
LAZY EVALUATION AND DELIMITED CONTROL

RONALD GARCIA^a, ANDREW LUMSDAINE^b, AND AMR SABRY^c

^a Carnegie Mellon University
e-mail address: rxg@cs.cmu.edu

^{b,c} Indiana University
e-mail address: {lums,sabry}@cs.indiana.edu

ABSTRACT. The call-by-need lambda calculus provides an equational framework for reasoning syntactically about lazy evaluation. This paper examines its operational characteristics.

By a series of reasoning steps, we systematically unpack the standard-order reduction relation of the calculus and discover a novel abstract machine definition which, like the calculus, goes “under lambdas.” We prove that machine evaluation is equivalent to standard-order evaluation.

Unlike traditional abstract machines, delimited control plays a significant role in the machine’s behavior. In particular, the machine replaces the manipulation of a heap using store-based effects with disciplined management of the evaluation stack using control-based effects. In short, state is replaced with control.

To further articulate this observation, we present a simulation of call-by-need in a call-by-value language using delimited control operations.

1. INTRODUCTION

From early on, the connections between lazy evaluation (Friedman and Wise, 1976; Henderson and Morris, 1976) and control operations seemed strong. One of these seminal papers on lazy evaluation (Henderson and Morris, 1976) advocates laziness for its coroutine-like behavior. Specifically, it motivates lazy evaluation with a solution to the *same fringe* problem: how to determine if two trees share the same fringe without first flattening each tree and then comparing the resulting lists. A successful solution to the problem traverses just enough of the two trees to tell that they do not match. The same fringe problem is also addressed in Sussman and Steele’s original exposition of the Scheme programming

1998 ACM Subject Classification: D.3.1.

Key words and phrases: call-by-need, reduction semantics, abstract machines, delimited continuations, lambda calculus.

^a This work was supported by the National Science Foundation under Grant #0937060 to the Computing Research Association for the CIFellows Project.

^{b,c} This work was supported by NSF awards CSR/EHS 0720857 and CCF 0702717.

language (Sussman and Steele Jr., 1975). One of their solutions uses a continuation passing-style representation of coroutines. More recently, Biernacki et al. (2005) explores a number of continuation-based solutions to some fringe variants.

Same fringe is not the only programming problem that can be solved using either lazy evaluation or continuations. For instance, lazy streams and continuations are also used to implement and reason about backtracking (Wand and Vaillancourt, 2004; Kiselyov et al., 2005). Strong parallels in the literature have long suggested that lazy evaluation elegantly embodies a stylized use of coroutines. Indeed, we formalize this connection.

Call-by-need evaluation combines the equational reasoning capabilities of call-by-name with a more efficient implementation technology that systematically shares the results of some computations. However, call-by-need's evaluation strategy makes it difficult to reason about the operational behavior and space usage of programs. In particular, call-by-need evaluation obscures the control flow of evaluation. To facilitate reasoning, semantic models (Launchbury, 1993; Sestoft, 1997; Friedman et al., 2007; Nakata and Hasegawa, 2009), simulations (Okasaki et al., 1994), and tracing tools (Gibbons and Wansbrough, 1996) for call-by-need evaluation have been developed. Many of these artifacts use an explicit store or store-based side-effects (Wang, 1990) to represent values that are shared between parts of a program. Stores, being amorphous structures, make it difficult to establish program properties or analyze program execution. This representation of program execution loses information about the control structure of evaluation.

The call-by-need lambda calculus was introduced by Ariola et al. (1995) as an alternative to store-based formalizations of lazy evaluation. It is an equational framework for reasoning about call-by-need programs and languages. Following Plotkin (1975), these authors present a calculus and prove a standardization theorem that links the calculus to a complete and deterministic (i.e. standard order) reduction strategy. The calculus can be used to formally justify transformations, particularly compiler optimizations, because any terms it proves equal are also contextually equivalent under call-by-need evaluation.

Call-by-need calculi were investigated by two groups (Maraist et al., 1998; Ariola and Felleisen, 1997). The resulting two calculi are quite similar but their subtle differences yield trade-offs that are discussed in the respective papers. Nonetheless, both papers connect their calculi to similar standard-order reduction relations.

One notable feature of Ariola and Felleisen's calculus (and both standard-order reduction relations) is the use of evaluation contexts within the notions of reduction. It has been observed that evaluation contexts correspond to continuations in some presentations of language semantics (Felleisen and Friedman, 1986; Biernacka and Danvy, 2007). However, in these systems evaluation contexts are used to model variable references and demand-driven evaluation, not first-class continuations.

This paper exposes how Ariola et al.'s call-by-need evaluation relates to continuations. By systematically unpacking the standard-order reduction relation of the calculus, we discover a novel abstract machine that models call-by-need style laziness and sharing without using a store. Instead, the machine manipulates its evaluation context in a manner that corresponds to a stylized use of delimited control operations. The machine's behavior reveals a connection between control operations and laziness that was present but hidden in the reduction semantics.

To directly interpret this connection in the terminology of delimited control, we construct a simulation of call-by-need terms in the call-by-value language of Dybvig et al.

(2007), which provides a general framework for delimited continuations with first-class generative prompts.

Our concrete specifications of the relationship between call-by-need and delimited control firmly establish how lazy evaluation relates to continuations and other control-oriented language constructs and effects. Implementations of both the machine and the simulation are available at the following url:

<http://osl.iu.edu/~garcia/call-by-need.tgz>.

2. THE CALL-BY-NEED LAMBDA CALCULUS

The remainder of this paper examines Ariola and Felleisen’s formalization of call-by-need (Ariola and Felleisen, 1997). The terms of the calculus are standard:

$$t ::= x \mid \lambda x.t \mid t t$$

The call-by-need calculus, in direct correspondence with the call-by-value and call-by-name lambda calculi, distinguishes lambda abstractions as values:

$$v ::= \lambda x.t$$

Call-by-need is characterized by two fundamental properties: a computation is only performed when its value is needed, and the result of any computation is remembered and shared so that it only needs to be computed once. This calculus distinguishes two additional subsets of the term language to help represent these properties.

To capture the notion of *necessary* computations, the calculus distinguishes the set of lambda abstractions that immediately need the value of their argument. To define this set, the calculus appeals to a notion of evaluation contexts, a set of terms that each have a single hole (\square) in them:

$$E ::= \square \mid E t \mid (\lambda x.E[x]) E \mid (\lambda x.E) t$$

Given the evaluation contexts, the set of lambda abstractions in question is defined syntactically as $(\lambda x.E[x])$, the set of lambda abstractions whose bodies can be decomposed into an evaluation context and a free instance of the abstracted variable.

The intuition for this definition is as follows. Evaluation contexts are essentially terms with a single *hole* in them. When used in the style of Felleisen and Hieb (1992), evaluation contexts indicate those locations in a program that are subject to evaluation. As such, a lambda abstraction of the form $(\lambda x.E[x])$ will immediately refer to the value of its argument when its body is evaluated. Here, $E[x]$ is not a language construct, but rather metalinguistic notation for a side condition on the term in the body of the lambda abstraction.

The syntactic structure of call-by-need evaluation contexts give some hint to the nature of computation captured by this calculus. First, the context production $E t$ indicates that evaluation can focus on the operator position of an application regardless of the structure of the operand. This property is also true of call-by-name and call-by-value and reflected in their respective evaluation contexts. Second, the $(\lambda x.E[x]) E$ production indicates the operand of an application expression can be evaluated only if the operator is a lambda abstraction that immediately needs its argument. This restriction does not hold for call-by-value¹, where any lambda abstraction in operator position justifies evaluating the operand.

¹assuming left-to-right evaluation of application expressions

This restriction on evaluation corresponds with our intuitive understanding of call-by-need. Third, the $(\lambda x.E) t$ production indicates that evaluation can proceed *under a lambda abstraction* when it is in operator position. Though not immediately obvious, this trait is used by the calculus to capture the process of sharing computations among subterms.

To capture the notion of *shared* computations, the call-by-need calculus distinguishes lambda abstractions with explicit bindings for some variables, calling them *answers*:

$$a ::= v \mid (\lambda x.a) t$$

Answers are a syntactic representation of (partial) closures. An answer takes the form of a lambda term nested inside some applications. The surrounding applications simulate environment bindings for free variables in the nested lambda term. This representation makes it possible for the calculus to explicitly account for variable binding and to syntactically model how call-by-need evaluation shares lazily computed values.

The calculus has three notions of reduction:

$$\begin{array}{lll} (\lambda x.E[x]) v & \rightarrow_{\text{need}} & (\lambda x.E[v]) v \\ (\lambda x.a) t_1 t_2 & \rightarrow_{\text{need}} & (\lambda x.a t_2) t_1 \\ (\lambda x_1.E[x_1]) ((\lambda x_2.a) t_1) & \rightarrow_{\text{need}} & (\lambda x_2.(\lambda x_1.E[x_1]) a) t_1 \end{array}$$

The first reduction rule substitutes a value for a single variable instance in an abstraction. The rule retains the binding and abstraction so as to share its computation with other variable references as needed. The second and third reduction rules commute an application with an answer binding to expose opportunities for reduction without duplicating not-yet-needed computations. These two rules help to ensure that computations will be shared among references to a common variable.

As popularized by Barendregt (1981), each reduction assumes a hygiene convention. When combined with the evaluation contexts, the notions of reduction yield a deterministic standard order reduction relation (\mapsto_{sr}) and its reflexive-transitive closure (\mapsto_{sr}^*).

Definition 2.1. $t_1 \mapsto_{sr} t_2$ if and only if $t_1 \equiv E[t_r]$, $t_2 \equiv E[t_c]$ and $t_r \rightarrow_{\text{need}} t_c$.

Terms of the calculus satisfy unique decomposition, meaning that any program (i.e. closed term) that is not an answer can be decomposed exactly one way into a context E and redex t_r . This property ensures that \mapsto_{sr} is deterministic. Standard order reduction is an effective specification of call-by-need evaluation: if t is a program (i.e. closed term), then t call-by-need evaluates to an answer if and only if $t \mapsto_{sr}^* a$ for some answer a .

3. FROM REDUCTION SEMANTICS TO MACHINE SEMANTICS

Some reduction semantics have been shown to correspond directly to abstract machine semantics, thereby establishing the equivalence of a reducer and a tail-recursive abstract machine implementation (Felleisen and Friedman, 1986; Felleisen et al., 2009). In particular, Danvy and Nielsen (2004) introduces a method and outlines criteria for mechanically transforming reduction semantics into abstract machine semantics. However, proceeding directly from the reduction semantics for call-by-need to a tail-recursive abstract machine semantics poses some challenges that do not arise with other reduction semantics like call-by-name and call-by-value.

A straightforward call-by-need reducer implementation naïvely decomposes a term into a context and a redex. Any application could be one of three different redexes, each of

which is nontrivial to detect, so whenever the decompose function detects an application, it sequentially applies several recursive predicates to the term in hopes of detecting a redex. If the term is a redex, it returns; if not, it recursively decomposes the operator position.

Call-by-need redexes require more computational effort to recognize than either call-by-name or call-by-value. For instance, given a term t , only a fixed number of terminal operations are required to detect whether t is a call-by-name redex: one to check if the term is an application, one to access the operator position, and one to check if the operator is an abstraction.

Contrast this with the call-by-need redex $(\lambda x.E[x])(\lambda y.a)t$. Given a call-by-need term t_x , testing whether it matches this redex form requires an unknown number of operations: check if t_x is an application; check if its operator position is a lambda abstraction; check, in an unknown number of steps, if the operator’s body can be decomposed into $E[x]$, where x is both free in E and bound by the operator; and finally check, in an unknown number of steps, if the operand has the inductive structure of an answer.

To make matters worse, some terms can be decomposed into the form $E[x]$ in more than one way. For instance, consider the term $(\lambda x.(\lambda y.y) x)$. It can be decomposed as both $(\lambda x.E_1[y])$ and $(\lambda x.E_2[x])$ where $E_1 \equiv (\lambda y.\square) x$ and $E_2 \equiv (\lambda y.y) \square$. As such, a recursive procedure for decomposing a term cannot stop at the first variable it finds: it must be able to backtrack in a way that guarantees it will find the right decomposition $E[x]$ —and in turn the right redex—if there is one.

Recall that one of the evaluation contexts has the form $(\lambda x.E)t$. This means that redex evaluation can occur “under binders” (Moggi and Sabry, 2004; Kameyama et al., 2008). All three call-by-need notions of reduction shuffle lambda abstractions about in unusual ways. Furthermore, while reducing a recursive routine, a call-by-need evaluator may end up performing reductions under multiple copies of the same lambda abstraction. Call-by-name and call-by-value evaluators can address hygiene concerns by using environments and closures, but a call-by-need evaluator must prevent its evaluation context from incorrectly capturing free variable references. Any evaluator that goes under lambdas must pay particular attention to hygiene (Xi, 1997).

Since this work was originally published, Danvy et al. (2010) have adapted and extended the method of Danvy and Nielsen (2004) to produce a related abstract machine for call-by-need.

3.1. Towards an Abstract Machine. To find call-by-need redexes tail-recursively, we apply an insight from the CK abstract machine (Felleisen and Friedman, 1986; Felleisen et al., 2009). The CK machine implements an evaluation strategy for call-by-value based on a reduction semantics using the (inside-out) evaluation contexts \square , $E[\square t]$ and $E[(\lambda x.t) \square]$. To find a redex, the machine iteratively examines the outermost constructor of a term and uses the evaluation context to remember what has been discovered. Since call-by-value satisfies unique decomposition, this approach will find the standard redex if there is one.

To illustrate this in action, we walk through an example. Consider the program $(\lambda x.x) \lambda y.y$. Evaluation begins with configuration $\langle [], (\lambda x.x) \lambda y.y \rangle$. Underlining indicates subterms that the machine knows nothing about; at the beginning of evaluation, it knows nothing about the entire term. On the first step of the reduction, the machine detects that the term is an application $(\lambda x.x)$ $\lambda y.y$. To examine the term further, the machine must move its focus to either the operator or operand of this application. Since the machine is tail recursive, it must also push an evaluation context to store the as-yet uncovered structure

of the term. The only context it can reliably push at this point is $[\square \lambda y.y]$: it cannot push $[(\lambda x.x) \square]$ because it has not yet discovered that the operator position is a value. So the machine pushes the $[\square \lambda y.y]$ context, which serves as a reminder that it is focused on the operator of an application.

On the second step of the reduction, the machine detects that the operator is an abstraction $\lambda x.x$, and observes that the innermost context is $[\square \lambda y.y]$. In response, the machine pops the context, focuses on $\lambda y.y$, and pushes the context $[(\lambda x.x) \square]$, since the operator is now known to be a value. This context serves as a reminder that evaluation is currently focused on the operand of an application that can be reduced once that operand becomes a value.

On the third step, the machine detects the abstraction $(\lambda y.y)$, and remembers that the innermost context is $[(\lambda x.x) \square]$. At this point, the machine has deduced enough information to recognize the redex $(\lambda x.x) \lambda y.y$. This example illustrates how the CK machine uses a depth-first left-to-right search strategy to detect call-by-value redexes.

Now consider the same term under call-by-need using a similar strategy. As with call-by-value, the top-level application can be detected, the operand can be pushed onto the evaluation context, and the operator can be exposed as the abstraction $\lambda x.x$. At this point behavior must diverge from call-by-value because the body of the abstraction is still unknown and call-by-need does not have $[(\lambda x.t) \square]$ contexts for arbitrary t . However, call-by-need does have contexts of the form $[(\lambda x.\square) \lambda y.y]$. Therefore, it is possible to proceed under the first lambda abstraction, push the context, and focus on x .

The term is exposed as a variable x , which combines with the context $[(\lambda x.\square) \lambda y.y]$ to form the term $(\lambda x.E[x]) \lambda y.y$ (where $E \equiv \square$). At this point, enough information has been uncovered to push the context $[(\lambda x.\square[x]) \square]$ and focus on $\lambda y.y$. The abstraction $\lambda y.y$ is recognized, and with that a call-by-need redex $(\lambda x.\square[x]) \lambda y.y$ has been found. Success with this example suggests a promising strategy for implementing call-by-need reduction tail-recursively.

3.2. An Initial Abstract Machine. In this section, we elaborate the above search strategy into a simple but inefficient tail-recursive abstract machine. We present it without proof and then by a series of correct transformations we derive an efficient machine that we prove correct.

This abstract machine uses the same terms, values, and answers as the calculus. However, it introduces two alternate notions. First, the machine uses a more versatile representation of evaluation contexts. As observed in Danvy and Nielsen (2004), evaluation contexts can be mathematically specified in more than one way. For optimal flexibility, we define evaluation contexts as lists of frames, where the empty list $[\]$ and single-frame lists $[f]$ are our simple units, and the operator \circ stands for list concatenation.

$$\begin{array}{l}
 f ::= \square t \mid (\kappa x.E) \square \mid (\lambda x.\square) t \\
 E ::= [\] \mid [f] \circ E \mid E \circ [f] \\
 \text{where } E \circ [\] = [\] \circ E = E \\
 \text{and } E_1 \circ (E_2 \circ E_3) = (E_1 \circ E_2) \circ E_3
 \end{array}$$

When two contexts are composed, the second context is plugged into the hole of the first context: for example $[\square t_2] \circ [\square t_1] = [(\square t_1) t_2]$.

We call the frame $[(\lambda x. \square) t]$ a *binder* frame. It represents a variable binding in the context. It can be read as $[\text{let } x = t \text{ in } \square]$, but we use the former notation to emphasize that call-by-need evaluation proceeds under lambdas. This observation motivates our analysis of hygiene in Section 4.5.

We call the frame $[(\kappa x. E) \square]$ a *cont* frame, in reference to continuations. The construction $(\kappa x. E)$ is called a cont and replaces the metalinguistic term notation $(\lambda x. E[x])$ from the calculus. We use a different notation for conts than lambda terms to indicate that in the machine conts are distinct from terms (they are of type **Cont** rather than type **Term** in an implementation). Conts indicate nontrivial structural knowledge that the machine retains as it searches for a redex. This distinction matters when we establish continuation semantics for machine states. As we shall see, a cont frame represents a suspended variable reference.

Finally we call the frame $[\square t]$ an *operand* frame, and it represents a term waiting for an abstraction.

The abstract machine also introduces a notion of *redexes*:

$$r ::= a \ t \mid (\kappa x. E) \ a$$

Redexes are distinguished from terms in the machine, meaning that in an implementation, the type **Redex** is distinct from the type **Term**. This distinction suggests that conts $(\kappa x. E)$ are neither terms nor first-class entities in the call-by-need language: they only appear in evaluation contexts and in redexes. As we discuss below, the machine distinguishes one more kind of redex than the calculus.

The transition rules for the machine are staged into four distinct groups: refocus, rebuild, need, and reduce. Each machine configuration can be related to a term in the language of the calculus. The *refocus* rules examine the current term and push as many operand frames $[\square t]$ as possible. A refocus configuration $\langle E, t \rangle_f$ represents the term $E[t]$.

$$\boxed{\langle E, t \rangle_f} \quad (\text{Refocus})$$

$$\begin{aligned} \langle E, x \rangle_f &\mapsto \langle E, [\], x \rangle_n \\ \langle E, \lambda x. t \rangle_f &\mapsto \langle E, \lambda x. t \rangle_b \\ \langle E, t_1 \ t_2 \rangle_f &\mapsto \langle E \circ [\square t_2], t_1 \rangle_f \end{aligned}$$

Upon reaching a variable, refocus transitions to the need rules; upon reaching a lambda abstraction, it transitions to the rebuild rules.

The *rebuild* rules search up into the context surrounding an answer for the next applicable redex. A rebuild configuration $\langle E, a \rangle_b$ represents the term $E[a]$.

$$\boxed{\langle E, a \rangle_b} \quad (\text{Rebuild})$$

$$\begin{aligned} \langle [\], a \rangle_b &\mapsto a \\ \langle E \circ [\square t_1], a \rangle_b &\mapsto \langle E, a \ t_1 \rangle_d \\ \langle E \circ [(\lambda x. \square) t_1], a \rangle_b &\mapsto \langle E, (\lambda x. a) \ t_1 \rangle_b \\ \langle E_1 \circ [(\kappa x. E_2) \square], a \rangle_b &\mapsto \langle E_1, (\kappa x. E_2) \ a \rangle_d \end{aligned}$$

These rules examine the current context and proceed to build a maximal answer-shaped term, progressively wrapping binder frames around the current answer. If the entire context is consumed then evaluation has completed and the entire program is an answer. Upon

reaching an operand or cont frame, a redex has been found, and rebuild transitions to the reduce rules. These rules resemble the $\text{refocus}_{\text{aux}}$ rules of Danvy and Nielsen (2004).

The *need* rules also examine the context, but they search for the binder frame that corresponds to the variable under focus. A need configuration $\langle E_1, E_2, x \rangle_n$ represents the term $E_1[E_2[x]]$.

$\langle E, E, x \rangle_n$ (Need)
$\langle E_1 \circ [(\lambda x. \square) t], E_2, x \rangle_n \mapsto \langle E_1 \circ [(\kappa x. E_2) \square], t \rangle_f$
$\langle E_1 \circ [f], E_2, x \rangle_n \mapsto \langle E_1, [f] \circ E_2, x \rangle_n$
where, $[f] \neq [(\lambda x. \square) t]$

Since input programs are closed, the associated binder must be somewhere in the context. Upon finding the right binder frame, a cont frame $[(\kappa x. E) \square]$ is pushed onto the context and evaluation proceeds to refocus on the operand from the associated binder frame.

The *reduce* rules simulate the notions of reduction from the calculus. A reduce configuration $\langle E, r \rangle_d$ represents the term $E[r]$ where a cont $\kappa x. E$ represents the term $\lambda x. E[x]$.

$\langle E, r \rangle_d$ (Reduce)
$\langle E_1, (\kappa x. E_2) v \rangle_d \mapsto \langle E_1, (\lambda x. E_2[v]) v \rangle_f$
$\langle E_1, (\kappa x_1. E_2) ((\lambda x_2. a) t) \rangle_d \mapsto \langle E_1, (\lambda x_2. (\lambda x_1. E_2[x_1]) a) t \rangle_f$
$\langle E, (\lambda x. a) t_1 t_2 \rangle_d \mapsto \langle E, (\lambda x. a t_2) t_1 \rangle_f$
$\langle E, (\lambda x. t_1) t_2 \rangle_d \mapsto \langle E \circ [(\lambda x. \square) t_2], t_1 \rangle_f$

Each of the first two reduce rules transforms a cont into a lambda abstraction by plugging its context with a term and abstracting its variable. As such, each reduce rule transforms a redex into a pure term of the calculus and transitions to a refocus configuration, which searches for the next redex.

The reduce rules also handle terms of the form $(\lambda x. t_1) t_2$, even though such terms are not call-by-need redexes. Including this rule gives the set of redexes greater uniformity: all terms of the form $a t$ are redexes, just like the terms of the form $(\kappa x. E) a$. This symmetry is not exhibited in the call-by-need calculus. However, Ariola and Felleisen (1997) defines and uses an auxiliary *let calculus* that adds the reduction

$$(\lambda x. t_1) t_2 \rightarrow_{\text{need}} \text{let } x = t_2 \text{ in } t_1$$

to the calculus and defines the other reductions in terms of the let expressions. The fourth reduce rule corresponds to this reduction rule. However, our presentation shows that an auxiliary let term, though compatible with this model, is not needed to specify call-by-need: the syntax of pure calculus terms suffices. Furthermore, treating this rule as a reduction here anticipates a change we make to the machine (Section 4.5) that adds explicit variable renaming to that rule. Finally, the reduce rules are improved in the next section so that all reduce rules change their representative terms nontrivially.

Machine evaluation of a program t begins with the refocus configuration $\langle [], t \rangle_f$ and terminates if it arrives at an answer a . Its behavior in between can be summarized as follows: search downwards until a value or variable reference is reached. If a variable reference is reached, store a cont in the context to remember the variable reference and proceed to

evaluate its binding. If an abstraction is reached, accumulate an answer up to the innermost redex, or the top of the evaluation context if none is found. In short, the machine performs a depth-first, left-to-right traversal in search of a call-by-need redex. Along the way it uses the evaluation context to store and retrieve information about program structure, particularly the location of variable bindings (using binder frames) and variable references (using cont frames). The refocus, rebuild, and need rules leave the term representation of their configurations unchanged (e.g. if $\langle E_1, t_1 \rangle_f \mapsto \langle E_2, t_2 \rangle_f$ then $E_1[t_1] \equiv E_2[t_2]$), and the reduce rules embody the notions of reduction from the calculus.

Our strategy for producing this machine builds on the strategy of Danvy and Nielsen (2004), which mechanizes the direct transformation of reduction semantics into abstract machine semantics. That report introduces and verifies a general method for using reduction semantics that meet certain criteria to construct a function that “refocuses” an arbitrary term-context pair to a redex-context pair. The resulting function can then be used to induce an abstract machine semantics. Unfortunately that refocus function construction does not apply to the call-by-need lambda calculus because the calculus does not meet the required criteria. In particular, the construction requires that a maximally-decomposed closed term (i.e. program) will focus on a value or a redex. However, call-by-need evaluation contexts can decompose down to variable references, which are neither redexes nor values under call-by-need. There are however other ways to produce an abstract machine from a reduction semantics which may apply to the call-by-need calculus studied here. For instance, Danvy et al. (2010) devise a variant of the let-based call-by-need reduction semantics, implement it, and use a program-transformation based approach to produce a refocus function and abstract machine implementation.

The following partial trace demonstrates how the initial abstract machine discovers the first redex for our running example $(\lambda x.x) \lambda y.y$:

$$\begin{array}{l}
\langle [], (\lambda x.x) \lambda y.y \rangle_f \quad \mapsto \quad \langle [\square \lambda y.y], \lambda x.x \rangle_f \quad \mapsto \quad \langle [\square \lambda y.y], \lambda x.x \rangle_b \\
\mapsto \langle [], (\lambda x.x) \lambda y.y \rangle_d \quad \mapsto \quad \langle [(\lambda x.\square) \lambda y.y], x \rangle_f \quad \mapsto \quad \langle [(\lambda x.\square) \lambda y.y], [], x \rangle_n \\
\mapsto \langle [(\kappa x.[]) \square], \lambda y.y \rangle_f \quad \mapsto \quad \langle [(\kappa x.[]) \square], \lambda y.y \rangle_b \quad \mapsto \quad \langle [], (\kappa x.[]) \lambda y.y \rangle_d
\end{array}$$

4. REFINING THE MACHINE

In this section we study the behavior of the abstract machine and make some improvements based on our observations. These changes lead us from the initial machine above to our final machine specification.

4.1. Grabbing and Pushing Conts. The need rules linearly search the nearest context for a binder frame that matches the variable under question. This process can be specified as one step:

$$\begin{array}{c}
\langle E_1 \circ [(\lambda x.\square) t] \circ E_2, x \rangle_n \mapsto \langle E_1 \circ [(\kappa x.E_2) \square], t \rangle_f \\
\text{where } [(\lambda x.\square) t] \notin E_2
\end{array}$$

This evaluation step accumulates a segment of the current evaluation context and stores it. In general, abstract machines that model control operators represent control capture in a similar manner. In this particular case, only part of the evaluation context is captured, and the amount of context captured depends on the dynamic location in the context of a certain

frame. As such, the need rules seem to perform some kind of *delimited* control capture. This analogy becomes stronger upon further analysis of the first reduce rule from Section 3.2. The machine uses its structural knowledge of $\kappa x.E$ to construct the abstraction $\lambda x.E[v]$. However, the resulting machine configuration no longer retains any of the structure that had previously been discovered. Recall our example execution trace from Section 3.2. The machine reduces the redex found at the end of that trace as follows:

$$\langle [], (\kappa x.[]) \lambda y.\underline{y} \rangle_d \mapsto \langle [], (\lambda x.\lambda y.y) \lambda y.\underline{y} \rangle_f$$

By returning to refocus following the reduction, the machine loses all structural knowledge of the term. To continue execution, it must examine the structure of the contractum from scratch. Fortunately, the evaluator can be safely improved so that it retains knowledge of the contractum's structure:

Proposition 4.1.

$$\langle E_1, (\lambda x.E_2[v]) v \rangle_f \mapsto \langle E_1 \circ [(\lambda x.\square) v] \circ E_2, v \rangle_b$$

Proof. Corollary of $\langle E_1, E_2[v] \rangle_f \mapsto \langle E_1 \circ E_2, v \rangle_b$, which is proven by induction on E_2 . \square

This proposition justifies replacing the first reduce rule with one that pushes the evaluation context embedded in the cont and proceeds to rebuild an answer:

$$\langle E_1, (\kappa x.E_2) v \rangle_d \mapsto \langle E_1 \circ [(\lambda x.\square) v] \circ E_2, v \rangle_b$$

This short-circuit rule extends the current evaluation context with a binder frame and the context E_2 that was inside the cont. The rule is suggestive of delimited control because machine models of control operators generally represent the reinstatement of delimited continuations by extending the current context with a piece of captured evaluation context. Of more immediate interest, though, is how reduction of our example now proceeds:

$$\langle [], (\kappa x.[]) \lambda y.\underline{y} \rangle_d \mapsto \langle [(\lambda x.\square) (\lambda y.\underline{y})], \lambda y.\underline{y} \rangle_b$$

All knowledge of the contractum's structure is retained, though much of it is now stored in the evaluation context.

4.2. Shifting Binder Frames. The second and third reduce rules from Section 3.2 also discard structural information. Specifically, they both transition to the forgetful refocus rule. However their information can be preserved.

Proposition 4.2.

$$\langle E, (\lambda x.a t_2) t_1 \rangle_f \mapsto \langle E \circ [(\lambda x.\square) t_1], a t_2 \rangle_d.$$

Proof. Corollary of $\langle E_1, a \rangle_f \mapsto \langle E_1, a \rangle_b$, which is proven by induction on a . \square

Proposition 4.3. *If E_2 does not capture x_1 (Section 4.5), then*

$$\langle E_1, (\lambda x_2.(\lambda x_1.E_2[x_1]) a) t \rangle_f \mapsto \langle E_1 \circ [(\lambda x_2.\square) t], (\kappa x_1.E_2) a \rangle_d.$$

Proof. Corollary of $\langle E_1, a \rangle_f \mapsto \langle E_1, a \rangle_b$ and $\langle E_1, (\lambda x_1.E_2[x_1]) t \rangle_f \mapsto \langle E_1 \circ [(\kappa x_1.E_2) \square], t \rangle_f$, which is proven by case analysis and induction on E_2 . \square

These propositions justify short-circuiting the respective evaluation rules. The new rules improve the behavior of the abstract machine.

$$\begin{aligned} \langle E_1, (\kappa x_1.E_2) ((\lambda x_2.a) t) \rangle_d &\mapsto \langle E_1 \circ [(\lambda x_2.\square) t], (\kappa x_1.E_2) a \rangle_d \\ \langle E, (\lambda x.a) t_1 t_2 \rangle_d &\mapsto \langle E \circ [(\lambda x.\square) t_1], a t_2 \rangle_d \end{aligned}$$

By fast-forwarding to reduce, the rules retain the discovered term structure and thereby avoid retracing the same terms again.

4.3. Answer = Binders \times Value. The transition rules repeatedly create binder frames out of terms and reabsorb those frames into answers. In this section we simplify this protocol. We distinguish answers from terms by providing them a separate representation:

$$a ::= \llbracket E, v \rrbracket, \text{ where } E = \overline{[(\lambda x_i.\square) t_i]}$$

An answer is now represented as a tuple containing the nested lambda abstraction and the sequence of binder frames that are wrapped around it in the old presentation (we use overlines to indicate sequences). This presentation bears strong similarity to calculi with explicit substitutions (Abadi et al., 1991) in that each binder frame $\overline{[(\lambda x.\square) t]}$ corresponds to a substitution $[t/x]$. An answer can be seen as a lambda term nested inside a sequence of explicit substitutions, $v \overline{[t_i/x_i]}$.

The rebuild rules could be reformulated as a three place configuration, $\langle E, E, v \rangle_b$, but instead we immediately apply the same improvement that we applied to the need rules in Section 4.1. For instance, the new transition rule for rebuilding to a cont frame is:

$$\begin{aligned} \langle E_1 \circ [(\kappa x.E_2) \square] \circ E_3, v \rangle_b &\mapsto \langle E_1, (\kappa x.E_2) \llbracket E_3, v \rrbracket \rangle_d \\ \text{where } E_3 &= \overline{[(\lambda x_i.\square) t_i]} \end{aligned}$$

Returning to our running example, reduction from its most recent state (at the end of Section 4.1) transitions to a final answer, signaling the end of execution:

$$\langle \overline{[(\lambda x.\square) \lambda y.y]}, \lambda y.y \rangle_b \mapsto \langle \llbracket \overline{[(\lambda x.\square) \lambda y.y]}, \lambda y.y \rrbracket \rangle$$

4.4. Aggregate Reduction. Now that answers explicitly carry their binders in aggregate, the reduce rules can be substantially consolidated. Currently, the second and third reduce rules iteratively remove the top binder frame from an answer and push it onto the evaluation context. This process repeats until the answer is just a lambda abstraction. At that point, the second and third reduce rules defer to the first and fourth reduce rules respectively. This corresponds exactly with standard-order reduction (cf. Definition 2.1):

Proposition 4.4.

$$\begin{aligned} E[\overline{[(\lambda x_n. \dots ((\lambda x_1. ((\lambda x_0.v) t_0)) t_1) \dots) t_n] t}] &\mapsto_{sr} \\ E[\overline{[(\lambda x_n. \dots ((\lambda x_1. ((\lambda x_0.v t) t_0)) t_1) \dots) t_n]}. \end{aligned}$$

$$\begin{aligned} E[(\lambda x.E[x]) \overline{[(\lambda x_n. \dots ((\lambda x_1. ((\lambda x_0.v) t_0)) t_1) \dots) t_n]}] &\mapsto_{sr} \\ E[\overline{[(\lambda x_n. \dots ((\lambda x_1. ((\lambda x_0.(\lambda x.E[x]) v) t_0)) t_1) \dots) t_n]}. \end{aligned}$$

Proof. By induction on the structure of the answer term, using the unique decomposition lemma of Ariola and Felleisen (1997). \square

Using the new answer representation, each pair of associated reduce rules can be merged into one omnibus rule that moves all the binder frames at once and simultaneously performs a reduction using the underlying value.

$$\begin{aligned} \langle E_1, (\kappa x.E_2) \llbracket E_3, v \rrbracket \rangle_d &\mapsto \langle E_1 \circ E_3 \circ [(\lambda x.\square) v] \circ E_2, v \rangle_b \\ \langle E_1, \llbracket E_2, (\lambda x.t_1) \rrbracket t_2 \rangle_d &\mapsto \langle E_1 \circ E_2 \circ [(\lambda x.\square) t_2], t_1 \rangle_f \end{aligned}$$

As a result of these transformations, both conts and answers contain evaluation contexts. Furthermore, conts and answers are not terms of the calculus, and the machine never reverts a cont or answer to a term. The rules that create them, rebuild for answers and need for conts, capture part of the evaluation context, and the rules that consume them, the reduce rules, reinstate the captured contexts.

4.5. Variable Hygiene. Presentations of calculi often invoke a hygiene convention and from then on pay little attention to bound or free variables. In this manner, calculi do not commit to any of the numerous ways that hygiene can be enforced. Many abstract machines, however, use environments or explicit sources of fresh names to guarantee hygiene and thereby provide a closer correspondence to concrete implementations. In this section, we augment the call-by-need machine with simple measures to enforce hygiene.

Our primary hygiene concerns are that evaluation occurs under binders and binders are shifted about in unusual ways. In order to ensure that binding structure is preserved throughout evaluation, we need to be able to reason locally, within each machine configuration, about bound variables. To make this possible, we make one last change to the machine. We add a list of names to each machine configuration.

$$X ::= \bar{x}_i$$

Most of the machine rules simply pass the list of names along to the next transition. One of the reduce rules manipulates the list of names.

$$\langle X \mid E_1, \llbracket E_2, \lambda x.t_1 \rrbracket t_2 \rangle_d \mapsto_{nam} \langle X, x' \mid E_1 \circ E_2 \circ [(\lambda x'.\square) t_2], t_1[x'/x] \rangle_f \quad x' \notin X \quad (\text{D.2})$$

When this rule goes under a lambda, it adds the name of its bound variable to the list X of variables. The notation X, x expresses adding a new name x to X . If the bound variable x on the left hand side of the rule is already a member of X , then the variable is renamed as part of the transition. As such, X can be considered a set.

Now each machine configuration has one of five forms:

$$\begin{aligned} \langle X \mid E, ? \rangle &::= \langle X \mid \llbracket E, v \rrbracket \rangle \mid \langle X \mid E, r \rangle_d \mid \langle X \mid E, t \rangle_f \\ &\mid \langle X \mid E, v \rangle_b \mid \langle X \mid E, x \rangle_n \end{aligned}$$

We use the notation $\langle X \mid E, ? \rangle$ below to uniformly discuss all configuration types, where X refers to the list of names, E refers to the context, and $?$ refers to the term or redex. For a final configuration $\langle X \mid \llbracket E, v \rrbracket \rangle$, $?$ refers to the answer's underlying value v , and E corresponds to the answer's binder frames E . We use the metavariable C to range over configurations when the internal structure does not matter.

The call-by-need abstract machine uses the set X of names to keep track of *active variables*: any variable x whose binding instance has been pushed into a binder frame $[(\lambda x. \square) t]$:

$AV([\])$	$= \emptyset$
$AV([\!(\lambda x. \square) t\!] \circ E)$	$= \{x\} \cup AV(E)$
$AV([\!\square t\!] \circ E)$	$= AV(E)$
$AV([\!(\kappa x. E_1) \square\!] \circ E)$	$= AV(E_1) \cup \{x\} \cup AV(E)$

Cont-bound variables are counted among the active variables because machine evaluation must have gone under a binding to construct the cont frame.

The renaming condition on the (D.2) reduce rule ensures that active variables are mutually distinguishable. This guarantees that the machine's need rule can never capture the wrong evaluation context and thus execute the wrong bound expression.

Renaming is not obviously sufficient to ensure bound variable hygiene because of how the machine manipulates evaluation contexts. For instance, even though the need rule is guaranteed to only match a variable with the right binder frame, we have no guarantee that the right binder frame could never be trapped inside a cont frame and hidden from view while a need transition searches for it. Were this to happen, the machine would get stuck. Furthermore, the reduction rules flip and shift evaluation contexts that might contain binder frames. If a binder frame were to end up below another context frame that contains references to its bound variable, then those references would no longer be bound in the context; the need rule would exhaust the evaluation context if it attempted to resolve any of these references.

To verify that machine evaluation is well-formed, we establish well-formedness conditions that suffice to ensure safe evaluation and we show that they hold for evaluation of all programs. The well-formedness conditions rely on straightforward notions of captured context variables (CV) and free context variables (FV):

$CV([\])$	$= \emptyset$
$CV(E \circ [\!(\lambda x. \square) t\!])$	$= \{x\} \cup CV(E)$
$CV(E \circ [\!\square t\!])$	$= CV(E)$
$CV(E \circ [\!(\kappa x. E_1) \square\!])$	$= CV(E)$
$FV([\])$	$= \emptyset$
$FV([\!(\lambda x. \square) t\!] \circ E)$	$= FV(t) \cup (FV(E) - \{x\})$
$FV([\!\square t\!] \circ E)$	$= FV(E) \cup FV(t)$
$FV([\!(\kappa x. E_1) \square\!] \circ E)$	$= FV(E) \cup (FV(E_1) - \{x\})$

As expected, these two notions are related.

Lemma 4.5.

- (1) $FV(E_1) \subseteq FV(E_1 \circ E_2)$.
- (2) $CV(E_1 \circ E_2) = CV(E_1) \cup CV(E_2)$.
- (3) $FV(E_1 \circ E_2) = FV(E_1) \cup (FV(E_2) \setminus CV(E_1))$.

Proof. By induction on the length of E_1 , E_2 , and E_1 respectively. □

To establish that binder frames remain properly positioned, we define a notion of well-formed evaluation contexts:

Definition 4.6. A Machine context/name pair is well formed, notation $X \mid E \mathbf{wf}$, iff:

- (1) $FV(E) = \emptyset$;
- (2) $AV(E) = \{x \mid x \text{ occurs in } X\}$
- (3) Each active variable of E is distinct: if $E_1 \circ E_2 \in E$ then $AV(E_1) \cap AV(E_2) = \emptyset$, and if $[(\kappa x.E_1) \square] \in E$ then $x \notin AV(E_1)$.

These well-formedness criteria ensure that a context has no unbound variable references, that active variables cannot interfere with each other, and that X is simply a particular ordering of the E 's active variables. The captured variables of E are also distinct since every captured variable is an active variable.

Furthermore, machine configurations also have a notion of well-formedness:

$$\boxed{\begin{array}{c} \frac{X \mid E \mathbf{wf} \quad FV(?) \subseteq CV(E)}{\langle X \mid E, ? \rangle_c \mathbf{wf}} \quad c \neq d \\ \\ \frac{X \mid (E_1 \circ [(\kappa x.E_2) \square] \circ E_3) \mathbf{wf} \quad FV(v) \subseteq CV(E_1 \circ [(\kappa x.E_2) \square] \circ E_3)}{\langle X \mid E_1, (\kappa x.E_2) \llbracket E_3, v \rrbracket \rangle_d \mathbf{wf}} \\ \\ \frac{X \mid (E_1 \circ [\square t] \circ E_2) \mathbf{wf} \quad FV(v) \subseteq CV(E_1 \circ [\square t] \circ E_2)}{\langle X \mid E_1, \llbracket E_2, v \rrbracket t_2 \rangle_d \mathbf{wf}} \end{array}}$$

These well-formedness conditions ensure that the evaluation contexts E are well-formed, that the list of names X matches the active variables of E , and that the free variables of the term under focus are captured by the context. To account for redexes, the well-formedness conditions for each reduce configuration reflect the well-formedness conditions for the corresponding rebuild configuration. As shown below, well-formed reduce configurations $\langle X \mid E, r \rangle_d \mathbf{wf}$ ensure that reduce rules can be safely performed without implicit renaming.

Well-formedness of the reduce configurations ensures that the reduce transitions require no implicit bound-variable renaming to preserve hygiene. Well-formedness of the need configuration guarantees that it cannot be “stuck”: since $x \in CV(E)$, a well-formed need configuration always has a binder frame $[(\lambda x.\square) t]$ that matches the variable under focus, so the configuration can transition.

Well-formedness of configurations combined with rule $D.2$'s name management ensures that machine evaluation respects variable binding structure.

Theorem 4.7. *If t is a closed term of the calculus, then $\langle \emptyset \mid [], t \rangle_f \mathbf{wf}$.*

Proof. $\emptyset \mid [] \mathbf{wf}$ and $FV(t) \subseteq CV([]) = \emptyset$. □

Theorem 4.8. *Let C_1 and C_2 be configurations. If $C_1 \mathbf{wf}$ and $C_1 \mapsto_{nam} C_2$ then $C_2 \mathbf{wf}$.*

Proof. By cases on \mapsto_{nam} . The cases are immediate except when $C_1 \mapsto_{nam} C_2$ by a reduce rule. The proofs for both kinds of reduce configurations are similar, so we present only one of the cases:

Case ($C_1 = \langle X \mid E_1, (\kappa x.E_2) \llbracket E_3, v \rrbracket \rangle_d$).

By definition, $C_1 \mapsto_{nam} C_2 = \langle X \mid E, v \rangle_b$, where $E = E_1 \circ E_3 \circ [(\lambda x.\square) v] \circ E_2$. Since the

transition rule introduces no new active variables, X should remain the same. Furthermore, all the active variables remain distinct, though x is now introduced by the binder frame $[(\lambda x. \square) v]$ rather than the cont frame $[(\kappa x. E_2) \square]$. It remains to show that $FV(v) \subseteq CV(E)$ and that $FV(E) = \emptyset$.

First, since C_1 **wf**, it follows by inversion that $X \mid (E_1 \circ [(\kappa x. E_2) \square] \circ E_3)$ **wf** and

$$FV(v) \subseteq CV(E_1 \circ [(\kappa x. E_2) \square] \circ E_3) = CV(E_1 \circ E_3) \subseteq CV(E).$$

By Lemma 4.5, $FV(E_1) = FV(E_1 \circ [(\kappa x. E_2) \square]) = FV(E_1 \circ [(\kappa x. E_2) \square] \circ E_3) = \emptyset$.

Since $CV([(\kappa x. E_2) \square]) = \emptyset$, it follows from Lemma 4.5 that $FV(E_3) \subseteq CV(E_1)$ and from the definition of FV that $FV(E_2) \subseteq CV(E_1) \cup \{x\}$. From these it follows that $FV(E) = \emptyset$.

□

In short, well-formedness of the reduce configurations ensures that the reduce rules can be safely performed without any implicit renaming. Since the machine preserves well-formedness, this property persists throughout evaluation. The rest of this paper only considers well-formed configurations.

4.6. An Abstract Machine for Call-by-need. Putting together our observations from the previous section, we now present the specification of the abstract machine. Figure 1 presents its \mapsto_{nam} transitions rules. We have derived a heap-less abstract machine for call-by-need evaluation. It replaces the traditional manipulation of a heap using store-based effects with disciplined management of the evaluation stack using control-based effects. In short, state is replaced with control.

Machine evaluation of a program t begins with $\langle \emptyset \mid [], t \rangle_f$ and terminates at $\langle X \mid \llbracket E, v \rrbracket \rangle$.

5. CORRECTNESS OF THE MACHINE

The previous section proves that the machine manipulates terms in a manner that preserves variable binding. In this section, we prove that those manipulations correspond to standard-order call-by-need evaluation.

To proceed, we first establish correspondences between abstract machine configurations and call-by-need terms. As we have alluded to previously, abstract machine contexts correspond directly to calculus contexts:

$\mathcal{C}[\square]$	$=$	\square
$\mathcal{C}[\square t \circ E]$	$=$	$\mathcal{C}[E] t$
$\mathcal{C}[(\kappa x. E_1) \square \circ E_2]$	$=$	$(\lambda x. \mathcal{C}[E_1][x]) \mathcal{C}[E_2]$
$\mathcal{C}[(\lambda x. \square) t \circ E]$	$=$	$(\lambda x. \mathcal{C}[E]) t$

Redexes also map to call-by-need terms:

$\mathcal{C}[\llbracket E, v \rrbracket t]$	$=$	$(\mathcal{C}[E][v]) t$
$\mathcal{C}[(\kappa x. E_1) \llbracket E_2, v \rrbracket]$	$=$	$(\lambda x. \mathcal{C}[E_1][x]) (\mathcal{C}[E_2][v])$

Given that terms map identically to terms, configuration mapping is defined uniformly:

$\mathcal{C}[\langle X \mid E, ? \rangle]$	$=$	$\mathcal{C}[E][\mathcal{C}[?]]$
--	-----	----------------------------------

$\langle X \mid E, r \rangle_d$	(Reduce)
(D.1)	$\langle X \mid E_1, (\kappa x.E_2) \llbracket E_3, v \rrbracket \rangle_d \mapsto_{nam} \langle X \mid E_1 \circ E_3 \circ [(\lambda x.\square) v] \circ E_2, v \rangle_b$
(D.2)	$\langle X \mid E_1, \llbracket E_2, \lambda x.t_1 \rrbracket t_2 \rangle_d \mapsto_{nam} \langle X, x' \mid E_1 \circ E_2 \circ [(\lambda x'.\square) t_2], t_1[x'/x] \rangle_f \quad x' \notin X$
$\langle X \mid E, t \rangle_f$	(Refocus)
(F.1)	$\langle X \mid E, x \rangle_f \mapsto_{nam} \langle X \mid E, x \rangle_n$
(F.2)	$\langle X \mid E, \lambda x.t \rangle_f \mapsto_{nam} \langle X \mid E, \lambda x.t \rangle_b$
(F.3)	$\langle X \mid E, t_1 t_2 \rangle_f \mapsto_{nam} \langle X \mid E \circ [\square t_2], t_1 \rangle_f$
$\langle X \mid E, v \rangle_b$	(Rebuild)
(B.1)	$\langle X \mid E_b, v \rangle_b \mapsto_{nam} \langle X \mid \llbracket E_b, v \rrbracket \rangle$
(B.2)	$\langle X \mid E_1 \circ [\square t] \circ E_b, v \rangle_b \mapsto_{nam} \langle X \mid E_1, \llbracket E_b, v \rrbracket t \rangle_d$
(B.3)	$\langle X \mid E_1 \circ [(\kappa x.E_2) \square] \circ E_b, v \rangle_b \mapsto_{nam} \langle X \mid E_1, (\kappa x.E_2) \llbracket E_b, v \rrbracket \rangle_d$
where $E_b = \llbracket (\lambda x_i.\square) t_i \rrbracket$	
$\langle X \mid E, x \rangle_n$	(Need)
(N.1)	$\langle X \mid E_1 \circ [(\lambda x.\square) t] \circ E_2, x \rangle_n \mapsto_{nam} \langle X \mid E_1 \circ [(\kappa x.E_2) \square], t \rangle_f$
where $\llbracket (\lambda x.\square) t \rrbracket \notin E_2$	

Figure 1: Call-by-need Machine

Since the calculus is defined over alpha equivalence classes, we reason up to alpha equivalence when relating terms to machine configurations.

We now state our fundamental correctness theorems. First we guarantee soundness, the property that every step of the abstract machine respects standard-order reduction.

Theorem 5.1. *If $t_1 = \mathcal{C}[\llbracket C_1 \rrbracket]$ and $C_1 \mapsto_{nam} C_2$, then $t_1 \mapsto_{sr} t_2$, for some $t_2 = \mathcal{C}[\llbracket C_2 \rrbracket]$.*

Proof. By cases on \mapsto_{nam} . Only rules *D.1* and *D.2* are not immediate. The other rules preserve equality under $\mathcal{C}[\llbracket C \rrbracket]$. \square

Corollary 1 (Soundness).

If $t = \mathcal{C}[\llbracket C \rrbracket]$ and $C \mapsto_{nam} \langle X \mid \llbracket E, v \rrbracket \rangle$, then $t \mapsto_{sr} a$, for some $a = \mathcal{C}[\llbracket \langle X \mid \llbracket E, v \rrbracket \rrbracket \rrbracket]$.

Proof. By induction on the length of the \mapsto_{nam} sequence. \square

We also prove completeness, namely that abstract machine reduction subsumes standard order reduction.

Theorem 5.2 (Completeness).

If $t = \mathcal{C}[\llbracket C \rrbracket]$ and $t \mapsto_{sr} a$, then $C \mapsto_{nam} \langle X \mid \llbracket E, v \rrbracket \rangle$, with $a = \mathcal{C}[\llbracket \langle X \mid \llbracket E, v \rrbracket \rrbracket \rrbracket]$.

Proof. This proof proceeds by induction on the length of \mapsto_{sr} sequences. It utilizes Proposition 4.4 to accelerate the \mapsto_{sr} rules in accordance with \mapsto_{nam} . It also relies on a number of lemmas to establish that \mapsto_{nam} will find the unique redex of a term from any decomposition of a term into a context E and a subterm t . \square

Theorem 5.3 (Correctness). *If $t = \mathcal{C}[\![C]\!]$, then*

$$t \mapsto_{sr} a \quad \text{if and only if} \quad C \mapsto_{nam} \langle X \mid \llbracket E, v \rrbracket \rangle$$

with $a = \mathcal{C}[\langle X \mid \llbracket E, v \rrbracket \rangle]$.

5.1. Discussion. This abstract machine has nice inductive properties. The refocus rules always dispatch on the outermost term constructor. The rebuild and need rules dispatch on a prefix of the context, though each has different criteria for bounding the prefix.

The abstract machine's evaluation steps should not be seen as merely a desperate search for a redex. Rather, the machine exposes the fine-grain structure of call-by-need evaluation, just as the CK machine and the Krivine machine (Krivine, 2007) model evaluation for call-by-value and call-by-name respectively. Answers are the partial results of computations, and the rebuild rules represent the process of reconstructing and returning a result to a reduction site. Furthermore, the need rules can be viewed as a novel form of variable-lookup combined with lazy evaluation. The evaluation context captures the rest of computation, but not in order: variable references cause evaluation to skip around in a manner that is difficult to predict.

The way that variables behave in these semantics reveals a connection to coroutines. The reduction rule *D.2* binds a variable to a delayed computation; referencing that variable suspends the current computation and jumps to its associated delayed computation. Upon completion of that computation, any newly delayed computations (i.e. binder frames) are added to the evaluation context and the original computation is resumed.

The standard-order reduction relation of the call-by-need lambda calculus defines an evaluator concisely but abstractly. Surely unique decomposition, standardization, and hygiene ensure the existence of a deterministic evaluator, but these properties do not spell out the details or implications. Based on a reasoned inspection of standard-order reduction, we expose its computational behavior and capture it in a novel abstract machine that has no store. The improvements to the initial machine produce a variant that effectively assimilates computational information, explicitly accounts for variable hygiene and thereby reveals coarse-grained operational structure implicit in call-by-need standard-order evaluation.

5.2. Extensions. The machine presented above describes evaluation only for the pure lambda calculus. In this subsection, we introduce some features that are typical of pragmatic programming languages.

5.2.1. Let binding. To help with a proof that the call-by-need calculus can simulate call-by-name, Ariola and Felleisen introduce a *let*-based calculus. The let-calculus adds the construction $\text{let } x = t \text{ in } t$ to the set of terms and considers a new axiom:

$$(\lambda x.t_1) t_2 = \text{let } x = t_2 \text{ in } t_1$$

The let-calculus is formulated by taking this axiom as a reduction rule running from left to right, and reformulating the original three axioms of the call-by-need calculus in terms of let expressions. These new axioms can be justified by proving that they are derivable from the original three axioms and the new let axiom.

One approach to producing an abstract machine that supports let is to repeat the entire process described in this section, but focusing on the let-calculus instead of the

lambda calculus. However, let-binding can be retro-fitted to the current lambda calculus much more simply by reading the new let axiom as a *right-to-left* reduction rule:

$$\text{let } x = t_2 \text{ in } t_1 \rightarrow_{\text{need}} (\lambda x.t_1) t_2$$

Essentially, this rule indicates that let bindings in the calculus can be understood as equivalent to an immediate application of an abstraction to a term.

The consequences for the naïve machine are the addition of a new kind of redex:

$$r ::= \dots \mid \text{let } x = t_2 \text{ in } t_1$$

and a new refocus rule:

$$\langle E, \text{let } x = t_2 \text{ in } t_1 \rangle_f \mapsto \langle E, \text{let } x = t_2 \text{ in } t_1 \rangle_d$$

and a new reduce rule:

$$\langle E, \text{let } x = t_2 \text{ in } t_1 \rangle_d \mapsto \langle E, (\lambda x.t_1) t_2 \rangle_f$$

Then, in the process of improving our machine, the reduce rule can be fast-forwarded to reduce the introduced abstraction and application immediately:

$$\langle X \mid E, \text{let } x = t_2 \text{ in } t_1 \rangle_d \mapsto \langle X, x' \mid E \circ [(\lambda x'.\square) t_2], t_1[x'/x] \rangle_f$$

This process confirms that the let expression form can be comfortably treated as a conservative extension of the call-by-need lambda calculus. Its addition does not force a radical reconstruction of the abstract machine. As we show in Section 5.4, adding a circular letrec construct to the language fundamentally alters the system.

5.2.2. Constants. Plotkin (1975) augments the lambda calculus with two sets of constants, the basic constants b and the function constants f . In that language, the basic constants are observables and serve as placeholders for real programming language constants like numbers, strings, and so on. The function constants serve as placeholders for real primitive functions, generally over basic constants. The function constants are not first-class values in the language, but instead appear as operators in primitive function expressions.

$$\boxed{\begin{array}{l} t ::= \dots \mid b \mid f t \\ v ::= \dots \mid b \end{array}}$$

Function expressions always take a single term argument.

The calculus is parametrized on a partial function δ that maps function-value pairs to values. As is standard, the notions of reduction for the calculus are augmented to handle function constant applications.

$$\boxed{f v_1 \rightarrow_{\text{need}} v_2 \text{ if } \delta(f, v_1) = v_2}$$

Since the notions of reduction for function constants require them to be applied to values, the calculus must account for answers. To handle this, the calculus can be extended with a rule to commute answer bindings with function expressions, as is done for application expressions:

$$\boxed{f ((\lambda x.a) t) \rightarrow_{\text{need}} (\lambda x.f a) t}$$

If function expressions were only defined for basic constant arguments, then answer bindings could be discarded instead of commuted. However this approach imposes an ad hoc

limitation on the semantics. It would be cleaner and more orthogonal to uniformly handle garbage collection with a specific notion of reduction (Ariola and Felleisen, 1997).

Function constants must be considered in the evaluation contexts. A function constant must force its argument in order to produce a value, so upon encountering a function constant application, evaluation must continue in its argument:

$$E ::= \dots \mid f E$$

The other notions of reduction remain exactly the same, but since basic constants are values, they are subject to rules that manipulate values and answers. For example, if the set of basic constants includes numbers and the set of function constants includes operations on numbers, then by the deref rule:

$$(\lambda x.\text{add1 } x) 5 \rightarrow_{\text{need}} (\lambda x.\text{add1 } 5) 5$$

The abstract machine requires few changes to accommodate these additions. The set of redexes is extended to include function expressions:

$$r ::= \dots \mid f a$$

The refocus rules are generalized to create context frames for function constant applications and to rebuild when any value, lambda abstraction or basic constant, is encountered:

$$\begin{aligned} \langle X \mid E, f t \rangle_f &\mapsto_{\text{nam}} \langle X \mid E \circ [f \square], t \rangle_f \\ \langle X \mid E, v \rangle_f &\mapsto_{\text{nam}} \langle X \mid E, v \rangle_b \end{aligned}$$

Furthermore, rebuild must account for function constant applications:

$$\langle X \mid E \circ [f \square] \circ E_b, v \rangle_f \mapsto_{\text{nam}} \langle X \mid E, f \llbracket E_b, v \rrbracket \rangle_d$$

Finally, the reduce rule for $f a$, pushes binder frames upwards and appeals to the delta rule $\delta(f, v)$ for its result:

$$\langle X \mid E, f \llbracket E_b, v \rrbracket \rangle_d \mapsto_{\text{nam}} \langle X \mid E \circ E_b, \delta(f, v) \rangle_b$$

Since $\delta(f, v)$ yields only values, the reduce configuration for $f a$ can immediately transition to rebuild from this result.

5.3. Lazy Constructors. As pointed out by Ariola and Felleisen (1997), the call-by-need lambda calculus can be easily extended with support for lazy constructors. The rules for constructors and destructors can be inferred from the church encoding for pairs:

$$\begin{aligned} \text{cons} &\equiv \lambda x_1. \lambda x_2. \lambda d. d x_1 x_2 \\ \text{car} &\equiv \lambda p. p (\lambda x_1. \lambda x_2. x_1) \\ \text{cdr} &\equiv \lambda p. p (\lambda x_1. \lambda x_2. x_2) \end{aligned}$$

To add support for lazy pairs, we first extend the syntax of the language:

$$\begin{aligned} t &::= \dots \mid \text{cons } t t \mid \text{car } t \mid \text{cdr } t \mid \langle x, x \rangle \\ v &::= \dots \mid \langle x, x \rangle \\ E &::= \dots \mid \text{car } E \mid \text{cdr } E \end{aligned}$$

The $\text{cons } t_1 t_2$ expression creates a lazy pair of its two arguments, t_1 and t_2 . The $\text{car } t$ and $\text{cdr } t$ expressions evaluate their respective arguments and extract the first or second element of the resulting pair. As with function constants in the previous section, the lazy pair constructor and destructors are second-class. To express first-class constructors, these

may be eta-expanded into, e.g., $\lambda x_1. \lambda x_2. \text{cons } x_1 x_2$. The $\langle x_1, x_2 \rangle$ value is a representation of a lazy pair. It contains variables that refer to shared computations.

According to the evaluation contexts, evaluation may focus on the argument of a `car` or `cdr` expression. However, evaluation never directly operates on the arguments to a constructor. This is why the evaluation contexts for the language do not examine the arguments to `cons`. As a direct consequence, the standard call-by-need reduction rules will never substitute a value for a variable inside a pair. A variable inside a pair can only be evaluated after decomposing the pair using `car` or `cdr`.

The Church encoding of pairs motivates the following notions of reduction:

$$\begin{array}{l} \text{cons } t_1 t_2 \rightarrow_{\text{need}} (\lambda x_1. (\lambda x_2. \langle x_1, x_2 \rangle) t_2) t_1 \\ \text{car } \langle x_1, x_2 \rangle \rightarrow_{\text{need}} x_1 \\ \text{cdr } \langle x_1, x_2 \rangle \rightarrow_{\text{need}} x_2 \end{array}$$

These rules model `cons` as it is applied to two arguments. The `car` and `cdr` operations each expose a previously inaccessible variable reference to evaluation.

To accommodate lazy pairs in the abstract machine, we extend the set of redexes:

$$r ::= \dots \mid \text{cons } t_1 t_2 \mid \text{car } a \mid \text{cdr } a$$

and we introduce several new rules. The first set of rules extends the refocus stage of the machine to handle our extensions.

$$\begin{array}{l} \langle X \mid E, \text{cons } t_1 t_2 \rangle_f \mapsto_{\text{nam}} \langle X \mid E, \text{cons } t_1 t_2 \rangle_d \\ \langle X \mid E, \text{car } t \rangle_f \mapsto_{\text{nam}} \langle X \mid E \circ [\text{car } \square], t \rangle_f \\ \langle X \mid E, \text{cdr } t \rangle_f \mapsto_{\text{nam}} \langle X \mid E \circ [\text{cdr } \square], t \rangle_f \end{array}$$

A `cons` expression is immediately ready for reduction. The `car` and `cdr` expressions proceed to evaluate their respective arguments. Another set of rules returns an answer to `car` or `cdr`.

$$\begin{array}{l} \langle X \mid E \circ [\text{car } \square] \circ E_b, v \rangle_b \mapsto_{\text{nam}} \langle X \mid E, \text{car } \llbracket E_b, v \rrbracket \rangle_d \\ \langle X \mid E \circ [\text{cdr } \square] \circ E_b, v \rangle_b \mapsto_{\text{nam}} \langle X \mid E, \text{cdr } \llbracket E_b, v \rrbracket \rangle_d \end{array}$$

Finally, a set of rules corresponds to the notions of reduction.

$$\begin{array}{l} \langle X \mid E, \text{cons } t_1 t_2 \rangle_d \mapsto_{\text{nam}} \langle X, x_1, x_2 \mid E \circ [(\lambda x_1. \square) t_1] \circ [(\lambda x_2. \square) t_2], \langle x_1, x_2 \rangle \rangle_b \\ \langle X \mid E_1, \text{car } \llbracket E_2, \langle x_1, x_2 \rangle \rrbracket \rangle_d \mapsto_{\text{nam}} \langle X \mid E_1 \circ E_2, x_1 \rangle_n \\ \langle X \mid E_2, \text{cdr } \llbracket E_2, \langle x_1, x_2 \rangle \rrbracket \rangle_d \mapsto_{\text{nam}} \langle X \mid E_1 \circ E_2, x_2 \rangle_n \end{array}$$

Lazy pair construction creates new binder frames for the two terms and produces a pair that references them. To evaluate a destructor, the binder frames associated with the answer are pushed upwards and the corresponding variable reference is extracted from the pair and the value of its associated computation is demanded.

5.4. Circularity. As pointed out by Ariola et al. (1995), the presence of constructors makes it possible to observe duplicated constructors when the `Y` combinator is used to express recursion.

Consider the expression:

$$Y(\lambda y. \text{cons } 1 y) \equiv (\lambda f. (\lambda x. f (x x)) (\lambda x. f (x x))) (\lambda y. \text{cons } 1 y)$$

Standard-order reduction in the calculus proceeds as follows, underlining either the active variable reference or the current redex:

$$\begin{aligned}
& (\lambda f. (\lambda x. \underline{f} (x x)) (\lambda x. f (x x))) (\lambda y. \mathbf{cons} \ 1 \ y) \\
\rightarrow_{\text{need}} & (\lambda f. (\lambda x. (\lambda y. \underline{\mathbf{cons}} \ 1 \ y) (x x)) (\lambda x. f (x x))) (\lambda y. \mathbf{cons} \ 1 \ y) \\
\rightarrow_{\text{need}} & (\lambda f. (\lambda x. (\lambda y. (\lambda x_1. (\lambda x_2. \langle x_1, x_2 \rangle) y) 1) (x x)) (\lambda x. f (x x))) (\lambda y. \mathbf{cons} \ 1 \ y)
\end{aligned}$$

Because of the two references to f in Y , the term $(\lambda y. \mathbf{cons} \ 1 \ y)$ is copied and ultimately recomputed each time the term is needed. Ideally a recursive $\langle x, x \rangle$ value could refer to itself and not recompute its value whenever its \mathbf{cdr} is demanded. Duplicated computation does not correspond to how recursion and lazy constructors interact in typical semantics for lazy evaluation (Henderson and Morris, 1976). An implementation would create a single self-referencing \mathbf{cons} cell.

Because of the scoping rules for λ , there is no way to explicitly define truly circular structures in this calculus. To address this, Ariola and Felleisen introduce a *letrec*-based calculus, inspired by the circular calculus of Ariola and Klop (1994, 1997). The syntax of the *letrec*-calculus follows:

t	$::=$	$x \mid \lambda x. t \mid t t \mid \mathbf{letrec} \ D \ \mathbf{in} \ t$
D	$::=$	$x_1 \ \mathbf{be} \ t_1, \dots, x_n \ \mathbf{be} \ t_n$
v	$::=$	$\lambda x. t$
a	$::=$	$v \mid \mathbf{letrec} \ D \ \mathbf{in} \ a$
E	$::=$	$\square \mid E t \mid \mathbf{letrec} \ D \ \mathbf{in} \ E$ $\mid \mathbf{letrec} \ D, x \ \mathbf{be} \ E \ \mathbf{in} \ E[x]$ $\mid \mathbf{letrec} \ x_n \ \mathbf{be} \ E, D[x, x_n] \ \mathbf{in} \ E[x]$
$D[x, x_n]$	$::=$	$x \ \mathbf{be} \ E[x_1], x_1 \ \mathbf{be} \ E[x_2], \dots, x_{n-1} \ \mathbf{be} \ E[x_n], D$

The *letrec*-calculus is similar to the *let*-calculus in how it adds an explicit binding form, but each **letrec** expression can contain an unordered set D of mutually recursive bindings $x \ \mathbf{be} \ t$ for distinct variables.

The evaluation contexts for the *letrec*-calculus resemble those for the *let*-calculus. The contexts **letrec** $D, x \ \mathbf{be} \ E \ \mathbf{in} \ E[x]$ and **letrec** $x_n \ \mathbf{be} \ E, D[x, x_n] \ \mathbf{in} \ E[x]$ represent evaluation taking place in the binding position of a **letrec** expression. The first form expresses that a variable has been referenced in the body of the **letrec** and that variable's definition is currently under evaluation. The second form expresses that some variables bound in the **letrec** form depend on each other, and a variable reference in the body of the *letrec* has forced evaluation of this chain of variable references. This chain of dependencies is denoted by the syntax $D[x, x_n]$. The definition of the last referenced variable in the chain of dependencies, x_n , is currently being evaluated. For evaluation to proceed, all the variables in a chain of dependencies must be disjoint. The *letrec* calculus regards a cyclic dependency chain as a diverging computation, a stuck expression that is not a valid answer.

The notions of reduction for the letrec calculus follow:

(1)	$(\lambda x.t_1) t_2$	$\rightarrow_{\text{need}}$	letrec x be t_2 in t_1
(2)	letrec D, x be v in $E[x]$	$\rightarrow_{\text{need}}$	letrec D, x be v , in $E[v]$
(3)	letrec x_n be $v, D[x, x_n]$ in $E[x]$	$\rightarrow_{\text{need}}$	letrec x_n be $v, D[x, v]$ in $E[x]$
(4)	(letrec D in a) t	$\rightarrow_{\text{need}}$	letrec D in a t
(5)	letrec x_n be (letrec D in a), $D[x, x_n]$ in $E[x]$	$\rightarrow_{\text{need}}$	letrec D, x_n be $a, D[x, x_n]$ in $E[x]$
(6)	letrec D_1, x be (letrec D_2 in a) in $E[x]$	$\rightarrow_{\text{need}}$	letrec D_2, D_1, x be a in $E[x]$

The rules operate as follows. Rule (1) is analogous to the equivalent let-calculus rule. Rule (2) is analogous to the basic dereference rule. Rule (3) resolves the last variable reference in a chain of dependency. The syntax $D[x, v]$ expresses replacing x_{n-1} **be** $E[x_n]$ with x_{n-1} **be** $E[v]$ in the list of bindings. Rules (4) through (6) are associativity rules that percolate bindings upward to ensure proper sharing. Rules (5) and (6) are important for recursion. They lift recursive bindings from the definition of a letrec-bound variable x and incorporate them into the same **letrec** expression that binds x . In essence they expand the set of variables bound by the outer **letrec** expression.

Using the letrec-calculus, the corresponding example using **cons** is as follows:

$$\begin{aligned}
 & \mathbf{letrec} \ y \ \mathbf{be} \ \mathbf{cons} \ 1 \ y \ \mathbf{in} \ y \\
 \rightarrow_{\text{need}} & \mathbf{letrec} \ y \ \mathbf{be} \ (\mathbf{letrec} \ x_1 \ \mathbf{be} \ 1 \ \mathbf{in} \ (\mathbf{letrec} \ x_2 \ \mathbf{be} \ y \ \mathbf{in} \ \langle x_1, x_2 \rangle)) \ \mathbf{in} \ y \\
 \rightarrow_{\text{need}} & \mathbf{letrec} \ x_1 \ \mathbf{be} \ 1, y \ \mathbf{be} \ (\mathbf{letrec} \ x_2 \ \mathbf{be} \ y \ \mathbf{in} \ \langle x_1, x_2 \rangle) \ \mathbf{in} \ y \\
 \rightarrow_{\text{need}} & \mathbf{letrec} \ x_1 \ \mathbf{be} \ 1, x_2 \ \mathbf{be} \ y, y \ \mathbf{be} \ \langle x_1, x_2 \rangle \ \mathbf{in} \ y \\
 \rightarrow_{\text{need}} & \mathbf{letrec} \ x_1 \ \mathbf{be} \ 1, x_2 \ \mathbf{be} \ y, y \ \mathbf{be} \ \langle x_1, x_2 \rangle \ \mathbf{in} \ \langle x_1, x_2 \rangle
 \end{aligned}$$

In this reduction, the recursive **cons** cell has been computed once and for all. No reference to the original **cons** lazy constructor remains.

5.4.1. *The letrec machine.* Now we express the call-by-need letrec-calculus as an abstract machine. Development of this machine proceeds along the same lines as the non-circular machine.

As described in the calculus, **letrec** bindings need not be ordered, but dependency chains have a natural order imposed by the order of dependencies. Furthermore, cyclic dependency chains are provable divergences in the letrec-calculus. In the face of such a cycle, a program is no longer subject to reduction. We can model this by letting the machine get “stuck”.

The basic terms for the machine remain close to those of the calculus:

t	$::=$	$x \mid \lambda x.t \mid tt \mid \mathbf{letrec} \ D^+ \ \mathbf{in} \ t$
D	$::=$	$\overline{x_i \ \mathbf{be} \ t_i}$
v	$::=$	$\lambda x.t$

As with the calculus, D refers to sets of recursive bindings. As needed, we distinguish possibly empty sets of bindings, D^* , from nonempty sets D^+ . This precision is needed to discuss machine behavior. The values of this machine are the lambda abstractions, but now answers are defined as values wrapped in zero or more tiers of **letrec** bindings.

As we did for the prior machine, we introduce some representation changes that help with presenting the letrec-machine. The evaluation context is once again replaced with a

list of context frames.

$$\begin{array}{l}
P_{x_0}^{x_n} ::= \{\langle x_n, E_n \rangle :: \langle x_{n-1}, E_{n-1} \rangle :: \cdots :: \langle x_0, E_0 \rangle\} \\
f ::= \square t \mid \text{LR } D^+ \text{ in } \square \\
\quad \mid \text{LR } x \text{ be } \square, D^* \text{ in } E \\
\quad \mid \text{LR } x_{n+1} \text{ be } \square, P_{x_0}^{x_n}, D^* \text{ in } E \\
E ::= [] \mid [f] \circ E \mid E \circ [f] \\
\quad \text{where } E \circ [] = [] \circ E = E \\
\quad \text{and } E_1 \circ (E_2 \circ E_3) = (E_1 \circ E_2) \circ E_3
\end{array}$$

The operand frame $[\square t]$ and binding frame $[\text{LR } D^+ \text{ in } \square]$ are directly analogous to the corresponding frames in the original abstract machine. The original cont frame, on the other hand, splits into two variants. The frame $[\text{LR } x_{n+1} \text{ be } \square, P_{x_0}^{x_n}, D^* \text{ in } E[x_0]]$ captures chain of dependencies that is currently under evaluation. The expression D^* stands for the inactive bindings in the frame, while the expression $P_{x_0}^{x_n}$ indicates a chain of dependencies. While no particular ordering is imposed on the bindings in D^* , the chain of dependencies represented by $P_{x_0}^{x_n}$ is ordered. Our machine representation separates dependencies from other bindings. In the calculus, $D[x_0, x_n]$ also includes any other bindings D . In the machine, $P_{x_0}^{x_n}$ only captures the dependencies, and the other bindings D^* are explicitly indicated. Machine dependencies $P_{x_0}^{x_n} \equiv \{\langle x_n, E_n \rangle :: \langle x_{n-1}, E_{n-1} \rangle :: \cdots :: \langle x_0, E_0 \rangle\}$ correspond to calculus dependencies $D[x_0, x_n] \equiv x \text{ be } E[x_1], x_1 \text{ be } E[x_2], \dots, x_{n-1} \text{ be } E[x_n]$. We allow $P_{x_0}^{x_0}$ to denote an empty chain. Observe that the order of the machine dependencies is reversed. This ordering expresses that dependencies are resolved in last-in first-out order. When the value of x_{n+1} is computed, its value will be used to compute the value of x_n and so on until the value of x_0 is computed and its value is returned to the body of the **letrec** expression. Machine dependencies form a stack of computations that reference one another. This behavior is clarified in the behavior of the cyclic abstract machine.

The frame $[\text{LR } x \text{ be } \square, D^* \text{ in } E]$ closely resembles the old cont frame $[(\kappa x.E) \square]$, except that the LR form may bind other variables as well. The absence of $P_{x_0}^{x_n}$ in this frame explicitly indicates that it has no chain of dependencies to be evaluated before substituting into the body of the LR.

Machine answers are still binding-value pairs, but now each binding is a recursive binding of multiple variables.

$$a ::= \llbracket E, v \rrbracket \quad \text{where } E = \overline{[\text{LR } D_i^+ \text{ in } \square]}$$

The set of redexes for the abstract machine follows.

$$\begin{array}{l}
r ::= a t \\
\quad \mid \text{LR } x \text{ be } a, D^* \text{ in } E \\
\quad \mid \text{LR } x_{n+1} \text{ be } a, P_{x_0}^{x_n}, D^* \text{ in } E
\end{array}$$

The first redex form is the same as a form from the original machine. The second form is analogous to the $(\kappa x.E) a$ form from the original machine. The final redex form captures the case where one link in the chain of dependencies is about to be resolved.

Figure 2 presents the transition rules of the cyclic abstract machine. The machine relies on an operator $\mathcal{F}(E_b)$, which given a list of binder frames $[\text{LR } D_i^+ \text{ in } \square]$, flattens them into a single binder frame $[\text{LR } \overline{D_i^+} \text{ in } \square]$. This operation captures in aggregate the treatment of bindings by rules (5) and (6) of the calculus.

The first three refocus rules are the same as the lambda calculus, while the fourth rule is analogous to the equivalent rule for the let-calculus. The rebuild rules are also analogous to the lambda calculus, though now the binder frames have the form $[\text{LR } D_i^+ \text{ in } \square]$.

The letrec-machine has two need rules, (N.1) for a variable reference in the body of the corresponding **letrec**, and (N.2) for a variable reference that extends a (possible empty) chain of dependencies. The need rule does not consider variable bindings that are currently in a dependency chain, so evaluation will get stuck upon arriving at a cycle.

There are four reduce rules. The (D.1) rule, substitutes a value into the body of a letrec after its value has been computed. The (D.2) rule resolves the most recent reference in a chain of dependencies. Having computed the value of x_{n+1} , it returns to computing the value of x_n , which needed x_{n+1} 's value. If the chain of dependencies has only one element (i.e. $P_{x_0}^{x_0}$), then the chain is fully resolved. The (D.3) rule handles when an answer is applied to an expression. It combines calculus rules (4) and (1). Finally, the (D.4) rule handles **letrec** occurrences in the source program. In order to address hygiene, this rule must simultaneously substitute for every bound variable in each binding as it focuses on the body of the **letrec** expression.

The following is a machine trace of the cyclic **cons** example:

$$\begin{aligned}
& \langle \emptyset \mid [], \text{letrec } y \text{ be cons } 1 \ y \text{ in } y \rangle_f \\
\mapsto & \langle \emptyset \mid [], \text{letrec } y \text{ be cons } 1 \ y \text{ in } y \rangle_d \\
\mapsto & \langle \{y\} \mid [\text{LR } y \text{ be cons } 1 \ y \text{ in } \square], y \rangle_f \\
\mapsto & \langle \{y\} \mid [\text{LR } y \text{ be cons } 1 \ y \text{ in } \square], y \rangle_n \\
\mapsto & \langle \{y\} \mid [\text{LR } y \text{ be } \square \text{ in } []], \text{cons } 1 \ y \rangle_f \\
\mapsto & \langle \{y\} \mid [\text{LR } y \text{ be } \square \text{ in } []], \text{cons } 1 \ y \rangle_d \\
\mapsto & \langle \{y, x_1, x_2\} \mid [\text{LR } y \text{ be } \square \text{ in } []] \circ [\text{LR } x_1 \text{ be } 1 \text{ in } \square] \circ [\text{LR } x_2 \text{ be } y \text{ in } \square], \langle x_1, x_2 \rangle \rangle_b \\
\mapsto & \langle \{y, x_1, x_2\} \mid [], \text{LR } y \text{ be } \llbracket [\text{LR } x_1 \text{ be } 1 \text{ in } \square] \circ [\text{LR } x_2 \text{ be } y \text{ in } \square] \rrbracket \text{ in } [] \rangle_d \\
\mapsto & \langle \{y, x_1, x_2\} \mid [\text{LR } y \text{ be } \langle x_1, x_2 \rangle, x_1 \text{ be } 1, x_2 \text{ be } y \text{ in } \square], \langle x_1, x_2 \rangle \rangle_b \\
\mapsto & \langle \{y, x_1, x_2\} \mid \llbracket [\text{LR } y \text{ be } \langle x_1, x_2 \rangle, x_1 \text{ be } 1, x_2 \text{ be } y \text{ in } \square] \rrbracket \rangle
\end{aligned}$$

6. SIMULATING CALL-BY-NEED USING CONTROL

As we allude to above, call-by-need machine evaluation is highly suggestive of delimited control operations, but the connection is indirect and mixed with the other details of lazy evaluation. In this section, we directly interpret this connection in the terminology of delimited control.

Based on the operational behavior of the abstract machine from Figure 1, we derive a simulation of call-by-need execution under call-by-value augmented with delimited control operators. In particular, we translate call-by-need terms into the framework of Dybvig et al. (2007). First we overview the language of delimited control operations. Then we describe how the abstract machine performs delimited control operations. Next we present the simulation of call-by-need using delimited control. Finally we show its correctness.

$\langle X \mid E, r \rangle_d$	(Reduce)	
(D.1)	$\langle X \mid E_1, (\text{LR } x \text{ be } \llbracket E_2, v \rrbracket, D^* \text{ in } E_3) \rangle_d$	$\mapsto_{nam} \langle X \mid E_1 \circ [\text{LR } x \text{ be } v, \mathcal{F}(E_2), D^* \text{ in } \square] \circ E_3, v \rangle_b$
(D.2)	$\langle X \mid E_1, (\text{LR } x_{n+1} \text{ be } \llbracket E_2, v \rrbracket, P_{x_0}^{x_n}, D^* \text{ in } E_3) \rangle_d$	$\mapsto_{nam} \langle X \mid E_1 \circ [\text{LR } x_n \text{ be } \square, P_{x_0}^{x_{n-1}}, x_{n+1} \text{ be } v, \mathcal{F}(E_2), D^* \text{ in } E_3] \circ E_n, v \rangle_b$
(D.3)	$\langle X \mid E_1, \llbracket E_2, \lambda x.t_1 \rrbracket t_2 \rangle_d$	$\mapsto_{nam} \langle X, x' \mid E_1 \circ E_2 \circ [\text{LR } x' \text{ be } t_2 \text{ in } \square], t_1[x'/x] \rangle_f \quad x' \notin X$
(D.4)	$\langle X \mid E, \text{letrec } \overline{x_i} \text{ be } t_i \text{ in } t \rangle_d$	$\mapsto_{nam} \langle X, \overline{x'_i} \mid E \circ [\text{LR } \overline{x'_i} \text{ be } t_i[\overline{x'_i}/x_i] \text{ in } \square], t \rangle_f \quad \overline{x'_i} \cap X = \emptyset$
$\langle X \mid E, t \rangle_f$	(Refocus)	
(F.1)	$\langle X \mid E, x \rangle_f$	$\mapsto_{nam} \langle X \mid E, x \rangle_n$
(F.2)	$\langle X \mid E, \lambda x.t \rangle_f$	$\mapsto_{nam} \langle X \mid E, \lambda x.t \rangle_b$
(F.3)	$\langle X \mid E, t_1 t_2 \rangle_f$	$\mapsto_{nam} \langle X \mid E \circ [\square t_2], t_1 \rangle_f$
(F.4)	$\langle X \mid E, \text{letrec } D^+ \text{ in } t \rangle_f$	$\mapsto_{nam} \langle X \mid E, \text{letrec } D^+ \text{ in } t \rangle_d$
$\langle X \mid E, v \rangle_b$	(Rebuild)	
(B.1)	$\langle X \mid E_b, v \rangle_b$	$\mapsto_{nam} \langle X \mid \llbracket E_b, v \rrbracket \rangle$
(B.2)	$\langle X \mid E_1 \circ [\square t] \circ E_b, v \rangle_b$	$\mapsto_{nam} \langle X \mid E_1, \llbracket E_b, v \rrbracket t \rangle_d$
(B.3)	$\langle X \mid E_1 \circ [(\text{LR } x \text{ be } \square, D^* \text{ in } E_2)] \circ E_b, v \rangle_b$	$\mapsto_{nam} \langle X \mid E_1, (\text{LR } x \text{ be } \llbracket E_b, v \rrbracket, D^* \text{ in } E_2) \rangle_d$
(B.4)	$\langle X \mid E_1 \circ [(\text{LR } x_{n+1} \text{ be } \square, P_{x_0}^{x_n}, D^* \text{ in } E_2)] \circ E_b, v \rangle_b$	$\mapsto_{nam} \langle X \mid E_1, (\text{LR } x_{n+1} \text{ be } \llbracket E_b, v \rrbracket, P_{x_0}^{x_n}, D^* \text{ in } E_2) \rangle_d$
where $E_b = [\text{LR } D_i^+ \text{ in } \square]$		
$\langle X \mid E, x \rangle_n$	(Need)	
(N.1)	$\langle X \mid E_1 \circ [\text{LR } (x \text{ be } t, D^*) \text{ in } \square] \circ E_2^\dagger, x \rangle_n$	$\mapsto_{nam} \langle X \mid E_1 \circ [\text{LR } x \text{ be } \square, D^* \text{ in } E_2], t \rangle_f$
(N.2)	$\langle X \mid E_1 \circ [\text{LR } x_n \text{ be } \square, P_{x_0}^{x_{n-1}}, (x \text{ be } t, D^*) \text{ in } E_2] \circ E_n^\dagger, x \rangle_n$	$\mapsto_{nam} \langle X \mid E_1 \circ [\text{LR } x \text{ be } \square, P_{x_0}^{x_n}, D^* \text{ in } E_2], t \rangle_f$
where $(x \text{ be } t) \notin E^\dagger$		

Figure 2: Letrec Machine

6.1. Delimited Control Operators. Dybvig et al. (2007) define a language with delimited control operators. We explain these operators using a simplified variant of the defining machine semantics.

t	$::= x \mid v \mid t t \mid \text{newPrompt} \mid \text{pushPrompt } t t$
	$\mid \text{withSubCont } t t \mid \text{pushSubCont } t t$
v	$::= \lambda x.t \mid p \mid \langle M \rangle$
E	$::= \square \mid E[\square t] \mid E[(\lambda x.t) \square] \mid E[\text{pushPrompt } \square t]$
	$\mid E[\text{withSubCont } \square t] \mid E[\text{withSubCont } p \square]$
	$\mid E[\text{pushSubCont } \square t]$
M	$::= [] \mid E : M \mid p : M$
p	$\in \mathbb{N}$

The language extends the call-by-value untyped lambda calculus with the four operators *newPrompt*, *pushPrompt*, *withSubCont*, and *pushSubCont* as well as two new values: first-class *prompts* p , and first-class delimited continuations $\langle M \rangle$. Its control structure is defined using evaluation contexts E , and metacontexts M , which are lists that interleave prompts and contexts. Metacontexts use Haskell list notation. Prompts are modeled using natural numbers.

A program state comprises an expression t , continuation E , metacontinuation M , and fresh prompt source p . The initial state for a program t is $\square[t], [], 0$.

$E[(\lambda x.t) v], M, p$	$\mapsto E[t[v/x]], M, p$
$E[\text{newPrompt}], M, p$	$\mapsto E[p], M, p + 1$
$E[\text{pushPrompt } p_1 t], M, p_2$	$\mapsto \square[t], p_1 : E : M, p_2$
$E[\text{withSubCont } p_1 \lambda x.t], M_1 ++ (p_1 : M_2), p_2$	$\mapsto \square[t[\langle E : M_1 \rangle / x]], M_2, p_2$
where $p_1 \notin M_1$	
$E[\text{pushSubCont } \langle M_1 \rangle t], M_2, p$	$\mapsto \square[t], M_1 ++ (E : M_2), p$
$\square[v], E : M, p$	$\mapsto E[v], M, p$
$\square[v], p_1 : M, p_2$	$\mapsto \square[v], M, p_2$

The four operators manipulate delimited continuations, or *subcontinuations*, which are part of an execution context. The *withSubCont* operator takes a prompt and a function; it captures the smallest subcontinuation that is delimited by the prompt and passes it to the function. The non-captured part of the continuation becomes the new continuation. The prompt instance that delimited the captured subcontinuation is discarded: it appears in neither the captured subcontinuation nor the current continuation. This operator generalizes \mathcal{F} (Felleisen, 1988) and *shift* (Danvy and Filinski, 1990).

The *pushSubCont* operator takes a subcontinuation and an expression; it composes the subcontinuation with the current continuation and proceeds to evaluate its second argument in the newly extended continuation.

The *pushPrompt* operator takes a prompt and an expression; it extends the current continuation with the prompt and evaluates the expression in the newly extended continuation. The *newPrompt* operator returns a distinguished fresh prompt each time it is called. These two operators generalize the delimiting operators $\#$ (Felleisen, 1988) and

`reset` (Danvy and Filinski, 1990), which extend a continuation with a single common delimiter.

To illustrate these operators in action, we consider a program that uses arithmetic and conditionals:

```

let p = newPrompt
in 2 + pushPrompt p
    if (withSubCont p
        (λk.(pushSubCont k False)+
         (pushSubCont k True)))
    then 3
    else 4

```

A fresh prompt is bound to p and pushed onto the continuation just prior to evaluation of the *if* expression. *withSubCont* captures the subcontinuation $[if \square then 3 else 4]$, which was delimited by p , and binds it to k . The subcontinuation k is pushed twice, given the value *False* the first time and *True* the second. The result of evaluation is the expression $2 + 4 + 3$ which yields 9.

6.2. Delimited Control Naïvely Simulates the Machine. The call-by-need abstract machine performs two different kinds of partial control capture. To review, the rebuild and need rules of the abstract machine both capture some portion of the evaluation context. In particular, the rebuild rules capture binder frames. If only binder frames remain, then execution is complete. When either of the other frames is found, then a reduction is performed. On the other hand, the need rule captures the evaluation context up to the binder that matches the variable whose value is needed.

These actions of the abstract machine can be recast in the language of delimited control capture. First, the need rule uses the identity of its variable, which must be an active variable, to delimit the context it captures. The well-formedness conditions from Section 4.5 guarantee that each binder frame binds a unique variable, so each active variable acts as a unique delimiter. Second, the rebuild rule uses the nearest non-binder frame to delimit the context it captures. This means that rebuild operates as though the operand frames, the cont frames, and the top of the evaluation context share a common delimiter. This guarantees that only binder frames are captured (as is stipulated in the rules).

In short, call-by-need evaluation captures partial evaluation contexts. These partial evaluation contexts correspond to delimited continuations, and there are two different kinds of delimitation, redex-based (for rebuild) and binder-based (for need).

It is useful to also consider how the machine manipulates these delimited continuations. Each reduce rule in Figure 1 immediately pushes the context associated with an answer onto the current evaluation context. In this manner, binders are consistently moved above the point of evaluation. The reduce rule then operates on the value part of the answer and the associated cont (for *D.1*) or term (for *D.2*).

Although each reduce rule pushes binders onto the evaluation context, only the *D.2* rule creates new binders. The variable bound by the answer's underlying lambda abstraction may already be a member of the set X , in which case it must be alpha-converted to a fresh name with respect to the set X . Also note that if $\lambda x.t$ is alpha converted to $\lambda x'.t[x'/x]$, the body under call-by-value satisfies the equation $t[x'/x] = (\lambda x.t) x'$. Since we are using the identifiers x' as delimiters, and we never turn the binder frame $[(\lambda x'.\square) t]$ back into a term, we can replace fresh variables x' with fresh *prompts* (Balat et al., 2004).

<p>Let s be a distinguished identifier:</p> $\mathcal{N}^P[[t]] = \text{runCC} (\text{let } s = \text{newPrompt} \text{ in } \text{pushPrompt } s \mathcal{N}[[t]])$ $\mathcal{N}[[x]] = \text{need } x$ $\mathcal{N}[[t_1 t_2]] = \text{do } v_a \leftarrow \mathcal{N}[[t_1]]$ $\quad \text{in } \text{let } x_p = \text{newPrompt}$ $\quad \text{in } \text{delay } \mathcal{N}[[t_2]] \text{ as } x_p \text{ in } (v_a x_p)$ $\mathcal{N}[[\lambda x.t]] = \text{return } \lambda x. \mathcal{N}[[t]]$ <hr style="border: 1px solid black;"/> $\text{return } v_a \equiv \text{withSubCont } s \lambda k_a. \langle k_a, v_a \rangle$ $\text{do } x \leftarrow t_1 \text{ in } t_2 \equiv \text{let } \langle k_a, x \rangle = \text{pushPrompt } s t_1$ $\quad \text{in } \text{pushSubCont } k_a t_2$ $\text{delay } t_1 \text{ as } x \text{ in } t_2 \equiv \text{let } f_k = \text{pushPrompt } x t_2$ $\quad \text{in } f_k \lambda().t_1$ $\text{force } f \equiv f ()$ $\text{need } x \equiv \text{withSubCont } x \lambda k.$ $\quad \lambda f_{th}. \text{do } v_a \leftarrow \text{force } f_{th}$ $\quad \text{in } \text{delay } (\text{return } v_a) \text{ as } x$ $\quad \text{in } \text{pushSubCont } k (\text{return } v_a)$

Figure 3: Translating CBN to CBV+Control

From these observations, we construct the simulation in Figure 3. The simulation can be understood as a direct encoding of the abstract machine semantics for call-by-need. To execute a program, $\mathcal{N}^P[[t]]$, the transformation uses *runCC* to initiate a control-based computation, acquires a fresh prompt, and binds it to a distinguished variable s . This prompt is the *redex prompt*, which is used to delimit every continuation that denotes a redex.

To expose the conceptual structure of the simulation, we define five syntactic macros, *do*, *return*, *delay*, *force*, and *need*. We accord no formal properties to them: they merely simplify the presentation. The *return* macro captures the nearest subcontinuation that is delimited by the redex prompt s . Since the s delimiter appears before every reduction, the captured continuation is guaranteed to contain only code equivalent to binder frames. The translation returns a tuple containing the subcontinuation and the argument to *return*, which must be a value; the tuple represents an answer. So the translation rule for lambda abstractions, $\mathcal{N}[[\lambda x.t]]$, literally simulates the rebuild rules.

The *do* macro executes a term t_1 under the current continuation extended with the redex prompt. If the term returns an answer $\langle k_a, x \rangle$ it immediately pushes the subcontinuation part and continues execution, binding the value part to the variable x . As such, the translation rule for applications, $\mathcal{N}[[t_1 t_2]]$, executes $[[t_1]]$ and binds the resulting operator to

v_a . The answer binders are pushed by the `do` macro, which starts the simulation of the $D.2$ rule.

The remainder of the $D.2$ rule is subtle. In the abstract machine, binder frame variables delimit the need rules. Since the delimited continuation framework relies on prompts to delimit continuations, fresh prompts literally substitute for variables (Kiselyov et al., 2006). The translation uses *newPrompt* to acquire a fresh prompt x_p and then uses the `delay` macro to simulate pushing a binder frame: the context `delay t as x in \square` is analogous to the binder frame $[(\lambda x. \square) t]$. The `delay` macro anticipates that its body returns a function f_k that expects the delayed argument, so it applies f_k to a suspension of t . As we see below, the function f_k is a cont $(\kappa x.E)$.

In the context of `delay`, the simulation executes $v_a x_p$. Since alpha conversion of $\lambda x.t$ can be written $(\lambda x_p.t[x_p/x])$, the term $v_a x_p$ is analogous to $(\lambda x.t) x_p = t[x_p/x]$: it substitutes a fresh prompt for a fresh variable.

The `need` macro, which defines the translation rule for variables, $\mathcal{N}[\![x]\!]$, captures the continuation delimited by x (which had better be a prompt!) and returns a function $\lambda f_{th} \dots$ that closes over both x and the captured continuation k . This function is the cont $\kappa x.E$, with x modeling the bound variable of the same name, and continuation k modeling E . The function expects the binder frame $[(\lambda x. \square) t]$, which is now at the top of the current continuation, to pass it the suspension $\lambda(). \mathcal{N}[\![t]\!]$. The simulation forces the suspension, and the `do` macro pushes the resulting answer binders and binds v_a to the underlying value. Pushing the answer binders begins the simulation of the $D.1$ rule.

The simulation of $D.1$ delays a computation that immediately returns the result v_a of evaluating the term t , pushes the continuation k associated with the cont, and returns v_a to the extended continuation. Now any subsequent evaluation of x immediately returns the memoized value v_a instead of recomputing t . This yields an answer $\langle k_a, v_a \rangle$ where k_a is an empty subcontinuation. The value v_a is delayed exactly as before and is also returned from the properly extended continuation. This part of the translation bears close resemblance to the paradoxical Y combinator (Curry and Feys, 1958), suggesting that the simulation requires recursive types (Shan, 2007).

7. CORRECTNESS OF THE SIMULATION

We prove correctness of the simulation relative to the machine semantics. Since we already proved correctness of the machine semantics relative to standard-order reduction, the result is a proof that our simulation provides a continuation semantics for call-by-need.

The previous discussion provides an informal justification for the structure of the call-by-need simulation. To prove the correctness of the simulation, we appeal to the continuation semantics for delimited control (Dybvig et al., 2007). This semantics is completely standard for the terms of the lambda calculus. Figure 4 presents the interesting parts of the semantics. All CPS terms take a standard continuation κ , but the control combinators also take a *metacontinuation* γ , which is a list of continuations and prompts, and a global prompt counter q . The base continuation κ_0 delimits each proper continuation and directs evaluation up the metacontinuation, discarding any intervening prompts. Given a CPS program t , the expression $t \kappa_0 [\] 0$ runs it.

To prove correctness, we compose $\mathcal{N}[\![\cdot]\!]$ with the delimited continuation semantics to produce a translation $\Lambda[\![\cdot]\!]$ to the $\lambda_{\beta\eta}$ calculus augmented with arithmetic, lists, and the

$ \begin{aligned} newPrompt_c &= \lambda\kappa.\lambda\gamma.\lambda q.\kappa\ \gamma\ (q+1) \\ withSubCont_c &= \lambda p.\lambda f.\lambda\kappa.\lambda\gamma.f\ (\kappa : \gamma^p)\ \kappa_0\ \gamma^p_{\downarrow} \\ pushPrompt_c &= \lambda p.\lambda t.\lambda\kappa.\lambda\gamma.t\ \kappa_0\ (p : \kappa : \gamma) \\ pushSubCont_c &= \lambda\gamma'.\lambda t.\lambda\kappa.\lambda\gamma.t\ \kappa_0\ (\gamma'^{++}(\kappa : \gamma)) \\ \kappa_0 &= \lambda v.\lambda\gamma.\lambda q.\mathcal{K}(v, \gamma, q) \\ \mathcal{K}(v, [], q) &= v \\ \mathcal{K}(v, p : \gamma, q) &= \mathcal{K}(v, \gamma, q) \\ \mathcal{K}(v, \kappa : \gamma, q) &= \kappa\ v\ \gamma\ q \end{aligned} $

Figure 4: Delimited Control Combinators

operator \mathcal{K} defined in Figure 4. We also give each abstract machine configuration a denotation, defined in terms of name-indexed denotations for its constituents $\mathcal{D}[\cdot]_X$ (see Figures 5 through 8).

$ \begin{aligned} \mathcal{D}[[t]]_X &= \Lambda[[t]][\overline{u(x_i, X)/x_i}] \\ \Lambda^P[[t]] &= \Lambda[[t]]\ \kappa_0\ (0 : [])\ 1 \\ \Lambda[[x]] &= \\ withSubCont_c\ x & \\ \lambda k_x.\lambda k_1.k_1 & \\ \lambda f_{th}.\lambda k_2. & \\ pushPrompt_c\ 0\ (f_{th}\ ()) & \\ (\lambda \langle k_a, v_a \rangle. & \\ pushSubCont_c\ k_a & \\ (\lambda k_3. & \\ pushPrompt_c\ x & \\ (pushSubCont_c\ k_x & \\ (withSubCont_c\ 0\ \lambda k_a.\lambda k.k\ \langle k_a, v_a \rangle)) & \\ (\lambda f_k.f_k\ (\lambda().withSubCont_c\ 0 & \\ \lambda k_a.\lambda k.k\ \langle k_a, v_a \rangle) & \\ k_3)) & \\ k_2) & \\ \Lambda[[t_1\ t_2]] &= \lambda k_1.pushPrompt_c\ 0\ \Lambda[[t_1]] \\ (\lambda \langle k_a, v_a \rangle. & \\ pushSubCont_c\ k_a & \\ (\lambda k_2. & \\ newPrompt_c & \\ \lambda x_p.pushPrompt_c\ x_p\ (v_a\ x_p) & \\ (\lambda f_k.f_k\ (\lambda().\Lambda[[t_2]])\ k_2)) & \\ k_1) & \\ \Lambda[[\lambda x.t]] &= withSubCont_c\ 0\ \lambda k_a.\lambda k.k\ \langle k_a, \lambda x.\Lambda[[t]] \rangle \end{aligned} $
--

Figure 5: Denotations for Terms

$$\begin{array}{l}
\iota(x_i, X) = \iota(x_i, [x_1, x_2, \dots, x_i, \dots, x_n]) = i \\
|X| = |[x_1, x_2, \dots, x_i, \dots, x_n]| = n \\
\Lambda[\langle X | E, ? \rangle_c] = \mathcal{D}[\cdot]_X \kappa_0 (\mathcal{D}[E]_{X^{++}}(0 : [])) (|X| + 1) \quad c \in \{d, f, b, n\} \\
\Lambda[\langle X | \llbracket E, \lambda x.t \rrbracket \rangle] = \langle \mathcal{D}[E]_X, \lambda x. \mathcal{D}[t]_X \rangle
\end{array}$$

Figure 6: Denotations for Names and Configurations

Denotations of machine configurations are constructed from their components: the configuration's focus $?$, context E , and list of names X . A machine configuration denotes the translation of its focus applied to three arguments: the base continuation κ_0 as its starting continuation, the denotation of its context, bounded by the redex delimiter 0 , as the metacontinuation, and the size $|X|$ of X plus 1 as its initial prompt. The redex delimiter attached to the metacontinuation handles the case when an answer subsumes the entire context by returning the answer as the result. The denotation of the terminal machine configuration $\langle X | \llbracket E, v \rrbracket \rangle$ is treated separately to show how it corresponds directly to a final answer.

Our semantic translation takes advantage of X being a proper list of unique names. Free active variables denote prompts in our translation, and since 0 is the redex delimiter, we assign to each variable its 1-based index in X . We use $|X|+1$ as the global prompt counter to ensure that no future prompts conflict with the current active variable denotations, thereby guaranteeing hygiene (see Section 4.5).

Each evaluation context frame denotes a two-element metacontinuation consisting of a prompt and a proper continuation. The prompt for a binder frame is the prompt translation $\iota(x, X)$ of the bound variable x . The cont and operand frames have redex prompts 0 . These prompts guarantee that answer building operations will arrive at the innermost redex. Each continuation function specializes a subexpression of the CPS translation for terms $\mathcal{D}[\cdot]_X$ with the denotations of the context frame's parts. Compare, for instance, the denotation of an application, $t_1 t_2$, to that of an operand frame, $[\square t_2]$. The application term pushes the global prompt, and executes t_1 in the context of a continuation that receives an answer $\langle k_a, v_a \rangle$. The denotation of the operand frame is a metacontinuation containing the same prompt and continuation.

A redex denotes a CPS'ed term that closes over the denotations of its constituents and implements the corresponding reduction step.

To facilitate our proof of correctness, we make a slight change to the machine semantics. In the machine, composing an empty context with the current context is an identity operation. The continuation semantics do not share this property. During execution, an empty continuation is denoted by the base continuation κ_0 . If a continuation is captured or pushed in the context of an empty continuation, then the empty continuation will be captured as part of the metacontinuation or pushed onto the current metacontinuation before reinstating the pushed continuation. In short, the call-by-need machine semantics guarantees that $E \circ [] = E$, but the continuation semantics do not prove that $\kappa_0 : \gamma = \gamma$. Dybvig et al. discuss the notion of proper tail recursion for delimited continuations. Their operational characterization of proper tail recursion corresponds to the latter equation.

$$\begin{aligned}
\mathcal{D}[[E \circ [f]]]_X &= \mathcal{D}[[f]]_X \mathbin{++} \mathcal{D}[[E]]_X \\
\mathcal{D}[[\square]]_X &= [] \\
\mathcal{D}[[\#]]_X &= \kappa_0 : [] \\
\mathcal{D}[[\square t_2]]_X &= 0 : k' : [] \\
\text{where } k' &= \lambda \langle k_a, v_a \rangle. \\
&\quad \text{pushSubCont}_c k_a \\
&\quad (\lambda k_2. \\
&\quad \quad \text{newPrompt}_c \\
&\quad \quad \lambda x_p. \text{pushPrompt}_c x_p (v_a x_p) \\
&\quad \quad (\lambda f_k. f_k (\lambda(). \mathcal{D}[[t_2]]_X) k_2)) \\
&\quad \kappa_0 \\
\mathcal{D}[[\lambda x. \square t_2]]_X &= \iota(x, X) : k' : [] \\
\text{where } k' &= \lambda f_k. f_k (\lambda(). \mathcal{D}[[t_2]]_X) \kappa_0 \\
\mathcal{D}[[\lambda x. E] \square]_X &= 0 : k' : [] \\
\text{where } k' &= \\
&\lambda \langle k_a, v_a \rangle. \\
&\quad \text{pushSubCont}_c k_a \\
&\quad (\lambda k_3. \\
&\quad \quad \text{pushPrompt}_c \iota(x, X) \\
&\quad \quad (\text{pushSubCont}_c \mathcal{D}[[E]]_X \\
&\quad \quad \quad (\text{withSubCont}_c 0 \lambda k_a. \lambda k. k \langle k_a, v_a \rangle)) \\
&\quad \quad (\lambda f_k. f_k (\lambda(). \text{withSubCont}_c 0 \lambda k_a. \lambda k. k \langle k_a, v_a \rangle) \\
&\quad \quad \quad k_3)) \\
&\quad \kappa_0
\end{aligned}$$

Figure 7: Denotations for Evaluation Contexts

To remove this mismatch, we add a *ghost* frame $[\#]$ to our definition of evaluation contexts. The ghost frame denotes the metacontinuation $\kappa_0 : []$. We also extend the unplug operation on evaluation contexts such that it discards ghost frames: $\mathcal{C}[[E \circ [\#]]] = \mathcal{C}[[E]]$. Finally, we alter the right hand side of transition rules that grab and push continuations to pair ghost frames with composed evaluation contexts in a manner consistent with the continuation semantics. For instance, the updated $D.1$ rule is as follows²:

$$\langle X \mid E_1, (\kappa x. E_2) \llbracket E_3, v \rrbracket \rangle_d \longmapsto_{nam} \langle X \mid E_1 \circ [\#] \circ E_3 \circ [(\lambda x. \square) v] \circ [\#] \circ E_2, v \rangle_b \quad (D.1)$$

These modifications do not alter the observable behavior of the machine while modeling the property that pushing empty frames has meaning in the continuation semantics.

Given these denotations, it is straightforward to prove correctness of the simulation relative to the abstract machine.

Theorem 7.1. *If t is a closed term, then $\Lambda^P[[t]] = \Lambda[[\langle \emptyset \mid [], t \rangle_f]]$.*

Proof. $\Lambda^P[[t]] = \mathcal{D}[[t]]_{\emptyset} \kappa_0 (0 : []) 1 = \Lambda[[\langle \emptyset \mid [], t \rangle_f]]$. □

²The exact placement of ghost frames falls right out of the correctness proof.

$$\begin{array}{l}
\mathcal{D}[(\kappa x_1.E_1) \llbracket E_2, \lambda x_2.t \rrbracket]_X = \\
\text{pushSubCont}_c \mathcal{D}\llbracket E_2 \rrbracket_X \\
(\lambda k_3. \\
\text{pushPrompt}_c \iota(x_1, X) \\
(\text{pushSubCont}_c \mathcal{D}\llbracket E_1 \rrbracket_X \\
(\text{withSubCont}_c 0 \lambda k_a. \lambda k.k \langle k_a, \lambda x_2. \llbracket t \rrbracket_X \rangle)) \\
(\lambda f_k.f_k (\lambda(). \text{withSubCont}_c 0 \lambda k_a. \lambda k.k \langle k_a, \lambda x_2. \llbracket t \rrbracket_X \rangle) \\
k_3)) \\
\\
\mathcal{D}[\llbracket E, \lambda x.t_1 \rrbracket t_2]_X = \\
\text{pushSubCont}_c \mathcal{D}\llbracket E \rrbracket_X \\
(\lambda k_2. \\
\text{newPrompt}_c \\
\lambda x_p. \text{pushPrompt}_c x_p ((\lambda x. \mathcal{D}\llbracket t_1 \rrbracket_X) x_p) \\
(\lambda f_k.f_k (\lambda(). \mathcal{D}\llbracket t_2 \rrbracket_X) k_2))
\end{array}$$

Figure 8: Denotations for Redexes

Theorem 7.2. *If $C_1 \mapsto_{nam} C_2$ then $\Lambda[C_1] = \Lambda[C_2]$.*

Proof. By cases on \mapsto_{nam} . The proof utilizes beta, eta and \mathcal{K} equivalences to establish correspondences. \square

8. SIMULATING EXTENSIONS

Many straightforward language extensions also have straightforward simulations under the above model.

Simulating *let* bindings essentially performs the same operations as immediately applying a lambda abstraction.

$$\begin{aligned}
\mathcal{N}[\text{let } x = t_1 \text{ in } t_2] &= \text{let } x = \text{newPrompt} \\
&\quad \text{in } \text{delay } \mathcal{N}\llbracket t_1 \rrbracket \text{ as } x \\
&\quad \text{in } \mathcal{N}\llbracket t_2 \rrbracket
\end{aligned}$$

In contrast to the application rule, the variable x is directly assigned a fresh prompt, rather than binding it to an auxiliary variable x_p . The body of the *let* can be interpreted in place and substitution of the prompt is implicit since x is already free in t_2 .

The translation for basic constants is analogous to that for lambda abstractions: the constant must be returned to the next redex.

$$\mathcal{N}\llbracket c \rrbracket = \text{return } c$$

For this translation, we assume that the call-by-value language provides the same constants as the call-by-need language.

The translation for function constants is as follows:

$$\begin{aligned}
\mathcal{N}\llbracket f t \rrbracket &= \text{do } v_a \leftarrow \mathcal{N}\llbracket t \rrbracket \\
&\quad \text{in return } (f v_a)
\end{aligned}$$

Interpreting a function constant application forces its argument and then acts on the value that is ultimately produced. Since function expressions yield values, the result is immediately returned.

The translation for `cons` acquires two fresh prompts, uses them to delay the argument to `cons`, and stores them as a pair.

$$\begin{aligned} \mathcal{N}[\text{cons } t_1 \ t_2] = & \text{let } x_1 = \text{newPrompt in} \\ & \text{let } x_2 = \text{newPrompt in} \\ & \text{delay } \mathcal{N}[[t_1]] \text{ as } x_1 \text{ in} \\ & \text{delay } \mathcal{N}[[t_2]] \text{ as } x_2 \text{ in} \\ & \text{cons } x_1 \ x_2 \end{aligned}$$

The translations for `car` and `cdr` evaluate their argument, retrieve a prompt from the resulting pair, and demand its value.

$$\begin{aligned} \mathcal{N}[\text{car } t] = & \text{do } v_p \leftarrow \mathcal{N}[[t]] \\ & \text{in need } (\text{car } v_p) \end{aligned}$$

$$\begin{aligned} \mathcal{N}[\text{cdr } t] = & \text{do } v_p \leftarrow \mathcal{N}[[t]] \\ & \text{in need } (\text{cdr } v_p) \end{aligned}$$

9. CONCLUSIONS

In this paper, we expose and examine the operational structure of lazy evaluation as embodied in call-by-need semantics. We present this understanding in two ways: as an abstract machine whose operational behavior involves control capture, and as a simulation of call-by-need under call-by-value plus delimited control operations. Delimited control can be used to simulate a global heap, but our particular simulation uses delimited control operations to manage laziness locally, just like the calculus reduction rules.

The artifacts of this investigation provide new tools for increasing our understanding of lazy evaluation and its connections to control. The abstract machine could be used to establish connections to heap-based implementations of call-by-need, and possibly modern graph-reduction based formulations (Peyton Jones and Salkild, 1989). In fact it seems that the calculus and abstract machine may point out new structural and dynamic invariants that are inherent to call-by-need evaluation but are hidden in the unstructured representations of heaps.

The abstract machine and simulation might also provide new opportunities for reasoning about the correctness of transformations applied to call-by-need programs. Surely the calculus provides the same equational reasoning powers as the abstract machine. However the machine may enable researchers to more easily conceive transformations and justifications that are not as easy to recognize in the reduction semantics. Our simulation might be connected to that of Okasaki et al. (1994). The simulation might suggest new mechanisms by which to embed call-by-need evaluation within call-by-value programs.

One significant difference between the two formulations of call-by-need lambda calculi (Maraist et al., 1998; Ariola and Felleisen, 1997) is the status of variables. Maraist et al. consider variables to be values, whereas Ariola and Felleisen do not. Ultimately, Maraist et al. (1998) prove standardization against a standard-order relation that does not consider variables to be values. This paper sheds no light on the inclusion of variables among the values, however it demonstrates in stark detail the consequences of the latter design. In the abstract machine, the transition rules for lambda terms, namely the rebuild rules, differ significantly from the transition rules for variables, the need rules. A similar

distinction can be seen simply by observing the complexity of their respective translations. In short, our semantics interpret variables as memoized computations rather than values. Variables can be treated as values under deterministic call-by-value and call-by-name reduction; it remains an open question whether the same could be achieved for call-by-need and if so what its operational implications would be.

Our results reveal that a proliferation of semantic frameworks—reduction semantics, machine semantics, etc—is a boon and not a crisis. The reduction semantics of call-by-need elegantly and mysteriously encode a rich semantics whose broad implications can be seen in equivalent machine semantics and continuation semantics. As such, our work provides new perspectives from which to reason about call-by-need, delimited control, and their respective expressive powers.

10. ACKNOWLEDGEMENTS

We thank Daniel P. Friedman, Roshan James, William Byrd, Michael Adams, and the rest of the Indiana University Programming Languages Group, as well as Jeremy Siek, Zena Ariola, Phil Wadler, Olivier Danvy, and anonymous referees for helpful discussions and feedback on this work.

REFERENCES

- ABADI, M., CARDELLI, L., CURIEN, P.-L., AND LÉVY, J.-J. 1991. Explicit substitutions. *Journal of Functional Programming* 1, 4, 375–416.
- ARIOLA, Z. AND KLOP, J. W. 1994. Cyclic lambda graph rewriting. In *Logic in Computer Science (LICS '94)*. IEEE Computer Society Press, Los Alamitos, Ca., USA, 416–425.
- ARIOLA, Z. M. AND FELLEISEN, M. 1997. The call-by-need lambda calculus. *Journal of Functional Programming* 7, 3 (May), 265–301.
- ARIOLA, Z. M. AND KLOP, J. W. 1997. Lambda calculus with explicit recursion. *Information and Computation* 139, 2, 154–233.
- ARIOLA, Z. M., MARAIST, J., ODERSKY, M., FELLEISEN, M., AND WADLER, P. 1995. A call-by-need lambda calculus. In *POPL '95: Proceedings of the 22nd ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*. ACM Press, New York, NY, USA, 233–246.
- BALAT, V., DI COSMO, R., AND FIORE, M. 2004. Extensional normalisation and type-directed partial evaluation for typed lambda calculus with sums. In *POPL '04: Proceedings of the 31st ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*. ACM, New York, NY, USA, 64–76.
- BARENDREGT, H. P. 1981. *The Lambda Calculus, its Syntax and Semantics*. North-Holland, Amsterdam, NL. Studies in Logic and the Foundations of Mathematics.
- BIERNACKA, M. AND DANVY, O. 2007. A concrete framework for environment machines. *ACM Transactions on Computational Logic* 9, 1, 6.
- BIERNACKI, D., DANVY, O., AND CHIEH SHAN, C. 2005. On the dynamic extent of delimited continuations. *Information Processing Letters* 96, 1, 7 – 17.
- CURRY, H. B. AND FEYS, R. 1958. *Combinatory Logic, Volume I*. Studies in Logic and the Foundations of Mathematics. North-Holland, Amsterdam. Second printing 1968.

- DANVY, O. AND FILINSKI, A. 1990. Abstracting control. In *LFP '90: Proceedings of the 1990 ACM Conference on LISP and Functional Programming*. ACM, New York, NY, USA, 151–160.
- DANVY, O., MILLIKIN, K., MUNK, J., AND ZERNY, I. 2010. Defunctionalized interpreters for call-by-need evaluation. In *FLOPS '10: Proceedings of the Tenth International Symposium on Functional and Logic Programming*. Springer-Verlag, London, UK. To appear.
- DANVY, O. AND NIELSEN, L. R. 2004. Refocusing in reduction semantics. Tech. Rep. RS-04-26, BRICS, DAIMI, Department of Computer Science, University of Aarhus, Aarhus, Denmark. November.
- DYBVIK, R. K., PEYTON JONES, S., AND SABRY, A. 2007. A monadic framework for delimited continuations. *Journal of Functional Programming* 17, 6, 687–730.
- FELLEISEN, M. 1988. The theory and practice of first-class prompts. In *POPL '88: Proceedings of the 15th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*. ACM, New York, NY, USA, 180–190.
- FELLEISEN, M., FINDLER, R., AND FLATT, M. 2009. *Semantics Engineering with PLT Redex*. MIT Press, Cambridge, MA.
- FELLEISEN, M. AND FRIEDMAN, D. P. 1986. Control operators, the SECD-machine, and the λ -calculus. In *Formal Description of Programming Concepts*, M. Wirsing, Ed. North-Holland, Netherlands, 193–217.
- FELLEISEN, M. AND HIEB, R. 1992. A revised report on the syntactic theories of sequential control and state. *Theoretical Computer Science* 103, 2, 235–271.
- FRIEDMAN, D. P., GHULOUM, A., SIEK, J. G., AND WINEBARGER, O. L. 2007. Improving the lazy Krivine machine. *Higher-Order and Symbolic Computation* 20, 3, 271–293.
- FRIEDMAN, D. P. AND WISE, D. S. 1976. CONS should not evaluate its arguments. In *Automata, Languages and Programming*, S. Michaelson and R. Milner, Eds. Edinburgh University Press, Edinburgh, Scotland, 257–284.
- GIBBONS, J. AND WANSBROUGH, K. 1996. Tracing lazy functional languages. In *Proceedings of Conference on Computing: The Australian Theory Symposium*, M. E. Houle and P. Eades, Eds. Australian Computer Science Communications, Townsville, 11–20.
- HENDERSON, P. AND MORRIS, JR., J. H. 1976. A lazy evaluator. In *POPL '76: Proceedings of the 3rd ACM SIGACT-SIGPLAN Symposium on Principles of Programming Languages*. ACM, New York, NY, USA, 95–103.
- KAMEYAMA, Y., KISELYOV, O., AND SHAN, C. 2008. Closing the stage: From staged code to typed closures. In *PEPM '08: Proceedings of the 2008 ACM SIGPLAN Symposium on Partial Evaluation and Semantics-based Program Manipulation*. ACM, New York, NY, USA, 147–157.
- KISELYOV, O., SHAN, C., FRIEDMAN, D. P., AND SABRY, A. 2005. Backtracking, interleaving, and terminating monad transformers: (functional pearl). In *ICFP '05: Proceedings of the Tenth ACM SIGPLAN International Conference on Functional Programming*. ACM, New York, NY, USA, 192–203.
- KISELYOV, O., SHAN, C., AND SABRY, A. 2006. Delimited dynamic binding. In *ICFP '06: Proceedings of the eleventh ACM SIGPLAN international conference on Functional programming*. ACM, New York, NY, USA, 26–37.
- KRIVINE, J.-L. 2007. A call-by-name lambda-calculus machine. *Higher-Order and Symbolic Computation* 20, 3, 199–207.
- LAUNCHBURY, J. 1993. A natural semantics for lazy evaluation. In *POPL '93: Proceedings of the 20th ACM SIGPLAN-SIGACT Symposium on Principles of Programming*

- Languages*. ACM, New York, NY, USA, 144–154.
- MARAIST, J., ODERSKY, M., AND WADLER, P. 1998. The call-by-need lambda calculus. *Journal of Functional Programming* 8, 3 (May), 275–317.
- MOGGI, E. AND SABRY, A. 2004. An abstract monadic semantics for value recursion. *Theoretical Informatics and Applications* 38, 4, 375–400.
- NAKATA, K. AND HASEGAWA, M. 2009. Small-step and big-step semantics for call-by-need. *Journal of Functional Programming* 19, 06, 699–722.
- OKASAKI, C., LEE, P., AND TARDITI, D. 1994. Call-by-need and continuation-passing style. *Lisp and Symbolic Computation* 7, 1 (Jan.), 57–81.
- PEYTON JONES, S. L. AND SALKILD, J. 1989. The spineless tagless G-machine. In *FPCA '89: Proceedings of the Fourth International Conference on Functional Programming Languages and Computer Architecture*. ACM, New York, NY, USA, 184–201.
- PLOTKIN, G. D. 1975. Call-by-name, call-by-value and the λ -calculus. *Theoretical Computer Science* 1, 2 (Dec.), 125–159.
- SESTOFT, P. 1997. Deriving a lazy abstract machine. *Journal of Functional Programming* 7, 3, 231–264.
- SHAN, C. 2007. A static simulation of dynamic delimited control. *Higher-Order and Symbolic Computation* 20, 4, 371–401.
- SUSSMAN, G. J. AND STEELE JR., G. L. 1975. Scheme: An interpreter for extended lambda calculus. AI Memo 349, Artificial Intelligence Laboratory, Massachusetts Institute of Technology, Cambridge, Massachusetts. Dec. Reprinted in *Higher-Order and Symbolic Computation* 11(4):405–439, 1998, with a foreword (Sussman and Steele Jr., 1998).
- SUSSMAN, G. J. AND STEELE JR., G. L. 1998. The first report on Scheme revisited. *Higher-Order and Symbolic Computation* 11, 4, 399–404.
- WAND, M. AND VAILLANCOURT, D. 2004. Relating models of backtracking. In *ICFP '04: Proceedings of the Ninth ACM SIGPLAN International Conference on Functional Programming*. ACM, New York, NY, USA, 54–65.
- WANG, C. 1990. Obtaining lazy evaluation with continuations in Scheme. *Information Processing Letters* 35, 2, 93–97.
- XI, H. 1997. Evaluation under lambda abstraction. In *PLILP '97: Proceedings of the Ninth International Symposium on Programming Languages: Implementations, Logics, and Programs*. Springer-Verlag, London, UK, 259–273.