# LABORATORY FOR
# COMPUTER SCIENCE

MASSACHUSETTS
INSTITUTE OF
TECHNOLOGY

## Semantics of Barriers in a Non-Strict, Implicitly-Parallel Language

**Shail Aditya**
**Arvind**
**MIT Laboratory for Computer Science**
`{shail,arvind}@lcs.mit.edu`

**Joseph E. Stoy**
**Oxford University Computing Laboratory**
`joe.stoy@comlab.ox.ac.uk`

545 TECHNOLOGY SQUARE, CAMBRIDGE, MASSACHUSETTS 02139

# Semantics of Barriers in a Non-Strict, Implicitly-Parallel Language

Shail Aditya        Arvind
Laboratory for Computer Science
Massachusetts Institute of Technology
{shail,arvind}@lcs.mit.edu

Joseph E. Stoy
Computing Laboratory
Oxford University
joe.stoy@comlab.ox.ac.uk

## Abstract

Barriers in parallel languages may be used to schedule parallel activities, control memory usage and ensure proper sequentialization of side-effects. In this paper we present operational semantics of barriers in Id and pH, which are non-strict, implicitly-parallel, functional languages extended with side-effects. The semantics are presented as a translation from a source language with barriers into a kernel language without barriers where the termination properties of an expression are made explicit in the form of signals. These signals are generated using a monotonic, primitive operator (W) that detects weak head-normal forms of expressions and are used to control the execution of expressions *via* a strict application operator (Sap). We present two versions of the semantics — the first uses purely data-driven, eager evaluation and the second mixes eager evaluation with a demand-driven identifier dereference mechanism. We compare and contrast the two for their ability to do resource management and preserve useful semantic properties.

## 1  Introduction

Purely functional languages have the advantage of being non-procedural yet determinate: the result of a functional program is always the same regardless of the order in which its sub-expressions are evaluated. The introduction of side-effects in a language usually forces a strict, sequential order of evaluation in order to guarantee determinacy (*e.g.*, Scheme and ML). However, it is possible to specify only a partial order on side-effects and still retain an overall consistent picture of the computation. This paper studies the use and the semantics of barriers as a partial sequentialization construct in Id [12] and pH [13], which are non-strict, implicitly-parallel, functional languages extended with side-effects.

The traditional view of barriers comes from data-parallel languages such as *Lisp and CM-Fortran, where concurrently executing copies of the same program periodically synchronize using a global barrier mechanism. Some parallel machines like the CM5 and Cray T3D even provide such global barrier support in hardware. However, in this paper we study the semantics of barriers under a more general,

multi-threaded, MIMD programming model where barriers may be localized to specified subsets of parallel tasks. This allows the specification of a partial order among the various tasks while providing fine-grain control over inter-task synchronization, resource usage and scheduling of side-effects.

Id and pH support imperative operations on I-structure [4] and M-structure [6] objects. I-structures allow the creation of a data structure to be separated from the definition of its components: attempts to use the value of a component are automatically delayed until that component is defined; attempts to redefine a component lead to an inconsistent state. M-structures, on the other hand, are fully mutable data structures whose components can be redefined repeatedly: a M-take operation reads and empties a full component; a M-put operation (re)defines it; two successive M-put operations on the same component lead to an inconsistent state unless separated by a M-take operation.

Barriers were introduced in Id to provide partial sequencing of M-structure operations. However, their semantics were not given precisely at that stage [5]. Subsequently, pH was designed to follow the syntax and the type structure of Haskell [8] combined with the eager evaluation model and the non-functional extensions of Id. This prompted a formalization of barrier semantics for both Id and pH. The first such attempt was made in [15], both in terms of CSP [7] and (for the functional layer) denotationally. This paper presents a different formalization in terms of graph rewriting systems previously used to describe the semantics of Id [3].

The rest of the paper is organized as follows. In Section 2 we discuss some important characteristics of non-strict, eager computations in Id and pH and introduce the concept of control regions and barriers. Section 3 presents several examples of using barriers in controlling parallelism, managing program resources, and providing fine-grain control over scheduling of side-effects. In Section 4 we establish the formal framework for defining the semantics of barriers by introducing a simple kernel language with parallel rewriting semantics. In Section 5 we describe the concept of termination and propose semantics for barriers based on a "data-driven" eager evaluation model. This version is sufficient for both resource management and side-effect control, although, as shown in Section 6 it turns out to be a little fragile with respect to some program transformations. In Section 7 we present another formulation of the semantics for barriers which attempts to overcome these shortcomings by mixing eager evaluation with "demand-driven" evaluation. Finally, Section 8 presents the summary and directions for future

work. No knowledge of dataflow is required to read this paper.

## 2 Non-Strict, Eager Evaluation

Operationally, a strict computation of a function cannot return a result until the values of all its arguments are available; an expression needs the results of all its sub-expressions; and a data structure is not available until all its components are initialized. Non-strict computations, on the other hand, relax these constraints (except when implied by data-dependencies): a function may return a partial result before all its arguments are available, or parts of a data structure may be read before it is fully computed. Id and pH languages follow an *eager* evaluation strategy for non-strict computations: all tasks execute in parallel, restricted only by the data dependencies among them. This strategy automatically exposes large amounts of parallelism both within and across procedures. This is in contrast with a *lazy* evaluation strategy followed by the Haskell language: only those tasks are evaluated which are required to produce the result. This strategy imposes a strong sequential constraint on the overall computation, although the exact ordering of tasks is decided dynamically.

We discuss the eager evaluation model of Id and pH in more detail below.

### 2.1 The Fully Parallel Execution Model

A program in Id consists of an expression to be evaluated within the scope of a set of top-level function and type declarations. Each function is broken up into several threads of computation (the length of the threads is determined in part by the compiler's ability to identify strict regions and in part by the ability of the run-time system and the hardware to exploit the exposed parallelism efficiently).

The parallel execution model of Id and pH is shown pictorially in Figure 1. Each function application executes within the context of an *activation frame* which records function arguments and keeps temporary, local values. The program starts by allocating a root activation frame and initiating the main thread in order to evaluate the top-level expression. Function applications within a thread give rise to parallel child activations, while loop invocations give rise to multiple parallel iterations. Threads belonging to a function share its activation frame and may be active concurrently. A thread may enable other threads by sending data or synchronization information to their associated activation frame: this characterizes the "data-driven" nature of this execution model. Thus at any time the overall computation is represented by a tree of activation frames, exploiting both intra-procedural and inter-procedural parallelism.

In contrast, under lazy evaluation parallel computation is spawned only if it is already known to contribute towards the final result. Otherwise, every potentially concurrent task is suspended in a *thunk* immediately upon creation. A "demand-driven" evaluation of the suspended thunks exposes only a small part of the parallel execution tree at any given time. Of course, some of the thunks may never get evaluated under demand-driven evaluation.

As shown in Figure 1, all threads participating in the parallel computation share a globally addressable heap. An activation frame is deallocated when its associated function or the loop terminates, but the data structures allocated on
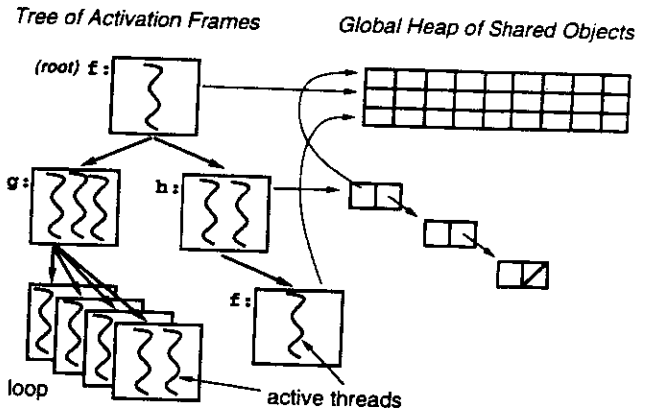


Figure 1: The Fully Parallel Execution Model of Id and pH.

the heap may continue to exist even after the function that allocated them has terminated. Such data structures either have to be explicitly deallocated or are garbage-collected when no more references to them remain.

### 2.2 Result vs. Termination

An important aspect of non-strict, eager evaluation is that the production of the result of a computation and its termination may not always coincide. For example, the function f shown below is allowed to execute and return the result 3 before the internal application g x terminates or even before its argument x becomes a value because the result does not depend on either of these events[1].

**Example 1:**
```
def f x = { y = g x; in 3 };
```

Under lazy evaluation a computation proceeds only as far as necessary to produce the requested result, so we never have to worry about the termination of the computation aside from the production of its result. On the other hand, under eager evaluation several activities may be spawned eagerly that do not immediately (or directly) contribute towards the result. It is therefore semantically important to separate the production of the result of a computation from its overall termination. From now on we will use the term *Value Semantics* to refer to the result of a computation, and *Termination Semantics* to identify its termination[2]. Semantically, a given computation may give rise to one of the following possible outcomes [1]:

1. **Result with proper Termination** — The computation produces a result and terminates.

2. **Result with improper Termination** — The computation produces a result but does not terminate. The non-termination could be either due to an infinite computation or due to a deadlocked computation.

---

[1] All our examples use the Id syntax [12]. In Id, *def* introduces function definitions, colon (:) is the infix list constructor, and block expressions (enclosed in { }) evaluate the expression following *in* within the scope of a set of parallel bindings.

[2] These correspond to what were called non-strict and strict semantics, respectively, in [15].

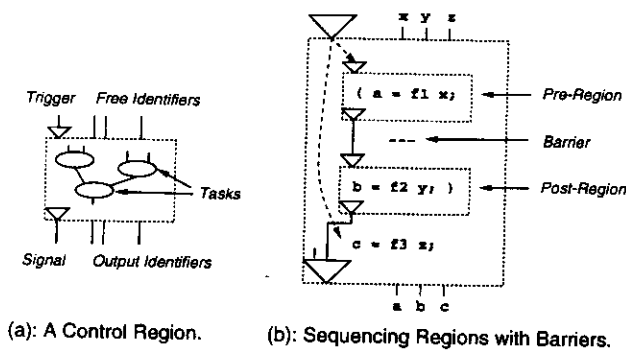(a): A Control Region.    (b): Sequencing Regions with Barriers.

Figure 2: Control Regions and Barriers.

3. **No Result or Termination** — The computation does not produce a result and does not terminate. Again, the non-termination could be due to infinite computation or deadlock[3].

For example, the function f1 below produces an infinite list of x's as the result and terminates properly, f2 produces the same result but does not terminate because of infinite computation, while f3 does not produce a result and deadlocks because the + operator is strict in its arguments.

**Example 2:**

```
def f1 x = { a = x:a in a };

def f2 x = x:(f2 x);

def f3 x = { a = a+x in a };
```

## 2.3 Control Regions

Another important aspect of eager evaluation is the notion of *control regions* (refer Figure 2 (a)). A control region is informally defined as a set of concurrent tasks that are under the same control dependence and therefore always execute together. Tasks within the same control region may execute in any order or in an interleaved manner as long as the data dependencies among them are respected. As Figure 2 (a) shows, a control region requires the values of its free identifiers as input and produces one or more result values as output. The combined termination of all tasks within a control region may be used as a synchronization event to initiate other control regions; a control region therefore also takes a trigger input and produces a signal output on termination.

As an example, the body of a function constitutes a control region which is triggered when that function is applied to all its arguments. The termination of the function body is signaled by the joint termination of each of its sub-computations which may be used to deallocate the activation frame associated with that function application. Control regions are always properly nested within one another and therefore form a "tree" of parallel activities. For example, the activation tree in Figure 1 shows the dynamic

nesting of function control regions. Similarly, each of the branches of a conditional expression forms a control region that is statically nested inside the enclosing region. In this case only one of the branches is initiated when the predicate is resolved, and the termination of the entire expression is determined by the termination of that branch.

## 2.4 Barrier Specification

*Barriers* provide a mechanism to detect the termination of a set of parallel activities enclosed within a control region. A barrier is specified in Id using three or more dashes (---) creating two sub-regions within a given control region — one above the barrier called the *pre-region* and the other below the barrier called the *post-region* (see Figure 2 (b)). Intuitively, no computation within the post-region is allowed to proceed until all computations within the pre-region have terminated[4]. The termination of control regions must also respect their nesting, making the barriers *hyper-strict*, i.e., the entire computation subtree spawned from the tasks within a pre-region must terminate before any computation in the post-region is allowed to proceed.

Id also provides the ability to delimit the scope of a control region to an arbitrary subset of the parallel tasks at hand. Unlike data-parallel computation, where all barriers are global, barriers in Id may be localized to a specified set of parallel activities by grouping those activities in parentheses. This is also illustrated in the example shown in Figure 2 (b). The activities (f1 x), and (f3 z) are initiated in parallel because the post-region of the barrier excludes the activity (f3 z) using parentheses. The activity (f2 y) is initiated upon termination of (f1 x) as prescribed by the barrier. The entire block terminates when all the activities within all sub-regions inside it have terminated. The power and usefulness of barriers in Id stem from the fact that the granularity of the regions being sequentialized is completely under user control.

Figure 2 illustrates another important point: the partitioning of parallel computation in Id into control regions is completely static and textual. The user always has a clear picture of which computations will evaluate and under what circumstances. This is possible only under eager evaluation. Under lazy evaluation, even if some computations are explicitly sequentialized or scheduled using barriers, it is not clear whether they will ever evaluate until something demands their value: there is no natural notion of control region as a group of concurrent activities.

## 3 The Uses of Barriers

In this section we describe the various uses of barriers under the parallel execution model discussed in Section 2.

### 3.1 Controlling Parallelism

As shown in Figure 1, the eager evaluation strategy unfolds the parallel computation into a tree of activation frames. In some "embarrassingly parallel" computations the potential parallelism is limited not by the data dependencies but purely by the size of the available frame or heap memory in the machine. An unrestricted unfolding of the activation tree may cause the machine to run out of memory very

---

[3]The fourth possibility, of proper termination but without producing a result, is not considered to be meaningful and is lumped with the third outcome. Some graph-reduction based semantic definitions [3] do, however, distinguish between these two cases.

[4]Actually, all we care about is that the *effect* (computed values and side-effects) of the post-region is not visible from outside until the pre-region has completely terminated.

quickly. In such cases it is important to "throttle" the parallelism to match the available resources. Barriers allow the programmer (or the compiler) to control the parallelism by sequentializing some parallel activities. For instance, without the barrier in the following example the `fib` function would unfold into an exponentially large tree of activation frames. The barrier reduces the resource requirement to be linear in the size of the problem, albeit also reducing the parallelism.

**Example 3:**
```
def fib n = if n < 2 then n
            else { x = fib (n-1);
                   ---
                   y = fib (n-2);
                in x+y };
```

## 3.2 Controlling Frame Storage

Sometimes a group of activities that are initiated simultaneously may not exhibit a lot of parallelism because of pre-existing data dependencies among them. But under eager evaluation all the initiated activities can grab memory resources and hold on to them for a long time without doing any useful work. Barriers help in scheduling the activities in the appropriate partial order so that the overall resource requirement is reduced. This behavior is illustrated by the following example.

**Example 4:**
```
def repeat1 f n x = if n==0 then x
                    else repeat1 f (n-1) (f x);

def repeat2 f n x = if n==0 then x
                    else { y = f x;
                           ---
                        in repeat2 f (n-1) y };
```

The `repeat1` function iterates the function `f` over its argument `x` $n$ times. Under eager evaluation it rapidly expands into a sequence of frames to compute the application `(f x)` at each iteration, using up to $O(n)$ frames. If the function `f` is known to be strict, then all except the first of these frames would wait for the previous iteration to produce a result before they can compute anything useful. In this case we may use a barrier as shown in function `repeat2` to wait for the current iteration to finish completely before moving on to the next iteration, using at most two frames during the entire computation. Note that, unlike the case of the `fib` function in Example 3, no parallelism is lost by inserting this barrier since the computation is already sequential.

## 3.3 Controlling Heap Storage

Barriers are also useful for scheduling the automatic (or programmer-specified) deallocation of heap resources once they are no longer in use. This strategy may be used to manage the heap without paying any extra overhead for garbage collection, as shown in the following example. Local data structures allocated and used during a large computation may be released at its end.

**Example 5:**
```
def compute graph =
{ ( note = make_notebook ();
    result = traverse graph note;
    ---
```

```
    deallocate note; )
in result };
```

## 3.4 Sequentializing Side-Effects

Barriers may be used in parallel M-structure computations for scheduling side-effects in the desired partial order. Under eager evaluation all side-effects present within a block are initiated in parallel. A barrier may therefore be required to ensure proper interleaving of M-take (read-and-lock) and M-put (write-and-unlock) operations. For instance, a barrier is necessary in the following example to ensure that the mutable array slot is taken before the new value is put back, because there is no explicit data dependency to control the order otherwise[5].

**Example 6:**
```
def replace a i v =
  { x = a![i];
    ---
    a![i] = v;
  in x };
```

Barriers provide the ability to detect the combined termination of a group of parallel update operations without sequentializing them completely. A typical example of this sort is the construction of a histogram from a tree of elements.

**Example 7:**
```
def histogram t n =
  { a = { def zero i = 0;
              in make_m_array (1,n) zero };
    _ = traverse a t;
    ---
    in a };

def traverse a (Leaf x) =
  { i = hash x (bounds a);
    a![i] = a![i]+1; }
|   traverse a (Node l r) =
  { _ traverse a l;
    _ = traverse a r; };
```

The `traverse` function is allowed to unfold automatically, accumulating the frequencies in parallel into the histogram array `a` using M-structure operations. No sequentialization is imposed over the accumulations, but the barrier in the function `histogram` guarantees that all the accumulations have terminated before the histogram array is returned. Note that the barrier ensures that all the computations above the barrier have terminated before returning the array, although actually we are interested only in the termination of all side-effects on the array itself.

## 3.5 Controlling Access to Shared External State

Barriers are also useful in managing activities that share a common external resource or state. Shared peripheral devices, I/O streams, and protected kernel data structures may constitute some of the external state that user programs may need to control or gain access to. Most functional languages allow only implicit control of such computation *via*

---

[5]In Id, *M-take* and *M-put* operations are represented using a slightly modified array indexing syntax *a![i]*. Indexed expressions on the right-hand side of a binding are *take* operations while those on the left-hand side are *put* operations.

4

Figure 3:

| | | |
|---|---|---|
| $c$ | $\in$ | Constant |
| $F, t, x, y, z \ldots$ | $\in$ | Identifier |
| $SE, X, Y, Z, \ldots$ | $\in$ | Simple Expression |
| $E$ | $\in$ | Expression |
| $S$ | $\in$ | Statement |
| $PF^n$ | $\in$ | Primitive Fn. with $n$ arguments |
| Constant | ::= | *Integer* | *Float* | *Boolean* | $\odot$ | $\bullet$ |
| $SE$ | ::= | Identifier | Constant |
| $PF$ | ::= | $+$ | $-$ | $\cdots$ | hd | tl | Cons | $\cdots$ |
| | | | Alloc | I-fetch | M-take | W | & |
| $E$ | ::= | $SE$ | $PF^n(X_1, \ldots, X_n)$ |
| | | | Cond $(X, E_1, E_2)$ | $\lambda x_1 \cdots x_n. E$ |
| | | | $Ap(F, X)$ | $Fap^n(F^n, X_1, \ldots, X_n)$ |
| | | | $Sap^n(F^n, X_1, \ldots, X_n)$ |
| | | | $Pap^i(F^n, \underline{n-i}, X_1, \ldots, X_i)$ | Block |
| Block | ::= | $\{ S^{par} \text{ in } (X_1, \ldots, X_n) \}$ |
| $S^{par}$ | ::= | $S_1; \ldots; S_n$ | $\epsilon$ |
| $S^{seq}$ | ::= | $(S_1^{par} \text{ --- } S_2^{par})$ |
| $S$ | ::= | Binding | Command | $S^{seq}$ |
| Binding | ::= | $x_1, \ldots, x_n = E$ |
| Command | ::= | $I\text{-store}(X, I, Z)$ | $M\text{-put}(X, I, Z)$ |

Figure 3: The Syntax of Kernel Expression Language.

| | | |
|---|---|---|
| $WHNF$ | $\in$ | Weak Head-Normal Form |
| Signal | $\in$ | $\{\perp, \bullet\}$ |
| $\mathcal{W}$ | :: | Expression $\rightarrow$ Signal |
| $WHNF$ | ::= | Constant | $Cons(X, Y)$ | $Alloc(\underline{n})$ |
| | | | $\lambda x_1 \cdots x_n. E$ |
| | | | $Pap(F^n, \underline{n-i}, X_1, \ldots, X_i)$ |
| $\mathcal{W}(c)$ | = | $\bullet$ |
| $\mathcal{W}(x)$ | = | $\mathcal{W}(E)$ Where $x$ is bound to $E$ |
| $\mathcal{W}(Cons(X, Y))$ | = | $\bullet$ |
| $\mathcal{W}(Alloc(E))$ | = | $\mathcal{W}(E)$ |
| $\mathcal{W}(\lambda x_1 \cdots x_n. E)$ | = | $\bullet$ |
| $\mathcal{W}(Pap(F^n, \underline{n-i}, X_1, \ldots, X_i))$ | = | $\bullet$ |
| $\mathcal{W}(E)$ | = | $\perp$ Otherwise |

Figure 4: Definition of Weak Head-Normal Forms (Values).

monads, continuations, abstract datatypes, or state transformers [9, 11, 14]. In Id, on the other hand, such objects are represented explicitly using M-structures, and access to them is controlled using locks and barriers that provide mutual exclusion and explicit sequentialization as necessary. For instance, currently there exist two I/O libraries in Id, sequential I/O [16] and threaded I/O [12], both of which crucially rely on barriers for proper sequentialization.

## 4 A Kernel Language

In this section we provide a framework for discussing the semantics of barriers. We present a small parallel language and its rewrite rule semantics.

### 4.1 Language Syntax

Figure 3 shows the grammar of our kernel expression language. The grammar ensures that every intermediate result of a complex computation is explicitly named using an identifier. This is extremely desirable in order to be able to express and preserve sharing constraints among computations.

Aside from the usual arithmetic primitives, the strict primitives hd and tl are used to manipulate lists. Cons is a non-strict constructor: it returns the allocated Cons-cell before it is completely filled. The primitive function Alloc is used to allocate either an I-structure or an M-structure memory block of specified size. The type-checker ensures that this memory is used consistently. The W operator is discussed in the next section.

There are several forms of application expressions in Figure 3. Ap is the general apply operation, $Fap^n$ is a full application of a known function to all its arguments, $Sap^n$ denotes strict application, and $Pap^i$ denotes a partial application of a function to $i$ arguments accumulated within the closure[6].

---

[6]We will omit the arity superscripts when they are obvious.

The latter three forms are generated during compiler optimizations or as intermediate results during rewriting. The language also provides nested $\lambda$-expressions that may contain free identifiers, and a block construct that controls lexical scoping and enables precise sharing of sub-expression values.

The kernel language allows several identifiers to be bound at once in a multiple-value binding. The number of multiple-values being returned from the right-hand side expression must match the number of identifiers being bound on the left-hand side. A group of statements separated by semicolons forms a parallel control region: the order of bindings in a parallel statement group is not significant. Finally, two parallel control regions may be sequentialized using a barrier.

### 4.2 Weak Head-Normal Forms

Typically a non-strict program computes a *weak head-normal form* as its answer. An expression is said to be in normal form if it cannot be simplified any further, while an expression is said to be in weak head-normal form as soon as its top-level structure becomes "stable". Under eager evaluation, an expression has to be in weak head-normal form before it can be passed as an argument, stored into data structures, or returned as a result from a function. Thus weak head-normal forms provide an exact characterization of the run-time *values* under eager evaluation.

The weak head-normal forms or values recognized in our language are shown in Figure 4. Constants like integers, floats, booleans and underlined expressions (e.g., $\underline{n-i}$) constitute literal values. An allocated data structure such as a Cons-cell or an Alloc memory block is also considered to be a value even if its components are not yet filled. Note that the size of the memory block must be known for the allocation to take place. $\lambda$-expressions are considered to be in weak head-normal form because we do not reduce inside the $\lambda$-body unless it is fully applied. In addition, partial applications (Pap) are also considered to be in weak head-normal form because they also represent functions, produced *via* currying. Note that under eager evaluation the values are either scalars or pointers to data objects residing in the heap — unevaluated expressions (*thunks*) are never treated as values. Similarly, selecting components from heap objects directly produces data values instead of thunks that would in lazy evaluation remain unevaluated until their values were

demanded.

It is not possible to test algorithmically within a functional language whether or not a given expression is in weak head-normal form without destroying the confluence property. An expression may not be in weak head-normal form at one point, but may become so after some evaluation. However, it is possible without destroying confluence to add a monotonic primitive function W to the language that produces a signal value (represented by •) when applied to expressions that are already in weak head-normal form. The meaning of the W operator is precisely captured by the semantic function $\mathcal{W}(E)$ defined over the range $\{\perp, \bullet\}$ as shown in Figure 4. This function produces the signal • if the given expression is in weak head-normal form; otherwise it produces no signal, and we say that its value is $\perp$. We also use a primitive operator & that combines two signals to produce a single signal.

### 4.3 Rewriting Semantics

Figure 5 shows the rewrite rules that apply to kernel language terms without barriers. We make extensive use of intermediate identifiers in order to preserve sharing of subexpressions. This system therefore represents a Graph Rewriting System [3].

The semantic function $\mathcal{W}$ is used to define the W operator (the **WHNF** rule) as well as to express additional conditions under which other rewrite rules apply. For example, the **List** rules (hd and tl) fire only when the Conscell is available and the component being selected becomes a value. The latter condition reflects the data-driven nature of eager evaluation. In contrast, under lazy evaluation, selecting a component from a Cons-cell would immediately return a *thunk* which is subsequently evaluated to a value only if necessary. From the point of view of value semantics, these two evaluation strategies are equivalent, but they yield very different dynamic resource utilization and termination characteristics.

Note that all the **Apply** rules except the Sap rule are *non-strict*: they fire even if the applied arguments have not yet become a value. In the case of Sap we explicitly wait for all the arguments before proceeding with the application. In other cases, partial applications are created whenever a function is applied to a number of arguments less than its arity[7]. Once all the arguments have been syntactically applied, the function body is instantiated, replacing $\lambda$-bound parameters by actual arguments. The operation $\mathcal{RB}[E]$ simply creates a new copy of the $\lambda$-body by renaming all its bound identifiers to new identifiers.

The rules for I-structures and M-structures implement their respective protocols. The rules for I-structures show that all I-fetch operations to a known location must wait until an I-store operation succeeds in storing a value to the same location. Two I-store operations to the same location take the program into an inconsistent state (T). The M-structure rules are similar, except that each M-take operation non-deterministically matches and eliminates one M-put operation within the same control region. Two successive M-put operations take the program into an inconsistent state.

In order to apply the rewriting rules shown in Figure 5 properly, program terms must be brought into a canoni-

---

[7]When creating partial applications we assume that their previously accumulated arguments are not copied, otherwise their termination behavior may get affected, as it will become clear shortly.

**WHNF Rule:**

$$\frac{\mathcal{W}(X) = \bullet}{\mathtt{W}(X) \longrightarrow \bullet}$$

$\delta$ **Rules:**

$$\&(\bullet, \bullet) \longrightarrow \bullet$$
$$+(\underline{m}, \underline{n}) \longrightarrow +(m, n)$$
$$\cdots$$

**Cond Rules:**

$$\mathtt{Cond\ True}\ (E_1, E_2) \longrightarrow E_1$$
$$\mathtt{Cond\ False}\ (E_1, E_2) \longrightarrow E_2$$

**List Rules:**

$$\frac{X = \mathtt{Cons}(Y, Z) \qquad \mathcal{W}(Y) = \bullet}{\mathtt{hd}(X) \longrightarrow Y}$$

$$\frac{X = \mathtt{Cons}(Y, Z) \qquad \mathcal{W}(Z) = \bullet}{\mathtt{tl}(X) \longrightarrow Z}$$

**Apply Rules:**

$$\frac{F = \lambda X_1 \cdots X_n.\ E}{\mathtt{Ap}(F, X) \longrightarrow \mathtt{Pap}(F, \underline{n-1}, X)}$$

$$\frac{F = \mathtt{Pap}(F', \underline{n}, X_1, \ldots, X_i) \qquad n > 1}{\mathtt{Ap}(F, X) \longrightarrow \mathtt{Pap}(F', \underline{n-1}, X_1, \ldots, X_i, X)}$$

$$\frac{F = \mathtt{Pap}(F', 1, X_1, \ldots, X_i)}{\mathtt{Ap}(F, X) \longrightarrow \mathtt{Fap}(F', X_1, \ldots, X_i, X)}$$

$$\frac{F = \lambda X_1 \cdots X_n.\ E}{\mathtt{Fap}(F, Y_1, \ldots, Y_n) \longrightarrow (\mathcal{RB}[E])[Y_1/X_1, \ldots, Y_n/X_n]}$$

**Strict Apply Rule:**

$$\frac{\mathcal{W}(X_1) = \bullet \quad \cdots \quad \mathcal{W}(X_n) = \bullet}{\mathtt{Sap}(F, X_1, \ldots, X_n) \longrightarrow \mathtt{Fap}(F, X_1, \ldots, X_n)}$$

**I-structure Rules:**

$$\frac{\mathtt{I\text{-}store}(X, \underline{i}, Z) \qquad \mathcal{W}(X) = \bullet \quad \mathcal{W}(Z) = \bullet}{\mathtt{I\text{-}fetch}(X, \underline{i}) \longrightarrow Z}$$

$$\frac{\mathtt{I\text{-}store}(X, \underline{i}, Z) \quad \mathcal{W}(X) = \bullet \quad \mathcal{W}(Z) = \bullet \quad \mathcal{W}(Z') = \bullet}{\mathtt{I\text{-}store}(X, \underline{i}, Z') \longrightarrow \mathsf{T}}$$

**M-structure Rules:**

$$\frac{\mathcal{W}(X) = \bullet \quad \mathcal{W}(Z) = \bullet}{\mathtt{M\text{-}put}(X, \underline{i}, Z);\ Y = \mathtt{M\text{-}take}(X, \underline{i}) \longrightarrow Y = Z}$$

$$\frac{\mathcal{W}(X) = \bullet \quad \mathcal{W}(Z) = \bullet \quad \mathcal{W}(Z') = \bullet}{\mathtt{M\text{-}put}(X, \underline{i}, Z);\ \mathtt{M\text{-}put}(X, \underline{i}, Z') \longrightarrow \mathsf{T}}$$

Figure 5: Rewrite Rules for the Kernel Language.

---

cal form after each rewriting step using the canonicalization rules shown in Figure 6. The canonicalization rules cover block flattening, identifier and constant substitution, elimination of redundant bindings, and propagation of T. These rules cannot be expressed as rewriting rules in the traditional sense. Nested blocks are flattened by merging their sets of bindings: this cannot cause any name clashes since we always rename all bound identifiers of a $\lambda$-expression on application. Bindings of the form $x = y$ or $x = c$ cause substitution of $[y/x]$ or $[c/x]$ respectively within their entire

**Block Flattening Rule:**

$$\{\; x = \{\; y_1 = E'_1;$$
$$\qquad\qquad \cdots$$
$$\qquad\quad y_m = E'_m;$$
$$\qquad\qquad \text{in } y \;\};$$
$$\quad x_1 = E_1;$$
$$\qquad \cdots$$
$$\quad x_n = E_n;$$
$$\quad \text{in } z \;\}$$

$$\xrightarrow{\;cn\;}$$

$$\{\; x = y;$$
$$\quad y_1 = E'_1;$$
$$\qquad \cdots$$
$$\quad y_m = E'_m;$$
$$\quad x_1 = E_1;$$
$$\qquad \cdots$$
$$\quad x_n = E_n;$$
$$\quad \text{in } z \;\}$$

**Substitution Rules:**

$$\{\; x_1 = E_1; \cdots; x = c; \cdots; x_n = E_n;\; \text{in } z \;\} \xrightarrow{\;cn\;} \qquad (c \neq \odot)$$
$$\{\; x_1 = E_1[c/x]; \cdots; x_n = E_n[c/x];\; \text{in } z[c/x] \;\}$$

$$\{\; x_1 = E_1; \cdots; x = y; \cdots; x_n = E_n;\; \text{in } z \;\} \xrightarrow{\;cn\;} \qquad (x \neq y)$$
$$\{\; x_1 = E_1[y/x]; \cdots; x_n = E_n[y/x];\; \text{in } z[y/x] \;\}$$

**Degenerate Cycle Rule:**

$$\{\; x_1 = E_1; \cdots; x = x; \cdots; x_n = E_n;\; \text{in } z \;\} \xrightarrow{\;cn\;}$$
$$\{\; x_1 = E_1[\odot/x]; \cdots; x = \odot; \cdots; x_n = E_n[\odot/x];\; \text{in } z[\odot/x] \;\}$$

**Blowup Rule:**

$$\{\; x_1 = E_1; \cdots; x = \top; \cdots; x_n = E_n;\; \text{in } z \;\} \xrightarrow{\;cn\;} \top$$

Figure 6: Canonicalization Rules.

lexical scope[8]. Only identifiers or constants are allowed to be substituted, in order to preserve precise sharing of sub-expressions among terms. After the substitution all such simple bindings are eliminated from the program term, leaving it in a unique canonical form.

It is also possible to encounter degenerate bindings of the form $x = x$ during rewriting and substitution, because of circular dependencies present within the program. In this case, the semantic meaning of $x$ is bottom ($\perp$), because such bindings represent deadlocked computation. Such bindings are canonicalized by substituting a special constant $\odot$ for $x$ ($\mathcal{W}(\odot) = \perp$), retaining the binding $x = x$ in the form $x = \odot$ to distinguish this case from the properly terminating computations. Finally, the blowup rule shows that once an inconsistent state ($\top$) is created the whole program becomes inconsistent.

## 5 Semantics of Barriers

We define the semantics of barriers by translating a kernel program with barriers into another kernel program without them, in which the termination of each control region is represented explicitly as a signal. To achieve this, we must precisely define the meaning of termination for each construct in our kernel language and provide a way of combining it in order to produce a joint termination signal.

### 5.1 What is Termination?

Let us consider the termination semantics of a simple parallel block. We translate it into another block as follows:

---

$$\mathbf{TE}[\; \{\; x_1 = E_1;$$
$$\qquad \cdots$$
$$\quad x_n = E_n;$$
$$\quad \text{in } x \;\} \;]$$

$$=$$

$$\{\; x_1, s_1 = \mathbf{TE}[E_1];$$
$$\quad s'_1 = s_1 \;\&\; \mathbf{W}(x_1);$$
$$\qquad \cdots$$
$$\quad x_n, s_n = \mathbf{TE}[E_n];$$
$$\quad s'_n = s_n \;\&\; \mathbf{W}(x_n);$$
$$\quad \text{in } x, s'_1 \;\&\; \cdots \;\&\; s'_n \;\}$$

An expression $E_i$ is translated as $\mathbf{TE}[E_i]$ which dynamically returns a value (bound to $x_i$) along with an explicit termination signal $s_i$ defined over the range $\{\perp, \bullet\}$. We use the weak head-normal form operator ($\mathbf{W}$) to wait for $x_i$ to become a value and combine it with the signal $s_i$ using the signal combining operator $\&$. All such signals within the block are similarly combined into a single signal for the whole block.

Since every expression node is explicitly named in the kernel language, the signal produced in the above manner ensures the termination of the actual computation taking place within the block. But what about references to free identifiers of the block? Some of these references are used within the actual block computation and hence are automatically included in its termination signal, but others are not. For example, free identifiers present in a non-strict data constructor or the free identifiers of a conditional branch that is not taken do not directly contribute towards the result of the computation. Should we nonetheless explicitly include all such references as part of the overall termination of the block?

There are two choices before us at this point, leading to two different versions of the termination semantics. Under "data-driven" semantics, we account for all the syntactically free identifiers of a block as part of its computation, and hence always include them in the block's overall termination. Under "demand-driven" semantics, only those identifiers are included in a block's termination which are demanded by the actual computation taking place within the block. Correspondingly, we obtain two different translations for barriers, described below and in Section 7 respectively.

### 5.2 Data-Driven Barrier Translation

Figure 7 shows the rules for transforming Kernel Id language constructs to make their termination events explicit. We provide two sets of translation rules. $\mathbf{TE}[]$ translates a kernel language expression $E$ into a new composite expression that computes both the value of the original expression $E$ and a termination signal. $\mathbf{TS}[]$ translates a kernel language statement $S$ into a parallel statement group $S'$ and a signal identifier $s$ that is bound to a signal expression representing the termination of $S$. Additional syntactic properties of the given expression or the statement may also be used, such as the free identifiers ($FV$) or the bound identifiers ($BV$) of an expression or a statement.

In all the rules shown in Figure 7 we maintain the invariant that the weak head-normal forms of all the free identifiers statically present within a construct are implicitly or explicitly collected into its signal. This is reflected in the rules for identifiers, primitive function applications, conditionals, $\lambda$-expressions, and user-defined function applications. For strict primitives such as $+$ this may be redundant, but it keeps the semantics intuitive and *compositional*, in that the termination of a program term (expression or statement) depends solely on the termination of its sub-terms.

EXPRESSIONS

**TE** :: Expression → Expression

$\mathbf{TE}[c]$ = $c, \bullet$

$\mathbf{TE}[x]$ = $x, \mathtt{W}(x)$

$\mathbf{TE}[PF^n(X_1,\ldots,X_n)]$ = $PF^n(X_1,\ldots,X_n), \quad \mathtt{W}(X_1) \ \& \ \cdots \ \& \ \mathtt{W}(X_n)$

$\mathbf{TE}[\mathtt{Cond}\ (X,E_1,E_2)]$ = $\{\ y,s = \mathtt{Cond}\ (X, \mathbf{TE}[E_1], \mathbf{TE}[E_2]);$
$\quad\quad$ in $y,s \ \& \ \mathtt{W}(z_1) \ \& \ \cdots \ \& \ \mathtt{W}(z_n)\ \}$

where $z_1,\ldots,z_n = FV(E_1) \cup FV(E_2)$

$\mathbf{TE}[\lambda x_1 \cdots x_n.\ E]$ = $\lambda x_1 \cdots x_n.\ \{\ y,s = \mathbf{TE}[E]$
$\quad\quad\quad$ in $y,s \ \& \ \mathtt{W}(x_1) \ \& \ \cdots \ \& \ \mathtt{W}(x_n)\ \}, \quad \mathtt{W}(z_1) \ \& \ \cdots \ \& \ \mathtt{W}(z_n)$

where $z_1,\ldots,z_n = FV(\lambda x_1 \cdots x_n.\ E)$

$\mathbf{TE}[\mathtt{Ap}(F,X)]$ = $\{\ y,s = \mathtt{Ap}(F,X)$ in $y,s \ \& \ \mathtt{W}(F) \ \& \ \mathtt{W}(X)\ \}$

$\mathbf{TE}[\{\ S^{par}\ \text{in}\ (X_1,\ldots,X_n)\ \}]$ = $\{\ S'^{par}\ \text{in}\ (X_1,\ldots,X_n,s)\ \}$

where $S'^{par}, s = \mathbf{TS}[S^{par}]$

STATEMENTS

**TS** :: Statement → Statement × Signal Identifier

$\mathbf{TS}[x_1,\ldots,x_n = E]$ = $(x_1,\ldots,x_n,s_1 = \mathbf{TE}[E];$
$\quad\quad s = s_1 \ \& \ \mathtt{W}(x_1) \ \& \ \cdots \ \& \ \mathtt{W}(x_n)), \quad s$

$\mathbf{TS}[\epsilon]$ = $(s = \bullet), \quad s$

$\mathbf{TS}[S_1;\ldots;S_n]$ = $(S'_1;\ldots;S'_n; s = s_1 \ \& \ \cdots \ \& \ s_n), \quad s$

where $S'_i, s_i = \mathbf{TS}[S_i] \quad 1 \le i \le n$

$\mathbf{TS}[S_1^{par} \mathbin{\text{---}} S_2^{par}]$ = $(S'^{par}_1;$
$\quad\quad F = \lambda s.\{\ S'^{par}_2\ \text{in}\ y_1,\ldots,y_m,s_2\ \};$
$\quad\quad y_1,\ldots,y_m,s'_2 = \mathtt{Sap}(F,s_1)), \quad s'_2$

where $S'^{par}_i, s_i = \mathbf{TS}[S^{par}]\quad 1 \le i \le 2$
$\quad\quad y_1,\ldots,y_m = BV(S_2^{par})$

$\mathbf{TS}[\mathtt{I\text{-}store}(X_0,X_1,X_2)]$ = $(\mathtt{I\text{-}store}(X_0,X_1,X_2);$
$\quad\quad s = \mathtt{W}(X_0) \ \& \ \mathtt{W}(X_1) \ \& \ \mathtt{W}(X_2)), \quad s$

$\mathbf{TS}[\mathtt{M\text{-}put}(X_0,X_1,X_2)]$ = $(\mathtt{M\text{-}put}(X_0,X_1,X_2);$
$\quad\quad s = \mathtt{W}(X_0) \ \& \ \mathtt{W}(X_1) \ \& \ \mathtt{W}(X_2)), \quad s$

Figure 7: Barrier Translation Rules based on Data-Driven Eager Evaluation.

The translation rule for function application shows that we wait not only for the operator and the operand to become values, but also for the termination of the function application itself. The rule for $\lambda$-expressions correspondingly translates each function body to return a signal upon its termination. This body signal waits for all the arguments of the function, thereby ensuring that unused or non-strict arguments are also collected before signaling the termination of the function application.

It may seem redundant to collect the weak head-normal form of a function argument at the application site when we already collect it within the body of the function. Collection at the application site is necessary, however, because any of these sites could be a partial application. In that case the argument contributes towards building the partial application closure rather than executing a function body. Building a closure is a resource-consuming computation and its termination must therefore be tracked independently.

Finally, the translation of a barrier encloses its post-region in a $\lambda$-expression which is applied to the accumulated termination signal of its pre-region using the strict application primitive Sap. This prevents the initiation of the computation below the barrier until after the computation above the barrier has terminated.

It is easy to see that the data-driven termination semantics presented above provides precise control over both resource management and scheduling of side-effects, as discussed in Section 3. The signal generation and composition rules shown in Figure 7 ensure that a termination signal from a control region is not produced until all the region's inputs have become values and the entire computation within it, including side-effects, has terminated. Any resources used during the computation may therefore then be safely released.

### 5.3 Syntactic Properties of Barriers

It is easy to show that the barrier translation defined in Figure 7 satisfies the following two desirable syntactic properties:

**Property 1 (Eliminating Redundant Barriers)** *Given a statement group* $S^{par}$, *we have*

$$\left(\begin{array}{c}\text{---}\\S^{par}\end{array}\right) \equiv \left(\begin{array}{c}S^{par}\\\text{---}\end{array}\right) \equiv S^{par}$$

**Property 2 (Associativity of Barriers)** *Given statement groups* $S_1^{par}$, $S_2^{par}$, *and* $S_3^{par}$, *we have*

$$\left(\begin{array}{c}(S_1^{par}\\\text{---}\\S_2^{par})\\\text{---}\\S_3^{par}\end{array}\right) \equiv \left(\begin{array}{c}(S_1^{par}\\\text{---}\\(S_2^{par}\\\text{---}\\S_3^{par}))\end{array}\right)$$

8

## 6 Placement of Barrier Translation within a Compiler

The barrier translation presented above is a source-to-source transformation of a Kernel Id program. A natural question is when should this translation be performed within the compiler: should we do it before or after other program transformations such as λ-lifting or pattern-matching? Moreover, most of the rewrite rules (Figure 5) and the canonicalization rules (Figure 6) may be used at compile-time as compiler optimizations [2]. We may even want to postpone the barrier translation until after the optimization phase because the translation increases the complexity of the program. Does it matter (semantically) when this translation is performed?

### 6.1 Canonicalization vs. Barrier Translation

Unfortunately, the barrier translation shown in Figure 7 is very fragile even with respect to the simple canonicalization rules of Figure 6. Canonicalization may eliminate some free identifiers from a control region, which seriously affects its termination semantics as shown in the following example.

**Example 8:**
```
{ ( x = a;
    ---
    a = 1; )
  in a }
```

Without canonicalization this program deadlocks under the barrier translation because the value of a is required to discharge the barrier. However, the binding (x = a) is redundant and can be eliminated by canonicalization, and this would make the above program terminate normally.

### 6.2 Program Transformations vs. Barrier Translation

**The Role of Free Identifiers**

Compile-time rewriting and other program transformations can either increase or decrease the set of free identifiers of a control region, again altering its termination semantics. For instance, consider the following program fragment.

**Example 9:**
```
def f x = 1;
( y = f a;
  ---
  ... )
```

If the user (or the compiler) decided to inline the function f before performing barrier translation, it would eliminate the free identifier a from the pre-region of the barrier. The computation below the barrier would no longer wait for a, as it would have done if the inlining had been performed after barrier translation.

Most other rewrite rules shown in Figure 5 would also yield a different termination behavior if used at compile-time before performing barrier translation. For instance, applying Cond rules would eliminate the free identifiers of one of the branches that would otherwise be collected into the signal of the translated conditional expression (see Figure 7). Optimizing general applications Ap to partial applications Pap would introduce new free identifiers into a control region that would not be present if the optimization had been performed after barrier translation.

Another program transformation that changes the set of free identifiers of a control region is λ-lifting [10]. λ-lifting converts free identifiers of a function into additional parameters and then lifts the functions to the top-level as combinators. All occurrences of the lifted function identifier are replaced by partial applications of the transformed function to its previously free identifiers. This transformation would eliminate free identifiers from the region containing the lifted λ-expression, and introduce free identifiers into regions where the lifted function identifier occurs. Again, this would alter the termination semantics of those regions.

**The Role of Bound Identifiers**

The translation of a λ-expression in Figure 7 shows that the signal of the λ-body collects the weak head-normal forms of all incoming bound identifiers of the function. This is essentially a resource management requirement — if an incoming argument is not used within the body of the function then it must be collected into its termination signal. This ensures that no dangling argument values are left behind to be received after the production of the termination signal has allowed the function's activation frame to be deallocated.

This data-driven compilation scheme creates some subtle problems for transformations involving functions that do not use some of their arguments. Consider the following example.

**Example 10:**

```
    def f x y = y;                    def f x y = y;
    { f' = f b;                       { f' = f b;
      ( a = f' 1;       Inlining        ( a = 1;
        ---              ──→             ---
        b = a; );                        b = a; );
      in a }                           in a }
```

It is easy to detect that the application (f' 1) before the barrier is a full-arity application of the function f. If the compiler decided to inline this function, then the resulting code after canonicalization shown on the right would be able to discharge the barrier immediately and the program would terminate normally. If the compiler decided not to inline, then the data-driven implementation of the function f would wait for its first argument b at run-time before releasing the termination signal which in turn is supposed to discharge the barrier and allow the identifier b to become a value. Therefore, the program would deadlock instead.

The above discrepancy arises because the data-driven implementation deliberately waits for the unused bound identifiers of a function even though at compile-time we eliminate such bindings via canonicalization. On the other hand, this problem would not arise if the compiler was careful not to inline functions whose partial applications spanned multiple control regions. Whether this condition severely restricts the applicability of the inlining optimization is a topic for further investigation.

### 6.3 Multiple Argument Functions vs. Barrier Translation

Most functional language implementations make a distinction between the following two function definitions for pragmatic reasons; the function f1 on the left-hand side takes two arguments before expanding into its body while the function f1 on the right-hand side takes one argument and evaluates to another function of one argument.

```
                          def f2 x =
   def f1 x y = ...;   vs.   { def g y = ...;
                               in g };
```

However, this distinction is not semantically observable: these functions can be used interchangeably. Unfortunately, the termination semantics shown in Figure 7 exposes this distinction as shown by the following example.

**Example 11:**

```
{ f' = f1 b;              { f' = f2 b;
  ( a = f' 1;               ( a = f' 1;
    ---           vs.         ---
    b = a; );                 b = a; );
  in a }                   in a }
```

The left-hand side uses the function f1 with two arguments that is curried over two applications. f' gets bound to a partial application of the function f1 which behaves like the identity function in value semantics but has a slightly more complex behavior for termination semantics. The function application (f' 1) above the barrier does not produce a signal until it also receives the first argument of the function f1. This is because the barrier translation rules mandate that functions should collect their unused arguments into signals. But that can never happen, because the actual value for that argument b cannot be produced until the application itself terminates. So the program deadlocks.

The right-hand side computes the identifier f' using the function f2 which evaluates to a pure identity function that does not contain any free identifiers within its closure. The application (f' 1) above the barrier therefore successfully terminates, producing the value for b which allows the binding for f' to terminate. Thus the program successfully produces a value and terminates.

### 6.4 Making Data-Driven Barriers more Robust

Judging by the above discussion, it seems that the simplest way to guarantee that the source program semantics are kept intact is to perform barrier translation before any other program transformation is performed. This in itself is not problematic, and can be easily done within the compiler, although it shows that the barrier translation scheme presented in Figure 7 is not sufficiently robust. We can improve its robustness by performing the barrier translation in the following manner.

1. Record the set of free identifiers of every pre-region of a barrier in the source program.

2. Perform program transformations: λ-lifting, inlining, etc.

3. Canonicalize the program eliminating all redundant bindings.

4. Translate the barriers in a manner similar to Figure 7 using the weak head-normal forms of previously saved free identifiers for each pre-region.

Recording the set of free identifiers of every pre-region beforehand ensures that subsequent changes in this set due to canonicalization or program transformations would have no effect on the termination semantics of that control region. This set would be recorded at the barrier itself by inspecting the original source program[9]. These identifiers would then be used to control the discharge of the barrier at the time of translation irrespective of the actual set of free identifiers

---

[9]Identifier substitutions would still need to be performed on this set of free identifiers.

of the pre-region. Furthermore, canonicalizing the whole program before performing the barrier translation ensures that redundant bindings are consistently eliminated from everywhere within the program without interfering with the barrier translation. Note that this scheme would correctly handle Example 8 and Example 9 by permanently associating the free identifier a with the region above the barrier and collecting it within its termination signal during barrier translation.

The problem of preserving equivalence of multiple argument function implementations (Example 11) is a little more difficult to overcome, and this shows that the data-driven termination semantics are fundamentally at odds with desirable semantic properties enjoyed by the value semantics. In the next section we explore a different approach to this problem.

## 7 Mixed Eager and Lazy Evaluation

In this section we present an alternate translation for barriers based on a demand-driven identifier dereferencing mechanism within the context of non-strict, eager evaluation. This translation overcomes some of the semantic shortcomings of the data-driven barrier translation presented above and is also sufficient for resource management under an environment-based implementation. We will show the translation rules and discuss their properties, comparing them with the data-driven translation.

The concept of a control region remains unchanged under this mixed model of evaluation. We spawn all tasks within a control region eagerly, but each computation fetches the values of its free (input) identifiers as needed rather than being explicitly provided with them. Similarly, the values of the bound (output) identifiers computed within a region are not automatically sent to the computations that may need them; instead they are stored in a known place from where other computations may obtain them as needed. We now outline a possible implementation of this idea.

In order that each spawned computation would "know" where to get the data from if and when it needed it, every control region would keep an *environment* in the form of a table of values. Nested control regions would give rise to a chain of environments. The value for each bound identifier in a control region would now be stored at a fixed location specified by the compiler within its environment table, and all subsequent computations that need this value would fetch it from there. It is not safe to flatten or copy the environment tables in this model because reading an identifier that is not used in a control region can change the termination behavior of that region. The inability to flatten environments may have serious implications for the efficiency of the implementation.

### 7.1 Barrier Translation based on Mixed Evaluation

Figure 8 shows the new barrier translation rules under the mixed execution model. We no longer collect weak head-normal forms of the free identifiers of every sub-expression into its signal. Nevertheless, we still need to collect the weak head-normal forms of the bound identifiers of a computation into its signal, identifying the termination of that computation. Similarly, the weak head-normal forms of the identifiers participating in side-effect operations I-store and M-put signify the completion of those operations. Within

10

$$
\begin{array}{lll}
& & \text{EXPRESSIONS} \\
\mathbf{TE} & :: & \text{Expression} \to \text{Expression} \\[4pt]
\mathbf{TE}[c] & = & c, \quad \bullet \\
\mathbf{TE}[x] & = & x, \quad \bullet \\
\mathbf{TE}[PF^n(X_1,\ldots,X_n)] & = & PF^n(X_1,\ldots,X_n), \quad \bullet \\
\mathbf{TE}[\text{Cond }(X,E_1,E_2)] & = & \text{Cond }(X,\mathbf{TE}[E_1],\mathbf{TE}[E_2]) \\
\mathbf{TE}[\lambda x_1 \cdots x_n.\, E] & = & \lambda x_1 \cdots x_n.\, \mathbf{TE}[E], \quad \bullet \\
\mathbf{TE}[\text{Ap}(F,X)] & = & \text{Ap}(F,X) \\
\mathbf{TE}[\{\ S^{par}\ \text{in}\ (X_1,\ldots,X_n)\ \}] & = & \{\ S'^{par}\ \text{in}\ (X_1,\ldots,X_n,s)\ \} \\
\text{where}\quad S'^{par}, s = \mathbf{TS}[S^{par}] & &
\end{array}
$$

$$
\begin{array}{lll}
& & \text{STATEMENTS} \\
\mathbf{TS}[] & :: & \text{Statement} \to \text{Statement} \times \text{Signal Identifier} \\[4pt]
\mathbf{TS}[x_1,\ldots,x_n = E] & = & (x_1,\ldots,x_n,s_1 = \mathbf{TE}[E]; \\
& & \quad s = s_1\ \&\ W(x_1)\ \&\ \cdots\ \&\ W(x_n)),\quad s \\
\mathbf{TS}[\epsilon] & = & (s = \bullet),\quad s \\
\mathbf{TS}[S_1;\ldots;S_n] & = & (S_1';\ldots;S_n';s = s_1\ \&\ \cdots\ \&\ s_n),\quad s \\
\text{where}\quad S_i', s_i = \mathbf{TS}[S_i]\quad 1 \le i \le n & & \\
\mathbf{TS}[S_1^{par} \mathrel{---} S_2^{par}] & = & (S_1'^{par}; \\
& & \quad F = \lambda s.\{\ S_2'^{par}\ \text{in}\ y_1,\ldots,y_m,s_2\ \}; \\
& & \quad y_1,\ldots,y_m,s_2' = \text{Sap}(F,s_1)),\quad s_2' \\
\text{where}\quad S_i'^{par}, s_i\ =\ \mathbf{TS}[S^{par}]\quad 1 \le i \le 2 & & \\
\phantom{\text{where}}\quad y_1,\ldots,y_m\ =\ BV(S_2^{par}) & & \\
\mathbf{TS}[\text{I-store}(X_0,X_1,X_2)] & = & (\text{I-store}(X_0,X_1,X_2); \\
& & \quad s = W(X_0)\ \&\ W(X_1)\ \&\ W(X_2)),\quad s \\
\mathbf{TS}[\text{M-put}(X_0,X_1,X_2)] & = & (\text{M-put}(X_0,X_1,X_2); \\
& & \quad s = W(X_0)\ \&\ W(X_1)\ \&\ W(X_2)),\quad s
\end{array}
$$

Figure 8: Barrier Translation Rules based on Mixed Evaluation.

expressions, identifiers are expected to be fetched on demand and therefore do not need any separate termination signal.

Another interesting difference between this translation and that of Figure 7 is the rule for $\lambda$-expressions. Note that the rule presented in Figure 7 collects all $\lambda$-bound identifiers into its termination signal in order to capture the identifiers that are left unused. On the other hand, the rule presented in Figure 8 does not collect any such identifiers. This is because unused identifiers are never fetched under mixed evaluation. This property allows the equivalence of multiple argument function implementations (shown in Example 11) to be preserved. Under the new translation rule, the signal from the body of the function f1 on the left-hand side would no longer depend on its unused bound identifier x; so the left-hand side would also produce a value and terminate properly.

## 7.2 Properties of Mixed Barrier Semantics

The translation shown in Figure 8 is automatically robust with respect to compiler optimizations and program transformations, since termination of a computation depends only upon the values of identifiers that are actually required during that computation. No artificial data dependencies may be introduced or eliminated by program transformations. It also satisfies the desired equivalence of multiple argument function implementations, and is therefore a semantically cleaner alternative to purely eager, data-driven barrier translation.

As expected, the mixed evaluation model is much more relaxed in terms of the condition governing the discharge of a barrier. The termination of the pre-region of a barrier now solely depends upon the termination of the actual computation spawned within it and not on the weak headnormal forms of unused identifiers. This may cause some programs to terminate properly under this semantics that would have deadlocked under data-driven termination semantics, as shown below.

**Example 12:**
```
{ ( def f x = x + a ;
    ---
    a = 1; )
  in f 3 }
```

Under data-driven termination semantics the value of a is still required to discharge the barrier since it is a free identifier of the pre-region; but since a is bound in the post-region, it never gets defined and the program deadlocks. Under mixed barrier semantics function definitions always terminate immediately; so the barrier discharges, allowing the overall expression to produce a value and terminate properly.

Both data-driven and mixed termination semantics provide the power to control the termination of actual computations including those involving side-effects. The major difference between them shows up while using them for resource management. Data-driven termination semantics automatically provide a full resource management capability, by ensuring that even unused values have been received by the time a termination signal is produced from a control region. We may need to move to an environment-based implementation in order to provide a similar capability under the mixed termination semantics.

11

## 8 Conclusions and Future Work

In this paper we have presented two different termination semantics for barriers in an eagerly evaluated, implicitly parallel, imperative language. In each case, the semantics have been presented as a translation from a language with barriers into a kernel language without barriers where termination signals have been made explicit using a weak head-normal form operator ($W$) and a strict application operator ($Sap$).

The purely eager termination semantics were motivated by a data-driven execution model and have been implemented in the Id/pH compiler. The formalization of these semantics has helped us in uncovering some bugs in the current implementation. This model unfortunately lacks certain desirable semantic properties, such as the equivalence of multiple argument function implementations; this motivated us to devise alternate, mixed termination semantics that combine eager spawning of computation along with lazy handling of identifier references. It remains to be seen whether this mixture is adequate for resource management (for example, it might prove better to retain the rules of Figure 7 for $Ap$ and $PF^n$ in order to enforce a stronger view about the termination of function applications and non-strict data constructors, respectively).

At this point we do not have enough experience to judge whether an environment-based implementation of the mixed evaluation approach would be as efficient as the data-driven implementation of purely eager evaluation. It is possible that a clean solution would require us to separate barriers for resource management (such as the deallocation of frames and heap objects) from barriers that detect the completion of side-effects (such as writing into I-structures and M-structures). The resource management barrier semantics must capture some aspects of a particular implementation while the side-effects barrier semantics may be definable in an implementation independent manner.

## 9 Acknowledgments

## References

[1] Zena M. Ariola and Arvind. P-TAC: A Parallel Intermediate Language. In *Proc. Functional Programming Languages and Computer Architecture*, September 1989.

[2] Zena M. Ariola and Arvind. A Syntactic Approach to Program Transformations. In *Proc. Partial Evaluation and Semantics-Based Program Manipulation*, pages 116–129. ACM Press, September 1991.

[3] Zena M. Ariola and Arvind. Properties of a First-order Functional Language with Sharing. CSG Memo 347-1, Laboratory for Computer Science, MIT, Cambridge, MA 02139, June 1994. To appear in Theoretical Computer Science, September 1995.

[4] Arvind, Rishiyur S. Nikhil, and Keshav K. Pingali. I-Structures: Data Structures for Parallel Computing. *ACM Transactions on Programming Languages and Systems*, 11(4):598–632, 1989.

[5] Paul S. Barth. *Atomic Data Structures for Parallel Computing*. PhD thesis, Laboratory for Computer Science, MIT, Cambridge, MA 02139, March 1992. Available as Technical Report MIT/LCS/TR-532.

[6] Paul S. Barth, Rishiyur S. Nikhil, and Arvind. M-Structures: Extending a Parallel, Non-Strict, Functional Language with State. In *Proc. Functional Programming Languages and Computer Architecture*, pages 538–568. Springer-Verlag, 1991. LNCS 523.

[7] C.A.R. Hoare. *Communicating Sequential Processes*. Prentice Hall International Series in Computer Science. Prentice Hall, Englewood Cliffs, NJ, 1985.

[8] P. Hudak and P. Wadler (editors). Report on the programming language Haskell, a non-strict purely functional language (Version 1.0). Technical Report YALEU/DCS/RR777, Department of Computer Science, Yale University, April 1990.

[9] Paul Hudak. Mutable Abstract Datatypes -or- How to Have Your State and Munge It Too. Research Report YALEU/DCS/RR-914, Department of Computer Science, Yale University, New Haven, CT 06520, December 1992. Revised May 1993.

[10] Thomas Johnsson. Lambda Lifting: Transforming Programs to Recursive Equations. In *Proc. Functional Programming Languages and Computer Architecture*, pages 190–203. Springer-Verlag, 1985. LNCS 201.

[11] John Launchbury and Simon L. Peyton Jones. Lazy Functional State Threads. In *Proc. Programming Language Design and Implementation*, pages 24–35. ACM Press, June 1994.

[12] Rishiyur S. Nikhil. Id Language Reference Manual Version 90.1. Technical Report CSG Memo 284-2, Laboratory for Computer Science, MIT, Cambridge, MA 02139, July 15 1991.

[13] Rishiyur S. Nikhil, Arvind, James Hicks, Shail Aditya, Lennart Augustsson, Jan-Willem Maessen, and Yuli Zhou. pH Language Reference Manual, Version 1.0—preliminary. CSG Memo 369, Laboratory for Computer Science, MIT, Cambridge, MA 02139, January 1995.

[14] Simon L. Peyton Jones and Philip Wadler. Imperative Functional Programming. In *Proc. Principles of Programming Languages*, pages 71–84. ACM Press, January 1993.

[15] J.E. Stoy. The Semantics of Id. In A.W. Roscoe, editor, *A Classical Mind: Essays in Honour of C.A.R. Hoare*, pages 379–404. Prentice Hall, New York, 1994.

[16] Thorsten von Eicken, David Culler, and Klaus Erik Schauser. Berkeley Id90 I/O proposal. Draft Memo, University of California, Berkeley, August 1991.