

# The Call-by-Need Lambda Calculus, Revisited

Stephen Chang and Matthias Felleisen

College of Computer Science  
Northeastern University  
Boston, Massachusetts, USA  
{stchang,matthias}@ccs.neu.edu

**Abstract.** The existing call-by-need  $\lambda$  calculi describe lazy evaluation via equational logics. A programmer can use these logics to safely ascertain whether one term is behaviorally equivalent to another or to determine the value of a lazy program. However, neither of the existing calculi models evaluation in a way that matches lazy implementations.

Both calculi suffer from the same two problems. First, the calculi never discard function calls, even after they are completely resolved. Second, the calculi include re-association axioms even though these axioms are merely administrative steps with no counterpart in any implementation.

In this paper, we present an alternative axiomatization of lazy evaluation using a single axiom. It eliminates both the function call retention problem and the extraneous re-association axioms. Our axiom uses a grammar of contexts to describe the exact notion of a *needed computation*. Like its predecessors, our new calculus satisfies consistency and standardization properties and is thus suitable for reasoning about behavioral equivalence. In addition, we establish a correspondence between our semantics and Launchbury's natural semantics.

**Keywords:** call-by-need, laziness, lambda calculus.

## 1 A Short History of the $\lambda$ Calculus

Starting in the late 1950s, programming language researchers began to look to Church's  $\lambda$  calculus [6] for inspiration. Some used it as an analytic tool to understand the syntax and semantics of programming languages, while others exploited it as the basis for new languages. By 1970, however, a disconnect had emerged in the form of call-by-value programming, distinct from the notion of  $\beta$  and normalization in Church's original calculus. Plotkin [25] reconciled the  $\lambda$  calculus and Landin's SECD machine for the ISWIM language [16] with the introduction of a notion of correspondence and with a proof that two distinct variants of the  $\lambda$  calculus corresponded to two distinct variants of the ISWIM programming language: one for call-by-value and one for call-by-name.

In the early 1970s, researchers proposed call-by-need [12, 14, 28], a third kind of parameter passing mechanism that could be viewed as yet another variant of the ISWIM language. Call-by-need is supposed to represent the best of both worlds. While call-by-value ISWIM always evaluates the argument of a function,

the call-by-name variant evaluates the argument every time it is needed. Hence, if an argument (or some portion) is never needed, call-by-name wins; otherwise call-by-value is superior because it avoids re-evaluation of arguments. Call-by-need initially proceeds like call-by-name, evaluating a function's body before the argument—until the value of the argument is needed; at that point, the argument is evaluated and the resulting value is used from then onward. In short, call-by-need evaluates an argument at most once, and only if needed.

Since then, researchers have explored a number of characterizations of call-by-need [8, 11, 13, 15, 23, 24, 26]. Concerning this paper, three stand out. Launchbury's semantics [17] specifies the meaning of complete programs with a Kahn-style natural semantics. The call-by-need  $\lambda$  calculi of Ariola and Felleisen [2–4], and of Maraist, Odersky, and Wadler [4, 20, 21] are equational logics in the spirit of the  $\lambda$  calculus.

The appeal of the  $\lambda$  calculus has several reasons. First, a calculus is sound with respect to the observational (behavioral) equivalence relation [22]. It can therefore serve as the starting point for other, more powerful logics. Second, its axioms are rich enough to mimic machine evaluation, meaning programmers can reduce programs to values without thinking about implementation details. Finally, the  $\lambda$  calculus gives rise to a substantial meta-theory [5, 7] from which researchers have generated useful and practical results for its cousins.

Unfortunately, neither of the existing by-need calculi model lazy evaluation in a way that matches lazy language implementations. Both calculi suffer from the same two problems. First, unlike the by-name and by-value calculi, the by-need calculi never discard function calls, even after the call is resolved and the argument is no longer needed. Lazy evaluation does require some accumulation of function calls due to the delayed evaluation of arguments but the existing calculi adopt the extreme solution of retaining every call. Indeed, the creators of the existing calculi acknowledge that a solution to this problem would strengthen their work but they could not figure out a proper solution.

Second, the calculi include re-association axioms even though these axioms have no counterpart in any implementation. The axioms are mere administrative steps, needed to construct  $\beta$ -like redexes. Hence, they should not be considered computationally on par with other axioms.

In this paper, we overcome these problems with an alternative axiomatization. Based on a single axiom, it avoids the retention of function calls and eliminates the extraneous re-association axioms. The single axiom uses a grammar of contexts to describe the exact notion of a *needed computation*. Like its predecessors, our new calculus satisfies consistency and standardization properties and is thus suitable for reasoning about behavioral equivalence. In addition, we establish an intensional correspondence with Launchbury's semantics.

The second section of this paper recalls the two existing by-need calculi in some detail. The third section presents our new calculus, as well as a way to derive it from Ariola and Felleisen's calculus. Sections 4 and 5 show that our calculus satisfies the usual meta-theorems and that it is correct with respect to Launchbury's semantics. Finally, we discuss some possible extensions.

## 2 The Original Call-by-Need $\lambda$ Calculi

The original call-by-need  $\lambda$  calculi are independently due to two groups: Ariola and Felleisen [2, 3] and Maraist, et al. [20, 21]. They were jointly presented at POPL in 1995 [4]. Both calculi use the standard set of terms as syntax:

$$e = x \mid \lambda x.e \mid e e \quad (\text{Terms})$$

Our treatment of syntax employs the usual conventions, including Barendregt’s standard hygiene condition for variable bindings [5]. Figure 1 specifies the calculus of Maraist et al.,  $\lambda_{mow}$ , and  $\lambda_{af}$ , Ariola and Felleisen’s variant. Nonterminals in some grammar productions have subscript tags to differentiate them from similar sets elsewhere in the paper. Unsubscripted definitions have the same denotation in all systems.

$v_m = x \mid \lambda x.e$	$v = \lambda x.e$
$C = [ \ ] \mid \lambda x.C \mid C e \mid e C$	$a_{af} = v \mid (\lambda x.a_{af}) e$
$E_{af} = [ \ ] \mid E_{af} e \mid (\lambda x.E_{af}) e \mid (\lambda x.E_{af}[x]) E_{af}$	$E_{af} = [ \ ] \mid E_{af} e \mid (\lambda x.E_{af}) e \mid (\lambda x.E_{af}[x]) E_{af}$
$(\lambda x.C[x]) v_m = (\lambda x.C[v_m]) v_m \quad (\mathcal{V})$	$(\lambda x.E_{af}[x]) v = (\lambda x.E_{af}[v]) v \quad (\text{deref})$
$(\lambda x.e_1) e_2 e_3 = (\lambda x.e_1 e_3) e_2 \quad (\mathcal{C})$	$(\lambda x.a_{af}) e_1 e_2 = (\lambda x.a_{af} e_2) e_1 \quad (\text{lift})$
$(\lambda x.e_1)((\lambda y.e_2) e_3) =$ $\quad (\lambda y.(\lambda x.e_1) e_2) e_3 \quad (\mathcal{A})$	$(\lambda x.E_{af}[x]) ((\lambda y.a_{af}) e) =$ $\quad (\lambda y.(\lambda x.E_{af}[x]) a_{af}) e \quad (\text{assoc})$
$(\lambda x.e_1) e_2 = e_1, x \notin \text{fv}(e_1) \quad (\mathcal{G})$	$(\lambda y.(\lambda x.E_{af}[x]) a_{af}) e$

**Fig. 1.** Existing call-by-need  $\lambda$  calculi (left:  $\lambda_{mow}$ , right:  $\lambda_{af}$ )

In both calculi, the analog to the  $\beta$  axiom—also called a *basic notion of reduction* [5]—replaces variable occurrences, one at a time, with the value of the function’s argument. Value substitution means that there is no duplication of work as far as argument evaluation is concerned. The function call is retained because additional variable occurrences in the function body may need the argument. Since function calls may accumulate, the calculi come with axioms that re-associate bindings to pair up functions with their arguments. For example, re-associating  $(\lambda x.(\lambda y.\lambda z.z) v_y) v_x v_z$  in  $\lambda_{af}$  exposes a *deref* redex:

$$(\lambda x.(\lambda y.\lambda z.z) v_y) v_x v_z \xrightarrow{\text{lift}} (\lambda x.(\lambda y.\lambda z.z) v_y v_z) v_x \xrightarrow{\text{lift}} (\lambda x.(\lambda y.(\lambda z.z) v_z) v_y) v_x$$

The two calculi differ from each other in their timing of variable replacements. The  $\lambda_{mow}$  calculus allows the replacement of a variable with its value anywhere in the body of its binding  $\lambda$ . The  $\lambda_{af}$  calculus replaces a variable with its argument only if evaluation of the function body needs it, where “need” is formalized via so-called evaluation contexts ( $E_{af}$ ). Thus evaluation contexts in  $\lambda_{af}$  serve the

double purpose of specifying demand for arguments and the standard reduction strategy. The term  $(\lambda x.\lambda y.x) v$  illustrates this difference between the two calculi. According to  $\lambda_{mow}$ , the term is a  $\mathcal{V}$  redex and reduces to  $(\lambda x.\lambda y.v) v$ , whereas in  $\lambda_{af}$ , the term is irreducible because the  $x$  occurs in an inner, unapplied  $\lambda$ , and is thus not “needed.”

Also,  $\lambda_{mow}$  is more lenient than  $\lambda_{af}$  when it comes to re-associations. The  $\lambda_{af}$  calculus re-associates the left or right hand side of an application only if it has been completely reduced to an answer, but  $\lambda_{mow}$  permits re-association as soon as one nested function layer is revealed. In short,  $\lambda_{mow}$  proves more equations than  $\lambda_{af}$ , i.e.,  $\lambda_{af} \subset \lambda_{mow}$ .

In  $\lambda_{af}$ , programs reduce to answers:

$$\mathbf{eval}_{af}(e) = \mathbf{done} \text{ iff there exists an answer } a_{af} \text{ such that } \lambda_{af} \vdash e = a_{af}$$

In contrast, Maraist et al. introduce a “garbage collection” axiom into  $\lambda_{mow}$  to avoid answers and to use values instead. This suggests the following definition:

$$\mathbf{eval}_{mow}(e) = \mathbf{done} \text{ iff there exists a value } v_m \text{ such that } \lambda_{mow} \vdash e = v_m$$

This turns out to be incorrect, however. Specifically, let  $\mathbf{eval}_{name}$  be the analogous call-by-name evaluator. Then  $\mathbf{eval}_{af} = \mathbf{eval}_{name}$  but  $\mathbf{eval}_{mow} \neq \mathbf{eval}_{name}$ . Examples such as  $(\lambda x.\lambda y.x) \Omega$  confirm the difference.

In recognition of this problem, Maraist et al. use Ariola and Felleisen’s axioms and evaluation contexts to create their Curry-Feys-style standard reduction sequences. Doing so reveals the inconsistency of  $\lambda_{mow}$  with respect to Plotkin’s *correspondence criteria* [25]. According to Plotkin, a useful calculus *corresponds to* a programming language, meaning its axioms (1) satisfy the Church-Rosser and Curry-Feys Standardization properties, and (2) define a standard reduction function that is equal to the evaluation function of the programming language. Both the call-by-name and the call-by-value  $\lambda$  calculi satisfy these criteria with respect to call-by-name and call-by-value SECD machines for ISWIM, respectively. So does  $\lambda_{af}$  with respect to a call-by-need SECD machine, but some of  $\lambda_{mow}$ ’s axioms cannot be used as standard reduction relations.

Finally, the inclusion of  $\mathcal{G}$  is a brute-force attempt to address the function call retention problem. Because  $\mathcal{G}$  may discard arguments even before the function is called, both sets of authors consider it too coarse and acknowledge that a tighter solution to the function call retention issue would “strengthen the calculus and its utility for reasoning about the implementations of lazy languages” [4].

### 3 A New Call-by-Need $\lambda$ Calculus

Our new calculus,  $\lambda_{need}$ , uses a single axiom,  $\beta_{need}$ . The new axiom evaluates the argument when it is first demanded, replaces all variable occurrences with that result, and then discards the argument and thus the function call. In addition, the axiom performs the required administrative scope adjustments as part of the same step, rendering explicit re-association axioms unnecessary. In short, every reduction step in our calculus represents computational progress.

Informally, to perform a reduction, three components must be identified:

1. the next demanded variable,
2. the function that binds that demanded variable,
3. and the argument to that function.

In previous by-need calculi the re-association axioms rewrite a term so that the binding function and its argument are adjacent.

Without the re-association axioms, finding the function that binds the demanded variable and its argument requires a different kind of work. The following terms show how the demanded variable, its binding function, and its argument can appear at seemingly arbitrary locations in a program:

- $(\underline{\lambda x} . (\underline{\lambda y} . \underline{\lambda z} . \underline{x}) e_y) e_x e_z$
- $(\underline{\lambda x} . (\underline{\lambda y} . \underline{\lambda z} . \underline{y}) e_y) e_x e_z$
- $(\underline{\lambda x} . (\underline{\lambda y} . \underline{\lambda z} . \underline{z}) e_y) e_x e_z$

Our  $\beta_{need}$  axiom employs a grammar of contexts to describe the path from a demanded variable to its binding function and from there to its argument.

The first subsection explains the syntax and the contexts of  $\lambda_{need}$  in a gradual fashion. The second subsection presents the  $\beta_{need}$  axiom and also shows how to derive it from Ariola and Felleisen's  $\lambda_{af}$  calculus.

### 3.1 Contexts

Like the existing by-need calculi, the syntax of our calculus is that of Church's original calculus. In  $\lambda_{need}$ , calculations evaluate terms  $e$  to answers  $A[v]$ , which generalize answers from Ariola and Felleisen's calculus:

$$e = x \mid \lambda x . e \mid e e \quad (\text{Terms})$$

$$v = \lambda x . e \quad (\text{Values})$$

$$a = A[v] \quad (\text{Answers})$$

$$A = [ \ ] \mid A[\lambda x . A] e \quad (\text{Answer Contexts})$$

Following Ariola and Felleisen, the basic axiom uses evaluation contexts to specify the notion of demand for variables:

$$E = [ \ ] \mid E e \mid \dots \quad (\text{Evaluation Contexts})$$

The first two kinds, taken from  $\lambda_{af}$ , specify that a variable is demanded, and that a variable in the operator position of an application is demanded, respectively.

Since the calculus is to model program evaluation, we are primarily interested in demanded variables under a  $\lambda$ -abstraction. This kind of evaluation context is defined using an answer context  $A$ :

$$E = \dots \mid A[E] \mid \dots \quad (\text{Another Evaluation Context})$$

Using answer contexts, this third evaluation context dictates that demand exists under a  $\lambda$  if a corresponding argument exists for that  $\lambda$ . Note how *an answer context descends under the same number of  $\lambda$ s as arguments for those  $\lambda$ s*. In particular, for any term  $A[\lambda x.e_1] e_2$ ,  $e_2$  is always the argument of  $\lambda x.e_1$ . The third evaluation context thus generalizes the function-is-next-to-argument requirement found in both call-by-name and call-by-value. The generalization is needed due to the retention of function calls in  $\lambda_{need}$ .

Here are some example answer contexts that might be used:

$$\begin{aligned}
 A_0 &= \underbrace{(\lambda x.[ \ ] e_x)} \\
 A_1 &= (\lambda x. \underbrace{(\lambda y.[ \ ] e_y)} e_x) \\
 A_2 &= (\lambda x. (\lambda y. \underbrace{(\lambda z.[ \ ] e_z)} e_y) e_x)
 \end{aligned}$$

An underbrace matches each function to its argument. The examples all juxtapose functions and their arguments. In contrast, the next two separate functions from their arguments:

$$\begin{aligned}
 A_3 &= (\lambda x. \underbrace{\lambda y. \lambda z. [ \ ]} e_x) e_y e_z \\
 A_4 &= (\lambda x. \underbrace{(\lambda y. \lambda z. [ \ ] e_y)} e_x) e_z
 \end{aligned}$$

To summarize thus far, when a demanded variable is discovered under a  $\lambda$ , the surrounding context looks like this:

$$A[E[x]]$$

where both the function binding  $x$  and its argument are in  $A$ . The decomposition of the surrounding context into  $A$  and  $E$  assumes that  $A$  encompasses as many function-argument pairs as possible; in other words, it is impossible to merge the outer part of  $E$  with  $A$  to form a larger answer context.

To know which argument corresponds to the demanded variable, we must find the  $\lambda$  that binds  $x$  in  $A$ . To this end, we split answer contexts so that we can “highlight” a function-argument pair within the context:

$$\begin{aligned}
 \hat{A} &= [ \ ] \mid A[\hat{A}] e && \text{(Partial Answer Contexts–Outer)} \\
 \check{A} &= [ \ ] \mid A[\lambda x.\check{A}] && \text{(Partial Answer Contexts–Inner)}
 \end{aligned}$$

Using these additional contexts, any answer context can be decomposed into

$$\hat{A}[A[\lambda x.\check{A}[ \ ]] e]$$

where  $e$  is the argument of  $\lambda x.\check{A}[\ ]$ . For a fixed function-argument pair in an answer context, this partitioning into  $\hat{A}$ ,  $A$ , and  $\check{A}$  is unique. The  $\hat{A}$  subcontext represents the part of the answer context around the chosen function-argument pair; the  $\check{A}$  subcontext represents the part of the answer context in its body; and  $A$  here is the subcontext between the function and its argument. Naturally we must demand that  $\hat{A}$  composed with  $\check{A}$  is an answer context as well so that the overall context remains an answer context. The following table lists the various subcontexts for the example  $A_4$  for various function-argument pairs:

	$A_4 = (\lambda x.(\lambda y.\lambda z.[\ ] e_y) e_x e_z$		
$\hat{A} =$	$[\ ] e_z$	$(\lambda x.[\ ]) e_x e_z$	$[\ ]$
$A =$	$[\ ]$	$[\ ]$	$(\lambda x.(\lambda y.[\ ]) e_y) e_x$
$\check{A} =$	$(\lambda y.\lambda z.[\ ]) e_y$	$\lambda z.[\ ]$	$[\ ]$
$A_4 =$	$\hat{A}[A[\lambda x.\check{A}[E[x]]] e_x]$	$\hat{A}[A[\lambda y.\check{A}[E[x]]] e_y]$	$\hat{A}[A[\lambda z.\check{A}[E[x]]] e_z]$

Now we can define the fourth kind of evaluation context:

$$E = \dots \mid \hat{A}[A[\lambda x.\check{A}[E[x]]] E], \quad \text{where } \hat{A}[\check{A}] \in A \quad (\text{Final Eval. Context})$$

This final evaluation context shows how demand shifts to an argument when a function parameter is in demand within the function body.

### 3.2 The $\beta_{need}$ Axiom and a Derivation

Figure 2 summarizes the syntax of  $\lambda_{need}$  as developed in the preceding section.<sup>1</sup> In this section we use these definitions to formulate the  $\beta$  axiom for our calculus.

$$\begin{aligned}
e &= x \mid \lambda x.e \mid e e && (\text{Terms}) \\
v &= \lambda x.e && (\text{Values}) \\
a &= A[v] && (\text{Answers}) \\
A &= [\ ] \mid A[\lambda x.A] e && (\text{Answer Contexts}) \\
\hat{A} &= [\ ] \mid A[\hat{A}] e && (\text{Partial Answer Contexts--Outer}) \\
\check{A} &= [\ ] \mid A[\lambda x.\check{A}] && (\text{Partial Answer Contexts--Inner}) \\
E &= [\ ] \mid E e \mid A[E] \mid \hat{A}[A[\lambda x.\check{A}[E[x]]] E], && (\text{Evaluation Contexts}) \\
&\quad \text{where } \hat{A}[\check{A}] \in A
\end{aligned}$$

**Fig. 2.** The syntax and contexts of the new call-by-need  $\lambda$  calculus,  $\lambda_{need}$

Here is the single axiom of  $\lambda_{need}$ :

$$\begin{aligned}
\hat{A}[A_1[\lambda x.\check{A}[E[x]]] A_2[v]] &= \hat{A}[A_1[A_2[\check{A}[E[x]]\{x:=v\}]]], && (\beta_{need}) \\
&\quad \text{where } \hat{A}[\check{A}] \in A
\end{aligned}$$

<sup>1</sup> We gratefully acknowledge Casey Klein's help with the  $A$  production.

A  $\beta_{need}$  redex determines which parameter  $x$  of some function is “in demand” and how to locate the corresponding argument  $A_2[v]$ , which might be an answer not necessarily a value. The contexts from the previous section specify the path from the binding position ( $\lambda$ ) to the variable occurrence and the argument. A  $\beta_{need}$  reduction substitutes the value in  $A_2[v]$  for all free occurrences of the function parameter—just like in other  $\lambda$  calculi. In the process, the function call is discarded. Since the argument has been reduced to a value, there is no duplication of work, meaning our calculus satisfies the requirements of lazy evaluation. Lifting  $A_2$  to the top of the evaluation context ensures that its bindings remain intact and visible for  $v$ .

Here is a sample reduction in  $\lambda_{need}$ , where  $\longrightarrow$  is the one-step reduction:

$$\begin{aligned} & ((\lambda x. (\lambda y. \underline{\lambda z. \mathbf{z}} \ y \ x) \ \lambda y. y) \ \lambda x. x) \ \underline{\lambda z. z} & (1) \\ \longrightarrow & (\lambda x. (\lambda y. (\underline{\lambda z. \underline{z}}) \ \mathbf{y} \ x) \ \underline{\lambda y. y}) \ \lambda x. x & (2) \\ \longrightarrow & (\lambda x. ((\underline{\lambda z. \mathbf{z}}) \ \underline{\lambda y. y}) \ x) \ \lambda x. x & (3) \\ \longrightarrow & (\lambda x. (\underline{\lambda y. y}) \ \mathbf{x}) \ \underline{\lambda x. x} & (4) \end{aligned}$$

The “in demand” variable is in bold; its binding  $\lambda$  and argument are underlined. Line 1 is an example of a reduction that involves a non-adjoined function and argument pair. In line 2, the demand for the value of  $z$  (twice underlined) triggers a demand for the value of  $y$ ; line 4 contains a similar demand chain.

$$\begin{aligned} v &= \lambda x. e && \text{(Values)} \\ a_{af} &= A_{af}[v] && \text{(Answers)} \\ A_{af} &= [ \ ] \mid (\lambda x. A_{af}) e && \text{(Answer Contexts)} \\ E_{af} &= [ \ ] \mid E_{af} e \mid A_{af}[E_{af}] \mid (\lambda x. E_{af}[x]) E_{af} && \text{(Evaluation Contexts)} \\ & (\lambda x. E_{af}[x]) v = E_{af}[x] \{x := v\} && (\beta'_{need}) \\ & (\lambda x. A_{af}[v]) e_1 e_2 = (\lambda x. A_{af}[v e_2]) e_1 && (lift') \\ & (\lambda x. E_{af}[x]) ((\lambda y. A_{af}[v]) e) = (\lambda y. A_{af}[(\lambda x. E_{af}[x]) v]) e && (assoc') \end{aligned}$$

**Fig. 3.** A modified calculus,  $\lambda_{af-mod}$

To furnish additional intuition into  $\beta_{need}$ , we use the rest of the section to derive it from the axioms of  $\lambda_{af}$ . The  $\lambda_{af-mod}$  calculus in figure 3 combines  $\lambda_{af}$  with two insights. First, Garcia et al. [13] observed that when the answers in  $\lambda_{af}$ ’s *lift* and *assoc* redexes are nested deeply, multiple re-associations are performed consecutively. Thus we modify *lift* and *assoc* to perform all these re-associations in one step.<sup>2</sup> The modified calculus defines answers via answer contexts,  $A_{af}$ , and the modified *lift'* and *assoc'* axioms utilize these answer contexts to do the multi-step re-associations. Thus programs in this modified calculus reduce to answers

<sup>2</sup> The same modifications cannot be applied to  $\mathcal{C}$  and  $\mathcal{A}$  in  $\lambda_{mow}$  because they allow earlier re-association and thus not all the re-associations are performed consecutively.



$A_{af}[v]$ . Also, the  $A_{af}$  answer contexts are identical to the third kind of evaluation context in  $\lambda_{af-mod}$  and the new definition of  $E_{af}$  reflects this relationship.

Second, Maraist et al. [19] observed that once an argument is reduced to a value, all substitutions can be performed at once. The  $\beta'_{need}$  axiom exploits this idea and performs a full substitution. Obviously  $\beta'_{need}$  occasionally performs more substitutions than *deref*. Nevertheless, any term with an answer in  $\lambda_{af}$  likewise has an answer when reducing with  $\beta'_{need}$ .

Next an inspection of the axioms shows that the contractum of a *assoc'* redex contains a  $\beta'_{need}$  redex. Thus the *assoc'* re-associations and  $\beta'_{need}$  substitutions can be performed with one merged axiom:<sup>3</sup>

$$(\lambda x.E_{af}[x]) A_{af}[v] = A_{af}[E_{af}[x]\{x:=v\}] \quad (\beta''_{need})$$

The final step is to merge *lift'* with  $\beta''_{need}$ , which requires our generalized answer and evaluation contexts. A naïve attempt may look like this:

$$A_1[\lambda x.E[x]] A_2[v] = A_1[A_2[E[x]\{x:=v\}]] \quad (\beta'''_{need})$$

As the examples in the preceding subsection show, however, the binding occurrence for the “in demand” parameter  $x$  may not be the inner-most binding  $\lambda$  once the re-association axioms are eliminated. That is, in comparison with  $\beta_{need}$ ,  $\beta'''_{need}$  incorrectly assumes  $E$  is always next to the binder. We solve this final problem with the introduction of partial answer contexts.

## 4 Consistency, Determinism, and Soundness

If a calculus is to model a programming language, it must satisfy some essential properties, most importantly a Church-Rosser theorem and a Curry-Feys standardization theorem [25]. The former guarantees consistency of evaluation; that is, we can define an evaluator *function* with the calculus. The latter implies that the calculus comes with a *deterministic* evaluation strategy. Jointly these properties imply the calculus is *sound* with respect to observational equivalence.

### 4.1 Consistency: Church-Rosser

The  $\lambda_{need}$  calculus defines an evaluator for a by-need language:

$$\text{eval}_{need}(e) = \text{done} \text{ iff there exists an answer } a \text{ such that } \lambda_{need} \vdash e = a$$

To prove that the evaluator is indeed a (partial) function, we prove that the notion of reduction satisfies the Church-Rosser property.

**Theorem 1.** *eval<sub>need</sub> is a partial function.*

*Proof.* The theorem is a direct consequence of lemma 1 (Church-Rosser).

<sup>3</sup> Danvy et al. [8] dub a  $\beta''_{need}$  redex a “potential redex” in unrelated work.

Our strategy is to define a parallel reduction relation for  $\lambda_{need}$  [5]. Define  $\rightarrow$  to be the compatible closure of a  $\beta_{need}$  reduction, and  $\twoheadrightarrow$  to be the reflexive, transitive closure of  $\rightarrow$ . Additionally, define  $\Rightarrow$  to be the relation that reduces  $\beta_{need}$  redexes in parallel.

**Definition 1** ( $\Rightarrow$ ).

$$\begin{aligned}
 e &\Rightarrow e \\
 \hat{A}[A_1[\lambda x.\check{A}[E[x]]] A_2[v]] &\Rightarrow \hat{A}'[A'_1[A'_2[\check{A}'[E'[x]]\{x:=v'\}]]], \\
 &\text{if } \hat{A}[\check{A}] \in A, \hat{A}'[\check{A}'] \in A, \hat{A} \Rightarrow \hat{A}', A_1 \Rightarrow A'_1, \\
 &A_2 \Rightarrow A'_2, \check{A} \Rightarrow \check{A}', E \Rightarrow E', v \Rightarrow v' \\
 e_1 e_2 &\Rightarrow e'_1 e'_2, \text{ if } e_1 \Rightarrow e'_1, e_2 \Rightarrow e'_2 \\
 \lambda x.e &\Rightarrow \lambda x.e', \text{ if } e \Rightarrow e'
 \end{aligned}$$

The parallel reduction relation  $\Rightarrow$  relies on notion of parallel reduction for contexts; for simplicity, we overload the relation symbol to denote both relations.

**Definition 2** ( $\Rightarrow$  for Contexts).

$$\begin{aligned}
 [ ] &\Rightarrow [ ] \\
 A_1[\lambda x.A_2] e &\Rightarrow A'_1[\lambda x.A'_2] e', \text{ if } A_1 \Rightarrow A'_1, A_2 \Rightarrow A'_2, e \Rightarrow e' \\
 A[\hat{A}] e &\Rightarrow A'[\hat{A}'] e', \text{ if } A \Rightarrow A', \hat{A} \Rightarrow \hat{A}', e \Rightarrow e' \\
 A[\lambda x.\check{A}] &\Rightarrow A'[\lambda x.\check{A}'], \text{ if } A \Rightarrow A', \check{A} \Rightarrow \check{A}' \\
 E e &\Rightarrow E' e', \text{ if } E \Rightarrow E', e \Rightarrow e' \\
 A[E] &\Rightarrow A'[E'], \text{ if } A \Rightarrow A', E \Rightarrow E' \\
 \hat{A}[A[\lambda x.\check{A}[E_1[x]]] E_2] &\Rightarrow \hat{A}'[A'[\lambda x.\check{A}'[E'_1[x]]] E'_2], \\
 &\text{if } \hat{A}[\check{A}] \in A, \hat{A}'[\check{A}'] \in A, \hat{A} \Rightarrow \hat{A}', A \Rightarrow A', \\
 &\check{A} \Rightarrow \check{A}', E_1 \Rightarrow E'_1, E_2 \Rightarrow E'_2
 \end{aligned}$$

**Lemma 1 (Church-Rosser).** *If  $e \twoheadrightarrow e_1$  and  $e \twoheadrightarrow e_2$ , then there exists a term  $e'$  such that  $e_1 \twoheadrightarrow e'$  and  $e_2 \twoheadrightarrow e'$ .*

*Proof.* By lemma 2,  $\Rightarrow$  satisfies a diamond property. Since  $\Rightarrow$  extends  $\rightarrow$ ,  $\twoheadrightarrow$  is also the transitive-reflexive closure of  $\Rightarrow$ , so  $\twoheadrightarrow$  also satisfies a diamond property.

**Lemma 2 (Diamond Property of  $\Rightarrow$ ).** *If  $e \Rightarrow e_1$  and  $e \Rightarrow e_2$ , there exists  $e'$  such that  $e_1 \Rightarrow e'$  and  $e_2 \Rightarrow e'$ .*

*Proof.* The proof proceeds by structural induction on the derivation of  $e \Rightarrow e_1$ .

## 4.2 Deterministic Behavior: Standard Reduction

A language calculus should also come with a deterministic algorithm for applying the reductions to evaluate a program. Here is our *standard reduction*:

$$E[e] \mapsto E[e'], \text{ where } e \beta_{need} e'$$

Our standard reduction strategy picks exactly one redex in a term.

**Proposition 1 (Unique Decomposition).** *For all closed terms  $e$ ,  $e$  either is an answer or  $e = E[e']$  for a unique evaluation context  $E$  and  $\beta_{need}$  redex  $e'$ .*

*Proof.* The proof proceeds by structural induction on  $e$ .

Since our calculus satisfies the unique decomposition property, we can use the standard reduction relation to define a (partial) evaluator function:

$$\mathbf{eval}_{need}^{\mathbf{sr}}(e) = \mathbf{done} \text{ iff there exists an answer } a \text{ such that } e \mapsto^* a$$

where  $\mapsto^*$  is the reflexive, transitive closure of  $\mapsto$ . Proposition 1 shows  $\mathbf{eval}_{need}^{\mathbf{sr}}$  is a function. The following theorem confirms that it equals  $\mathbf{eval}_{need}$ .

**Theorem 2.**  $\mathbf{eval}_{need} = \mathbf{eval}_{need}^{\mathbf{sr}}$

*Proof.* The theorem follows from lemma 3, which shows how to obtain a standard reduction sequence for any arbitrary reduction sequence. The front-end of the former is a series of standard reduction steps.

**Definition 3 (Standard Reduction Sequences  $\mathcal{R}$ ).**

- $x \in \mathcal{R}$
- $\lambda x.e_1 \diamond \dots \diamond \lambda x.e_m \in \mathcal{R}$ , if  $e_1 \diamond \dots \diamond e_m \in \mathcal{R}$
- $e_0 \diamond e_1 \diamond \dots \diamond e_m \in \mathcal{R}$ , if  $e_0 \mapsto e_1$  and  $e_1 \diamond \dots \diamond e_m \in \mathcal{R}$
- $(e_1 e'_1) \diamond \dots \diamond (e_m e'_m) \in \mathcal{R}$ , if  $e_1 \diamond \dots \diamond e_m, e'_1 \diamond \dots \diamond e'_m \in \mathcal{R}$ .

**Lemma 3 (Curry-Feys Standardization).**  $e \twoheadrightarrow e'$  iff there exists  $e_1 \diamond \dots \diamond e_n \in \mathcal{R}$  such that  $e = e_1$  and  $e' = e_n$ .

*Proof.* Replace  $\twoheadrightarrow$  with  $\Rightarrow$ s, and the lemma immediately follows from lemma 4.

The key to the remaining proofs is a size metric for parallel reductions.

**Definition 4 (Size of  $\Rightarrow$  Reduction).**

$$\begin{aligned} |e \Rightarrow e| &= 0 \\ |(e_1 e_2) \Rightarrow (e'_1 e'_2)| &= |e_1 \Rightarrow e'_1| + |e_2 \Rightarrow e'_2| \\ |\lambda x.e \Rightarrow \lambda x.e'| &= |e \Rightarrow e'| \\ |r| &= 1 + |\hat{A} \Rightarrow \hat{A}'| + |A_1 \Rightarrow A'_1| + |\check{A}[E[x]] \Rightarrow \check{A}'[E'[x]]| + \\ &\quad |A_2 \Rightarrow A'_2| + \#(x, \check{A}'[E'[x]]) \times |v \Rightarrow v'| \\ \text{where } r &= \hat{A}[A_1[\lambda x.\check{A}[E[x]]] A_2[v] \Rightarrow \hat{A}'[A'_1[A'_2[\check{A}'[E'[x]]\{x := v'\}]]] \\ \#(x, e) &= \text{the number of free occurrences of } x \text{ in } e \end{aligned}$$

The size of a parallel reduction of a context equals the sum of the sizes of the parallel reductions of the subcontexts and subterms that comprise the context.

**Lemma 4.** If  $e_0 \Rightarrow e_1$  and  $e_1 \diamond \dots \diamond e_n \in \mathcal{R}$ , there exists  $e_0 \diamond e'_1 \diamond \dots \diamond e'_p \diamond e_n \in \mathcal{R}$ .

*Proof.* By triple lexicographic induction on (1) length  $n$  of the given standard reduction sequence, (2)  $|e_0 \Rightarrow e_1|$ , and (3) structure of  $e_0$ .<sup>4</sup>

<sup>4</sup> We conjecture that the use of Ralph Loader's technique [18] may simplify our proof.

### 4.3 Observational Equivalence

Following Morris [22] two expressions  $e_1$  and  $e_2$  are observationally equivalent,  $e_1 \simeq e_2$ , if they are indistinguishable in all contexts. Formally,  $e_1 \simeq e_2$  if and only if  $\text{eval}_{\text{need}}(C[e_1]) = \text{eval}_{\text{need}}(C[e_2])$  for all contexts  $C$ , where

$$C = [ \ ] \mid \lambda x.C \mid C e \mid e C \quad (\text{Contexts})$$

An alternative definition of the behavioral equivalence relation uses co-induction. In either case,  $\lambda_{\text{need}}$  is sound with respect to observational equivalence.

**Theorem 3 (Soundness).** *If  $\lambda_{\text{need}} \vdash e_1 = e_2$ , then  $e_1 \simeq e_2$ .*

*Proof.* Following Plotkin, a calculus is sound if it satisfies Church-Rosser and Curry-Feys theorems.

## 5 Correctness

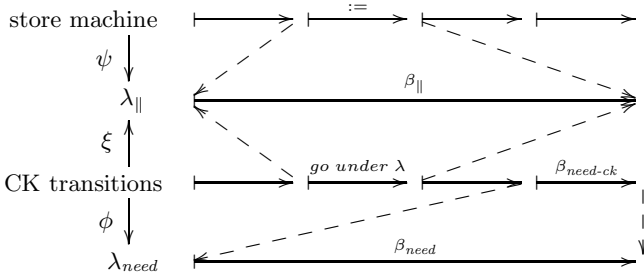
Ariola and Felleisen [3] prove that  $\lambda_{\text{af}}$  defines the same evaluation function as the call-by-name  $\lambda$  calculus. Nakata and Hasegawa [23] additionally demonstrate extensional correctness of the same calculus with respect to Launchbury’s natural semantics [17]. In this section, we show that  $\lambda_{\text{need}}$  defines the same evaluation function as Launchbury’s semantics. While our theorem statement is extensional, the proof illuminates the tight intensional relationship between the two systems.

### 5.1 Overview

The gap between the  $\lambda_{\text{need}}$  standard reduction “machine” and Launchbury’s natural semantics is huge. While the latter’s store-based natural semantics uses the equivalent of assignment statements to implement the “evaluate once, only when needed” policy, the  $\lambda_{\text{need}}$  calculus exclusively relies on term substitutions. To close the gap, we systematically construct a series of intermediate systems that makes comparisons easy, all while ensuring correctness at each step. A first step is to convert the natural semantics into a store-based machine [27].

To further bridge the gap we note that a single-use assignment statement is equivalent to a program-wide substitution of shared expressions [10]. A closely related idea is to reduce shared expressions simultaneously. This leads to a parallel program rewriting system, dubbed  $\lambda_{\parallel}$ . Equipped with  $\lambda_{\parallel}$  we get closer to  $\lambda_{\text{need}}$  but not all the way there because reductions in  $\lambda_{\text{need}}$  and  $\lambda_{\parallel}$  are too coarse-grained for direct comparison. Fortunately, it is easy to construct an intermediate transition system that eliminates the remainder of the gap. We convert  $\lambda_{\text{need}}$  to an equivalent CK transition system [9], where the program is partitioned into a control string (C) and an explicit context (K) and we show that there is a correspondence between this transition system and  $\lambda_{\parallel}$ .

Figure 4 outlines our proof strategy pictorially. The four horizontal layers correspond to the four rewriting systems. While  $\lambda_{\text{need}}$  and  $\lambda_{\parallel}$  use large steps to



**Fig. 4.** Summary of correctness proof technique

progress from term to term, the machine-like systems take several small steps. The solid vertical arrows between the layers figure indicate how mapping functions relate the rewriting sequences and the dashed arrows show how the smaller machine steps correspond to the larger steps of  $\lambda_{need}$  and  $\lambda_{\parallel}$ :

- The  $\psi$  function maps states from the store-based machine to terms in the  $\lambda_{\parallel}$  world. For every step in the natural-semantics machine, the resulting operation in  $\lambda_{\parallel}$  is either a no-op or a  $\beta_{\parallel}$  reduction, with assignment in the store-machine being equivalent to program-wide substitution in  $\lambda_{\parallel}$ .
- Similarly, the  $\xi$  function maps states of the CK transition system to the  $\lambda_{\parallel}$  space and for every CK transition, the resulting  $\lambda_{\parallel}$  operation is also either a no-op or a  $\beta_{\parallel}$  reduction, with the transition that descends under a  $\lambda$ -abstraction being equivalent to substitution in  $\lambda_{\parallel}$ .
- Finally, the  $\phi$  function maps states of the CK transition system to  $\lambda_{need}$  terms and is used to show that the CK system and  $\lambda_{need}$  are equivalent. For every CK transition, the equivalent  $\lambda_{need}$  operation is either a no-op or a  $\beta_{need}$  reduction.

Syntax

$$\begin{aligned}
 S_L &= \langle e, F_Ls, \Gamma \rangle && \text{(States)} \\
 F_Ls &= F_L, \dots && \text{(List of Frames)} \\
 F_L &= (\mathbf{arg} \ e) \mid (\mathbf{var} \ x) && \text{(Frames)} \\
 \Gamma &= (x \mapsto e, \dots) && \text{(Heaps)}
 \end{aligned}$$

Transitions

$$\begin{aligned}
 \langle e_1 \ e_2, F_Ls, \Gamma \rangle &\xrightarrow{ckh} \langle e_1, ((\mathbf{arg} \ e_2), F_Ls), \Gamma \rangle && (\text{push-arg-ckh}) \\
 \langle \lambda x. e_1, ((\mathbf{arg} \ e_2), F_Ls), \Gamma \rangle &\xrightarrow{ckh} \langle e_1 \{x := y\}, F_Ls, (\Gamma, y \mapsto e_2) \rangle, \ y \ \text{fresh} && (\text{descend-lam-ckh}) \\
 \langle x, F_Ls, (\Gamma, x \mapsto e) \rangle &\xrightarrow{ckh} \langle e, ((\mathbf{var} \ x), F_Ls), \Gamma \rangle && (\text{lookup-var-ckh}) \\
 \langle v, ((\mathbf{var} \ x), F_Ls), \Gamma \rangle &\xrightarrow{ckh} \langle v, F_Ls, (\Gamma, x \mapsto v) \rangle && (\text{update-heap-ckh})
 \end{aligned}$$

**Fig. 5.** The natural semantics as an abstract machine

Subsections 5.2 and 5.3 present the store-based machine and the parallel rewriting semantics, respectively, including a proof of equivalence. Subsection 5.4 presents the CK system and subsection 5.5 explains the rest of the proof.

### 5.2 Adapting Launchbury’s Natural Semantics

Figure 5 describes the syntax and transitions of the store machine.<sup>5</sup> It is dubbed CKH because it resembles a three-register machine [9]: a machine state  $S_L$  is comprised of a control string (C), a list of frames (K) that represents the control context in an inside-out manner, and a heap (H). The  $\dots$  notation means “zero or more of the preceding kind of element.” An  $(\mathbf{arg} e)$  frame represents the argument in an application and the  $(\mathbf{var} x)$  frame indicates that a heap expression is the current control string. Parentheses are used to group a list of frames when necessary. The initial machine state for a program  $e$  is  $\langle e, (), () \rangle$ . Computation terminates when the control string is a value and the list of frames is empty.

The *push-arg-ckh* transition moves the argument in an application to a new  $\mathbf{arg}$  frame in the frame list and makes the operator the next control string. When that operator is a  $\lambda$ -abstraction, the *descend-lam-ckh* transition adds its argument to the heap,<sup>6</sup> mapped to a fresh variable name, and makes the body of the operator the new control string. The *lookup-var-ckh* transition evaluates an argument from the heap when the control string is a variable. The mapping is removed from the heap and a new  $(\mathbf{var} x)$  frame remembers the variable whose corresponding expression is under evaluation. Finally, when the heap expression is reduced to a value, the *update-heap-ckh* transition extends the heap again.

### 5.3 Parallel Rewriting

The syntax of the parallel  $\lambda$ -rewriting semantics is as follows:

$$\begin{aligned}
 e_{\parallel} &= e \mid e_{\parallel}^x && \text{(Terms)} \\
 v_{\parallel} &= v \mid v_{\parallel}^x && \text{(Values)} \\
 E_{\parallel} &= [ \ ] \mid E_{\parallel} e_{\parallel} \mid E_{\parallel}^x && \text{(Evaluation Contexts)}
 \end{aligned}$$

This system expresses computation with a selective parallel reduction strategy. When a function application is in demand, the system substitutes the argument for all free occurrences of the bound variable, regardless of the status of the argument. When an instance of a substituted argument is reduced, however, all instances of the argument are reduced in parallel. Here is a sample reduction:

$$(\lambda x.x x) (I I) \multimap (I I)^x (I I)^x \multimap I^x I^x \multimap I$$

<sup>5</sup> To aid comparisons, we slightly alter Launchbury’s rules (and the resulting machine) to use pure  $\lambda$  terms. Thus we avoid Launchbury’s preprocessing and special syntax.

<sup>6</sup> The notation  $(\Gamma, x \mapsto e)$  is a heap  $\Gamma$ , extended with the variable-term mapping  $x \mapsto e$ .

The  $\lambda_{\parallel}$  semantics keeps track of arguments via labeled terms  $e_{\parallel}^x$ , where labels are variables. Values in  $\lambda_{\parallel}$  also include labeled  $\lambda$ -abstractions. Reducing a labeled term triggers the simultaneous reduction of all other terms with the same label. Otherwise, labels do not affect program evaluation.

We require that all expressions with the same label must be identical.

**Definition 5.** A program  $e_{\parallel}$  is consistently labeled (CL) when for any two sub-terms  $e_{\parallel 1}^{x_1}$  and  $e_{\parallel 2}^{x_2}$  of  $e_{\parallel}$ ,  $x_1 = x_2$  implies  $e_{\parallel 1} = e_{\parallel 2}$ .

In the reduction of  $\lambda_{\parallel}$  programs, evaluation contexts  $E_{\parallel}$  determine which part of the program to reduce next. The  $\lambda_{\parallel}$  evaluation contexts are the call-by-name evaluation contexts with the addition of the labeled  $E_{\parallel}^x$  context, which dictates that a redex search goes under labeled terms. Essentially, when searching for a redex, terms tagged with a label are treated as if they were unlabeled.

The parallel semantics can exploit simpler evaluation contexts than  $\lambda_{need}$  because substitution occurs as soon as an application is encountered:

$$E_{\parallel}[\lambda x.e_{\parallel 1}^{\bar{y}}] e_{\parallel 2} \mapsto \begin{cases} E_{\parallel}[e_{\parallel}], & \text{if } [ ] \text{ is not under a label in } E_{\parallel} \\ E_{\parallel}[e_{\parallel}]\{\{z \leftarrow E_{\parallel 2}[e_{\parallel}]\}\}, & \text{if } E_{\parallel}[ ] = E_{\parallel 1}[(E_{\parallel 2}[ ])^z] \\ & \text{and } [ ] \text{ is not under a label in } E_{\parallel 2} \end{cases} \quad (\beta_{\parallel})$$

where  $e_{\parallel} = e_{\parallel 1}\{x := e_{\parallel 2}^w\}$ ,  $w$  fresh

On the left-hand side of  $\beta_{\parallel}$ , the program is partitioned into a context and a  $\beta$ -like redex. A term  $e^{\bar{y}}$  may have any number of labels and possibly none. On the right-hand side, the redex is contracted to a term  $e_{\parallel 1}\{x := e_{\parallel 2}^w\}$  such that the argument is tagged with an unique label  $w$ . Obsolete labels  $\bar{y}$  are discarded.

There are two distinct ways to contract a redex: when the redex is not under any labels and when the redex occurs under at least one label. For the former, the redex is the only contracted part of the program. For the latter, all other instances of that labeled term are similarly contracted. In this second case, the evaluation context is further subdivided as  $E_{\parallel}[ ] = E_{\parallel 1}[(E_{\parallel 2}[ ])^z]$ , where  $z$  is the label nearest the redex, i.e.,  $E_{\parallel 2}$  contains no additional labels. A whole-program substitution function is used to perform the parallel reduction:

$$\begin{aligned} e_{\parallel 1}^x\{\{x \leftarrow e_{\parallel}\}\} &= e_{\parallel}^x \\ e_{\parallel 1}^x\{\{y \leftarrow e_{\parallel}\}\} &= (e_{\parallel 1}\{\{y \leftarrow e_{\parallel}\}\})^x, \quad x \neq y \\ (\lambda x.e_{\parallel 1})\{\{x \leftarrow e_{\parallel}\}\} &= \lambda x.(e_{\parallel 1}\{\{x \leftarrow e_{\parallel}\}\}) \\ (e_{\parallel 1} e_{\parallel 2})\{\{x \leftarrow e_{\parallel}\}\} &= (e_{\parallel 1}\{\{x \leftarrow e_{\parallel}\}\} e_{\parallel 2}\{\{x \leftarrow e_{\parallel}\}\}) \\ \text{otherwise, } e_{\parallel 1}\{\{x \leftarrow e_{\parallel}\}\} &= e_{\parallel 1} \end{aligned}$$

Rewriting terms with  $\beta_{\parallel}$  preserves the consistent labeling property.

**Proposition 2.** If  $e_{\parallel}$  is CL and  $e_{\parallel} \mapsto e_{\parallel}'$ , then  $e_{\parallel}'$  is CL.

The  $\psi$  function reconstructs a  $\lambda_{\parallel}$  term from a CKH machine configuration:

$$\begin{aligned} \psi(\langle e, ((\mathbf{var} \ x), F_{LS}), \Gamma \rangle) &= \psi(\langle x, F_{LS}, (x \mapsto e), \Gamma \rangle) & \boxed{\psi : S_L \rightarrow e_{\parallel}} \\ \psi(\langle e_1, ((\mathbf{arg} \ e_2), F_{LS}), \Gamma \rangle) &= \psi(\langle e_1 \ e_2, F_{LS}, \Gamma \rangle) \\ \psi(\langle e, (), \Gamma \rangle) &= e\{\!\{ \Gamma \}\!\} \end{aligned}$$

The operation  $e\{\!\{ \Gamma \}\!\}$ , using overloaded notation, replaces all free variables in  $e$  with their corresponding terms in  $\Gamma$  and tags them with appropriate labels.

Lemma 5 demonstrates the bulk of the equivalence of the store machine and  $\lambda_{\parallel}$ .<sup>7</sup> The rest of the equivalence proof is straightforward [9].

**Lemma 5.** *If  $\langle e, F_{LS}, \Gamma \rangle \xrightarrow{ckh} \langle e', F_{LS}', \Gamma' \rangle$ , then either:*

1.  $\psi(\langle e, F_{LS}, \Gamma \rangle) = \psi(\langle e', F_{LS}', \Gamma' \rangle)$
2.  $\psi(\langle e, F_{LS}, \Gamma \rangle) \xrightarrow{\parallel} \psi(\langle e', F_{LS}', \Gamma' \rangle)$

#### 5.4 A Transition System for Comparing $\lambda_{need}$ and $\lambda_{\parallel}$

The CK layer in figure 4 mediates between  $\lambda_{\parallel}$  and  $\lambda_{need}$ . The corresponding transition system resembles a two-register CK machine [9]. Figure 6 describes the syntax and the transitions of the system.<sup>8</sup>

##### Syntax

$$\begin{aligned} S &= \langle e, Fs \rangle && \text{(States)} \\ Fs &= F, \dots && \text{(List of Frames)} \\ F &= (\mathbf{arg} \ e) \mid (\mathbf{lam} \ x) \mid (\mathbf{bod} \ x \ Fs \ Fs) && \text{(Frames)} \end{aligned}$$

##### Transitions

$$\begin{aligned} \langle e_1 \ e_2, Fs \rangle &\xrightarrow{ck} \langle e_1, ((\mathbf{arg} \ e_2), Fs) \rangle && \text{(push-arg-ck)} \\ \langle \lambda x.e, Fs \rangle &\xrightarrow{ck} \langle e, ((\mathbf{lam} \ x), Fs) \rangle && \text{(descend-lam-ck)} \\ &\text{if } \mathbf{balance}(Fs) > 0 \\ \langle x, (Fs_1, (\mathbf{lam} \ x), Fs_2, (\mathbf{arg} \ e), Fs) \rangle &\xrightarrow{ck} \langle e, ((\mathbf{bod} \ x \ Fs_1 \ Fs_2), Fs) \rangle && \text{(lookup-var-ck)} \\ &\text{if } \phi_F(Fs_1) \in \hat{A}[E], \phi_F(Fs_2) \in A, \phi_F(Fs) \in E[\hat{A}], \hat{A}[\hat{A}] \in A \\ \langle v, (Fs_3, (\mathbf{bod} \ x \ Fs_1 \ Fs_2), Fs) \rangle &\xrightarrow{ck} \langle v, (Fs_1 \{x := v\}, Fs_3, Fs_2, Fs) \rangle && \text{(\beta}_{need}\text{-ck)} \\ &\text{if } \phi_F(Fs_3) \in A \end{aligned}$$

**Fig. 6.** A transition system for comparing  $\lambda_{need}$  and  $\lambda_{\parallel}$

States consist of a subterm and a list of frames representing the context. The first kind of frame represents the argument in an application and the second

<sup>7</sup> The lemma relies on an extension of the typical  $\alpha$ -equivalence classes of terms to include variables in labels as well.

<sup>8</sup> The CK transition system is a proof-technical device. Unlike the original CK machine, ours is ill-suited for an implementation.



frame represents a  $\lambda$ -abstraction with a hole in the body. The last kind of frame has two frame list components, the first representing a context in the body of the  $\lambda$ , and the second representing the context between the  $\lambda$  and its argument. The variable in this last frame is the variable bound by the  $\lambda$  expression under evaluation. The initial state for a program  $e$  is  $\langle e, () \rangle$ , where  $()$  is an empty list of frames, and evaluation terminates when the control string is a value and the list of frames is equivalent to an answer context.

The *push-arg-ck* transition makes the operator in an application the new control string and adds a new **arg** frame to the frame list containing the argument. The *descend-lam-ck* transition goes under a  $\lambda$ , making the body the control string, but only if that  $\lambda$  has a corresponding argument in the frame list, as determined by the **balance** function, defined as follows:

$$\boxed{\text{balance} : Fs \rightarrow \mathbb{Z}}$$

$$\begin{aligned} \text{balance}(Fs_3, (\text{bod } x \text{ } Fs_1 \text{ } Fs_2), Fs) &= \text{balance}(Fs_3) \\ \text{balance}(Fs) &= \#arg\text{-frames}(Fs) - \#lam\text{-frames}(Fs) \\ &\quad Fs \text{ contains no bod frames} \end{aligned}$$

The **balance** side condition for *descend-lam-ck* dictates that evaluation goes under a  $\lambda$  only if there is a matching argument for it, thus complying with the analogous evaluation context. The **balance** function employs **#arg-frames** and **#lam-frames** to count the number of **arg** or **lam** frames, respectively, in a list of frames. Their definitions are elementary and therefore omitted.

The *lookup-var-ck* transition is invoked if the control string is a variable, somewhere in a  $\lambda$  body, and the rest of the frames have a certain shape consistent with the corresponding parts of a  $\beta_{need}$  redex. With this transition, the argument associated with the variable becomes the next control string and the context around the variable in the  $\lambda$  body and the context between the  $\lambda$  and argument are saved in a new **bod** frame. Finally, when an argument is an answer, indicated by a value control string and a **bod** frame in the frame list—with the equivalent of an answer context in between—the value gets substituted into the body of the  $\lambda$  according to the  $\beta_{need-ck}$  transition. The  $\beta_{need-ck}$  transition uses a substitution function on frame lists,  $Fs\{x := e\}$ , which overloads the notation for regular term substitution and has the expected definition.

Figure 7 defines metafunctions for the CK transition system. The  $\phi$  function converts a CK state to the equivalent  $\lambda_{need}$  term, and uses  $\phi_F$  to convert a list of frames to an evaluation context.

Now we can show that an evaluator defined with  $\mapsto^{ck}$  is equivalent to  $\text{eval}_{need}^{sr}$ . The essence of the proof is a lemma that relates the shape of CK transition sequences to the shape of  $\lambda_{need}$  standard reduction sequences. The rest of the equivalence proof is straightforward [9].

**Lemma 6.** *If  $\langle e, Fs \rangle \mapsto^{ck} \langle e', Fs' \rangle$ , then either:*

1.  $\phi(\langle e, Fs \rangle) = \phi(\langle e', Fs' \rangle)$
2.  $\phi(\langle e, Fs \rangle) \mapsto \phi(\langle e', Fs' \rangle)$

<div style="border: 1px solid black; display: inline-block; padding: 2px; margin-bottom: 10px;"><math>\phi : S \rightarrow e</math></div> $\phi(\langle e, Fs \rangle) = \phi_F(Fs)[e]$ <div style="border: 1px solid black; display: inline-block; padding: 2px; margin-bottom: 10px;"><math>\phi_F : Fs \rightarrow E</math></div> $\begin{aligned} \phi_F(()) &= [ ] \\ \phi_F(\mathbf{lam} \ x, Fs) &= \phi_F(Fs)[\lambda x. [ ] ] \\ \phi_F(\mathbf{arg} \ e, Fs) &= \phi_F(Fs)[ [ ] e ] \\ \phi_F(\mathbf{bod} \ x \ Fs_1 \ Fs_2, Fs) &= \\ \phi_F(Fs)[\phi_F(Fs_2)[\lambda x. \phi_F(Fs_1)[x]] [ ] ] \end{aligned}$	<div style="border: 1px solid black; display: inline-block; padding: 2px; margin-bottom: 10px;"><math>\xi : S \rightarrow e_{\parallel}</math></div> $\xi(\langle e, Fs \rangle) = \xi_F(Fs, e)$ <div style="border: 1px solid black; display: inline-block; padding: 2px; margin-bottom: 10px;"><math>\xi_F : Fs \times e_{\parallel} \rightarrow e_{\parallel}</math></div> $\begin{aligned} \xi_F((), e_{\parallel}) &= e_{\parallel} \\ \xi_F(\langle \mathbf{arg} \ e_{\parallel 1}, Fs \rangle, e_{\parallel}) &= \xi_F(Fs, e_{\parallel} \ e_{\parallel 1}) \\ \xi_F(\langle \mathbf{bod} \ x \ Fs_1 \ Fs_2 \rangle, Fs, e_{\parallel}) &= \\ \xi_F(Fs_1, (\mathbf{lam} \ x), Fs_2, \langle \mathbf{arg} \ e_{\parallel} \rangle, Fs, x) &= \\ \xi_F(\langle \mathbf{lam} \ x \rangle, Fs_1, \langle \mathbf{arg} \ e_{\parallel 1} \rangle, Fs_2, e_{\parallel}) &= \\ \xi_F(Fs_1, Fs_2, e_{\parallel} \{x := e_{\parallel}^y\}) &= \\ \phi_F(Fs_1) \in A, \ y \ \text{fresh} \end{aligned}$
--	--

**Fig. 7.** Functions to map CK states to  $\lambda_{need}$  ( $\phi$ ) and  $\lambda_{\parallel}$  ( $\xi$ )

Finally, we show how the CK system corresponds to  $\lambda_{\parallel}$ . The  $\xi$  function defined in figure 7 constructs a  $\lambda_{\parallel}$  term from a CK configuration.

**Lemma 7.** *If  $\langle e, Fs \rangle \xrightarrow{ck} \langle e', Fs' \rangle$ , then either:*

1.  $\xi(\langle e, Fs \rangle) = \xi(\langle e', Fs' \rangle)$
2.  $\xi(\langle e, Fs \rangle) \xrightarrow{\parallel} \xi(\langle e', Fs' \rangle)$

## 5.5 Relating All Layers

In the previous subsections, we have demonstrated the correspondence between  $\lambda_{\parallel}$ , the natural semantics, and the  $\lambda_{need}$  standard reduction sequences via lemmas 5 through 7. We conclude this section with the statement of an extensional correctness theorem, where  $eval_{natural}$  is an evaluator defined with the store machine transitions. The theorem follows from the composition of the equivalences of our specified rewriting systems.

**Theorem 4.**  $eval_{need} = eval_{natural}$

## 6 Extensions and Variants

**Data Constructors** Real-world lazy languages come with data structure construction and extraction operators. Like function arguments, the arguments to a data constructor should not be evaluated until there is demand for their values [12, 14]. The standard  $\lambda$  calculus encoding of such operators [5] works well:

$$\mathbf{cons} = \lambda x. \lambda y. \lambda s. s \ x \ y, \quad \mathbf{car} = \lambda p. p \ \lambda x. \lambda y. x, \quad \mathbf{cdr} = \lambda p. p \ \lambda x. \lambda y. y$$

Adding true algebraic systems should also be straightforward.

**Recursion** Our  $\lambda_{need}$  calculus represents just a core  $\lambda$  calculus and does not include an explicit `letrec` constructor for cyclic terms. Since cyclic programming is an important idiom in lazy programming languages, others have extensively explored cyclic by-need calculi, e.g., Ariola and Blum [1], and applying their solutions to our calculus should pose no problems.

## 7 Conclusion

Following Plotkin’s work on call-by-name and call-by-value, we present a call-by-need  $\lambda$  calculus that expresses computation via a single axiom in the spirit of  $\beta$ . Our calculus is close to implementations of lazy languages because it captures the idea of by-need computation without retaining every function call and without need for re-associating terms. We show that our calculus satisfies Plotkin’s criteria, including an intensional correspondence between our calculus and a Launchbury-style natural semantics. Our future work will leverage our  $\lambda_{need}$  calculus to derive a new abstract machine for lazy languages.

**Acknowledgments.** We thank J. Ian Johnson, Casey Klein, Vincent St-Amour, Asumu Takikawa, Aaron Turon, Mitchell Wand, and the ESOP 2012 reviewers for their feedback on early drafts. This work was supported in part by NSF Infrastructure grant CNS-0855140 and AFOSR grant FA9550-09-1-0110.

## References

1. Ariola, Z., Blom, S.: Cyclic Lambda Calculi. In: Ito, T., Abadi, M. (eds.) TACS 1997. LNCS, vol. 1281, pp. 77–106. Springer, Heidelberg (1997)
2. Ariola, Z.M., Felleisen, M.: The call-by-need lambda-calculus. Tech. Rep. CIS-TR-94-23, University of Oregon (1994)
3. Ariola, Z.M., Felleisen, M.: The call-by-need lambda calculus. *J. Funct. Program.* 7, 265–301 (1997)
4. Ariola, Z.M., Maraist, J., Odersky, M., Felleisen, M., Wadler, P.: A call-by-need lambda calculus. In: Proc. 22nd Symp. on Principles of Programming Languages, pp. 233–246 (1995)
5. Barendregt, H.P.: *Lambda Calculus, Syntax and Semantics*. North-Holland (1985)
6. Church, A.: *The Calculi of Lambda Conversion*. Princeton University Press (1941)
7. Curry, H.B., Feys, R.: *Combinatory Logic, vol. I*. North-Holland (1958)
8. Danvy, O., Millikin, K., Munk, J., Zerny, I.: Defunctionalized Interpreters for Call-by-Need Evaluation. In: Blume, M., Kobayashi, N., Vidal, G. (eds.) FLOPS 2010. LNCS, vol. 6009, pp. 240–256. Springer, Heidelberg (2010)
9. Felleisen, M., Findler, R.B., Flatt, M.: *Semantics Engineering with PLT Redex*. MIT Press (2009)
10. Felleisen, M., Friedman, D.P.: A syntactic theory of sequential state. *Theor. Comput. Sci.* 69(3), 243–287 (1989)
11. Friedman, D.P., Ghuloum, A., Siek, J.G., Winebarger, O.L.: Improving the lazy Krivine machine. *Higher Order Symbolic Computation* 20, 271–293 (2007)

12. Friedman, D.P., Wise, D.S.: Cons should not evaluate its arguments. In: Proc. 3rd Intl. Colloq. on Automata, Languages and Programming. pp. 256–284 (1976)
13. Garcia, R., Lumsdaine, A., Sabry, A.: Lazy evaluation and delimited control. In: Proc. 36th Symp. on Principles of Programming Languages, pp. 153–164 (2009)
14. Henderson, P., Morris Jr., J.H.: A lazy evaluator. In: Proc. 3rd Symp. on Principles of Programming Languages, pp. 95–103 (1976)
15. Josephs, M.B.: The semantics of lazy functional languages. *Theor. Comput. Sci.* 68(1) (1989)
16. Landin, P.J.: The next 700 programming languages. *Comm. ACM* 9, 157–166 (1966)
17. Launchbury, J.: A natural semantics for lazy evaluation. In: Proc. 20th Symp. on Principles of Programming Languages, pp. 144–154 (1993)
18. Loader, R.: Notes on simply typed lambda calculus. Tech. Rep. ECS-LFCS-98-381, Department of Computer Science, University of Edinburgh (1998)
19. Maraist, J., Odersky, M., Turner, D.N., Wadler, P.: Call-by-name call-by-value, call-by-need, and the linear lambda calculus. In: Proc. 11th Conference on Mathematical Foundations of Programming Semantics, pp. 370–392 (1995)
20. Maraist, J., Odersky, M., Wadler, P.: The call-by-need lambda calculus (unabridged). Tech. Rep. 28/94, Universität Karlsruhe (1994)
21. Maraist, J., Odersky, M., Wadler, P.: The call-by-need lambda calculus. *J. Funct. Program.* 8, 275–317 (1998)
22. Morris, J.H.: Lambda Calculus Models of Programming Languages. Ph.D. thesis, MIT (1968)
23. Nakata, K., Hasegawa, M.: Small-step and big-step semantics for call-by-need. *J. Funct. Program.* 19(6), 699–722 (2009)
24. Peyton Jones, S.L., Salkild, J.: The spineless tagless g-machine. In: Proc. 4th Conf. on Functional Programming Lang. and Computer Architecture, pp. 184–201 (1989)
25. Plotkin, G.D.: Call-by-name, call-by-value and the lambda-calculus. *Theor. Comput. Sci.* 1, 125–159 (1975)
26. Purushothaman, S., Seaman, J.: An adequate operational semantics for sharing in lazy evaluation. In: Proc. 4th European Symp. on Program, pp. 435–450 (1992)
27. Sestoft, P.: Deriving a lazy abstract machine. *J. Func. Program.* 7(3), 231–264 (1997)
28. Wadsworth, C.P.: Semantics and Pragmatics of the Lambda Calculus. Ph.D. thesis, Oxford University (1971)