]\text{out} & \text{get } c & \triangleq c \langle x \rangle. [x]c. [x] \\
\text{read} & \triangleq \text{in} \langle x \rangle. [x] & \text{set } c & \triangleq \langle x \rangle. c \langle - \rangle. [x]c \\
\text{rand} & \triangleq \text{rnd} \langle x \rangle. [x] & (x=N); M & \triangleq [N]. \langle x \rangle. M
\end{aligned}$$

Constant operations such as the conditional **if** pop the required number of items from the main stack, and reinstate their result, as is standard for stack languages. For example:

$$\frac{(\ S_A; S_\lambda \cdot 2 \cdot 3 \ , \ +. M )}{(\ S_A; S_\lambda \cdot 5 \ , \ M )} \quad \frac{(\ S_A; S_\lambda \cdot P \cdot N \cdot \perp \ , \ \text{if. } M )}{(\ S_A; S_\lambda \cdot P \ , \ M )}$$

The FMC then operates similarly to a stack calculus for arithmetic: an expression  $1 + ((2 + 3) \times 4)$  is given as a term  $[4]. [3]. [2]. +. \times. [1]. +$  which indeed returns 21. This results in an imperative programming style similar to Haskell's *do*-notation, with the difference that terms may have *any* number of return values, and consume any number of previously returned values.

**Example 4.4** Consider the following example, where **rnd** is taken to randomly sample natural numbers.

$$(f = \text{rand} ; \text{set } c ; \text{get } c) ; f ; f ; + ; \text{print}$$

The term assigns  $f$  to be the function that draws a random number, stores it in cell  $c$ , and reads the value at  $c$  again as its return value. It then executes  $f$  twice, sums the results, and sends that to output. The overall actions should be to take two random inputs  $i$  and  $j$ , to update the cell  $c$  with the last value  $j$ , and to output  $i + j$ . In Figure 1 the term is first interpreted as an FMC-term and reduced to normal form, where each line is a beta-step, and then evaluated on the machine, where the initial memory provides the two expected inputs on **rnd** and one on  $c$ . (For compactness we give locations as a header and show only necessary stack elements.)

## 4.3 Confluence

In demonstrating confluence for the  $\lambda$ -calculus, the difficulty is reduction inside an argument: duplicating or deleting it creates converging reductions of different length. By contrast, *spine reduction*, which reduces in every context *except* argument position, is *diamond* (peaks converge in one step).

While the two extensions of the FMC, *locations* and *sequencing*, do create new configurations of overlapping redexes, the situation is fundamentally the same. Spine reduction, defined analogously as reduction in every context except argument position, is diamond, and the remaining problem is the same as for the

$$\begin{array}{l}
 [\text{rnd}\langle x \rangle. [\underline{x}]. \langle y \rangle. \underline{c}(-). [y]c. c(z). [z]c. [z]]. \langle f \rangle. f. f. + . \langle p \rangle. [p]out \\
 [\text{rnd}\langle x \rangle. \underline{c}(-). [\underline{x}]c. c(z). [z]c. [z]]. \langle f \rangle. f. f. + . \langle p \rangle. [p]out \\
 [\underline{\text{rnd}\langle x \rangle. \underline{c}(-). [\underline{x}]c. [x]}. \langle f \rangle. f. f. + . \langle p \rangle. [p]out} \\
 \text{rnd}\langle x \rangle. \underline{c}(-). [\underline{x}]c. [x] . \text{rnd}\langle y \rangle. \underline{c}(-). [y]c. [y] . + . \langle p \rangle. [p]out \\
 \text{rnd}\langle x \rangle. \underline{c}(-). [x] . \text{rnd}\langle y \rangle. [y]c. [y] . + . \langle p \rangle. [p]out
 \end{array}$$
  

out	rnd	c	$\lambda$
(	6 · 7	*	, $\text{rnd}\langle x \rangle. \underline{c}(-). [x]. \text{rnd}\langle y \rangle. [y]c. [y]. + . \langle p \rangle. [p]out$ )
(	6	*	, $\underline{c}(-). [7]. \text{rnd}\langle y \rangle. [y]c. [y]. + . \langle p \rangle. [p]out$ )
(	6		, $[7]. \text{rnd}\langle y \rangle. [y]c. [y]. + . \langle p \rangle. [p]out$ )
(	6		, $\text{rnd}\langle y \rangle. [y]c. [y]. + . \langle p \rangle. [p]out$ )
(			, $[6]c. [6]. + . \langle p \rangle. [p]out$ )
(		6 ; 7	, $[6]. + . \langle p \rangle. [p]out$ )
(		6 ; 7 · 6	, $+ . \langle p \rangle. [p]out$ )
(		6 ; 13	, $\langle p \rangle. [p]out$ )
(		6 ;	, $[13]out$ )
(	13 ;	6 ;	, $\star$ )

Fig. 1. Reduction followed by machine evaluation of the term in Example 4.4

$\lambda$ -calculus. The calculus thus remains confluent, which will be proved by the standard *parallel reduction* technique of Tait and Martin-Löf (see [2]) and Takahashi [35].

The new configurations are the following. *Sequencing* introduces terms of the form  $N ; x. M$ , with reduction in  $M$  or in  $N$ , where the latter may induce substitutions in  $x. M$ . But reduction in  $N$  may not duplicate a redex in  $M$ , and vice versa, so a peak of this kind converges immediately.

*Locations* create two new overlapping configurations,

$$\begin{array}{ll}
 \text{nested:} & [N]a. [P]b. b\langle y \rangle. a\langle x \rangle. M \\
 \text{interleaved:} & [N]a. [P]b. a\langle x \rangle. b\langle y \rangle. M
 \end{array}$$

but both resolve immediately: in each case the two reducts will converge in one step to  $\{N/x\}\{P/y\}M$ .

Because of this, the problem of confluence amounts to the problem of reduction in arguments, which may be duplicated or deleted, as it does in the regular  $\lambda$ -calculus. We formalize this observation in the following proposition (which is independent of the confluence result). *Spine reduction* is given by closing the  $\beta$ -step under all contexts except in argument position. That is, reduction in each of  $x. M$ ,  $[N]a. M$ , and  $\langle x \rangle. M$  may take place in  $M$ , but not in  $N$ .

**Proposition 4.5** *Spine reduction is diamond.*

The full confluence proof is a standard application of parallel reduction, and is given in Appendix A.

**Theorem 4.6** *Reduction  $\rightarrow$  is confluent.*

#### 4.4 Simple types for the FMC

As with sequential types, FMC types will represent the input/output behaviour of the machine. Since a memory is a family of stacks, types will use families of type vectors.

**Definition 4.7** *FMC-types*  $\rho, \sigma, \tau, \nu$  are given by:

$$\tau ::= \vec{\sigma}_A \Rightarrow \vec{\tau}_A \quad \vec{\tau}_A ::= \{\vec{\tau}_a \mid a \in A\} \quad \vec{\tau} ::= \tau_1 \dots \tau_n$$

$$\begin{array}{c}
\frac{}{\Gamma \vdash \star : \vec{\tau}_A \Rightarrow \vec{\tau}_A} \star \\
\frac{\Gamma \vdash N : \rho \quad \Gamma \vdash M : a(\rho) \vec{\sigma}_A \Rightarrow \vec{\tau}_A}{\Gamma \vdash [N]a. M : \vec{\sigma}_A \Rightarrow \vec{\tau}_A} \text{app} \\
\frac{\Gamma, x : \vec{\rho}_A \Rightarrow \vec{\sigma}_A \vdash M : \vec{\sigma}_A \vec{\tau}_A \Rightarrow \vec{v}_A}{\Gamma, x : \vec{\rho}_A \Rightarrow \vec{\sigma}_A \vdash x. M : \vec{\rho}_A \vec{\tau}_A \Rightarrow \vec{v}_A} \text{var} \\
\frac{\Gamma, x : \rho \vdash M : \vec{\sigma}_A \Rightarrow \vec{\tau}_A}{\Gamma \vdash a\langle x \rangle. M : a(\rho) \vec{\sigma}_A \Rightarrow \vec{\tau}_A} \text{abs}
\end{array}$$

Fig. 2. Typing rules for the Functional Machine Calculus

The typing rules for the FMC are in Figure 2. They use the following notation. *Concatenation* of two vector families is pointwise,  $\vec{\sigma}_A \vec{\tau}_A = \{\vec{\sigma}_a \vec{\tau}_a \mid a \in A\}$ , and the empty vector is denoted  $\varepsilon$ . A *slice*  $\vec{\tau}_A|_a$  of a family  $\vec{\tau}_A = \{\vec{\tau}_a \mid a \in A\}$  is the vector  $\vec{\tau}_a$ , and a slice  $\tau_a$  of a type  $\tau = \vec{\rho}_A \Rightarrow \vec{\sigma}_A$  is the type  $\vec{\rho}|_a \Rightarrow \vec{\sigma}|_a$  (i.e.  $\tau_a$  restricts  $\tau$  to a single location  $a$ ). *Composition* is slice-wise:  $\sigma \cdot \tau = \{\sigma|_a \cdot \tau|_a \mid a \in A\}$ . Finally, a *singleton*  $a(\vec{\tau})$  is a type vector at a single location, with all other locations empty, defined by  $a(\vec{\tau})|_a = \vec{\tau}$  and  $a(\vec{\tau})|_b = \varepsilon$  for  $a \neq b$ . A singleton  $\lambda(\vec{\tau})$  on the main location  $\lambda$  may be written as  $\vec{\tau}$ .

Poly-types embed as types of the form  $\vec{\tau}_A \Rightarrow$  by the definitions  $o \triangleq (\Rightarrow)$  and  $a(\sigma) \rightarrow (\vec{\tau}_A \Rightarrow) \triangleq a(\sigma) \vec{\tau}_A \Rightarrow$ , and sequential types embed directly as types over only the location  $\lambda$ . The properties proved of sequential types in Section 3 carry over straightforwardly: strict composition, expansion, composition, substitution, subject reduction, and termination of the machine.

**Example 4.8** The *singleton* construct  $a(\vec{\tau})$  gives a natural way of writing types in practice. The term from Example 4.4 may be typed as follows, where  $\mathbb{Z}$  is a base type of integers.

$$(f = \text{rand} ; \text{set } c ; \text{get } c) ; f ; f ; + ; \text{print} : \text{rnd}(\mathbb{Z} \ \mathbb{Z}) \ c(\mathbb{Z}) \Rightarrow c(\mathbb{Z}) \ \text{out}(\mathbb{Z})$$

The type expresses that the term pops two integers from `rnd` and one from `c`, and pushes one integer to `c` and one to `out`. (Note that there are other ways of writing the same type, since singleton types on different locations may permute.) Below left are the types of the defined subterms, and the full type is built up below right by composing these.

$$\begin{array}{ll}
+ : \mathbb{Z} \ \mathbb{Z} \Rightarrow \mathbb{Z} & \text{rand} ; \text{set } c : \text{rnd}(\mathbb{Z}) \ c(\mathbb{Z}) \Rightarrow c(\mathbb{Z}) \\
\text{rand} = \text{rnd}\langle x \rangle. [x] : \text{rnd}(\mathbb{Z}) \Rightarrow \mathbb{Z} & \text{rand} ; \text{set } c ; \text{get } c : \text{rnd}(\mathbb{Z}) \ c(\mathbb{Z}) \Rightarrow c(\mathbb{Z}) \ \mathbb{Z} \\
\text{print} = \langle x \rangle. [x] \text{out} : \mathbb{Z} \Rightarrow \text{out}(\mathbb{Z}) & (f = \text{rand} ; \text{set } c ; \text{get } c) : (\Rightarrow) \\
\text{set } c = \langle x \rangle. c\langle - \rangle. [x]c : \mathbb{Z} \ c(\mathbb{Z}) \Rightarrow c(\mathbb{Z}) & (f = \dots) ; f ; f : \text{rnd}(\mathbb{Z} \ \mathbb{Z}) \ c(\mathbb{Z}) \Rightarrow c(\mathbb{Z}) \ \mathbb{Z} \ \mathbb{Z} \\
\text{get } c = c\langle x \rangle. [x]c. [x] : c(\mathbb{Z}) \Rightarrow c(\mathbb{Z}) \ \mathbb{Z} & (f = \dots) ; f ; f ; + ; \text{print} : \text{rnd}(\mathbb{Z} \ \mathbb{Z}) \ c(\mathbb{Z}) \Rightarrow c(\mathbb{Z}) \ \text{out}(\mathbb{Z})
\end{array}$$

The third term is only a redex, which pushes and then pops a value, giving it an empty type  $(\Rightarrow)$ . The type of the fourth term is that of  $f ; f$ . Separating the self-composition of the type of  $f$  for each location exposes how the inputs on `rnd` and outputs on  $\lambda$  accumulate, while the output and input on `c` interact:

$$\begin{array}{lll}
\text{rnd}(\mathbb{Z}) \Rightarrow & \cdot & \text{rnd}(\mathbb{Z}) \Rightarrow & = & \text{rnd}(\mathbb{Z} \ \mathbb{Z}) \Rightarrow \\
c(\mathbb{Z}) \Rightarrow c(\mathbb{Z}) & \cdot & c(\mathbb{Z}) \Rightarrow c(\mathbb{Z}) & = & c(\mathbb{Z}) \Rightarrow c(\mathbb{Z}) \\
\Rightarrow \mathbb{Z} & \cdot & \Rightarrow \mathbb{Z} & = & \Rightarrow \mathbb{Z} \ \mathbb{Z}
\end{array}$$

By way of illustration, in Haskell the same example may be written as follows.

```

example :: RandT StdGen (StateT Int IO) ()
example = do
  let f = do x <- rand; lift (put x); lift get
        y <- f
        z <- f
        lift (lift (print (y+z)))

```

To combine the effects of I/O, state, and random generation, the Haskell example uses a stack of monad transformers. The effects are then layered in a fixed order, and to access each effect the function `lift` is used to move to the next layer. Thus, the `rand` action requires no lifting since the random generator is at the top of the transformer stack; the state actions `put` and `get` must be lifted once; and the `print` action must be lifted twice, since I/O is at the bottom of the transformer stack.

The example highlights the following differences between monad transformers and the FMC. Firstly, reader/writer effects in the FMC combine seamlessly, without requiring their organisation in a stack, and are accessed by their locations instead of a lifting function. Secondly, sequential composition in the FMC is a basic operation, with a close connection between the syntax and its execution, where Haskell `do`-notation is syntactic sugar for more involved monadic operations. Thirdly, the FMC type system accounts for the effectful operations that a term performs, where monad transformer types indicate which effects may be present, but not which operations they perform.

Of course, where monads are universal, the FMC is presently restricted to reader/writer effects. How to broaden the range of effects covered in the FMC is an important direction for future work.

## 5 Further work

We have given an exploratory overview of the Functional Machine Calculus with the most essential results: the natural capture of algebraic laws for effects by reductions and permutations; the encoding of related formalisms for controlling execution behaviour such as monadic constructs, `cbpv`,  $\kappa$ -calculus, and Arrows; confluence; and termination of the machine with simple types. A forthcoming paper will strengthen these results with domain-theoretic and categorical semantics, and strong normalization with simple types.

Present and future work aims to extend the FMC beyond reader/writer effects and with standard features. One direction is to include sum types, datatypes, and error handling, where it looks possible to capture all three in a uniform way. A second direction is to introduce parallel composition and explore the relation with process calculi, where our type system promises to give something closely related to session types. A third direction is local mutable store, by introducing a *new* construct for locations, and generalizing locations to *regions* to capture mutable data structures (arrays, graphs), which then leads naturally into an exploration of dependent types for the FMC. A fourth direction is to explore the connection with string diagrams, and to introduce constructs to capture diagrammatic reasoning, for example interaction nets or quantum diagrammatic systems.

An important theoretical challenge is type inference. This appears to be open also for the sequential  $\lambda$ -calculus: existing algorithms for concatenative languages are limited [33,6], and (we believe) unable to find the type for the self-application in Example 3.6.

### Acknowledgements

I am deeply grateful to Chris Barrett and Guy McCusker for valuable contributions to both presentation and content of this work. I would further like to thank the following people for constructive discussions and pointers to the literature: Ugo Dal Lago, Giulio Guerrieri, Paul Levy, John Power, Alex Simpson, and Vincent van Oostrom. This work was supported by EPSRC Project EP/R029121/1 *Typed lambda-calculi with sharing and unsharing*.

## References

- [1] Ahman, D. and S. Staton, *Normalization by evaluation and algebraic effects*, Electronic Notes in Theoretical Computer Science **298** (2013), pp. 51–69. DOI:10.1016/j.entcs.2013.09.007
- [2] Barendregt, H. P., “The Lambda Calculus – Its Syntax and Semantics,” Studies in Logic and the Foundations of Mathematics **103**, North-Holland, 1984.
- [3] Benton, N. and P. Wadler, *Linear logic, monads and the lambda calculus*, in: *Proceedings of the 11th Annual IEEE Symposium on Logic in Computer Science (LICS)* (1996), pp. 420–431. DOI:10.1109/LICS.1996.561458
- [4] Dal Lago, U., G. Guerrieri and W. Heijltjes, *Decomposing probabilistic lambda-calculi*, in: J. Goubault-Larrecq and B. König, editors, *Foundations of Software Science and Computation Structures*, LNCS **12077** (2020), pp. 136–156. DOI:10.1007/978-3-030-45231-5\_8



- [5] de Bruijn, N. G., *Algorithmic definition of lambda-typed lambda calculus*, in: G. Huet and G. Plotkin, editors, *Logical Environments*, Cambridge University Press, 1993 pp. 131–146.
- [6] Diggins, C., *Simple type inference for higher-order stack-oriented languages*, Technical Report Cat-TR-2008-001 (2008).
- [7] Douence, R. and P. Fradet, *A systematic study of functional language implementations*, ACM Transactions on Programming Languages and Systems **20** (1998), pp. 344–387. DOI:10.1145/276393.276397
- [8] Egger, J., R. E. Møgelberg and A. Simpson, *The enriched effect calculus: syntax and semantics*, Journal of Logic and Computation **24** (2014), pp. 615–654. DOI:10.1093/logcom/exs025
- [9] Ehrhard, T. and G. Guerrieri, *The bang calculus: An untyped lambda-calculus generalizing call-by-name and call-by-value*, in: *Proceedings of the 18th International Symposium on Principles and Practice of Declarative Programming (PPDP'16)*, 2016, pp. 174–187. DOI:10.1145/2967973.2968608
- [10] Filinski, A., “Controlling Effects,” Ph.D. thesis, School of Computer Science, Carnegie Mellon University (1996).
- [11] Hasegawa, M., *Decomposing typed lambda-calculus into a couple of categorical programming languages*, in: *Category Theory and Computer Science (CTCS 1995)*, Lecture Notes in Computer Science **953** (1995). DOI:10.1007/3-540-60164-3\_28
- [12] Herzberg, D. and T. Reichert, *Concatenative programming: An overlooked paradigm in functional programming*, in: *4th International Conference on Software and Data Technologies (ICSOFIT)*, 2009.
- [13] Hirschhoff, D., E. Prebet and D. Sangiorgi, *On the representation of references in the pi-calculus*, in: *31st International Conference on Concurrency Theory (CONCUR)*, LIPICs (2020), pp. 34:1–34:20. DOI:10.4230/LIPICs.CONCUR.2020.34
- [14] Hughes, J., *Generalising monads to arrows*, Science of Computer Programming **37** (2000), pp. 67–111. DOI:10.1016/S0167-6423(99)00023-4
- [15] Krivine, J.-L., *A call-by-name lambda-calculus machine*, Higher-Order and Symbolic Computation **20** (2007), pp. 199–207. DOI:10.1007/s10990-007-9018-9
- [16] Landin, P. J., *The mechanical evaluation of expressions*, The Computer Journal **6** (1964), pp. 308–320. DOI:10.1093/comjnl/6.4.308
- [17] Landin, P. J., *Correspondence between Algol 60 and Church’s lambda-notation: Part I*, Communications of the ACM **8** (1965), pp. 89–101. DOI:10.1145/363744.363749
- [18] Levy, P. B., “Call-by-push-value: A functional/imperative synthesis,” *Semantic Structures in Computation* **2**, Springer Netherlands, 2003.
- [19] Levy, P. B., *Call-by-push-value: Decomposing call-by-value and call-by-name*, Higher-Order and Symbolic Computation **19** (2006), pp. 377–414. DOI:10.1007/s10990-006-0480-6
- [20] Levy, P. B., J. Power and H. Thielecke, *Modelling environments in call-by-value programming languages*, Information and Computation **185** (2003), pp. 182–210. DOI:10.1016/S0890-5401(03)00088-9
- [21] Lynas, A. R. and B. Stoddart, *Adding lambda expressions to Forth*, in: *22nd EuroForth Conference*, 2006.
- [22] Maraist, J., M. Odersky, D. N. Turner and P. Wadler, *Call-by-name, call-by-value, call-by-need, and the linear lambda calculus*, Theoretical Computer Science **228** (1999), pp. 175–210. DOI:10.1016/S0304-3975(98)00358-2
- [23] Milner, R., *Functions as processes*, Mathematical Structures in Computer Science **2** (1992), pp. 119–141. DOI:10.1017/S0960129500001407
- [24] Moggi, E., *Computational lambda-calculus and monads*, in: *Fourth Annual Symposium on Logic in Computer Science* (1989), pp. 14–23. DOI:10.1109/LICS.1989.39155
- [25] Moggi, E., *Notions of computation and monads*, Information and Computation **93** (1991), pp. 55–92. DOI:10.1016/0890-5401(91)90052-4
- [26] Pestov, S., D. Ehrenberg and J. Groff, *Factor: A dynamic stack-based programming language*, ACM SIGPLAN Notices **45** (2010), pp. 43–58. DOI:10.1145/1899661.1869637
- [27] Plotkin, G. and J. Power, *Notions of computation determine monads*, in: *International Conference on Foundations of Software Science and Computation Structures (FoSSaCS)*, Lecture Notes in Computer Science **2303** (2002), pp. 342–356. DOI:10.1007/3-540-45931-6\_24
- [28] Plotkin, G. D., *Call-by-name, call-by-value and the  $\lambda$ -calculus*, Theoretical Computer Science **1** (1975), pp. 125–159. DOI:10.1016/0304-3975(75)90017-1
- [29] Plotkin, G. D. and M. Pretnar, *Handlers of algebraic effects*, in: *European Symposium on Programming (ESOP)*, Lecture Notes in Computer Science **5502** (2009), pp. 80–94. DOI:10.1007/978-3-642-00590-9\_7

- [30] Power, A. and H. Thielecke, *Closed Freyd- and  $\kappa$ -categories*, in: *International Colloquium on Automata, Languages, and Programming (ICALP)*, Lecture Notes in Computer Science **1644**, Springer, 1999, pp. 625–634. DOI:[10.1007/3-540-48523-6\\_59](https://doi.org/10.1007/3-540-48523-6_59)
- [31] Power, J. and E. Robinson, *Premonoidal categories and notions of computation*, Mathematical Structures in Computer Science **7** (1997), pp. 453–468. DOI:[10.1017/S0960129597002375](https://doi.org/10.1017/S0960129597002375)
- [32] Smetsers, J., E. Barendsen, M. v. Eekelen and M. Plasmeijer, *Guaranteeing safe destructive updates through a type system with uniqueness information for graphs*, in: *Workshop Graph Transformations in Computer Science*, Lecture Notes in Computer Science **776**, Schloss Dagstuhl (1993). DOI:[10.1007/3-540-57787-4\\_23](https://doi.org/10.1007/3-540-57787-4_23)
- [33] Stoddart, B. and P. J. Knaggs, *Type inference in stack based languages*, Formal Aspects of Computing **5** (1993), pp. 289–298. DOI:[10.1007/BF01212404](https://doi.org/10.1007/BF01212404)
- [34] Tait, W., *Intensional interpretations of functionals of finite type I*, The Journal of Symbolic Logic **32** (1967), pp. 198–212. DOI:[10.2307/2271658](https://doi.org/10.2307/2271658)
- [35] Takahashi, M., *Parallel reductions in lambda-calculus*, Information and Computation **118** (1995), pp. 120–127. DOI:[10.1006/inco.1995.1057](https://doi.org/10.1006/inco.1995.1057)
- [36] von Thun, M., *Joy: Forth’s functional cousin*, in: *Proc. 17th EuroForth Conference*, 2001.
- [37] Winskel, G., “The formal semantics of programming languages: An introduction,” MIT Press, Cambridge, Massachusetts, 1993.

## A The confluence proof

This section gives a complete proof of confluence using parallel reduction. The particular approach is to extend the syntax of terms with marked redexes, to define parallel reduction and to identify residuals.

**Definition A.1** A *marked redex* is one whose application and abstraction are marked by the symbols  $(\blacktriangleleft)$  and  $(\blacktriangleright)$ , as follows.

$$[N]_{a\blacktriangleleft} H.\blacktriangleright a\langle x \rangle.M$$

A *redex-marked* or *marked* term  $M^\blacktriangleright$  is one where a selection of redexes is marked. For a marked term consisting of a head context and a term  $(H.M)^\blacktriangleright$ , the marking of both components separately is indicated by  $H^\blacktriangleright$  and  $M^\blacktriangleright$ . The *marked reduct*  $(M^\blacktriangleright)_\blacktriangleleft$  of a marked term  $M^\blacktriangleright$  is defined as follows.

$$\begin{aligned} (\star)_\blacktriangleleft &= \star & ([N^\blacktriangleright]_{a\blacktriangleleft} M^\blacktriangleright)_\blacktriangleleft &= [(N^\blacktriangleright)_\blacktriangleleft]_{a\blacktriangleleft} (M^\blacktriangleright)_\blacktriangleleft \\ (x.M^\blacktriangleright)_\blacktriangleleft &= x.(M^\blacktriangleright)_\blacktriangleleft & (a\langle x \rangle.M^\blacktriangleright)_\blacktriangleleft &= a\langle x \rangle.(M^\blacktriangleright)_\blacktriangleleft \\ & & ([N^\blacktriangleright]_{a\blacktriangleleft} H^\blacktriangleright.\blacktriangleright a\langle x \rangle.M^\blacktriangleright)_\blacktriangleleft &= \{(N^\blacktriangleright)_\blacktriangleleft/x\}(H^\blacktriangleright.M^\blacktriangleright)_\blacktriangleleft \end{aligned}$$

A *parallel reduction step*  $M^\blacktriangleright \Rightarrow (M^\blacktriangleright)_\blacktriangleleft$  relates a marked term to its marked reduct, and an unmarked term to all its marked reducts:  $M \Rightarrow (M^\blacktriangleright)_\blacktriangleleft$  for each marking  $M^\blacktriangleright$  of  $M$ .

To reduce clutter, a marked reduct  $(M^\blacktriangleright)_\blacktriangleleft$  may be abbreviated as  $M_\blacktriangleleft$ .

**Proposition A.2** A *parallel reduction step* is a *beta-reduction*:  $M \twoheadrightarrow (M^\blacktriangleright)_\blacktriangleleft$  whenever  $M^\blacktriangleright$  is a marking of  $M$ .

**Proof.** By induction on the size of  $M^\blacktriangleright$ .

- Unit case: immediate by  $(\star)_\blacktriangleleft = \star$ .
- Variable case: if  $M^\blacktriangleright \twoheadrightarrow M_\blacktriangleleft$  then  $x.M^\blacktriangleright \twoheadrightarrow x.M_\blacktriangleleft$ .
- Unmarked application case: if  $M^\blacktriangleright \twoheadrightarrow M_\blacktriangleleft$  and  $N^\blacktriangleright \twoheadrightarrow N_\blacktriangleleft$  then  $[N^\blacktriangleright]_{a\blacktriangleleft} M^\blacktriangleright \twoheadrightarrow [N_\blacktriangleleft]_{a\blacktriangleleft} M_\blacktriangleleft$ .
- Unmarked abstraction case: if  $M^\blacktriangleright \twoheadrightarrow M_\blacktriangleleft$  then  $a\langle x \rangle.M^\blacktriangleright \twoheadrightarrow a\langle x \rangle.M_\blacktriangleleft$ .
- Redex case: if  $(H.M)^\blacktriangleright \twoheadrightarrow (H.M)_\blacktriangleleft$  and  $N^\blacktriangleright \twoheadrightarrow N_\blacktriangleleft$  then

$$\begin{aligned} [N^\blacktriangleright]_{a\blacktriangleleft} H^\blacktriangleright.\blacktriangleright a\langle x \rangle.M^\blacktriangleright &\twoheadrightarrow [N_\blacktriangleleft]_{a\blacktriangleleft} H^\blacktriangleright.\blacktriangleright a\langle x \rangle.M^\blacktriangleright \\ &\rightarrow H^\blacktriangleright.\{N_\blacktriangleleft/x\}M^\blacktriangleright \\ &= \{N_\blacktriangleleft/x\}(H^\blacktriangleright.M^\blacktriangleright) \\ &\twoheadrightarrow \{N_\blacktriangleleft/x\}(H.M)_\blacktriangleleft. \end{aligned}$$

where in the second line,  $a \in \text{loc}(H)$  by the definition of a redex, and for the equality on the third line,  $H$  does not bind in  $N$  and  $x$  is not free in  $H$ , by  $\alpha$ -equivalence.  $\square$

To show that parallel reduction is diamond, a term is reduced according to two markings,  $\blacktriangleleft$  and  $\diamond$ . These are then applied simultaneously and interchangeably, and may be considered a single marking  $\blacktriangleleft \diamond = \diamond \blacktriangleleft$ :

$$(M^\blacktriangleright)^\diamond = M^{\blacktriangleright \diamond} = (M^\diamond)^\blacktriangleright$$

If this term is reduced relative to one marking, the other is preserved:  $(M^{\blacktriangleright \diamond})_\diamond = ((M^\diamond)^\blacktriangleright)_\blacktriangleleft$ . The proof then amounts to showing that reducing by each marking is commutative and the same as reducing along both markings simultaneously:

$$(M_\blacktriangleleft)_\diamond = M_{\blacktriangleleft \diamond} = (M_\diamond)_\blacktriangleleft$$

**Lemma A.3** *Parallel reduction commutes with composition*:  $M_\blacktriangleleft ; N_\blacktriangleleft = (M ; N)_\blacktriangleleft$ .

**Proof.** By induction on the size of  $M$ .

- Unit case:  $(\star)_\blacktriangleleft ; N_\blacktriangleleft = \star ; N_\blacktriangleleft = N_\blacktriangleleft = (\star ; N)_\blacktriangleleft$ .

- Variable case: if  $M_\bullet; N_\bullet = (M; N)_\bullet$  then

$$(x.M)_\bullet; N_\bullet = x.(M_\bullet; N_\bullet) = x.(M; N)_\bullet = ((x.M); N)_\bullet.$$

- Unmarked application case: if  $M_\bullet; N_\bullet = (M; N)_\bullet$  then

$$\begin{aligned} ([P]a.M)_\bullet; N_\bullet &= [P_\bullet]a.(M_\bullet; N_\bullet) \\ &= [P_\bullet]a.(M; N)_\bullet \\ &= (([P]a.M); N)_\bullet. \end{aligned}$$

- Unmarked abstraction case: if  $M_\bullet; N_\bullet = (M; N)_\bullet$  then

$$\begin{aligned} (a\langle x \rangle.M)_\bullet; N_\bullet &= a\langle x \rangle.(M_\bullet; N_\bullet) \\ &= a\langle x \rangle.(M; N)_\bullet \\ &= ((a\langle x \rangle.M); N)_\bullet. \end{aligned}$$

- Marked redex case: if  $(H.M)_\bullet; N_\bullet = ((H.M); N)_\bullet$  and  $x \notin \text{fv}(N)$  then

$$\begin{aligned} ([P]a\blacktriangleleft.H.\blacktriangleright a\langle x \rangle.M)_\bullet; N_\bullet &= (\{P_\bullet/x\}(H.M)_\bullet); N_\bullet \\ &= \{P_\bullet/x\}((H.M); N)_\bullet \\ &= (([P]a\blacktriangleleft.H.\blacktriangleright a\langle x \rangle.M); N)_\bullet. \end{aligned}$$

□

**Lemma A.4** *Parallel reduction commutes with substitution:  $\{N_\bullet/x\}M_\bullet = (\{N/x\}M)_\bullet$ .*

**Proof.** By induction on the size of  $M^\bullet$ .

- Unit case: immediate by  $\{N/x\}\star = \star$ .
- Variable case: if  $\{N_\bullet/x\}M_\bullet = (\{N/x\}M)_\bullet$  then by Lemma A.3,

$$\begin{aligned} \{N_\bullet/x\}(x.M)_\bullet &= \{N_\bullet/x\}(x.M_\bullet) \\ &= N_\bullet; \{N_\bullet/x\}M_\bullet \\ &= N_\bullet; (\{N/x\}M)_\bullet \\ &= (N; \{N/x\}M)_\bullet \\ &= (\{N/x\}(x.M))_\bullet. \end{aligned}$$

- Unmarked application case: if  $\{N_\bullet/x\}M_\bullet = (\{N/x\}M)_\bullet$  and  $\{N_\bullet/x\}P_\bullet = (\{N/x\}P)_\bullet$  then

$$\begin{aligned} \{N_\bullet/x\}([P]a.M)_\bullet &= \{N_\bullet/x\}([P_\bullet]a.M_\bullet) \\ &= [\{N_\bullet/x\}P_\bullet]a.\{N_\bullet/x\}M_\bullet \\ &= [(\{N/x\}P)_\bullet]a.(\{N/x\}M)_\bullet \\ &= ((\{N/x\}P)a.\{N/x\}M)_\bullet. \end{aligned}$$

- Unmarked abstraction case: if  $\{N_\bullet/x\}M_\bullet = (\{N/x\}M)_\bullet$  then

$$\begin{aligned}
\{N_\bullet/x\}(a\langle y \rangle. M)_\bullet &= \{N_\bullet/x\}(a\langle y \rangle. M_\bullet) \\
&= a\langle y \rangle. \{N_\bullet/x\}M_\bullet \\
&= a\langle y \rangle. (\{N/x\}M)_\bullet \\
&= (a\langle y \rangle. \{N/x\}M)_\bullet .
\end{aligned}$$

- Marked redex case: if  $\{N_\bullet/x\}(H. M)_\bullet = (\{N/x\}(H. M))_\bullet$  and  $\{N_\bullet/x\}P_\bullet = (\{N/x\}P)_\bullet$  then

$$\begin{aligned}
&\{N_\bullet/x\}([P]a\leftarrow H. \blacktriangleright a\langle y \rangle. M)_\bullet \\
&= \{N_\bullet/x\}\{P_\bullet/y\}(H. M)_\bullet \\
&= \{\{N_\bullet/x\}P_\bullet/y\}\{N_\bullet/x\}(H. M)_\bullet \\
&= \{\{N/x\}P_\bullet/y\}(\{N/x\}(H. M))_\bullet \\
&= \{\{N/x\}P_\bullet/y\}(\{N/x\}H). \{N/x\}M)_\bullet \\
&= ([N/x]P]a\leftarrow (\{N/x\}H). \blacktriangleright a\langle y \rangle. \{N/x\}M)_\bullet \\
&= (\{N/x\}([P]a\leftarrow H. \blacktriangleright a\langle y \rangle. M))_\bullet .
\end{aligned}$$

□

**Lemma A.5** For a doubly marked term,  $M^{\diamond\circ}$ , reducing each marking in turn gives the same result as reducing both simultaneously:  $(M_\bullet)_\circ = M_{\circ\circ}$ .

**Proof.** By induction on the size of  $M$ .

- Unit case: immediate by  $((\star)_\bullet)_\circ = \star = (\star)_{\circ\circ}$ .
- Variable case: if  $(M_\bullet)_\circ = M_{\circ\circ}$  then

$$((x. M)_\bullet)_\circ = x. (M_\bullet)_\circ = x. M_{\circ\circ} = (x. M)_{\circ\circ} .$$

- Unmarked application case: if  $(M_\bullet)_\circ = M_{\circ\circ}$  and  $(N_\bullet)_\circ = N_{\circ\circ}$  then

$$\begin{aligned}
(([N]a. M)_\bullet)_\circ &= [(N_\bullet)_\circ]a. (M_\bullet)_\circ \\
&= [N_{\circ\circ}]a. M_{\circ\circ} \\
&= ([N]a. M)_{\circ\circ} .
\end{aligned}$$

- Unmarked abstraction case: if  $(M_\bullet)_\circ = M_{\circ\circ}$  then

$$((a\langle x \rangle. M)_\bullet)_\circ = a\langle x \rangle. (M_\bullet)_\circ = a\langle x \rangle. M_{\circ\circ} = (a\langle x \rangle. M)_{\circ\circ} .$$

- Doubly-marked redex case: if  $(H. M_\bullet)_\circ = H. M_{\circ\circ}$  and  $(N_\bullet)_\circ = N_{\circ\circ}$  then

$$\begin{aligned}
([N]a\leftarrow H. \blacktriangleright a\langle x \rangle. M)_\bullet)_\circ &= \{(N_\bullet)_\circ/x\}((H. M)_\bullet)_\circ \\
&= \{N_{\circ\circ}/x\}(H. M)_{\circ\circ} \\
&= ([N]a\leftarrow H. \blacktriangleright a\langle x \rangle. M)_{\circ\circ} .
\end{aligned}$$

- First singly-marked redex case: if  $(H. M_\bullet)_\diamond = (H. M)_{\bullet\bullet}$  and  $(N_\bullet)_\diamond = N_{\bullet\bullet}$  then, using Lemma A.4,

$$\begin{aligned}
 ([N]a\blacktriangleleft H.\blacktriangleright a\langle x \rangle. M)_\bullet)_\diamond &= (\{N_\bullet/x\}(H. M)_\bullet)_\diamond \\
 &= \{(N_\bullet)_\diamond/x\}((H. M)_\bullet)_\diamond \\
 &= \{N_{\bullet\bullet}/x\}(H. M)_{\bullet\bullet} \\
 &= ([N]a\blacktriangleleft H.\blacktriangleright a\langle x \rangle. M)_{\bullet\bullet} .
 \end{aligned}$$

- Second singly-marked redex case: the redex considered is

$$[N^{\blacklozenge}]a\blacktriangleleft H^{\blacktriangleleft}.\blacktriangleright a\langle x \rangle. M^{\blacktriangleright}$$

where  $H$  does not use the location  $a$  or a free variable  $x$ . First, it is established that  $(H^{\blacktriangleleft}.\blacktriangleright a\langle x \rangle. M^{\blacktriangleright})_\bullet$  is of the form  $K^{\blacktriangleleft}.\blacktriangleright a\langle x \rangle. N^{\blacktriangleright}$  where  $(H^{\blacktriangleleft}.\blacktriangleright M^{\blacktriangleright})_\bullet = K^{\blacktriangleleft}.\blacktriangleright N^{\blacktriangleright}$  for some head context  $K$  and term  $N$ . For readability the marking  $\diamond$  will be suppressed, except on the abstraction  $\blacktriangleright a\langle x \rangle$ . The proof is by induction on the size of  $H$ .

- Unit case: if  $H = \{\}$  then let  $K = \{\}$  and  $N = M_\bullet$ , which gives the following, as required.

$$\begin{aligned}
 (H.\blacktriangleright a\langle x \rangle. M)_\bullet &= (\blacktriangleright a\langle x \rangle. M)_\bullet & (H. M)_\bullet &= M_\bullet \\
 &= \blacktriangleright a\langle x \rangle. M_\bullet & &= N \\
 &= K.\blacktriangleright a\langle x \rangle. N & &= K. N .
 \end{aligned}$$

- Unmarked application case: if  $H = [P]b. H'$  then the inductive hypothesis for  $H'$  and  $M$  gives  $K'$  and  $N$ . Let  $K = [P]b. H'$ , which gives the following, as required.

$$\begin{aligned}
 ([P]b. H'.\blacktriangleright a\langle x \rangle. M)_\bullet & & ([P]b. H'. M)_\bullet \\
 &= [P_\bullet]b. (H'.\blacktriangleright a\langle x \rangle. M)_\bullet & &= [P_\bullet]b. (H'. M)_\bullet \\
 &= [P_\bullet]b. K'.\blacktriangleright a\langle x \rangle. N & &= [P_\bullet]b. K'. N \\
 &= K.\blacktriangleright a\langle x \rangle. N & &= K. N .
 \end{aligned}$$

- Unmarked abstraction case: if  $H = b\langle y \rangle. H'$  then the inductive hypothesis for  $H'$  and  $M$  gives  $K'$  and  $N$ . Let  $K = b\langle y \rangle. H'$ , which gives the following, as required.

$$\begin{aligned}
 (b\langle y \rangle. H'.\blacktriangleright a\langle x \rangle. M)_\bullet & & (b\langle y \rangle. H'. M)_\bullet \\
 &= b\langle y \rangle. (H'.\blacktriangleright a\langle x \rangle. M)_\bullet & &= b\langle y \rangle. (H'. M)_\bullet \\
 &= b\langle y \rangle. K'.\blacktriangleright a\langle x \rangle. N & &= b\langle y \rangle. K'. N \\
 &= K.\blacktriangleright a\langle x \rangle. N & &= K. N .
 \end{aligned}$$

- First marked redex case: if  $H = [P]b\blacktriangleleft H'.\blacktriangleright b\langle y \rangle. H''$  then the inductive hypothesis for  $H'. H''$  and  $M$  gives  $K'$  and  $N'$ . Let  $K = \{P_\bullet/x\}K'$  and  $N = \{P_\bullet/x\}N'$ , which gives the following, as required.

$$\begin{aligned}
 ([P]b\blacktriangleleft H'.\blacktriangleright b\langle y \rangle. H''.\blacktriangleright a\langle x \rangle. M)_\bullet & & ([P]b\blacktriangleleft H'.\blacktriangleright b\langle y \rangle. H''. M)_\bullet \\
 &= \{P_\bullet/y\}(H'. H''.\blacktriangleright a\langle x \rangle. M)_\bullet & &= \{P_\bullet/y\}(H'. H''. M)_\bullet \\
 &= \{P_\bullet/y\}(K'.\blacktriangleright a\langle x \rangle. N') & &= \{P_\bullet/y\}(K'. N') \\
 &= (\{P_\bullet/y\}K').\blacktriangleright a\langle x \rangle. \{P_\bullet/y\}N' & &= (\{P_\bullet/y\}K'). \{P_\bullet/y\}N' \\
 &= K.\blacktriangleright a\langle x \rangle. N & &= K. N .
 \end{aligned}$$

- Second marked redex case: if  $H = [P]b\blacktriangleleft H'$  and  $M = H''\blacktriangleright b\langle y \rangle.M'$  then the inductive hypothesis for  $H'$  and  $H''.M$  gives  $K'$  and  $N'$ . Let  $K = \{P_\blacktriangleright/x\}K'$  and  $N = \{P_\blacktriangleright/x\}N'$ , which gives the following, as required.

$$\begin{aligned}
 & ([P]b\blacktriangleleft H'.\blacktriangleright a\langle x \rangle.H''.\blacktriangleright b\langle y \rangle.M)_\blacktriangleright & ([P]b\blacktriangleleft H'.H''.\blacktriangleright b\langle y \rangle.M)_\blacktriangleright \\
 & = \{P_\blacktriangleright/x\}(H'.\blacktriangleright a\langle x \rangle.H''.M)_\blacktriangleright & = \{P_\blacktriangleright/x\}(H'.H''.M)_\blacktriangleright \\
 & = \{P_\blacktriangleright/x\}(K'.\blacktriangleright a\langle x \rangle.N') & = \{P_\blacktriangleright/x\}(K'.N') \\
 & = (\{P_\blacktriangleright/x\}K').\blacktriangleright a\langle x \rangle.\{P_\blacktriangleright/x\}N' & = (\{P_\blacktriangleright/x\}K').\{P_\blacktriangleright/x\}N' \\
 & = K.\blacktriangleright a\langle x \rangle.N & = K.N .
 \end{aligned}$$

This concludes the subproof. Returning to the main proof, let  $(H.\blacktriangleright a\langle x \rangle.M)_\blacktriangleright = K.\blacktriangleright a\langle x \rangle.N$  where  $(H.M)_\blacktriangleright = K.N$ . The case concludes as follows.

$$\begin{aligned}
 (([N]a\blacktriangleleft H.\blacktriangleright a\langle x \rangle.M)_\blacktriangleright)_\blacktriangleright & = ([N_\blacktriangleright]a\blacktriangleleft (H.\blacktriangleright a\langle x \rangle.M)_\blacktriangleright)_\blacktriangleright \\
 & = ([N_\blacktriangleright]a\blacktriangleleft K.\blacktriangleright a\langle x \rangle.N)_\blacktriangleright \\
 & = \{(N_\blacktriangleright)_\blacktriangleright/x\}(K.N)_\blacktriangleright \\
 & = \{N_{\blacktriangleright\blacktriangleright}/x\}(H.M)_{\blacktriangleright\blacktriangleright} \\
 & = ([N]a\blacktriangleleft H.\blacktriangleright a\langle x \rangle.M)_{\blacktriangleright\blacktriangleright}
 \end{aligned}$$

□

**Lemma A.6** *Parallel reduction is diamond: if  $N \leftarrow M \Rightarrow P$  then  $N \Rightarrow Q \leftarrow P$  for some  $Q$ .*

**Proof.** Let  $M^\blacktriangleright \Rightarrow M_\blacktriangleright = N$  and  $M^\blacktriangleright \Rightarrow M_\blacktriangleright = P$  for two separate markings  $\blacktriangleright$  and  $\blacktriangleright$  of  $M$ . Then with both markings on  $M$  the peak becomes

$$(M_\blacktriangleright)^\blacktriangleright \leftarrow M^{\blacktriangleright\blacktriangleright} \Rightarrow (M_\blacktriangleright)^\blacktriangleright .$$

Let  $Q = M_{\blacktriangleright\blacktriangleright}$  so that by Lemma A.5 the peak converges as

$$(M_\blacktriangleright)^\blacktriangleright \Rightarrow M_{\blacktriangleright\blacktriangleright} \leftarrow (M_\blacktriangleright)^\blacktriangleright .$$

□

**Theorem 4.6 (restatement)** Reduction  $\rightarrow$  is confluent.

**Proof.** A reduction step  $M \rightarrow N$  is a parallel step  $M^\blacktriangleright \Rightarrow M_\blacktriangleright = N$  by marking only the reduced redex. A peak  $M \leftarrow N \rightarrow P$  in regular reduction is then immediately one in parallel reduction,  $M \leftarrow N \Rightarrow P$ . By the diamond property, Lemma A.6, this converges with parallel reductions,  $M \Rightarrow Q \leftarrow P$ . By Proposition A.2 a parallel reduction ( $\Rightarrow$ ) is a regular reduction ( $\rightarrow$ ), so that the peak converges as  $M \rightarrow Q \leftarrow P$ . □