

# Normalization by Evaluation for the Computational Lambda-Calculus

Andrzej Filinski

BRICS\*, Department of Computer Science, University of Aarhus  
andrzej@brics.dk

**Abstract.** We show how a simple semantic characterization of normalization by evaluation for the  $\lambda_{\beta\eta}$ -calculus can be extended to a similar construction for normalization of terms in the computational  $\lambda$ -calculus. Specifically, we show that a suitable *residualizing* interpretation of base types, constants, and computational effects allows us to extract a syntactic normal form from a term's denotation. The required interpretation can itself be constructed as the meaning of a suitable functional program in an ML-like language, leading directly to a practical normalization algorithm. The results extend easily to product and sum types, and can be seen as a formal basis for call-by-value type-directed partial evaluation.

## 1 Introduction

The basic idea of normalization by evaluation is to extract the normal form (with respect to some notion of conversion) of a term from its interpretation in a suitably chosen, quasi-syntactic denotational model of the conversion relation [5].

For instance, let us consider the interpretation of a pure, simply typed lambda-term  $E$  in a model where all base types are interpreted as the set  $A$  of well-formed lambda-terms, and function types are interpreted as full set-theoretic function spaces. Then it is fairly simple to (at least informally) construct for any type  $\tau$ , a function  $nf_\tau \in \llbracket \tau \rrbracket \rightarrow A$ , such that for any closed term  $\tilde{E} : \tau$  in  $\beta\eta$ -long normal form,  $nf_\tau(\llbracket \tilde{E} \rrbracket) =_\alpha \tilde{E}$ . We proceed as follows:

Let  $\tau = \tau_1 \rightarrow \dots \rightarrow \tau_n \rightarrow b$  ( $n \geq 0$ ), where each  $\tau_i = \tau_{i1} \rightarrow \dots \rightarrow \tau_{im_i} \rightarrow b_i$ . Then  $\tilde{E} : \tau$  must be of the form  $\lambda x_1. \dots \lambda x_n. x_i \tilde{E}_1 \dots \tilde{E}_{m_i}$  where each  $\tilde{E}_j : \tau_{ij}$  is again in normal form. We can thus define  $nf_\tau$  inductively as:

$$nf_\tau = \lambda f. LAM(v_1, \dots, LAM(v_n, \\ f(\lambda a_1. \dots \lambda a_{m_1}. APP(\dots APP( VAR v_1, nf_{\tau_{11}} a_1) \dots, nf_{\tau_{1m_1}} a_{m_1})) \\ \vdots \\ (\lambda a_1. \dots \lambda a_{m_n}. APP(\dots APP( VAR v_n, nf_{\tau_{n1}} a_1) \dots, nf_{\tau_{nm_n}} a_{m_n}))))))$$

where the  $v_i$  are “fresh” variable names, and we use  $VAR$ ,  $LAM$ , and  $APP$  for constructing elements of  $A$ , to distinguish them from function abstraction and application in the set-theoretic model.

---

\* Basic Research in Computer Science ([www.brics.dk](http://www.brics.dk)),  
funded by the Danish National Research Foundation.

Moreover, it is easy to see that  $nf_\tau(\llbracket - \rrbracket)$  is in fact a normalization function: since  $\beta\eta$ -convertibility is sound for equality in all set-theoretic interpretations (and hence also in our chosen one), we have, for all terms  $E$  and  $E'$ , that  $E =_{\beta\eta} E'$  implies  $\llbracket E \rrbracket = \llbracket E' \rrbracket$ . So if  $\tilde{E}$  is the  $\beta\eta$ -long normal form of  $E$ , then  $nf_\tau(\llbracket E \rrbracket) = nf_\tau(\llbracket \tilde{E} \rrbracket) =_\alpha E$ .

Finally, if we can also construct a syntactic term  $nf_\tau : \tau \rightarrow \Lambda$  such that  $nf_\tau \llbracket E \rrbracket = \llbracket nf_\tau E \rrbracket$ , then  $nf_\tau E$  is a closed term of base type  $\Lambda$ , and can thus be executed as a functional program. This gives us a very efficient executable *algorithm* for computing normal forms, and was indeed one of the motivations behind the construction [4,3]: we are reducing the general problem of term normalization to a special case for which we already have a good solution.

A natural question arises whether this semantic technique for normalization of lambda-terms is inherently tied to  $\beta\eta$ -conversion. Somewhat surprisingly, it is not: in the following, we show how the same idea can be used to normalize terms with respect to the computational lambda-calculus [16], where it also extends to product and – more notably – sum types. In fact, we can systematically extract the computational normal form of any pure, typed lambda-term from only its observable behavior in an imperative functional language such as ML.

Despite the relative simplicity of the construction, there are still a few technical details to nail down, even in the purely functional case. Accordingly, we will first present in Section 2 the normalization algorithm for a call-by-name setting, then show in Section 3 how it can be refined to call-by-value. In Section 4 we show how to further extend the normalizer with product and sum types, and in Section 5 we consider the relationship between normalization by evaluation and type-directed partial evaluation. Finally, Section 6 concludes and points out some directions for further work.

## 2 Normalization by Evaluation for Call by Name

The normalization construction sketched above, essentially due to Berger and Schwichtenberg [4], has been studied in many formulations [8], including more syntactic variants [2] as well as category-theoretic ones [1,6]. In the following, we present it in a call-by-name functional-programming setting [11]; this formulation extends particularly naturally to the call-by-value variant in the next section.

### 2.1 Language and Semantic Framework

*Syntax.* A *signature*  $\Sigma$  includes, first, a collection of base types  $b$ . The set of well-formed  $\Sigma$ -types  $\tau$  is then given by the grammar

$$\tau ::= b \mid \tau_1 \rightarrow \tau_2$$

Further,  $\Sigma$  assigns  $\Sigma$ -types to a (possibly infinite) collection of constants  $c$ . Let  $x$  range over variable names, and  $\Gamma$  be a finite assignment of  $\Sigma$ -types to variables. Then the set of well-typed  $\Sigma$ -terms  $\Gamma \vdash_\Sigma E : \tau$  is again given by the usual rules:

$$\frac{\Sigma(c) = \tau}{\Gamma \vdash_\Sigma c : \tau} \quad \frac{\Gamma(x) = \tau}{\Gamma \vdash_\Sigma x : \tau} \quad \frac{\Gamma, x : \tau_1 \vdash_\Sigma E : \tau_2}{\Gamma \vdash_\Sigma \lambda x^{\tau_1}. E : \tau_1 \rightarrow \tau_2} \quad \frac{\Gamma \vdash_\Sigma E_1 : \tau_1 \rightarrow \tau_2 \quad \Gamma \vdash_\Sigma E_2 : \tau_1}{\Gamma \vdash_\Sigma E_1 E_2 : \tau_2}$$

Finally, a  $\Sigma$ -program is a closed  $\Sigma$ -term of base type.

*Semantics.* For concreteness, and to accommodate the refinements in Section 5, we consider only a specific, domain-theoretic framework, but the results also adapt easily to a set-theoretic setting, by forgetting the order structure.

We work in the setting of (bottomless) cpos and (total) continuous functions. We also use the concept of a *monad* (more precisely, a *Kleisli triple*)  $\mathcal{T} = (T, \eta, \star)$ , where  $T$  maps cpos to cpos,  $\eta_A \in A \rightarrow TA$  is the *unit* function, and  $\star_{A,B} \in TA \times (A \rightarrow TB) \rightarrow TB$  is the *extension* operation written backwards, i.e., with the function last; this makes longer sequences of extensions easier to read. We omit the subscripts on units and extensions where they are clear from the context. A particularly important instance is the *lifting monad*,  $\mathcal{T}_1$ , where  $T^1 A = A_\perp = \{\iota a \mid a \in A\} \cup \{\perp\}$  with the usual ordering,  $\eta^1 a = \iota a$ ,  $\perp \star^1 f = \perp$ , and  $(\iota a) \star^1 f = f a$ .

An *interpretation*  $\mathcal{I}$  of a signature  $\Sigma$  is a pair of functions  $(\mathcal{B}, \mathcal{C})$ .  $\mathcal{B}$  assigns to every base type  $b$  in  $\Sigma$ , a cpo  $\mathcal{B}(b)$ . This assignment determines for any  $\Sigma$ -type  $\tau$ , a pointed (i.e., containing a least element) cpo  $\llbracket \tau \rrbracket^{\mathcal{I}}$  as follows:

$$\llbracket b \rrbracket^{\mathcal{I}} = T^1 \mathcal{B}(b) \qquad \llbracket \tau_1 \rightarrow \tau_2 \rrbracket^{\mathcal{I}} = \llbracket \tau_1 \rrbracket^{\mathcal{I}} \rightarrow \llbracket \tau_2 \rrbracket^{\mathcal{I}}$$

where  $A \rightarrow B$  denotes the cpo of all continuous functions between  $A$  and  $B$ . We also give meaning to a type assignment  $\Gamma$  as a finite product:

$$\llbracket \Gamma \rrbracket^{\mathcal{I}} = \prod_{x \in \text{dom } \Gamma} \llbracket \Gamma(x) \rrbracket^{\mathcal{I}} = \{\rho \mid \forall x \in \text{dom } \Gamma. \rho x \in \llbracket \Gamma(x) \rrbracket^{\mathcal{I}}\}$$

The function  $\mathcal{C}$  assigns to every  $\Sigma$ -constant  $c$  an element  $\mathcal{C}(c) \in \llbracket \Sigma(c) \rrbracket^{\mathcal{I}}$ . Again this assignment extends to a full semantics of terms: for any  $\Gamma \vdash_{\Sigma} E : \tau$ , we define a continuous function  $\llbracket E \rrbracket^{\mathcal{I}} \in \llbracket \Gamma \rrbracket^{\mathcal{I}} \rightarrow \llbracket \tau \rrbracket^{\mathcal{I}}$  in the usual way:

$$\begin{aligned} \llbracket c \rrbracket^{\mathcal{I}} \rho &= \mathcal{C}(c) & \llbracket \lambda x^{\tau}. E \rrbracket^{\mathcal{I}} \rho &= \lambda a. \llbracket E \rrbracket^{\mathcal{I}} (\rho[x \mapsto a]) \\ \llbracket x \rrbracket^{\mathcal{I}} \rho &= \rho x & \llbracket E_1 E_2 \rrbracket^{\mathcal{I}} \rho &= \llbracket E_1 \rrbracket^{\mathcal{I}} \rho (\llbracket E_2 \rrbracket^{\mathcal{I}} \rho) \end{aligned}$$

*Equivalence and normal forms.* We say that two  $\Sigma$ -terms  $E$  and  $E'$  are semantically equivalent, written  $\models E = E'$ , if for all interpretations  $\mathcal{I}$  of  $\Sigma$ ,  $\llbracket E \rrbracket^{\mathcal{I}} = \llbracket E' \rrbracket^{\mathcal{I}}$ . It is easy to see that if  $E =_{\beta\eta} E'$  then  $\models E = E'$ . More generally, if  $\text{Int}$  is a subset of all possible interpretations of  $\Sigma$  (e.g., constraining the meanings of some of the constants), we write  $\models^{\text{Int}} E = E'$  iff for all  $\mathcal{I} \in \text{Int}$ ,  $\llbracket E \rrbracket^{\mathcal{I}} = \llbracket E' \rrbracket^{\mathcal{I}}$ ; we will return to this notion in Section 5.

Among the well-typed terms  $\Gamma \vdash_{\Sigma} E : \tau$ , we distinguish those in *normal* and *atomic* (also known as *neutral*) form:

$$\begin{aligned} \frac{\Gamma \Vdash^{\text{nf}} E : b}{\Gamma \Vdash^{\text{nf}} E : b} \quad & \frac{\Gamma, x : \tau_1 \Vdash^{\text{nf}} E : \tau_2}{\Gamma \Vdash^{\text{nf}} \lambda x^{\tau_1}. E : \tau_1 \rightarrow \tau_2} \\ \frac{\Gamma(x) = \tau}{\Gamma \Vdash^{\text{nf}} x : \tau} \quad & \frac{\Sigma(c) = \tau}{\Gamma \Vdash^{\text{nf}} c : \tau} \quad & \frac{\Gamma \Vdash^{\text{nf}} E_1 : \tau_1 \rightarrow \tau_2 \quad \Gamma \Vdash^{\text{nf}} E_2 : \tau_1}{\Gamma \Vdash^{\text{nf}} E_1 E_2 : \tau_2} \end{aligned}$$

A *normalization function*, in the sense of Coquand and Dybjer [5] (but with a semantic notion of equivalence), then maps any term  $E$  to a normal-form term  $\text{norm}(E)$ , such that  $\models \text{norm}(E) = E$ , and such that for all  $E'$  with  $\models E' = E$ ,  $\text{norm}(E') = \text{norm}(E)$ .

## 2.2 A Normalization Result

The traditional way of computing  $\text{norm}(E)$  is by repeated  $\beta$ -reductions, possibly followed by  $\eta$ -expansions. However, we can also compute  $\text{norm}$  by a subtler, semantic method, *reduction-free normalization*.

*Representing lambda-terms.* Let  $\mathbf{V}$  be a set (= discrete cpo) of explicitly typed variable names, and  $\mathbf{E}$  a set suitable for representing lambda-terms, i.e., allowing us to define injective functions with mutually disjoint ranges,

$$CST \in \text{dom } \Sigma \rightarrow \mathbf{E}, \quad VAR \in \mathbf{V} \rightarrow \mathbf{E}, \quad LAM \in \mathbf{V} \times \mathbf{E} \rightarrow \mathbf{E}, \quad APP \in \mathbf{E} \times \mathbf{E} \rightarrow \mathbf{E}$$

Then for any term  $E$  with variables from  $\mathbf{V}$ , we define its representation  $\lceil E \rceil \in \mathbf{E}$  in the obvious way, e.g.,  $\lceil \lambda x. E \rceil = LAM(x, \lceil E \rceil)$ . Because of the injectivity and disjointness assumptions, for any  $e \in \mathbf{E}$ , there is *at most one*  $E$  such that  $\lceil E \rceil = e$ ; we need not require that all elements of  $\mathbf{E}$  represent well-formed lambda-terms, let alone well-typed ones.

(We deliberately use a very concrete representation of terms, rather than a higher-level notion based on abstract syntax with binding constructs such as [12]. Our ultimate goal is to implement the normalization process as a simple functional program, without assuming potentially expensive operations, such as capture-avoiding substitution, as primitives.)

The task is now to construct an interpretation  $\mathcal{I}_\tau$  of  $\Sigma$  such that we can recover  $E$ 's normal form from  $\llbracket E \rrbracket^{\mathcal{I}_\tau}$ . We want to use the idea from the introduction, but need to account rigorously for “fresh” variable names. Freshness could be captured abstractly in a framework such as Fraenkel-Mostowski sets [14], but this again removes us a level from a direct implementation. Instead, we will explicitly generate non-clashing variable names. Perhaps the simplest scheme for doing so is through de Bruijn levels [2,11], but we adopt instead a scheme for generating “globally unique”, gensym-style names using a monad, as it scales better to the constructions of the next section.

*Auxiliary definitions.* We first define the name-generation monad  $\mathcal{T}_g$ . This is just a state-passing monad atop  $\mathcal{T}_1$ ; the state is the “next free index”:

$$T^g A = \mathbf{N} \rightarrow T^1(A \times \mathbf{N}) \quad \eta^g a = \lambda i. \eta^1(a, i) \quad t \star^g f = \lambda i. t i \star^1 \lambda(a, i'). f a i'$$

With respect to  $T^g$ , we can define an effectful computation that generates a fresh name, and one that initializes the index within a delimited subcomputation:

$$\begin{aligned} new_\tau &\in T^g \mathbf{V} & withct_A &\in T^g A \rightarrow T^1 A \\ new_\tau &= \lambda i. \eta^1(g_i^\tau, i + 1) & withct_A t &= t 0 \star^1 \lambda(a, i'). \eta^1 a \end{aligned}$$

where the  $g_i^\tau \in \mathbf{V}$  are assumed distinct for distinct  $i$ . Note that the codomain of  $withct_A$  is simply  $T^1 A$ , i.e., *withct*  $t$  represents a side-effect-free, purely functional computation, for any name-generating computation  $t$ .

*The residualizing interpretation.* We can now define a suitable residualizing interpretation  $\mathcal{I}_r = (\mathcal{B}_r, \mathcal{C}_r)$ . For  $\mathcal{B}_r$  we take, for all base types  $b$  in  $\Sigma$ ,

$$\mathcal{B}_r(b) = T^g \mathbf{E}$$

Formalizing the construction from the introduction, we further define, for any  $\Sigma$ -type  $\tau$ , a pair of functions commonly called *reification* and *reflection*:

$$\begin{aligned} \downarrow^\tau &\in \llbracket \tau \rrbracket^{\mathcal{I}_r} \rightarrow T^g \mathbf{E} \\ \downarrow^b &= \lambda \varepsilon. \varepsilon \\ \downarrow^{\tau_1 \rightarrow \tau_2} &= \lambda f. \text{new}_{\tau_1} \star^g \lambda v. \downarrow^{\tau_2} (f(\uparrow_{\tau_1}(\eta^g(\text{VAR } v)))) \star^g \lambda e. \eta^g(\text{LAM}(v, e)) \\ \uparrow_\tau &\in T^g \mathbf{E} \rightarrow \llbracket \tau \rrbracket^{\mathcal{I}_r} \\ \uparrow^b &= \lambda \varepsilon. \varepsilon \\ \uparrow_{\tau_1 \rightarrow \tau_2} &= \lambda \varepsilon. \lambda a. \uparrow_{\tau_2} (\varepsilon \star^g \lambda e. \downarrow^{\tau_1} a \star^g \lambda e'. \eta^g(\text{APP}(e, e')))) \end{aligned}$$

(It may be helpful, on a first reading, to think of  $\mathcal{T}_g$  as just the identity monad, and  $\text{new}_\tau$  as “magically” generating fresh variable names; then the reification function simplifies to precisely the construction of  $nf$  sketched in Section 1.) Finally, we define the residualizing interpretation of constants by

$$\mathcal{C}_r(c) \in \llbracket \Sigma(c) \rrbracket^{\mathcal{I}_r}, \quad \mathcal{C}_r(c) = \uparrow_{\Sigma(c)}(\eta^g(\text{CST } c))$$

*The normalization function.* To extract the syntactic normal form from the residualizing meaning of a term, we only need to supply a starting index for name generation. We can thus define an *extraction function*:

$$nf_\tau \in \llbracket \tau \rrbracket^{\mathcal{I}_r} \rightarrow T^1 \mathbf{E}, \quad nf_\tau = \lambda a. \text{withct}_{\mathbf{E}}(\downarrow^\tau a)$$

Finally, we define the (potentially partial) syntax-to-syntax function *norm* on closed terms  $\vdash_\Sigma E : \tau$  by

$$\text{norm}(E) = \tilde{E} \quad \text{iff} \quad nf_\tau(\llbracket E \rrbracket^{\mathcal{I}_r} \emptyset) = \eta^1 \ulcorner \tilde{E} \urcorner$$

(We can find normal forms of open terms by explicitly lambda-abstraction over their free variables. The closed-term formulation leads to a particularly natural implementation, as sketched below.)

**Theorem 1 (CBN semantic normalization).** *Let  $\vdash_\Sigma E : \tau$  be a closed  $\Sigma$ -term. Then (0)  $\tilde{E} = \text{norm}(E)$  is defined, (1)  $\models^{\text{nf}} \tilde{E} : \tau$ , (2)  $\models \tilde{E} = E$ , and (3) for all  $\vdash_\Sigma E' : \tau$  such that  $\models E' = E$ ,  $\text{norm}(E') = \tilde{E}$ .*

If we already know that any  $\Sigma$ -term has a unique (up to  $\alpha$ -conversion)  $\beta\eta$ -long normal form, the proof is fairly simple, using the argument sketched in the introduction. However, it is also possible to prove the theorem directly, using a suitable Kripke logical relation between the meanings of terms in the residualizing interpretation  $\mathcal{I}_r$  and in an arbitrary one  $\mathcal{I}$ , with the base relation taken as the denotational meaning function. The details (for a more general setting, as sketched in Section 5) can be found in [11].

*A normalization algorithm.* The normalization function described above can be effectively computed as a program in any PCF-like functional language. The key idea is to express the residualizing semantic interpretation as a syntactic *realization* of all base types and constants in  $\Sigma$  in terms of types and terms of the programming-language signature  $\Sigma_{\text{pl}}$ , giving a substitution  $\Phi^r$ , such that  $\llbracket E\{\Phi^r\} \rrbracket^{\mathcal{I}_{\text{pl}}} = \llbracket E \rrbracket^{\mathcal{I}_r}$ . Likewise, for any  $\Sigma$ -type  $\tau$ , we construct a term  $\text{nf}_\tau$  such that  $\llbracket \text{nf}_\tau \rrbracket^{\mathcal{I}_{\text{pl}}} \emptyset = \text{nf}_\tau$ . Then for any closed  $\Sigma$ -term  $\vdash_\Sigma E : \tau$ , we can compute its normal form by evaluating the  $\Sigma_{\text{pl}}$ -program “ $\text{nf}_\tau E\{\Phi^r\}$ ”. Again, the details can be found in [11].

### 3 Normalization by Evaluation for Call by Value with Effects

We now refine the normalization result to a language based on Moggi’s computational lambda-calculus  $\lambda_c$  [16], which provides a semantic framework for ML-like languages where “functions” may have effects such as mutating the state, performing input/output operations, or raising exceptions.

#### 3.1 Language and Semantic Framework

*Syntax and semantics.* The syntax of types is the same as before. For terms, we also add a let-construct, with the usual typing rule:

$$\frac{\Gamma \vdash_\Sigma E_1 : \tau_1 \quad \Gamma, x : \tau_1 \vdash_\Sigma E_2 : \tau_2}{\Gamma \vdash_\Sigma \text{let } x = E_1 \text{ in } E_2 : \tau_2}$$

Now an interpretation  $\mathcal{I}$  of a signature  $\Sigma$  consists a triple  $(\mathcal{B}, \mathcal{T}, \mathcal{C})$ . As before,  $\mathcal{B}$  assigns cpos to base types of  $\Sigma$ . The new component  $\mathcal{T} = (T, \eta, \star)$  is a monad used to model computational effects. These could be just divergence (modeled with the lifting monad), but also state, exceptions, continuations, etc.; the actual effectful operations are invoked through suitable constants from  $\Sigma$ .

We need to assume that  $\mathcal{T}$  is layered atop  $\mathcal{T}_!$  [10]; this amounts to requiring that  $TA$  is pointed for any  $A$ , and that  $\lambda t. t \star f \in TA \rightarrow TB$  is strict for any  $f \in A \rightarrow TB$ . The CBV semantics of types is then given by:

$$\llbracket b \rrbracket_v^\mathcal{I} = \mathcal{B}(b) \quad \llbracket \tau_1 \rightarrow \tau_2 \rrbracket_v^\mathcal{I} = \llbracket \tau_1 \rrbracket_v^\mathcal{I} \rightarrow T \llbracket \tau_2 \rrbracket_v^\mathcal{I}$$

The meaning of a typing environment,  $\llbracket \Gamma \rrbracket_v^\mathcal{I}$ , is a  $(\text{dom } \Gamma)$ -indexed product of the meanings of the individual types, as before.

For the semantic function  $\mathcal{C}$ , we again require that for any  $c \in \text{dom } \Sigma$ ,  $\mathcal{C}(c) \in \llbracket \Sigma(c) \rrbracket_v^\mathcal{I}$ . Then we define the meaning of a well-typed term  $\Gamma \vdash_\Sigma E : \tau$  as a continuous function  $\llbracket E \rrbracket_v^\mathcal{I} \in \llbracket \Gamma \rrbracket_v^\mathcal{I} \rightarrow T \llbracket \tau \rrbracket_v^\mathcal{I}$ , as follows:

$$\begin{aligned} \llbracket c \rrbracket_v^\mathcal{I} \rho &= \eta(\mathcal{C}(c)) & \llbracket \lambda x^\tau. E \rrbracket_v^\mathcal{I} \rho &= \eta(\lambda a. \llbracket E \rrbracket_v^\mathcal{I} (\rho[x \mapsto a])) \\ \llbracket x \rrbracket_v^\mathcal{I} \rho &= \eta(\rho x) & \llbracket E_1 E_2 \rrbracket_v^\mathcal{I} \rho &= \llbracket E_1 \rrbracket_v^\mathcal{I} \rho \star \lambda f. \llbracket E_2 \rrbracket_v^\mathcal{I} \rho \star \lambda a. f a \\ \llbracket \text{let } x = E_1 \text{ in } E_2 \rrbracket_v^\mathcal{I} \rho &= \llbracket E_1 \rrbracket_v^\mathcal{I} \rho \star \lambda a. \llbracket E_2 \rrbracket_v^\mathcal{I} (\rho[x \mapsto a]) \end{aligned}$$

Note that the **let**-construct appears redundant, because  $\llbracket \text{let } x = E_1 \text{ in } E_2 \rrbracket_v^\mathcal{I} \rho = \llbracket (\lambda x. E_2) E_1 \rrbracket_v^\mathcal{I} \rho$ , but including it enables a nicer syntactic characterization of normal forms.

*Equivalence and normal forms.* Analogously to the CBN case, we write  $\models_v E = E'$  if for all  $\mathcal{I}$ ,  $\llbracket E \rrbracket_v^{\mathcal{I}} = \llbracket E' \rrbracket_v^{\mathcal{I}}$ . The shape of normal forms is now somewhat different, however: instead of normal and atomic forms, we have *normal values* and *normal computations*:

$$\frac{\Sigma(c) = b}{\Gamma \models^v c : b} \quad \frac{\Gamma(x) = b}{\Gamma \models^v x : b} \quad \frac{\Gamma, x : \tau_1 \models^c E : \tau_2}{\Gamma \models^v \lambda x^{\tau_1}. E : \tau_1 \rightarrow \tau_2}$$

$$\frac{\Gamma \models^v E : \tau \quad \Sigma(c) = \tau_1 \rightarrow \tau_2 \quad \Gamma \models^v E : \tau_1 \quad \Gamma, x : \tau_2 \models^c E' : \tau}{\Gamma \models^c \text{let } x = cE \text{ in } E' : \tau}$$

$$\frac{\Gamma(x') = \tau_1 \rightarrow \tau_2 \quad \Gamma \models^v E : \tau_1 \quad \Gamma, x : \tau_2 \models^c E' : \tau}{\Gamma \models^c \text{let } x = x'E \text{ in } E' : \tau}$$

That is, a normal value is either a base-typed constant or variable, or of the form  $\lambda x. \text{let } x_1 = f_1 V_1 \text{ in } \dots \text{let } x_n = f_n V_n \text{ in } V$  where all the  $V$ s are normal values, and each  $f_i$  is a function-typed constant or variable.

The set of normal-form terms is similar to Flanagan et al.'s *A-normal forms* [13]. However, their notion of A-reduction does not include even restricted  $\beta$ -conversion, so a term such as  $(\lambda x. x)y$  is already A-normal. (Nor does it include  $\eta$ -like let-conversions: both  $fx$  and  $\text{let } y = fx \text{ in } y$  are A-normal.) A much closer match is Ohori's language of *cut-free A-normal sequent proofs* [17], but still with one important difference: since we also care about uniqueness of normal forms, a variable is only a normal value in our sense if it is of base type; function-typed normal values must always be syntactic lambda-abstractions.

### 3.2 A Normalization Result

*Term representations.* Corresponding to the extended source syntax, we also assume given an additional constructor function  $LET \in \mathbf{V} \times \mathbf{E} \times \mathbf{E} \rightarrow \mathbf{E}$ , injective and with range disjoint from the others. The representation function for terms is also extended in the obvious way.

*Residualizing monad.* For constructing the residualizing interpretation, we now also need to pick a residualizing monad  $\mathcal{T}_r$ . It is easy to see that we cannot simply use the lifting monad here, even if we only care about “purely functional” call-by-value languages, i.e., interpretations with  $\mathcal{T}$  taken as lifting. The reason is that the two normal-form terms  $E_1 = \lambda f^{b \rightarrow b}. \lambda x^b. \text{let } y = fx \text{ in } x$  and  $E_2 = \lambda f^{b \rightarrow b}. \lambda x^b. x$  are not semantically equivalent: for any  $\mathcal{I}$  where  $\mathcal{T}$  is the lifting monad, we only have  $\llbracket E_1 \rrbracket_v^{\mathcal{I}} \sqsubseteq \llbracket E_2 \rrbracket_v^{\mathcal{I}} \rho$  (with the strictness of the inequality demonstrated by application of both sides to  $\lambda a. \perp$ ). But  $\lceil E_1 \rceil \not\sqsubseteq \lceil E_2 \rceil$ , so there can be no *monotone* (let alone continuous) function  $nf : \llbracket (b \rightarrow b) \rightarrow b \rightarrow b \rrbracket^{\mathcal{I}_r} \rightarrow \mathbf{E}_{\perp}$  such that  $nf(\llbracket E_1 \rrbracket^{\mathcal{I}_r} \emptyset) = \eta^{\lceil E_1 \rceil}$  and  $nf(\llbracket E_2 \rrbracket^{\mathcal{I}_r} \emptyset) = \eta^{\lceil E_2 \rceil}$ , if we require  $\mathcal{I}_r$  to use only the lifting monad to interpret computational effects.

Indeed, looking at the shape of normal computations, we see that the residualizing monad must allow us to register exactly where and when a function-typed constant or variable was applied, even if its return value is never used. We will show that it suffices that  $\mathcal{T}_r$  can be equipped with operations

$$\text{bind}_r \in \mathbf{E} \rightarrow T^r \mathbf{V} \quad \text{and} \quad \text{collect} \in T^r \mathbf{E} \rightarrow T^g \mathbf{E}$$

satisfying the equational constraints

$$\begin{aligned} \text{collect}(\eta^r e) &= \eta^g e \\ \text{collect}(\text{bind}_\tau e \star^r f) &= \text{new}_\tau \star^g \lambda v. \text{collect}(f v) \star^g \lambda e'. \eta^g(\text{LET}(v, e, e')) \end{aligned}$$

These equations ensure that a  $T^r \mathbf{E}$ -computation consisting of a sequence of calls to  $\text{bind}$  followed by returning a term has the effect of wrapping that term in a corresponding sequence of  $\text{LET}$ s:

$$\begin{aligned} \text{collect}(\text{bind}_{\tau_1} e_1 \star^r \lambda v_1. \dots \text{bind}_{\tau_n} (e_n v_1 \dots v_{n-1}) \star^r \lambda v_n. \eta^r(e v_1 \dots v_n)) = \\ \text{new}_{\tau_1} \star^g \lambda v_1. \dots \text{new}_{\tau_n} \star^g \lambda v_n. \eta^g(\text{LET}(v_1, e_1, \dots \text{LET}(v_n, e_n v_1 \dots v_{n-1}, e v_1 \dots v_n))) \end{aligned}$$

To view  $T^g$ -computations as special cases of  $T^r$ -computations, we will also need a *monad morphism* from  $\mathcal{T}_g$  to  $\mathcal{T}_r$ , i.e., a collection of functions  $\gamma_A^{g,r} \in T^g A \rightarrow T^r A$ , such that  $\gamma^{g,r}(\eta^g a) = \eta^r a$  and  $\gamma^{g,r}(t \star^g f) = \gamma^{g,r} t \star^r \lambda a. \gamma^{g,r}(f a)$ . We can construct a monad with these operations in several ways, notably including the following two:

*The continuation monad with answer domain  $T^g \mathbf{E}$ .* We take:

$$\begin{aligned} T^r A &= (A \rightarrow T^g \mathbf{E}) \rightarrow T^g \mathbf{E} & \gamma^{g,r} t &= \lambda \kappa. t \star^g \kappa \\ \eta^r a &= \lambda \kappa. \kappa a & \text{bind}_\tau e &= \lambda \kappa. \text{new}_\tau \star^g \lambda v. \kappa v \star^g \lambda e'. \eta^g(\text{LET}(v, e, e')) \\ t \star^r f &= \lambda \kappa. t(\lambda a. f a \kappa) & \text{collect } t &= t \eta^g \end{aligned}$$

This continuation-based wrapping of syntactic bindings was originally used for an “administrative-reduction free” continuation-passing transformation [9], and later adapted for a similar purpose in type-directed partial evaluation [7]. It is notable for also allowing an extension of CBV NBE to sum types (see Section 4).

*The accumulation monad over the monoid of  $(\mathbf{V} \times \mathbf{E})$ -lists.* Writing  $[]$  for the empty list,  $[-]$  for a singleton list, and  $@$  for list concatenation, we take:

$$\begin{aligned} T^r A &= T^g(A \times (\mathbf{V} \times \mathbf{E})^*) & \gamma^{g,r} t &= t \star^g \lambda a. \eta^g(a, []) \\ \eta^r a &= \eta^g(a, []) & \text{bind}_\tau e &= \text{new}_\tau \star^g \lambda v. \eta^g(v, [(v, e)]) \\ t \star^r f &= t \star^g \lambda(a, l). f a \star^g \lambda(b, l'). \eta^g(b, l @ l') & \text{collect } t &= t \star^g \lambda(e, l). \eta^g(\text{wrap } l e) \end{aligned}$$

with the auxiliary function  $\text{wrap} : (\mathbf{V} \times \mathbf{E})^* \rightarrow \mathbf{E} \rightarrow \mathbf{E}$  defined inductively as

$$\text{wrap}[] e = e \quad \text{wrap}([(v, e')] @ l) e = \text{LET}(v, e', \text{wrap } l e)$$

This choice is a refinement of state-based TDPE [18]; see the end of this section for a brief account of the relationship between accumulation and state. Other constructions are also possible, such as accumulation with respect to the monoid  $(\mathbf{E} \rightarrow \mathbf{E}, \text{id}, \circ)$ .

*Residualizing interpretation.* For the residualizing interpretation, we again interpret all base types of  $\Sigma$  as syntactic lambda-terms; this time, however, we do not need to involve the name-generation monad yet, but simply take:

$$\mathcal{B}_r(b) = \mathbf{E}$$



For any definition of  $\mathcal{T}_r$  satisfying the equational constraints on *bind* and *collect*, we can then define new reification and reflection functions:

$$\begin{aligned}
\downarrow^\tau &\in \llbracket \tau \rrbracket_{\mathbf{V}}^{\mathcal{T}_r} \rightarrow T^{\mathbf{g}} \mathbf{E} \\
\downarrow^b &= \lambda e. \eta^{\mathbf{g}} e \\
\downarrow^{\tau_1 \rightarrow \tau_2} &= \lambda f. \text{new}_{\tau_1} \star^{\mathbf{g}} \lambda v. \text{collect} \left( \uparrow_{\tau_1} (\text{VAR } v) \star^{\mathbf{r}} \lambda a. f a \star^{\mathbf{r}} \lambda b. \gamma^{\mathbf{g}, \mathbf{r}} (\downarrow^{\tau_2} b) \right) \star^{\mathbf{g}} \lambda e. \\
&\quad \eta^{\mathbf{g}} (\text{LAM } (v, e)) \\
\uparrow_\tau &\in \mathbf{E} \rightarrow T^{\mathbf{r}} \llbracket \tau \rrbracket_{\mathbf{V}}^{\mathcal{T}_r} \\
\uparrow_b &= \lambda e. \eta^{\mathbf{r}} e \\
\uparrow_{\tau_1 \rightarrow \tau_2} &= \lambda e. \eta^{\mathbf{r}} \left( \lambda a. \gamma^{\mathbf{g}, \mathbf{r}} (\downarrow^{\tau_1} a) \star^{\mathbf{r}} \lambda e'. \text{bind}_{\tau_2} (\text{APP } (e, e')) \star^{\mathbf{r}} \lambda v. \uparrow_{\tau_2} (\text{VAR } v) \right)
\end{aligned}$$

(The codomain of  $\uparrow_\tau$  is  $T^{\mathbf{r}} \llbracket \tau \rrbracket_{\mathbf{V}}^{\mathcal{T}_r}$ , rather than simply  $\llbracket \tau \rrbracket_{\mathbf{V}}^{\mathcal{T}_r}$ , to accommodate the extensions in Section 4.) Note in particular how every construction of an *APP*-term is wrapped in a *bind*.

Finally, as for call by name, we interpret all  $\Sigma$ -constants as reflected *CST*-constructors:

$$\mathcal{C}_r(c) \in \llbracket \Sigma(c) \rrbracket_{\mathbf{V}}^{\mathcal{T}_r}, \quad \mathcal{C}_r(c) = a, \text{ where } \uparrow_{\Sigma(c)} (\text{CST } c) = \eta^{\mathbf{r}} a$$

( $\mathcal{C}_r(c)$  is well defined, because the reflection function factors through the injective  $\eta^{\mathbf{r}}$ .) We also define the extraction function essentially as before:

$$nf_\tau \in \llbracket \tau \rrbracket_{\mathbf{V}}^{\mathcal{T}_r} \rightarrow T^1 \mathbf{E}, \quad nf_\tau = \lambda a. \text{withct}_{\mathbf{E}} (\downarrow^\tau a)$$

and the CBV normalization function for a closed value (constant or lambda-abstraction)  $E$  as

$$\text{norm}_{\mathbf{V}}(E) = \tilde{E} \quad \text{iff} \quad nf_\tau a = \eta^1 \ulcorner \tilde{E} \urcorner, \text{ where } \llbracket E \rrbracket_{\mathbf{V}}^{\mathcal{T}_r} \emptyset = \eta^{\mathbf{r}} a.$$

(We can find the normal form of a non-value term by wrapping a dummy lambda-abstraction around it.)

**Theorem 2 (CBV semantic normalization).** *Let  $\vdash_\Sigma E : \tau$  be a value, and take  $\tilde{E} = \text{norm}_{\mathbf{V}}(E)$ . Then (0)  $\tilde{E}$  is defined, (1)  $\models^{\text{bv}} \tilde{E} : \tau$ , (2)  $\models_{\mathbf{V}} \tilde{E} = E$ , and (3) if  $\models_{\mathbf{V}} E' = E$  then  $\text{norm}_{\mathbf{V}}(E') = \tilde{E}$ .*

The proof is similar to the CBN case, but using a pair of mutually inductively defined logical relations, one for values and one for computations. Very roughly, one again establishes that the residualizing and the arbitrary interpretations of all terms are related, and that for a pair of related values  $(a, a') \in \llbracket \tau \rrbracket_{\mathbf{V}}^{\mathcal{T}_r} \times \llbracket \tau \rrbracket_{\mathbf{V}}^{\mathcal{T}_r}$  the  $\mathcal{I}$ -meaning of  $\downarrow^\tau a$  equals  $a'$ .

*A normalization algorithm.* Phrasing the normalization function as a functional program is a bit more complicated than for the CBN case. A typical CBV host language will have its own notion of effects, modeled by some monad  $\mathcal{T}_{\text{pl}}$ , which is not likely to be exactly our residualizing monad  $\mathcal{T}_r$ . However, much as we can embed the normalization algorithm into a host language with a signature much

larger than what we need for the construction, we can realize both the residualizing and the name-generating monad through a uniform *effect-embedding* into a more general notion of effect in the host language [10].

For example, the accumulation-based choice of  $\mathcal{T}_r$  can be easily (and more efficiently) implemented by passing around the bindings accumulated so far in a mutable state cell, rather than appending the bindings from both subcomputations in  $\star$ . The current name-generation index is naturally kept in another cell. An analogous, but somewhat more involved, construction also allows us to simulate  $\mathcal{T}_r$  taken as a continuation monad, provided the host language provides both (higher-typed) state and first-class continuations, as found in Scheme or SML/NJ. The resulting implementation forms the basis of the CBV normalization algorithm used in the context of type-directed partial evaluation [7].

## 4 Structured Data Types

In this section, we consider normalization for the call-by-value language extended with product and sum types. (Adding products to the call-by-name language is trivial, but it does not appear possible to add even weak sum types, at least in the domain-theoretic semantics.)

*Syntax.* We extend the set of types by two new type constructors:

$$\tau ::= \dots \mid \tau_1 \times \tau_2 \mid \tau_1 + \tau_2$$

(The generalizations to  $n$ -ary ( $n \geq 0$ ) products and sums are completely straightforward and thus omitted.) The associated new terms are:

$$\frac{\Gamma \vdash_{\Sigma} E_1 : \tau_1 \quad \Gamma \vdash_{\Sigma} E_2 : \tau_2}{\Gamma \vdash_{\Sigma} (E_1, E_2) : \tau_1 \times \tau_2} \quad \frac{\Gamma \vdash_{\Sigma} E : \tau_1 \times \tau_2 \quad \Gamma, x_1 : \tau_1, x_2 : \tau_2 \vdash_{\Sigma} E' : \tau}{\Gamma \vdash_{\Sigma} \mathbf{split}(E, x_1. x_2. E') : \tau}$$

$$\frac{\Gamma \vdash_{\Sigma} E : \tau_1}{\Gamma \vdash_{\Sigma} \mathbf{inl}(E) : \tau_1 + \tau_2} \quad \frac{\Gamma \vdash_{\Sigma} E : \tau_2}{\Gamma \vdash_{\Sigma} \mathbf{inr}(E) : \tau_1 + \tau_2}$$

$$\frac{\Gamma \vdash_{\Sigma} E : \tau_1 + \tau_2 \quad \Gamma, x_1 : \tau_1 \vdash_{\Sigma} E_1 : \tau \quad \Gamma, x_2 : \tau_2 \vdash_{\Sigma} E_2 : \tau}{\Gamma \vdash_{\Sigma} \mathbf{case}(E, x_1. E_1, x_2. E_2) : \tau}$$

(Instead of **split**, we could have used explicit projections, but the characterization of normal forms becomes more uniform with **split**. In practice, the separate split-construct above is usually folded into pattern-matching let- and lambda-bindings.)

*Semantics.* The semantics of the type constructors is standard:

$$\llbracket \tau_1 \times \tau_2 \rrbracket_v^{\mathcal{I}} = \llbracket \tau_1 \rrbracket_v^{\mathcal{I}} \times \llbracket \tau_2 \rrbracket_v^{\mathcal{I}} = \{(a_1, a_2) \mid a_1 \in \llbracket \tau_1 \rrbracket_v^{\mathcal{I}}, a_2 \in \llbracket \tau_2 \rrbracket_v^{\mathcal{I}}\}$$

$$\llbracket \tau_1 + \tau_2 \rrbracket_v^{\mathcal{I}} = \llbracket \tau_1 \rrbracket_v^{\mathcal{I}} + \llbracket \tau_2 \rrbracket_v^{\mathcal{I}} = \{\iota_1 a \mid a \in \llbracket \tau_1 \rrbracket_v^{\mathcal{I}}\} \cup \{\iota_2 a \mid a \in \llbracket \tau_2 \rrbracket_v^{\mathcal{I}}\}$$

as is the semantics of the associated terms:

$$\llbracket (E_1, E_2) \rrbracket_v^{\mathcal{I}} \rho = \llbracket E_1 \rrbracket_v^{\mathcal{I}} \rho \star \lambda a_1. \llbracket E_2 \rrbracket_v^{\mathcal{I}} \rho \star \lambda a_2. \eta(a_1, a_2)$$

$$\llbracket \mathbf{split}(E, x_1. x_2. E') \rrbracket_v^{\mathcal{I}} \rho = \llbracket E \rrbracket_v^{\mathcal{I}} \rho \star \lambda(a_1, a_2). \llbracket E' \rrbracket_v^{\mathcal{I}} (\rho[x_1 \mapsto a_1, x_2 \mapsto a_2])$$

$$\begin{aligned}
\llbracket \mathbf{inl}(E) \rrbracket_v^\tau \rho &= \llbracket E \rrbracket_v^\tau \rho \star \lambda a. \eta(\iota_1 a) \\
\llbracket \mathbf{inr}(E) \rrbracket_v^\tau \rho &= \llbracket E \rrbracket_v^\tau \rho \star \lambda a. \eta(\iota_2 a) \\
\llbracket \mathbf{case}(E, x_1. E_1, x_2. E_2) \rrbracket_v^\tau \rho &= \llbracket E \rrbracket_v^\tau \rho \star \lambda s. \begin{cases} \llbracket E_1 \rrbracket_v^\tau (\rho[x_1 \mapsto a_1]) & \text{if } s = \iota_1 a_1 \\ \llbracket E_2 \rrbracket_v^\tau (\rho[x_2 \mapsto a_2]) & \text{if } s = \iota_2 a_2 \end{cases}
\end{aligned}$$

*Normal forms.* With the addition of product and sum types, **CBV normal forms exhibit a striking similarity with cut-free proofs in Gentzen-style intuitionistic sequent calculus**, as also noted by Ohori [17]. In fact, we also get the usual sequent-calculus inconvenience of having to make arbitrary choices about the order in which we apply left-rules to decompose the types of variables. To keep normal forms unique, we choose to eliminate structured-type variables immediately as they are introduced, in a stack-like manner. (Immediate elimination leads to a slight anomaly for constants introduced by the signature: we can only allow  $\Sigma$  to declare constants of base and top-level-functional types. In the rare cases where we need, e.g., a sum-typed constant  $c : \tau_1 + \tau_2$  in  $\Sigma$ , it can be provided as a function  $c' : 1 \rightarrow \tau_1 + \tau_2$  where 1 is the zero-ary product type.)

We now have three mutually recursive notions of normality: **normal values**  $\Gamma \Vdash^v E : \tau$ , **normal computations**  $\Gamma \Vdash^c E : \tau$ , and **normal bodies**  $\Gamma \mid \Theta \Vdash^b E : \tau$ , where  $\Theta$  is an *ordered* list of typing assumptions:

$$\begin{array}{c}
\frac{\Sigma(c) = b}{\Gamma \Vdash^v c : b} \quad \frac{\Gamma(x) = b}{\Gamma \Vdash^v x : b} \quad \frac{\Gamma \mid x : \tau_1 \Vdash^b E : \tau_2}{\Gamma \Vdash^v \lambda x. \tau_1. E : \tau_1 \rightarrow \tau_2} \\
\\
\frac{\Gamma \Vdash^v E_1 : \tau_1 \quad \Gamma \Vdash^v E_2 : \tau_2}{\Gamma \Vdash^v (E_1, E_2) : \tau_1 \times \tau_2} \quad \frac{\Gamma \Vdash^v E : \tau_1}{\Gamma \Vdash^v \mathbf{inl}(E) : \tau_1 + \tau_2} \quad \frac{\Gamma \Vdash^v E : \tau_2}{\Gamma \Vdash^v \mathbf{inr}(E) : \tau_1 + \tau_2} \\
\\
\frac{\Gamma \Vdash^v E : \tau}{\Gamma \Vdash^c E : \tau} \quad \frac{\Sigma(c) = \tau_1 \rightarrow \tau_2 \quad \Gamma \Vdash^v E : \tau_1 \quad \Gamma \mid x : \tau_2 \Vdash^b E' : \tau}{\Gamma \Vdash^c \mathbf{let } x = c E \mathbf{ in } E' : \tau} \\
\\
\frac{\Gamma(x') = \tau_1 \rightarrow \tau_2 \quad \Gamma \Vdash^v E : \tau_1 \quad \Gamma \mid x : \tau_2 \Vdash^b E' : \tau}{\Gamma \Vdash^c \mathbf{let } x = x' E \mathbf{ in } E' : \tau} \\
\\
\frac{\Gamma \Vdash^c E : \tau}{\Gamma \mid \Vdash^b E : \tau} \quad \frac{\Gamma, x : b \mid \Theta \Vdash^b E : \tau}{\Gamma \mid \Theta, x : b \Vdash^b E : \tau} \quad \frac{\Gamma, x : \tau_1 \rightarrow \tau_2 \mid \Theta \Vdash^b E : \tau}{\Gamma \mid \Theta, x : \tau_1 \rightarrow \tau_2 \Vdash^b E : \tau} \\
\\
\frac{\Gamma \mid \Theta, x_1 : \tau_1, x_2 : \tau_2 \Vdash^b E : \tau}{\Gamma \mid \Theta, x : \tau_1 \times \tau_2 \Vdash^b \mathbf{split}(x, x_1. E) : \tau} \quad \frac{\Gamma \mid \Theta, x_1 : \tau_1 \Vdash^b E_1 : \tau \quad \Gamma \mid \Theta, x_2 : \tau_2 \Vdash^b E_2 : \tau}{\Gamma \mid \Theta, x : \tau_1 + \tau_2 \Vdash^b \mathbf{case}(x, x_1. E_1, x_2. E_2) : \tau}
\end{array}$$

Note how newly-introduced variables are put into the quarantined context  $\Theta$ , where their types are decomposed and the pieces migrate back into the ordinary context  $\Gamma$ . In particular, without the rules for product and sum types, the new definitions of normal values and computations agree exactly with the original ones.

*Normalization by evaluation.* We show only how the continuation-based residualizing interpretation can be extended. Products could be added to an accumulation-based interpretation without too much trouble, but sums apparently require the full power of applying a single continuation multiple times.

We assume the syntax-constructor functions for **E** are extended with functions *PAIR*, *SPLIT*, *INL*, *INR*, and *CASE* with the obvious types. Then we can

define additional helper functions analogous to  $bind_\tau$  from before:

$$\begin{aligned}
bindp_{\tau_1, \tau_2} &\in \mathbf{E} \rightarrow T^r(\mathbf{V} \times \mathbf{V}) \\
&= \lambda e. \lambda \kappa. new_{\tau_1} \star^g \lambda v_1. new_{\tau_2} \star^g \lambda v_2. \kappa(v_1, v_2) \star^g \lambda e'. \eta^g(SPLIT(e, v_1, v_2, e')) \\
binds_{\tau_1, \tau_2} &\in \mathbf{E} \rightarrow T^r(\mathbf{V} + \mathbf{V}) \\
&= \lambda e. \lambda \kappa. new_{\tau_1} \star^g \lambda v_1. new_{\tau_2} \star^g \lambda v_2. \kappa(\iota_1 v_1) \star^g \lambda e_1. \kappa(\iota_2 v_2) \star^g \lambda e_2. \\
&\quad \eta^g(CASE(e, v_1, e_1, v_2, e_2))
\end{aligned}$$

and the corresponding cases for the reification and reflection functions, extending the ones from Section 2.2:

$$\begin{aligned}
\downarrow^{\tau_1 \times \tau_2} &= \lambda(a_1, a_2). \downarrow^{\tau_1} a_1 \star^g \lambda e_1. \downarrow^{\tau_2} a_2 \star^g \lambda e_2. \eta^g(PAIR(e_1, e_2)) \\
\downarrow^{\tau_1 + \tau_2} &= \lambda s. \begin{cases} \downarrow^{\tau_1} a_1 \star^g \lambda e_1. \eta^g(INLe_1) & \text{if } s = \iota_1 a_1 \\ \downarrow^{\tau_2} a_2 \star^g \lambda e_2. \eta^g(INRe_2) & \text{if } s = \iota_2 a_2 \end{cases} \\
\uparrow_{\tau_1 \times \tau_2} &= \lambda e. bindp_{\tau_1, \tau_2} e \star^r \lambda(v_1, v_2). \uparrow_{\tau_1}(VAR v_1) \star^r \lambda a_1. \uparrow_{\tau_2}(VAR v_2) \star^r \lambda a_2. \eta^r(a_1, a_2) \\
\uparrow_{\tau_1 + \tau_2} &= \lambda e. binds_{\tau_1, \tau_2} e \star^r \lambda s. \begin{cases} \uparrow_{\tau_1}(VAR v_1) & \text{if } s = \iota_1 v_1 \\ \uparrow_{\tau_2}(VAR v_2) & \text{if } s = \iota_2 v_2 \end{cases}
\end{aligned}$$

The residualizing interpretation is as before. (Note that we cannot give a residualizing interpretation to constants of top-level product or sum type, since the reflection function does *not* factor through  $\eta^r$  in these cases.) The normalization function and Theorem 2 remain the same. For the implementation, the realizations of  $bindp$ ,  $binds$ , and the new clauses for reification and reflection in terms of continuation-manipulating primitives are straightforward.

## 5 Type-Directed Partial Evaluation

A primary application of normalization by evaluation is for type-directed partial evaluation (TDPE) [7]. The goal is to simplify a partially applied function of multiple arguments by propagating the values of the known arguments throughout the body of the function. Here, in addition to eliminating  $\beta$ -redexes, we also want to simplify occurrences of constants, such as arithmetic operations, when they are applied to literal values. In other words, we now want to normalize terms with *interpreted* base types and constants. There are in fact two natural ways to achieve this, both expressible in terms of the notion of constrained interpretations:

*Offline TDPE.* For offline TDPE [7, Section 3], we say that a program is *binding-time separated* if it is expressed over a signature  $\Sigma$  partitioned into a static and a dynamic part,  $\Sigma_s$  and  $\Sigma_d$ , each containing some type and term constants. The interpretation is likewise partitioned into  $\mathcal{I}_s$  and  $\mathcal{I}_d$ . We then constrain the allowable interpretations so that  $\mathcal{I}_s$  is always the standard interpretation (i.e.,  $\text{int}$  as  $\mathbf{Z}$ ,  $+$  as addition,  $\text{fix}$  as the domain-theoretic least fixed point, etc.), while the dynamic part remains completely unconstrained. That is, we consider the notion of *static equivalence*,  $\models^{Int_s} E = E'$  where  $Int_s = \{\mathcal{I} \mid \mathcal{I}|_{\Sigma_s} = \mathcal{I}_s\}$ . When the static normal form  $\bar{E}$  of a term  $E$  exists, it is then equivalent to  $E$  with

respect to all interpretations of  $\Sigma_d$ . The residualizing interpretation of  $\Sigma$ , like any interpretation in  $Int_s$ , also uses the standard interpretation of  $\Sigma_s$ , and the syntax-reconstructing interpretation from Section 3.2 for  $\mathcal{I}_d$ .

The separation allows us to realize the static part of the signature completely natively in terms of the corresponding construct of the programming language, and in fact we can use syntactic conveniences of the host language, such as pattern matching or **letrec**-forms directly for the static computations. It also becomes possible to self-apply the partial evaluator (the so-called second Futamura projection) [15].

*Online TDPE.* In the online variant [7, Section 4], like in general online partial evaluation, we do not annotate types, nor most occurrences of constants, with their binding times. Instead, the partial evaluator “opportunistically” propagates statically known data and performs reductions such as  $2 + 3 \rightarrow 5$ .

This corresponds to normalizing with respect to the set  $Int_c$  of interpretations  $(\mathcal{B}, \mathcal{C})$  that satisfy constraints such as  $\mathcal{C}(+) (\mathcal{C}(2), \mathcal{C}(3)) = \mathcal{C}(5)$ , and possibly also additional ones, such as  $\mathcal{C}(+) (x, \mathcal{C}(0)) = x$  for all  $x$  in  $\mathcal{B}(\text{int})$ . Again, the standard interpretation of  $\mathcal{C}(+)$  satisfies these constraints automatically. The residualizing one includes explicit checks for the reducible cases, to avoid constructing the corresponding redexes in the generated code. This formulation is similar to recent work on merging the reduction-free normalization of function abstraction and application with explicit reduction rules for constants [3].

The advantage of the online approach is that we do not have to explicitly separate the binding times in the source program, but correspondingly it becomes less predictable how much of the source program can be simplified at partial-evaluation time. An online normalizer also requires all primitive operations to explicitly check whether their arguments are literals (so the operation can be eliminated) or more general expressions (so the operation must remain in the normal form), slowing down the specialization process somewhat. Finally, fixed-point operators must still either be explicitly classified as static or dynamic, or need some ad hoc mechanism for deciding whether their unfolding equation should be applied.

In both cases, the normalization function may be partial, i.e., without part (0) of Theorems 1 and 2. This is unavoidable, since in the presence of recursion, some terms simply have no normal form. However, for the CBN case, one can still show that when it is defined,  $\tilde{E}$  satisfies parts (1-3), and also that whenever a  $\tilde{E}$  satisfying (1-3) exists,  $norm(E)$  is defined [11]. The situation for CBV TDPE has not been fully analyzed yet, although it seems reasonable to conjecture an analogous result.

In any case, the semantic treatment of TDPE allows us to uniformly analyze the construction and state its correctness criterion independently of the details of its implementation. That is, we can think about partial evaluation in terms of normalization with respect to a class of interpretations, without worrying about whether the normalization is achieved through repeated reductions, or through reduction-less normalization by evaluation.

## 6 Conclusions and Future Work

We have seen that the same basic idea that allows us to compute normal forms of lambda-terms with respect to purely functional interpretations also allows us to compute such normal forms with respect to general computational interpretations. In both cases, we chose a “quasi-syntactic” interpretation of the types and constants, and in the latter, also a binding-accumulating monad as the interpretation of computational effects. Both variants of semantic normalization can be phrased as functional-program evaluation, although the construction is significantly more involved in the computational case.

An important application of the normalization construction is type-directed partial evaluation, which seeks to compute normal forms not with respect to *all* interpretations of a signature, but only with respect to a subset of those; different choices of such subsets lead to offline and online partial evaluation. The offline, call-by-name case is analyzed in isolation in an earlier paper [11], but the more general constrained-interpretation formulation presented here seems worth investigating further, especially since it also leads to natural TDPE formulations of other partial-evaluation concepts, such as polyvariant program-point specialization expressed as suitable constraints on the recursion operator.

Additionally, it should be possible to obtain a syntactic analog of the CBV normalization result, showing that the computed normal form is not only equivalent to the original term with respect to arbitrary set-, or domain-theoretic interpretations, but is provably equal to it using the axioms of the computational lambda-calculus [16]. However, it seems that the semantic characterization may be ultimately more convenient for reasoning about programs, since it seems to scale more directly to partially constrained interpretations, and especially static recursion.

**Acknowledgments.** The author wishes to thank Olivier Danvy and Peter Dybjer for many fruitful discussions of the topics presented here, and the anonymous TLCA’01 reviewers for numerous helpful suggestions (of which regrettably not all could be followed here, due to space constraints.)

## References

1. Thorsten Altenkirch, Martin Hofmann, and Thomas Streicher. Categorical reconstruction of a reduction free normalization proof. In *Category Theory and Computer Science, 6th International Conference*, number 953 in Lecture Notes in Computer Science, 1995.
2. Ulrich Berger. Program extraction from normalization proofs. In M. Bezem and J. F. Groote, editors, *Typed Lambda Calculi and Applications*, number 664 in Lecture Notes in Computer Science, pages 91–106, Utrecht, The Netherlands, March 1993.
3. Ulrich Berger, Matthias Eberl, and Helmut Schwichtenberg. Normalization by evaluation. In *Prospects for Hardware Foundations (NADA)*, number 1546 in Lecture Notes in Computer Science, pages 117–137, 1998.

4. Ulrich Berger and Helmut Schwichtenberg. An inverse of the evaluation functional for typed  $\lambda$ -calculus. In *Proceedings of the Sixth Annual IEEE Symposium on Logic in Computer Science*, pages 203–211, Amsterdam, The Netherlands, July 1991.
5. Thierry Coquand and Peter Dybjer. Intuitionistic model constructions and normalization proofs. *Mathematical Structures in Computer Science*, 7:75–94, 1997.
6. Djordje Čubrić, Peter Dybjer, and Philip Scott. Normalization and the Yoneda embedding. *Mathematical Structures in Computer Science*, 8:153–192, 1998.
7. Olivier Danvy. Type-directed partial evaluation. In J. Hatchcliff, T. Æ. Mogensen, and P. Thieman, editors, *Partial Evaluation – Practice and Theory; Proceedings of the 1998 DIKU Summer School*, number 1706 in Lecture Notes in Computer Science, pages 367–411. Springer-Verlag, Copenhagen, Denmark, July 1998.
8. Olivier Danvy and Peter Dybjer, editors. *Preliminary Proceedings of the APPSEM Workshop on Normalization by Evaluation*. Department of Computer Science, University of Aarhus, May 1998. BRICS Note NS-98-1.
9. Olivier Danvy and Andrzej Filinski. Representing control: A study of the CPS transformation. *Mathematical Structures in Computer Science*, 2(4):361–391, December 1992.
10. Andrzej Filinski. Representing layered monads. In *Proceedings of the 26th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, pages 175–188, San Antonio, Texas, January 1999.
11. Andrzej Filinski. A semantic account of type-directed partial evaluation. In G. Nadathur, editor, *International Conference on Principles and Practice of Declarative Programming*, number 1702 in Lecture Notes in Computer Science, pages 378–395, Paris, France, September 1999.
12. Marcelo Fiore, Gordon Plotkin, and Daniele Turi. Abstract syntax and variable binding. In *Proceedings of the 14th Annual IEEE Symposium on Logic in Computer Science*, pages 193–202, Trento, Italy, July 1999.
13. Cormac Flanagan, Amr Sabry, Bruce F. Duba, and Matthias Felleisen. The essence of compiling with continuations. In *Proceedings of the SIGPLAN '93 Conference on Programming Language Design and Implementation*, pages 237–247, Albuquerque, New Mexico, June 1993.
14. Murdoch Gabbay and Andrew Pitts. A new approach to abstract syntax involving binders. In *Proceedings of the 14th Annual IEEE Symposium on Logic in Computer Science*, pages 214–224, Trento, Italy, July 1999.
15. Bernd Grobauer and Zhe Yang. The second Futamura projection for type-directed partial evaluation. In *ACM SIGPLAN Workshop on Partial Evaluation and Semantics-Based Program Manipulation*, pages 22–32. ACM Press, January 2000.
16. Eugenio Moggi. Computational lambda-calculus and monads. In *Proceedings of the Fourth Annual Symposium on Logic in Computer Science*, pages 14–23, Pacific Grove, California, June 1989. IEEE.
17. Atsushi Ohori. A Curry-Howard isomorphism for compilation and program execution. In J.-Y. Girard, editor, *Typed Lambda-Calculi and Applications*, volume 1581 of *Lecture Notes in Computer Science*, L'Aquila, Italy, April 1999.
18. Eijiro Sumii and Naoki Kobayashi. Online-and-offline partial evaluation: A mixed approach. In *ACM SIGPLAN Workshop on Partial Evaluation and Semantics-Based Program Manipulation*, pages 12–21. ACM Press, January 2000.