

# Witnessing Side-Effects \*

Tachio Terauchi

EECS Department  
University of California, Berkeley  
tachio@cs.berkeley.edu

Alex Aiken

Computer Science Department  
Stanford University  
aiken@cs.stanford.edu

## Abstract

We present a new approach to the old problem of adding side effects to purely functional languages. Our idea is to extend the language with “witnesses,” which is based on an arguably more pragmatic motivation than past approaches. We give a semantic condition for correctness and prove it is sufficient. We also give a static checking algorithm that makes use of a network flow property equivalent to the semantic condition.

**Categories and Subject Descriptors** D.3.1 [Programming Languages]: Formal Definitions and Theory; D.3.3 [Programming Languages]: Language Constructs and Features

**General Terms** Design, Languages, Theory

**Keywords** Functional languages, side-effects

## 1. Introduction

Adding side-effects to a purely functional language is a well-known problem with a number of solutions [7, 9, 11, 15, 6, 1] with monads [10, 11, 7, 9] being arguably the most popular. In this paper, we propose a new approach to this old problem by attacking it from a different angle. Instead of starting from a language theoretic point of view, we start by introducing a simple programming feature called *witnesses* so that programs can explicitly order side-effects. This feature is motivated by a pragmatic observation and is straightforward. The catch is that, because it is so simple, it actually does not guarantee that a program is *correct* (i.e., that it can be viewed as a purely functional program). Instead, we argue that the feature makes it easy for programmers to write correct programs. We then formally state a natural semantic condition that is sufficient to guarantee correctness and give a static checking algorithm. The result is a new framework for guaranteeing correctness of side-effects in purely functional programs.

Besides arguably being more intuitive to programmers, our approach is more expressive than previous approaches. In particular, our approach does not force side-effects to occur in a sequential

\* This research was supported in part by NASA Grant No. NNA04CI57A; NSF Grant Nos. CCR-0234689, CCR-0085949, and CCR-0326577. The information presented here does not necessarily reflect the position or the policy of the Government and no official endorsement should be inferred.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

ICFP'05 September 26–28, 2005, Tallinn, Estonia.  
Copyright © 2005 ACM 1-59593-064-7/05/0009...\$5.00.

order. For example, a program is allowed to read from a reference cell in two unordered contexts as well as write to two different cells in two unordered contexts.

Besides providing new insights into the old problem of fitting side effects into functional languages for conventional von Neumann architectures, our work is motivated by the emergence of commercial parallel computer architectures (e.g., chip-multiprocessors or “multi-core” chips) that encourage parallel programming. It is well-known that the “explicit dependence” property of functional languages makes parallelization easier for both programmers and compilers.<sup>1</sup> However, frequent use of side-effects, namely manual destructive memory/resource updates, are believed to be important for programming high-performance parallel applications in practice. Hence a functional way to add side-effects without imposing parallelism-destroying sequentiality may be of practical interest for exploiting parallelism within these new architectures.

### 1.1 Contributions and Overview

This paper makes the following contributions:

- A simple language feature called *witnesses* that can be used to order side-effects. (Section 2)
- A semantic condition called *witness race freedom* for correct usage of witnesses and a proof of its sufficiency. (Section 3)
- An automatic algorithm for checking the afore-mentioned semantic condition that makes use of a network flow property. (Section 4)

The semantic condition is intuitive in the sense that it is directly motivated by the implications of writing race-free programs. The automatic algorithm is derived as a type inference algorithm for a substructural type system. The type system and its inference problem are somewhat subtle and interesting in their own right. Section 5 discusses related work. Section 6 concludes.

## 2. Preliminaries

We need a precise definition of what it means for side-effects to be “correct” within a functional language. A helpful idea is to show that a program’s semantics is independent of a class of “functional” program transformation rules. However, there is no consensus on the right set of transformations. For example, the transformation  $(\text{let } x = e \text{ in } e') \equiv (e'[x/e])$  for  $x \notin \text{free}(e)$  is not always true in systems based on linear-types. (Here  $\text{free}(e)$  is the set of free variables of  $e$ .)

To define correctness, we fix a set of program transformations expressive enough to model different functional reduction strate-

<sup>1</sup> But not easy, since there are other challenging issues such as selecting the right granularity of parallelism, but these issues apply equally to other languages and many solutions such as data-parallel operators and thread annotations already exist.

$$e ::= x \mid i \mid \lambda x. e \mid e e' \mid \text{let } x = e \text{ in } e' \mid e \otimes e' \mid \pi_i(e) \mid \text{write } e_1 e_2 e_3 \mid \text{read } e e' \mid \text{ref } e \mid \text{join } e e' \mid \bullet$$

Figure 1: The syntax of the language  $\lambda_{wit}$ .

gies, including call-by-value, call-by-need (i.e., lazy-evaluation), and parallel evaluation.<sup>2</sup> So, for example, a program that is invariant under this set of transformations evaluates to the same result regardless of whether the evaluation order is call-by-value or call-by-need. In parallel evaluation, invariance implies that a program is deterministic under any evaluation schedule.

For ease of exposition, we include the afore-mentioned program transformations directly in the semantics as non-deterministic reduction rules. The correctness criteria then reduces to showing that a program is confluent with respect to this semantics. Also, for the purpose of exposition, we restrict side-effects to imperative operations on first class references.

Figure 1 gives the syntax of  $\lambda_{wit}$ , a simple functional language with side-effects and witnesses. Note  $\lambda_{wit}$  has the usual features of a functional language: variables  $x$ , integers  $i$ , function abstractions  $\lambda x. e$ , function applications  $e e'$ , variable bindings  $\text{let } x = e \text{ in } e'$ , pairs  $e \otimes e'$ , and projections  $\pi_i(e)$  where  $i = 1$  or  $i = 2$ . Bindings  $\text{let } x = e \text{ in } e'$  can be recursive, i.e.,  $x \in \text{free}(e)$ . Three expression kinds work with references: reference writes  $\text{write } e_1 e_2 e_3$ , reference reads  $\text{read } e e'$ , and reference creations  $\text{ref } e$ . A read  $\text{read } e e'$  has a *witness* parameter  $e'$  along with a reference parameter  $e$  such that it does not read the reference  $e$  until it sees the witness  $e'$ . (Section 3 defines the formal meaning of “seeing the witness.”) Similarly, a write  $\text{write } e_1 e_2 e_3$  writes the expression  $e_2$  to the reference  $e_1$  after it sees the witness  $e_3$ . After completion of the read,  $\text{read } e e'$  returns a pair of the read value and a witness. Similarly,  $\text{write } e_1 e_2 e_3$  returns a witness after the write. In general, any side-effect primitive returns a witness of performing the corresponding side-effect; in the case of  $\lambda_{wit}$ , the side-effect primitives are just  $\text{write } e_1 e_2 e_3$  and  $\text{read } e e'$ .

Before describing the formal semantics of  $\lambda_{wit}$ , we describe novel properties of  $\lambda_{wit}$  informally by examples.

Programs in  $\lambda_{wit}$  can use witnesses to order side-effects. For example, the following program returns 2 regardless of the evaluation order because the read requires a witness of the write:

$$\text{let } x = (\text{ref } 1) \text{ in let } w = (\text{write } x \ 2 \ \bullet) \text{ in read } x \ w$$

(The symbol  $\bullet$  is used for dummy witnesses.) On the other hand,  $\lambda_{wit}$  does not guarantee correctness. For example, the following  $\lambda_{wit}$  program has no ordering between the read and the write and hence may return 1 or 2 depending on the evaluation order:

$$\text{let } x = (\text{ref } 1) \text{ in let } w = (\text{write } x \ 2 \ \bullet) \text{ in read } x \ \bullet$$

The expression kind  $\text{join } e e'$  joins two witnesses by waiting until it sees the witness  $e$  and the witness  $e'$  and returning a witness. For example, the following program returns the pair  $1 \otimes 1$  regardless of the evaluation order because the write waits until it sees witnesses of both reads:

$$\begin{aligned} &\text{let } x = (\text{ref } 1) \text{ in} \\ &\quad \text{let } y = (\text{read } x \ \bullet) \text{ in let } z = (\text{read } x \ \bullet) \text{ in} \\ &\quad \quad \text{let } w = (\text{write } x \ 2 \ (\text{join } \pi_2(y) \ \pi_2(z))) \text{ in} \\ &\quad \quad \quad \pi_1(y) \otimes \pi_1(z) \end{aligned}$$

Note that the two reads may be evaluated in any order. In general, witnesses are first class values and hence they can be passed to and

$$E ::= D \cup \{a \mapsto E\} \mid [] \mid E e \mid e E \mid E \otimes e \mid e \otimes E \mid \pi_i(E) \mid \text{write } E e e' \mid \text{write } e E e' \mid \text{write } e e' E \mid \text{read } e E \mid \text{read } E e \mid \text{ref } E \mid \text{join } E e \mid \text{join } e E$$

$$\begin{aligned} \text{App} & (S, E[(\lambda x. e) e']) \Rightarrow (S, E[e[a/x]] \uplus \{a \mapsto e'\}) \\ \text{Let} & (S, E[\text{let } x = e \text{ in } e']) \\ & \Rightarrow (S, E[e'[a/x]] \uplus \{a \mapsto e[a/x]\}) \\ \text{Pair} & (S, E[\pi_i(e_1 \otimes e_2)]) \Rightarrow (S, E[e_i]) \\ \text{Write} & (S, E[\text{write } \ell e \bullet]) \Rightarrow (S[\ell \leftarrow a], E[\bullet] \uplus \{a \mapsto e\}) \\ \text{Read} & (S, E[\text{read } \ell \bullet]) \Rightarrow (S, E[S(\ell) \otimes \bullet]) \\ \text{Ref} & (S, E[\text{ref } e]) \Rightarrow (S \uplus \{\ell \mapsto a\}, E[\ell] \uplus \{a \mapsto e\}) \\ \text{Join} & (S, E[\text{join } \bullet \bullet]) \Rightarrow (S, E[\bullet]) \\ \text{Arrive} & (S, E[a] \uplus \{a \mapsto e\}) \Rightarrow (S, E[e] \uplus \{a \mapsto e\}) \\ & \text{where } e \in V \\ \text{GC} & (S, D \uplus D') \Rightarrow (S, D) \\ & \text{where } \diamond \notin \text{dom}(D') \wedge \text{dom}(D') \cap \text{free}(D) = \emptyset \end{aligned}$$

Figure 2: The semantics of  $\lambda_{wit}$ .

returned from a function, captured in function closures, and even written to and read from a reference. Witnesses are a simple feature that can be used to order side-effects in a straightforward manner.

In the rest of this section, we describe the semantics of  $\lambda_{wit}$  so that we can formally define when a  $\lambda_{wit}$  program is correct, i.e., when it is confluent. Figure 2 shows the semantics of  $\lambda_{wit}$ , which is defined via reduction rules of the form  $(S, D) \Rightarrow (S', D')$  where  $S, S'$  are *reference stores* and  $D, D'$  are *expression stores*. A reference store is a function from a set of *reference locations*  $\ell$  to *ports*  $a$ , and an expression store is a function from a set of ports to expressions. Here, expressions include any expression from the source syntax extended with reference locations and ports. Given a program  $e$ , evaluation of  $e$  starts from the initial state  $(\emptyset, \{\diamond \mapsto e\})$  where the symbol  $\diamond$  denotes the special *root port*. Ports are used for evaluation sharing.<sup>3</sup> The reduction rules are parametrized by the evaluation contexts  $E$ . For an expression  $e$ ,  $\text{free}(e)$  is the set of free variables, ports, and reference locations of  $e$ . For an expression store  $D$ ,  $\text{free}(D) = \text{dom}(D) \cup \bigcup_{e \in \text{ran}(D)} \text{free}(e)$ .

We briefly describe the reduction rules from top-to-bottom. The rule **App** corresponds to a function application. For functions  $F$  and  $F'$ ,  $F \uplus F'$  denotes  $F \cup F'$  if  $\text{dom}(F) \cap \text{dom}(F') = \emptyset$  and is undefined otherwise. Thus **App** creates a fresh port  $a$  and stores  $e'$  at  $a$ . The rule **Let** is similar. The rule **Pair** projects the  $i$ th element of the pair. The rule **Write** creates a fresh port  $a$ , stores the expression  $e'$  at the port  $a$ , and stores the port  $a$  at the reference location  $\ell$ . We use  $S[\ell \leftarrow a]$  as a shorthand for  $\{\ell' \mapsto S(\ell') \mid \ell' \in \text{dom}(S) \wedge \ell' \neq \ell\} \cup \{\ell \mapsto a\}$ . Note that we use the dummy witness symbol  $\bullet$  as the run-time representation of any witness because only its presence is important to the semantics, i.e., operationally, a witness is like a dataflow token in dataflow machines. The rule **Read** reads from the reference location  $\ell$  and, as noted above, returns the value paired with a witness. The rule **Ref** creates a fresh location  $\ell$  and a fresh port  $a$ , initializes  $a$  to the expression  $e$  and  $\ell$  to the port  $a$ . The rule **Join** takes two witnesses and returns one witness.

The rule **Arrive** might look somewhat unfamiliar. Here  $V$  is the set of “safe to duplicate” expressions. Partly for the sake of the static-checking algorithm to be presented later, we fix  $V$  to *values* generated by the following grammar:

$$v ::= x \mid i \mid a \mid \bullet \mid \ell \mid v \otimes v' \mid \lambda x. e$$

<sup>2</sup>The results in Section 3 are general enough for most other functional transformations too, but the static checking algorithm requires a more stringent definition.

<sup>3</sup>In the literature, top-level `let`-bound variables often double as variables and ports.

**Arrive** says that if  $e$  is safe to duplicate, then we can replace  $a$  by  $e$ ; we say a safe to duplicate expression has *arrived* at port  $a$ . In essence, while standard operational semantics for functional languages [12, 8] implicitly combine **Arrive** with other rules, we separate **Arrive** for increased freedom in the evaluation order. Lastly, the rule **GC** garbage-collects unreachable (from the root port  $\diamond$ ) portions of the expression store.

Here is an example of a  $\lambda_{wit}$  evaluation:

$$\begin{aligned}
& (\emptyset, \{\diamond \mapsto (\lambda x. \text{read } x \bullet) \text{ ref } 1\}) \\
& \Rightarrow (\{\ell \mapsto a\}, \{\diamond \mapsto (\lambda x. \text{read } x \bullet) \ell, a \mapsto 1\}) & \mathbf{Ref} \\
& \Rightarrow (\{\ell \mapsto a\}, \{\diamond \mapsto \text{read } a' \bullet, a \mapsto 1, a' \mapsto \ell\}) & \mathbf{App} \\
& \Rightarrow (\{\ell \mapsto a\}, \{\diamond \mapsto \text{read } \ell \bullet, a \mapsto 1, a' \mapsto \ell\}) & \mathbf{Arrive} \\
& \Rightarrow (\{\ell \mapsto a\}, \{\diamond \mapsto a \otimes \bullet, a \mapsto 1, a' \mapsto \ell\}) & \mathbf{Read} \\
& \Rightarrow (\{\ell \mapsto a\}, \{\diamond \mapsto 1 \otimes \bullet, a \mapsto 1, a' \mapsto \ell\}) & \mathbf{Arrive} \\
& \Rightarrow (\{\ell \mapsto a\}, \{\diamond \mapsto 1 \otimes \bullet\}) & \mathbf{GC}
\end{aligned}$$

The semantics is non-deterministic and therefore also allows other reduction sequences for the same program. For example, we may take an **App** step immediately instead of first creating a new reference by a **Ref** step:

$$\begin{aligned}
& (\emptyset, \{\diamond \mapsto (\lambda x. \text{read } x \bullet) \text{ ref } 1\}) \\
& \Rightarrow (\emptyset, \{\diamond \mapsto \text{read } a \bullet, a \mapsto \text{ref } 1\}) & \mathbf{App}
\end{aligned}$$

Before defining confluence, we point out several important properties of this semantics. Firstly, note that the evaluation contexts  $E$  do not extend to subexpressions of a  $\lambda$  abstraction  $\lambda x.e$ , i.e., we do not reduce under  $\lambda$ . The evaluation contexts also do not extend to subexpressions of an expression  $\text{let } x = e \text{ in } e'$ , but  $e$  and  $e'$  may become available for evaluation via applications of the **Let** rule. As with call-by-value evaluation or call-by-need evaluation, evaluation of an expression is shared. For example, in the program  $(\lambda x.x \otimes x) e$ , the expression  $e$  is evaluated at most once.

The semantics of  $\lambda_{wit}$  has strictly more freedom in evaluation order than both call-by-value and call-by-need. In particular, call-by-need evaluation can be obtained by using the same reduction rules but restricting the evaluation contexts to the following

$$\begin{aligned}
E & := D \cup \{\diamond \mapsto E\} \mid [] \mid E e \mid \pi_i(E) \mid \text{write } E e e' \\
& \mid \text{write } \ell e E \mid \text{read } E e \mid \text{read } \ell E \\
& \mid \text{join } E e \mid \text{join } \bullet E
\end{aligned}$$

Call-by-value evaluation can be obtained by adding the following contexts to the evaluation contexts of the call-by-need evaluation

$$\begin{aligned}
E & := \dots \mid (\lambda x.e) E \mid E \otimes e \mid v \otimes E \mid \text{write } \ell E \bullet \\
& \mid \text{ref } E \mid \text{let } x = E \text{ in } e
\end{aligned}$$

in addition to restricting the rule **App** to the case  $e' \in V$ , the rule **Let** to the case  $e \in V$ , the rule **Pair** to the case  $e_1, e_2 \in V$ , the rule **Write** to the case  $e' \in V$ , and the rule **Ref** to the case  $e \in V$ .<sup>4</sup> Note that both lazy writes and strict writes are possible in  $\lambda_{wit}$ .

It is important to understand that we are only concerned with side-effects via references, and hence we are not concerned about issues like the number of ports that are created during an evaluation.

Having defined the semantics, we can formally define when a  $\lambda_{wit}$  program is confluent. To this end, we define *observational equivalence* as the smallest reflexive and transitive relation  $D \approx D'$  on expression stores satisfying:

- $D \approx D[a/a']$  where  $a \notin \text{free}(D)$
- $D \approx D[\ell/\ell']$  where  $\ell \notin \text{free}(D)$

That is, expression stores are observationally equivalent if they are equivalent up to consistent renaming of free ports and reference locations. Let  $\Rightarrow^*$  be a sequence of zero or more  $\Rightarrow$ .

<sup>4</sup> Strictly speaking, the context  $\text{let } x = E \text{ in } e$  is not in the semantics of  $\lambda_{wit}$ . But  $\lambda_{wit}$  can simulate the behavior via a **Let** step and then reducing  $e[a/x]$  which is now in an evaluation context.

**DEFINITION 1 (Confluence).** A program state  $(S, D)$  is *confluent* if for any two states  $(S_1, D_1)$  and  $(S_2, D_2)$  such that  $(S, D) \Rightarrow^* (S_1, D_1)$  and  $(S, D) \Rightarrow^* (S_2, D_2)$ , there exist two states  $(S'_1, D'_1)$  and  $(S'_2, D'_2)$  such that  $(S_1, D_1) \Rightarrow^* (S'_1, D'_1)$ ,  $(S_2, D_2) \Rightarrow^* (S'_2, D'_2)$ , and  $D'_1 \approx D'_2$ . A program  $e$  is *confluent* if its initial state  $(\emptyset, \{\diamond \mapsto e\})$  is confluent.

Note that the definition does not require any relation between reference location stores  $S'_1$  and  $S'_2$ . So, for example, a program that writes but never reads would be confluent. As shown before,  $\lambda_{wit}$  contains programs that are not confluent. Indeed, the difference between call-by-need and call-by-value is enough to demonstrate non-confluence:

$$(\lambda x. \text{read } x \bullet) (\text{let } x = (\text{ref } 1) \text{ in let } y = (\text{write } x \ 2 \bullet) \text{ in } x)$$

The above program evaluates to the pair  $1 \otimes \bullet$  under call-by-need and to the pair  $2 \otimes \bullet$  under call-by-value. No further reductions can make the two states observationally equivalent. (Here we implicitly read back the top-level expression from the root port instead of showing the actual expression stores for brevity.)

We have shown earlier that witnesses can aid in writing correct programs by directly ordering side-effects. Witnesses are first class values and hence can be treated like other expressions. For example, the program below captures a witness in a function which itself returns a witness to ensure that reads and writes happen in a correct order:

$$\begin{aligned}
& \text{let } x = \text{ref } 1 \text{ in} \\
& \quad \text{let } w = \text{write } x \ 2 \bullet \text{ in} \\
& \quad \quad \text{let } f = \lambda y. \text{read } x \ w \text{ in} \\
& \quad \quad \quad \text{let } z = (f \ 0) \otimes (f \ 0) \text{ in} \\
& \quad \quad \quad \quad \text{let } w = \text{write } x \ 3 \ \text{join } \pi_2(\pi_1(z)) \ \pi_2(\pi_2(z)) \text{ in } z
\end{aligned}$$

Note that a witness of the first write is captured in the function  $f$ . Hence both reads from the two calls to  $f$  see a witness of the write. A witness of each read is returned by  $f$ , and the last write waits until it sees witnesses from both reads. Therefore, the result of the program is  $(2 \otimes \bullet) \otimes (2 \otimes \bullet)$  regardless of the evaluation order. Note that the two calls to  $f$ , and thus the reads in the calls, can occur in either order.

### 3. Witness Race Freedom

As discussed in Section 2, witnesses aid in writing correct programs in the presence of side-effects but do not enforce correctness. In this section, we give a sufficient condition for guaranteeing confluence.

The guideline for writing a correct program should be intuitively clear at this point: we ensure that reads and writes happen in some race-free order by partially ordering them via witnesses. We now make this intuition more precise. First, we formally define what we mean by the phrase “side-effect  $A$  sees a witness of side effect  $B$ ” that we have used informally up to this point.

Intuitively, a *trace graph* is a program trace with all information other than reads, writes, and witnesses elided. There are three kinds of nodes: read nodes  $\text{read}(\ell)$ , write nodes  $\text{write}(\ell)$ , and the join node  $\text{join}$ . Read and write nodes are parametrized by a reference location  $\ell$ . There is a directed edge  $(A, B)$  from node  $A$  to node  $B$  if  $B$  directly sees a witness of  $A$ . A trace graph  $(\mathbb{V}, \mathbb{E})$  is

constructed as the program evaluates a modified semantics:

- Write**  $(S, E[\text{write } \ell e A]) \Rightarrow (S[\ell \leftarrow a], E[B] \uplus \{a \mapsto e\})$   
 $\mathbb{V} := \mathbb{V} \cup \{B\}$  where  $B$  is a new  $\text{write}(\ell)$  node  
 $\mathbb{E} := \mathbb{E} \cup \{(A, B)\}$
- Read**  $(S, E[\text{read } \ell A]) \Rightarrow (S, E[S(\ell) \otimes B])$   
 $\mathbb{V} := \mathbb{V} \cup \{B\}$  where  $B$  is a new  $\text{read}(\ell)$  node  
 $\mathbb{E} := \mathbb{E} \cup \{(A, B)\}$
- Join**  $(S, E[\text{join } A B]) \Rightarrow (S, E[C])$   
 $\mathbb{V} := \mathbb{V} \cup \{C\}$  where  $C$  is a new  $\text{join}$  node  
 $\mathbb{E} := \mathbb{E} \cup \{(A, C), (B, C)\}$

Note that we now use nodes as witnesses instead of  $\bullet$ . The line below each reduction rule shows the graph update action corresponding to that rule. The other rules remain unmodified and hence have no graph update actions. An evaluation starts with  $\mathbb{V} = \mathbb{E} = \emptyset$  and performs the corresponding graph update when taking a **Write** step, a **Read** step, or a **Join** step. Notice that a trace graph and the annotated semantics are only needed to state the semantic condition for correctness and are not needed in the actual execution of a  $\lambda_{wit}$  program.

We can now define what it means for a node  $A$  to see a witness of a node  $B$ , a notion we have used informally until now.

**DEFINITION 2.** Given a trace graph, we say that a node  $A$  sees a witness of a node  $B$  if there is a path from  $B$  to  $A$  in the trace graph. We write  $B \rightsquigarrow A$ .

The following is a trivial observation:

**THEOREM 1.** If  $B \rightsquigarrow A$  in a trace graph, then the side effect corresponding to  $B$  must have happened before the side effect corresponding to  $A$  in the evaluation that generated the trace graph.

Clearly, any trace graph is acyclic.

Having defined trace graphs and the  $\rightsquigarrow$  relation, we are now ready to state the semantic condition for correctness. Before showing the condition formally, we informally motivate it by making analogies to the conventional programming guideline for writing correct concurrent programs: prevent race conditions.

We first note that a program could produce different trace graphs depending on the choice of reductions, even when those trace graphs are from terminating evaluations. Furthermore, it is not necessarily the case that such a program is non-confluent. Therefore, instead of trying to argue about confluence by comparing different trace graphs, we shall define a condition that can be checked by observing each individual trace graph in isolation.

Let us write  $A : \text{nodetype}$  as a shorthand for a node  $A$  of type  $\text{nodetype}$ . If we have  $A : \text{read}(\ell)$  and  $B : \text{write}(\ell)$ , then we want either  $A \rightsquigarrow B$  or  $B \rightsquigarrow A$  to ensure that  $A$  always happens before  $B$  or  $B$  always happens before  $A$  because otherwise we may get a read-write race condition due to non-determinism. Also, for any  $A : \text{read}(\ell)$ , if there are two nodes  $B_1, B_2 : \text{write}(\ell)$  such that neither  $B_1 \rightsquigarrow B_2$  nor  $B_2 \rightsquigarrow B_1$  (so we do not know which write occurs first) and  $A$  could happen after both  $B_1$  and  $B_2$ , then we want  $C : \text{write}(\ell)$  such that  $C \rightsquigarrow A$ ,  $B_1 \rightsquigarrow C$  and  $B_2 \rightsquigarrow C$ , because otherwise the read at  $A$  might depend on whether the evaluation chose to do  $B_1$  first or  $B_2$  first, i.e., we have another kind of race-condition. Perhaps somewhat surprisingly, satisfying these two conditions turns out to be sufficient to ensure confluence.

We now formalizes this discussion. For  $B : \text{write}(\ell)$ , we use the shorthand  $B \rightsquigarrow^! A$  if for any  $C : \text{write}(\ell)$  such that  $C \rightsquigarrow A$  and  $C \neq B$ , we have  $C \rightsquigarrow B$ . Now, for any  $A : \text{read}(\ell)$ , there exists at most one  $B : \text{write}(\ell)$  such that  $B \rightsquigarrow^! A$ . Note that the second condition above is equivalent to requiring that for any  $A : \text{read}(\ell)$ , either there is no  $B : \text{write}(\ell)$  such that  $B \rightsquigarrow A$  or there is a  $B : \text{write}(\ell)$  such that  $B \rightsquigarrow^! A$ .

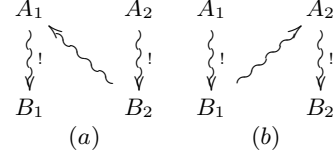


Figure 3: Possible orderings between pairs of reads and writes in a witness race free trace graph.

**DEFINITION 3 (Witness Race Freedom).** We say that a trace graph  $(\mathbb{V}, \mathbb{E})$  is witness race free if for every location  $\ell$ ,

- for every  $A : \text{read}(\ell) \in \mathbb{V}$  and  $B : \text{write}(\ell) \in \mathbb{V}$ , either  $A \rightsquigarrow B$  or  $B \rightsquigarrow A$ , and
- for every  $A : \text{read}(\ell) \in \mathbb{V}$ , either there is no  $B : \text{write}(\ell) \in \mathbb{V}$  such that  $B \rightsquigarrow A$  or there is a  $B : \text{write}(\ell) \in \mathbb{V}$  such that  $B \rightsquigarrow^! A$ .

We say that a program  $e$  is witness race free if every trace graph of  $e$  is witness race free.

**THEOREM 2.** If  $e$  is witness race free then  $e$  is confluent.

The proof appears in our companion technical report [13].

While witness race freedom is a sufficient condition, it is not necessary. For example, if for each reference location writes happen to never change the location's value, then the program is trivially confluent regardless of the order of reads and writes. Another example is a program using implicit order of evaluation: e.g., in  $\lambda_{wit}$ , expressions are not reduced under  $\lambda$  so a function body is evaluated only after a call. Hence a program that stores a function in a reference location, reads the reference location to call the function, and then writes in the same reference location from the body of the called function is confluent because the write always happens after the read despite the write not seeing the witness of the read.

Nevertheless, witness race freedom is “almost complete” in a sense that if the only way to order two side-effects is to make one see a witness of the other, and if we cannot assume anything about what expressions are written and how the contents are used, then it is the weakest condition guaranteeing correctness. In particular, if the trace graphs are the only information available about a program, then witness race freedom becomes a necessary condition.

The result in this section can be extended to most other functional program transformations, because witness race freedom is an entirely semantic condition. However, the static checking algorithm described in Section 4 is not as forgiving, which is why we have restricted the set of program transformations to that of the semantics of  $\lambda_{wit}$ . For example, the algorithm is unsound for general call-by-name reduction.

## 4. Types for Statically Checking Witness Race Freedom

While the concept of witnesses is straightforward, it may nevertheless be desirable to have an automated way of checking whether an arbitrary  $\lambda_{wit}$  program is witness race free. Witness race freedom may be checked directly by checking every program trace, which is computationally infeasible. Instead, we exploit a special property of witness-race-free trace graphs to design a sound algorithm that can efficiently verify a large subset of witness-race-free  $\lambda_{wit}$  programs.

The key observation is that any witness-race-free trace graph contains for each reference location  $\ell$  a subgraph that we shall call



Figure 4: A read-write pipeline with bottlenecks for a reference location  $\ell$ .

a *read-write pipeline with bottlenecks*. We shall design an algorithm that detects these subgraphs instead of directly checking the witness race freedom condition. Consider a witness-race-free trace graph. Suppose there are  $A_1, A_2 : \text{write}(\ell)$  and  $B_1, B_2 : \text{read}(\ell)$  such that  $A_1 \neq A_2$ ,  $A_1 \rightsquigarrow B_1$  and  $A_2 \rightsquigarrow B_2$ . Due to witness race freedom, it must be the case that  $B_2 \rightsquigarrow A_1$  or  $A_1 \rightsquigarrow B_2$ . If the former is the case, we have the relation as depicted in Figure 3 (a). Suppose that the latter is the case. Then, since  $A_2 \rightsquigarrow B_2$ , it must be the case that  $A_1 \rightsquigarrow A_2$ . Consider  $A_2$  and  $B_1$ . Due to witness race freedom again, it must be the case that either  $A_2 \rightsquigarrow B_1$  or  $B_1 \rightsquigarrow A_2$ . But if  $A_2 \rightsquigarrow B_1$ , then since  $A_1 \rightsquigarrow B_1$ , it must be the case that  $A_2 \rightsquigarrow A_1$ . But this is impossible since  $A_2 \rightsquigarrow A_1 \rightsquigarrow A_2$  forms a cycle. So it must be the case that  $B_1 \rightsquigarrow A_2$ , and we have the relation as depicted in Figure 3 (b).

Further reasoning along this line of thought reveals that for a witness-race-free trace graph, for any reference location  $\ell$ , the nodes in the set  $X = \{A : \text{write}(\ell) \mid \exists B : \text{read}(\ell). A \rightsquigarrow B\}$  are totally ordered (with  $\rightsquigarrow$  as the ordering relation), and that these nodes partition all  $\text{read}(\ell)$  nodes and  $\text{write}(\ell)$  nodes in a way depicted in Figure 4 where  $X = \{A_1, \dots, A_n\}$ . In the figure, each  $\mathbb{R}_i$  and  $\mathbb{W}_i$  is a collection of nodes. No  $\mathbb{R}_i$  contains a  $\text{write}(\ell)$  node and no  $\mathbb{W}_i$  contains a  $\text{read}(\ell)$  nodes. Each  $A_i$  is one  $\text{write}(\ell)$  node. An arrow from  $X$  to  $Y$  means that there is a path from each  $\text{write}(\ell)$  node or  $\text{read}(\ell)$  node in  $X$  to each  $\text{write}(\ell)$  node or  $\text{read}(\ell)$  node in  $Y$ , except if a  $\mathbb{W}_i$  contains no such node, then there is a path from each  $\text{read}(\ell) \in \mathbb{R}_i$  to  $A_i$ . Each  $\mathbb{R}_i$  for  $i \neq 1$  must contain at least one  $\text{read}(\ell)$ . Arrows just imply the presence of paths, and hence there can be more paths than the ones implied by the arrows, e.g., paths to/from nodes that are not in the diagram, paths to and from nodes in the same collection, and even paths relating the collections in the diagram such as one that goes directly from  $\mathbb{R}_i$  to  $A_i$ , bypassing  $\mathbb{W}_i$ .

The graph in Figure 4 can be described formally as a subgraph of the trace graph satisfying certain properties.

**DEFINITION 4.** Given a trace graph  $(\mathbb{V}, \mathbb{E})$  and a reference location  $\ell$ , we call its subgraph  $G_\ell$  a *read-write pipeline with bottlenecks* if  $G_\ell$  consists of collections of nodes  $\mathbb{R}_1, \mathbb{R}_2, \dots, \mathbb{R}_n$  and  $\mathbb{W}_1, \mathbb{W}_2, \dots, \mathbb{W}_n$  with the following properties:

- $\{A \mid A : \text{read}(\ell) \in \mathbb{V}\} \subseteq \bigcup_{i=1}^n \mathbb{R}_i$ ,
- $\{A \mid A : \text{write}(\ell) \in \mathbb{V}\} \subseteq \bigcup_{i=1}^n \mathbb{W}_i$ ,
- $\mathbb{R}_1, \dots, \mathbb{R}_n, \mathbb{W}_1, \dots, \mathbb{W}_n$ , restricted to  $\text{write}(\ell)$  nodes and  $\text{read}(\ell)$  nodes are pairwise disjoint,
- for each  $A : \text{read}(\ell) \in \mathbb{R}_i$  and  $B : \text{write}(\ell) \in \mathbb{W}_i$ ,  $A \rightsquigarrow B$ ,
- for each  $\mathbb{R}_i$  such that  $i \neq 1$ , there exists at least one  $A : \text{read}(\ell) \in \mathbb{R}_i$ , and
- there exists  $A : \text{write}(\ell) \in \mathbb{W}_i$  for all  $i \neq n$  such that for all  $B : \text{read}(\ell) \in \mathbb{R}_{i+1}$  and all  $C : \text{write}(\ell) \in \mathbb{W}_i$ ,  $A \rightsquigarrow B$  and  $C \rightsquigarrow A$ .

Note that each collection  $\mathbb{R}_i$  and  $\mathbb{W}_i$  corresponds to the collection of nodes marked by the same name in Figure 4 but with each

$$e ::= x \mid i \mid \lambda x.e \mid e e' \mid \text{let } x = e \text{ in } e' \mid e \otimes e' \mid \pi_i(e) \mid \text{write } e_1 e_2 e_3 \mid \text{read } e e' \mid \text{ref } e e' \mid \text{join } e e' \mid \bullet \mid \text{letreg } x e$$

Figure 5: The syntax of  $\lambda_{wit}^{reg}$ .

$$E ::= D \cup \{a \mapsto E\} \mid [] \mid E e \mid e E \mid E \otimes e \mid e \otimes E \mid \pi_i(E) \mid \text{write } E e e' \mid \text{write } e E e' \mid \text{write } e e' E \mid \text{read } e E \mid \text{read } E e \mid \text{ref } E e \mid \text{ref } e E \mid \text{join } E e \mid \text{join } e E$$

$$\begin{array}{ll} \text{App} & (R, S, E[(\lambda x.e) e']) \Rightarrow (R, S, E[e[a/x]] \uplus \{a \mapsto e'\}) \\ \text{Let} & (R, S, E[\text{let } x = e \text{ in } e']) \Rightarrow (R, S, E[e'[a/x]] \uplus \{a \mapsto e[a/x]\}) \\ \text{Pair} & (R, S, E[\pi_i(e_1 \otimes e_2)]) \Rightarrow (R, S, E[e_i]) \\ \text{Write} & (R, S, E[\text{write } \ell e \bullet]) \Rightarrow (R, S[\ell \leftarrow a], E[\bullet] \uplus \{a \mapsto e\}) \\ \text{Read} & (R, S, E[\text{read } \ell \bullet]) \Rightarrow (R, S, E[S(\ell) \otimes \bullet]) \\ \text{Ref} & (R, S, E[\text{ref } e r]) \Rightarrow (R, S \uplus \{\ell \mapsto a\}, E[\ell] \uplus \{a \mapsto e\}) \\ \text{Join} & (R, S, E[\text{join } \bullet \bullet]) \Rightarrow (R, S, E[\bullet]) \\ \text{LetReg} & (R, S, E[\text{letreg } x e]) \Rightarrow (R \uplus \{r\}, S, E[e[a/x]] \uplus \{a \mapsto r\}) \\ \text{Arrive} & (R, S, E[a] \uplus \{a \mapsto e\}) \Rightarrow (R, S, E[e] \uplus \{a \mapsto e\}) \\ & \text{where } e \in V \\ \text{GC} & (R, S, D \uplus D') \Rightarrow (R, S, D) \\ & \text{where } \diamond \notin \text{dom}(D') \wedge \text{dom}(D') \cap \text{free}(D) = \emptyset \end{array}$$

Figure 6: The semantics of  $\lambda_{wit}^{reg}$ .

node  $A_i$  included in the collection  $\mathbb{W}_i$ . The “bottlenecks” are the  $A_i$ ’s. Note that a trace graph  $(\mathbb{V}, \mathbb{E})$  actually contains a read-write pipeline with bottlenecks per each reference location  $\ell$  as a subgraph  $G_\ell$  (but the subgraphs may not be disjoint because the paths may involve other locations and share join nodes).

The following theorem formalizes our earlier informal discussion.

**THEOREM 3.** A trace graph  $(\mathbb{V}, \mathbb{E})$  is witness race free if and only if it has a read-write pipeline with bottlenecks for every  $\ell$ .

The proof appears in our companion technical report [13]. The following is an immediate consequence:

**COROLLARY 1.** A  $\lambda_{wit}$  program  $e$  is witness race free if and only if every trace graph of  $e$  has a read-write pipeline with bottlenecks for every reference location  $\ell$ .

#### 4.1 Regions

Corollary 1 reduces the problem of deciding whether a program  $e$  is witness race free to the problem of deciding if every trace graph of  $e$  has a read-write pipeline with bottlenecks for every reference location  $\ell$ . Therefore it suffices to design an algorithm for solving the latter problem. But before we do so, we make a slight change to  $\lambda_{wit}$  to make the problem more tractable. In  $\lambda_{wit}$ , there is a read-write pipeline with bottlenecks for each reference location  $\ell$ , but distinguishing dynamically allocated reference locations individually is difficult for a compile-time algorithm. Therefore, we add *regions* to the language so that programs can explicitly group reference locations that are to be tracked together.

Figure 5 shows  $\lambda_{wit}^{reg}$ ,  $\lambda_{wit}$  extended with regions. The syntax contains two new expression kinds:  $\text{letreg } x e$  which creates a

new region and  $\text{ref } e e'$  which places the newly created reference in region  $e'$ ;  $\text{ref } e e'$  replaces  $\text{ref } e$ . Figure 6 gives the semantics of  $\lambda_{wit}^{reg}$  which differs from  $\lambda_{wit}$  in two small ways. First, a state now contains a set of regions  $R$ . We use symbols  $r, r', r_i$ , etc. to denote regions. Regions are safe to duplicated, i.e.,  $r \in V$ . The  $R$ 's are used only for ensuring that the newly created region  $r$  at a **LetReg** step is fresh. (We overload the symbol  $\uplus$  such that  $R \uplus R' = R \cup R'$  if  $R \cap R' = \emptyset$  and, is undefined otherwise.) Note that evaluation contexts  $E$  do not extend to the subexpressions of  $\text{letreg } x e$ . The second difference is that a **Ref** step now takes a region  $r$  along with the initializer  $e$  to indicate that the newly created reference location  $\ell$  belongs to the region  $r$ . Note that the semantics does not actually associate the reference location  $\ell$  and the region  $r$ , and therefore grouping of reference locations via regions is entirely conceptual.<sup>5</sup>

Regions force programs to abide by witness race freedom at the granularity of regions instead of at the granularity of individual reference locations. That is, instead of  $\text{read}(\ell)$  nodes and  $\text{write}(\ell)$  nodes, we use  $\text{read}(r)$  nodes and  $\text{write}(r)$  nodes. Formally, a trace graph for  $\lambda_{wit}^{reg}$  is constructed by the following graph construction semantics:

$$\begin{array}{l}
\mathbf{Write} \quad (R, K, S, E[\text{write } \ell e A]) \\
\quad \Rightarrow (R, K, S[\ell \leftarrow a], E[B] \uplus \{a \mapsto e\}) \\
\quad \mathbb{V} := \mathbb{V} \cup \{B\} \text{ where } B \text{ is a new } \text{write}(K(\ell)) \text{ node} \\
\quad \mathbb{E} := \mathbb{E} \cup \{(A, B)\} \\
\mathbf{Read} \quad (R, K, S, E[\text{read } \ell A]) \Rightarrow (R, K, S, E[S(\ell) \otimes B]) \\
\quad \mathbb{V} := \mathbb{V} \cup \{B\} \text{ where } B \text{ is a new } \text{read}(K(\ell)) \text{ node} \\
\quad \mathbb{E} := \mathbb{E} \cup \{(A, B)\} \\
\mathbf{Join} \quad (R, K, S, E[\text{join } A B]) \Rightarrow (R, K, S, E[C]) \\
\quad \mathbb{V} := \mathbb{V} \cup \{C\} \text{ where } C \text{ is a new } \text{join} \text{ node} \\
\quad \mathbb{E} := \mathbb{E} \cup \{(A, C), (B, C)\} \\
\mathbf{Ref} \quad (R, K, S, E[\text{ref } e r]) \\
\quad \Rightarrow (R, K \uplus \{\ell \mapsto r\}, S \uplus \{\ell \mapsto a\}, E[\ell] \uplus \{a \mapsto e\})
\end{array}$$

Note that these reduction rules use an additional function  $K$  which is a mapping from reference locations to regions. The mapping  $K$  starts empty at the beginning of evaluation. Other reductions rules are unmodified except that the function  $K$  is passed from left to right in the obvious way.

Since there is less information available in a  $\lambda_{wit}^{reg}$  trace graph than in a  $\lambda_{wit}$  trace graph, the witness race freedom condition is more conservative for  $\lambda_{wit}^{reg}$ . That is, we still need the condition that for any  $A : \text{write}(r)$  and  $B : \text{read}(r)$ , either  $A \rightsquigarrow B$  or  $B \rightsquigarrow A$ . But we need to tighten the second condition so that for any  $A : \text{read}(r)$  if there are  $B_1, B_2 : \text{write}(r)$  such that  $B_1 \rightsquigarrow A$  and  $B_2 \rightsquigarrow A$ , then either  $B_1 \rightsquigarrow B_2$  or  $B_2 \rightsquigarrow B_1$ . This condition is strictly more conservative than for  $\lambda_{wit}$ , which only requires some  $C : \text{write}(r)$  such that  $C \rightsquigarrow A$  in such a situation. The reason for this conservativeness is that we do not know from a trace graph of  $\lambda_{wit}^{reg}$  whether  $B_1$  and  $B_2$  both write to the same reference location.

Formally, witness race freedom for  $\lambda_{wit}^{reg}$  can be defined as follows.

**DEFINITION 5 (Witness Race Freedom for  $\lambda_{wit}^{reg}$ ).** *We say that a  $\lambda_{wit}^{reg}$  trace graph  $(\mathbb{V}, \mathbb{E})$  is witness race free if for every region  $r$ ,*

- for every  $A : \text{read}(r) \in \mathbb{V}$  and  $B : \text{read}(r) \in \mathbb{V}$ , either  $A \rightsquigarrow B$  or  $B \rightsquigarrow A$ , and
- for every  $A : \text{read}(r) \in \mathbb{V}$  and  $B_1, B_2 : \text{write}(r) \in \mathbb{V}$  such that  $B_1 \rightsquigarrow A$  and  $B_2 \rightsquigarrow A$ , we have  $B_1 \rightsquigarrow B_2$  or  $B_2 \rightsquigarrow B_1$ .

**THEOREM 4.** *If a  $\lambda_{wit}^{reg}$  program  $e$  is witness race free, then  $e$  is confluent.*

<sup>5</sup>Regions are traditionally coupled with some semantic meaning such as memory management [14, 5]. It is possible to extend  $\lambda_{wit}^{reg}$  to do similar things with its regions.

**Proof (Sketch):** For any evaluation of  $e$ , carry out the same reduction sequence with the trace graph building action of  $\lambda_{wit}$ , i.e., the trace graph  $G$  generated is at the granularity of reference locations. Then it is easy to see that if  $G$  satisfies the above two conditions,  $G$  also satisfies the two conditions of Theorem 2.  $\square$

It is easy to see that Definition 5 is the weakest possible restriction to the original witness race freedom under the region abstraction because for any  $\lambda_{wit}^{reg}$  trace graph that is not witness race free, one can easily find a non-confluent program that produces the graph.

In a witness-race-free trace graph for  $\lambda_{wit}^{reg}$ , the read-write pipeline with bottlenecks for a region  $r$  consisting of the collections  $(\mathbb{R}_1, \dots, \mathbb{R}_n, \mathbb{W}_1, \dots, \mathbb{W}_n)$  has the following property: each set  $\{A \mid A : \text{write}(r) \in \mathbb{W}_i\}$  for  $i \neq n$  can be totally ordered (with  $\rightsquigarrow$  as the ordering relation). The theorem below is immediate from Corollary 1 under this additional property.

**THEOREM 5.** *A  $\lambda_{wit}^{reg}$  program  $e$  is witness race free if and only if every trace graph of  $e$  has a read-write pipeline with bottlenecks for every region  $r$ .*

This additional property helps in designing a static checking algorithm.

## 4.2 From Network Flow to Types

Now our goal is to design an algorithm for statically checking if every trace graph of a  $\lambda_{wit}^{reg}$  program  $e$  has a read-write pipeline with bottlenecks for every region  $r$ . Our approach exploits a *network flow* property of read-write pipelines with bottlenecks. Consider a trace graph as a network of nodes with each edge  $(A, B)$  able to carry any non-negative flow from  $A$  to  $B$ . (Recall edges are directed.) As usual with network flow, we require that the total incoming flow equal the total outgoing flow for every node in the graph. Now, let us add a *virtual source* node  $A_S$  and connect it to every node  $B$  by adding an edge  $(A_S, B)$ . We assign incoming flow 1 to  $A_S$ . Then it is not hard to see that if there exists a read-write pipeline with bottlenecks for the region  $r$  then there exists flow assignments such that every  $\text{read}(r)$  node and  $\text{write}(r)$  node gets a positive flow and every  $A : \text{write}(r) \in \mathbb{W}_i$  for  $i \neq n$  gets a flow equal to 1.

It turns out that the converse also holds. That is, given a trace graph, if there is a flow assignment such that each  $\text{read}(r)$  node and  $\text{write}(r)$  node gets a positive flow and each  $A : \text{write}(r)$  such that there exists  $B : \text{read}(r)$  with  $A \rightsquigarrow B$  gets a flow equal to 1, then there is a read-write pipeline with bottlenecks for the region  $r$ . By Theorem 5, this implies that there exists such a flow assignment for every region  $r$  if and only if the trace graph is witness race free. Because edges in a trace graph are traces of witnesses, our idea is to assign a type to a witness such that the type contains flow assignments for each (static) region. We use this idea to design a type system such that a well-typed program is guaranteed to be witness race free.

Formally, a witness type  $W$  is a function from the set of *static region identifiers* **RegIDs** to rational numbers in the range  $[0, 1]$ , i.e.,  $W : \mathbf{RegIDs} \rightarrow [0, 1]$ . The rational number  $W(\rho)$  indicates the flow amount for the static region identifier  $\rho$  in the witness type  $W$ . We use the notation  $\{\rho_1 \mapsto q_1, \dots, \rho_n \mapsto q_n\}$  to mean a witness type  $W$  such that  $W(\rho) = q_i$  if  $\rho = \rho_i$  for some  $1 \leq i \leq n$  and  $W(\rho) = 0$  otherwise. (We use the symbols  $q, q_i, q'$ , etc for non-negative rational numbers, including those larger than 1.)

The rest of the types are defined in Figure 7, including types include integer types  $\text{int}$ , function types  $\tau \xrightarrow{q} \tau'$ , pair types  $\tau \otimes \tau'$ , reference types  $\text{ref}(\tau, \tau', \rho)$ , and region types  $\text{reg}(\rho)$ . The non-negative rational number  $q$  in  $\tau \xrightarrow{q} \tau'$  represents the number of times the function can be called. We allow the symbols  $q, q'$ , etc

$$\tau := \text{int} \mid \tau \stackrel{q}{\rightarrow} \tau' \mid \tau \otimes \tau' \mid \text{ref}(\tau, \tau', \rho) \mid \text{reg}(\rho) \mid W$$

Figure 7: The type language.

$$\begin{array}{c}
\frac{\Gamma; W \vdash e : \tau \quad \tau \geq \tau'}{\Gamma; W \vdash e : \tau'} \text{Sub} \quad \frac{\Gamma(x) = \tau}{\Gamma; W \vdash x : \tau} \text{Var} \\
\frac{}{\Gamma; W \vdash i : \text{int}} \text{Int} \quad \frac{}{\Gamma; W \vdash \bullet : \emptyset} \text{Dummy} \\
\frac{\Gamma; W_1 \vdash e : W_2}{\Gamma; W_1 + W_3 \vdash e : W_2 + W_3} \text{Source} \\
\frac{\Gamma, x : \tau; W \vdash e : \tau'}{\Gamma \times q; W \times q \vdash \lambda x. e : \tau \stackrel{q}{\rightarrow} \tau'} \text{Abs} \\
\frac{\Gamma; W \vdash e : \tau \stackrel{q}{\rightarrow} \tau' \quad \Gamma'; W' \vdash e' : \tau \quad q \geq 1}{\Gamma + \Gamma'; W + W' \vdash e e' : \tau'} \text{App} \\
\frac{\Gamma; W \vdash e : \tau \quad \Gamma'; W' \vdash e' : \tau'}{\Gamma + \Gamma'; W + W' \vdash e \otimes e' : \tau \otimes \tau'} \text{Pair} \\
\frac{\Gamma; W \vdash e : \tau_1 \otimes \tau_2}{\Gamma; W \vdash \pi_i(e) : \tau_i} \text{Proj} \\
\frac{\Gamma; W \vdash e : \tau \quad \Gamma'; W' \vdash e' : \text{reg}(\rho)}{\Gamma + \Gamma'; W + W' \vdash \text{ref } e e' : \text{ref}(\tau, \tau', \rho)} \text{Ref} \\
\frac{\Gamma_1; W_1 \vdash e_1 : \text{ref}(\tau, \tau', \rho) \quad \Gamma_2; W_2 \vdash e_2 : \tau' \quad W(\rho) \geq 1}{\Gamma_1 + \Gamma_2 + \Gamma_3; W_1 + W_2 + W_3 \vdash \text{write } e_1 e_2 e_3 : W} \text{Write} \\
\frac{\Gamma; W_1 \vdash e : \text{ref}(\tau, \tau', \rho) \quad \Gamma'; W_2 \vdash e' : W \quad W(\rho) > 0}{\Gamma + \Gamma'; W_1 + W_2 \vdash \text{read } e e' : W \otimes \tau} \text{Read} \\
\frac{\Gamma, x : \text{reg}(\rho); W + \{\rho \mapsto q\} \vdash e : \tau \quad q \leq 1 \quad \rho \notin \text{free}(\Gamma) \cup \text{free}(W) \cup \text{free}(\tau)}{\Gamma; W \vdash \text{letreg } x e : \tau} \text{LetRegion} \\
\frac{\Gamma; W_1 \vdash e : W \quad \Gamma; W_2 \vdash e' : W'}{\Gamma + \Gamma'; W_1 + W_2 \vdash \text{join } e e' : W + W'} \text{Join} \\
\frac{\Gamma; W \vdash e'[(\text{let } x = e \text{ in } x)/x] : \tau \quad e \in V^+}{\Gamma; W \vdash \text{let } x = e \text{ in } e' : \tau} \text{LetA} \\
\frac{\Gamma, x : \tau; W \vdash e : \tau \quad \Gamma', x : \tau; W' \vdash e' : \tau' \quad \tau \geq \tau \times \infty \text{ if } x \in \text{free}(e)}{\Gamma + \Gamma'; W + W' \vdash \text{let } x = e \text{ in } e' : \tau'} \text{LetB}
\end{array}$$

Figure 8: Type judgment rules.

to take the valuation  $\infty$  to imply that the function can be called arbitrarily many times. We use the following arithmetic relation:  $q + \infty = \infty$ ,  $q \times \infty = \infty$  for  $q \neq 0$ , and  $0 \times \infty = 0$ .

Figure 8 shows the main type judgment rules. Our type system belongs to the family of substructural type systems, which includes linear types. We discuss the rules from top-to-bottom and left-to-

$$\begin{array}{l}
\text{reg}(\rho) + \text{reg}(\rho) = \text{reg}(\rho) \\
\text{int} + \text{int} = \text{int} \\
\tau \stackrel{q}{\rightarrow} \tau' + \tau \stackrel{q'}{\rightarrow} \tau' = \tau \stackrel{q+q'}{\rightarrow} \tau' \\
\tau_1 \otimes \tau_2 + \tau_3 \otimes \tau_4 = (\tau_1 + \tau_3) \otimes (\tau_2 + \tau_4) \\
\text{ref}(\tau_1, \tau, \rho) + \text{ref}(\tau_2, \tau, \rho) = \text{ref}(\tau_1 + \tau_2, \tau, \rho) \\
W + W' = \{\rho \mapsto W(\rho) + W'(\rho) \mid \rho \in \mathbf{RegIDs}\} \\
\text{reg}(\rho) \times q = \text{reg}(\rho) \\
\text{int} \times q = \text{int} \\
\tau \stackrel{q'}{\rightarrow} \tau' \times q = \tau \stackrel{q' \times q}{\rightarrow} \tau' \\
\tau \otimes \tau' \times q = (\tau \times q) \otimes (\tau' \times q) \\
\text{ref}(\tau, \tau', \rho) \times q = \text{ref}(\tau \times q, \tau', \rho) \\
W \times q = \{\rho \mapsto W(\rho) \times q \mid \rho \in \mathbf{RegIDs}\}
\end{array}$$

Figure 9: Arithmetic over types.

right, except for the rule **Sub** which we defer to the end. The rules **Var** and **Int** are standard. The rule **Dummy** gives a dummy witness  $\bullet$  an empty witness type; note that  $\emptyset(\rho) = 0$  for any static region identifier  $\rho$ .

The rule **Source** uses additive arithmetic over types defined in Figure 9. The rule adds  $W_3$  amount of flow from the virtual source nodes (i.e., nodes  $A_S$  from the first paragraph of this section) to  $W_2$ . In the type judgment relation  $\Gamma; W \vdash e : \tau$ , the witness type  $W$  represents the flow the expression  $e$  receives from the virtual source nodes. Therefore, the rule **Source** says that assuming that we took  $W_1$  flow from the virtual source nodes in the precondition, we are now taking  $W_3$  more.

In the rule **Abs**, we multiply the left hand side of the judgments by the number of times that the function can be used. Multiplication over type environments  $\Gamma$  is defined as follows:

$$(\Gamma, x : \tau) \times q = (\Gamma \times q), x : (\tau \times q)$$

So for example, if the  $\lambda$  abstraction  $\lambda x. e$  captures a witness as a free variable  $y$  and that  $\Gamma(y) = W$ , then  $(\Gamma \times q)(y) = W \times q$ . Thus if the function body requires  $W$  amount of flow in the witness, then we actually require  $W \times q$  amount of flow because the function may be called  $q$  times.

In the rule **App**, the precondition  $q \geq 1$  says that the number of times the function can be used must be at least 1. The left hand side of the two judgments in the precondition are added so that we can compute the combined flow required for the expressions  $e$  and  $e'$ . Addition over type environments is defined as follows:

$$(\Gamma, x : \tau) + (\Gamma', x : \tau') = (\Gamma + \Gamma'), x : (\tau + \tau')$$

The rules **Pair** and **Proj** are self-explanatory.

In a reference type  $\text{ref}(\tau, \tau', \rho)$ , the static region identifier  $\rho$  identifies the region where the reference belongs while the type  $\tau$  is the read type of the reference and the type  $\tau'$  is the write type of the reference. Initially the read and write types are the same as seen in the rule **Ref**. The rule **Write** matches the type of the to-be-assigned expression  $e_2$  with the write type of the reference while the rule **Read** uses the read type of the reference type. Note that we require  $W(\rho) \geq 1$  at **Write** and  $W(\rho) > 0$  at **Read**; both correspond to the flow requirement for writes and reads. The reason read-type/write-type separation is subtle. Consider the following expression where the expressions  $e_1$  and  $e_3$  are witnesses and the expression  $e_2$  is a region:

$$\text{let } x = (\text{ref } e_1 e_2) \text{ in let } w = (\text{write } x e_3 \bullet) \text{ in read } x w$$

Suppose we just have read types so that the type system uses read types at instances **Write** as well as at instances **Read**. Then the type system is unsound (even without **Sub**) for the following reason. The type system may assign some flow  $W$  to the occurrence of

$$\begin{array}{c}
\frac{}{\tau \geq \tau} \quad \frac{\tau_1 \leq \tau'_1 \quad \tau_2 \geq \tau'_2 \quad q \geq q'}{\tau_1 \xrightarrow{q} \tau_2 \geq \tau'_1 \xrightarrow{q'} \tau'_2} \\
\frac{\tau_1 \geq \tau'_1 \quad \tau_2 \geq \tau'_2}{\tau_1 \otimes \tau_2 \geq \tau'_1 \otimes \tau'_2} \\
\frac{\tau_1 \geq \tau'_1 \quad \tau_2 \leq \tau'_2}{\text{ref}(\tau_1, \tau_2, \rho) \geq \text{ref}(\tau'_1, \tau'_2, \rho)} \\
\frac{W(\rho) \geq W(\rho) \text{ for all } \rho \in \mathbf{RegIDs}}{W \geq W'}
\end{array}$$

Figure 10: Subtyping.

the variable  $x$  at the write and some flow  $W'$  to the occurrence of the variable  $x$  at the read. But there is no constraint to force  $W = W'$ , so the type system can let  $W' > W$  while keeping the sum  $W + W'$  fixed, i.e., we get more flow from a reference than what was assigned to the reference. Separating read and write types prevents this problem because addition and multiplication do not act on write types.

The rule **LetRegion** introduces a fresh static region identifier  $\rho$ . The witness type  $\{\rho \mapsto q\}$  represents the virtual source node for the new region. We constrain  $q \leq 1$  to ensure that we do not use more than 1 unit total from the source.

The rule **Join** combines two witnesses by adding their types.

There are two rules, **LetA** and **LetB**, for the expression kind  $\text{let } x = e \text{ in } e'$ . **LetA** is less conservative and should be used whenever  $x$  occurs more than once in  $e'$  and  $e \in V^+$  where  $V^+$  is the smallest set such that  $V^+ = V \cup \{\text{let } x = e \text{ in } x \mid e \in V^+\}$ . This rule corresponds to the usual substitution interpretation of let-based predicative polymorphism with the value restriction. **LetB** is used if  $e \notin V^+$  or  $x$  occurs at most once in  $e'$ . Here,  $\text{free}(\tau)$  is the set of static region identifiers in the type  $\tau$  where  $\text{free}(W) = \{\rho \mid W(\rho) \neq 0\}$ , and  $\text{free}(\Gamma) = \bigcup_{\tau \in \text{ran}(\Gamma)} \text{free}(\tau)$ .

Finally, we come return to **Sub**. The subtyping relation is defined in Figure 10. As usual, argument types of function types are contravariant. Write types of reference types are also contravariant; this treatment of reference subtyping is identical to that of a type-based formulation of Andersen's points-to analysis [4]. Intuitively, the rule **Sub** expresses the observation that the flow graph property may be relaxed so that the sum of the outgoing flow can be less than the sum of the incoming flow, i.e., if we could find a flow assignment satisfying the required flow constraints at reads and writes under this relaxed condition, then we still have a read-write pipeline with bottlenecks.

We say that a  $\lambda_{\text{wit}}^{\text{reg}}$  program  $e$  is well-typed if  $\emptyset; \emptyset \vdash e : \tau$  for some type  $\tau$ . The following theorem states that the type system is sound.

**THEOREM 6.** *If a  $\lambda_{\text{wit}}^{\text{reg}}$  program  $e$  is well-typed, then it is witness race free.*

The proof, which uses the network flow property, appears in the companion technical report [13].

We point out a few of the positive properties of this type system. If a program contains no reads or writes and can be typed by a standard Hindley-Milner polymorphic type system, then it can also be typed by our type system; for example, we may use the qualifier  $\infty$  for all function types and use 0 for all flows. In general, we can give the  $\infty$  qualifier to the function type of any function that does not capture a witness (directly or indirectly). We can also assign flow 0 to any flow for a region  $r$  that does not flow into a side-effect primitive operating on the region  $r$ .

$$\begin{array}{l}
\text{Fresh}(int) = int \\
\text{Fresh}(\sigma \rightarrow \sigma) = \text{Fresh}(\sigma) \xrightarrow{\beta} \text{Fresh}(\sigma) \\
\quad \text{where } \beta \text{ is fresh} \\
\text{Fresh}(\sigma \otimes \sigma') = \text{Fresh}(\sigma) \otimes \text{Fresh}(\sigma') \\
\text{Fresh}(\text{ref}(\sigma, \sigma', \rho)) = \text{ref}(\text{Fresh}(\sigma), \text{Fresh}(\sigma'), \rho) \\
\text{Fresh}(\text{reg}(\rho)) = \text{reg}(\rho) \\
\text{Fresh}(I) = \{\rho \mapsto \alpha \mid \rho \in I\} \\
\quad \text{where } \alpha \text{ is fresh}
\end{array}$$

Figure 11: *Fresh*.

$$\begin{array}{c}
\frac{\Gamma, W \vdash_b e^I : \tau, \mathcal{C}}{\Gamma, W + \text{Fresh}(I) \vdash_a e^I : \tau + \text{Fresh}(I), \mathcal{C}} \\
\frac{\Gamma, W \vdash_b e^\sigma : \tau, \mathcal{C} \quad \sigma \notin \text{type } I}{\Gamma, W \vdash_a e^\sigma : \tau, \mathcal{C}}
\end{array}$$

Figure 12: Type inference  $\vdash_a$ .

The type system is quite expressive. In particular, it is able to type all of the examples that were used as correct programs up to this point in the paper (with straightforward modification to translate  $\lambda_{\text{wit}}$  programs into  $\lambda_{\text{wit}}^{\text{reg}}$ ). In fact, assuming that for each region  $r$ ,  $\text{write}(r)$  nodes in each collection  $\mathbb{W}_i$  are totally ordered, the type system is “complete” for the first-order fragment (i.e., no higher order functions) with no recursion and no storing of witnesses in references. That is, such a program  $e$  is witness race free and well-typed by a standard Hindley-Milner type system if and only if it is well-typed by our type system. We also show later in Section 5 that the type system is more expressive than past approaches.

The limitations of the type system are the standard ones: let-based predicative polymorphism, flow-insensitivity of reference types, and unsoundness under call-by-name semantics; the latter is a typical limitation of a non-linear substructural type system. Another limitation is that the type system enforces for each region  $r$  that  $\text{write}(r)$  nodes in every collection  $\mathbb{W}_i$  are totally ordered whereas witness race freedom permits an absence of ordering for the case  $i = n$ ; we believe that this is a minor limitation.

### 4.3 Inference

We next present a type inference algorithm. By Theorem 6, this results in an automatic algorithm for statically checking witness race freedom.

At a high-level, our type system is a standard Hindley-Milner type system with some additional rational arithmetic constraints. Therefore we could perform inference by employing a standard type inference technique to solve all type-structural constraints while generating rational arithmetic constraints on the side, and then solving the generated arithmetic constraints separately. Unfortunately, the arithmetic constraints may be non-linear since they involve the multiplication of variables. Because there is no efficient algorithm for solving general non-linear rational arithmetic constraints, we need to dive into lower-level details of the type system.

Let us separate type inference into two phases. The first phase carries out type inference after erasing all rational numbers from the type system. That is, the types inferred in this phase are:

$$\sigma := int \mid \sigma \rightarrow \sigma' \mid \sigma \otimes \sigma' \mid \text{ref}(\sigma, \sigma', \rho) \mid \text{reg}(\rho) \mid I$$

where a type  $I$  is a subset of **RegIDs**. Intuitively, each type  $I$  represents the non-0 domain of some witness type  $W$ . The first phase can be carried out by a standard Hindley-Milner type inference, al-



$$\begin{array}{c}
\frac{\tau = \mathit{Fresh}(\sigma)}{\{x:\tau\}, \emptyset \vdash_b x^\sigma : \tau, \emptyset} \\
\\
\frac{\frac{\emptyset, \emptyset \vdash_b i : \mathit{int}, \emptyset}{\Gamma, W \vdash_a e : \tau, \mathcal{C}} \quad \frac{\emptyset, \emptyset \vdash_b \bullet : \emptyset, \emptyset}{\beta \text{ is fresh}}}{\Gamma \times \beta, W \times \beta \vdash_b \lambda x.e : \Gamma(x) \xrightarrow{\beta} \tau, \mathcal{C}} \\
\\
\frac{\frac{\Gamma_1, W_1 \vdash_a e_1 : \tau_1, \mathcal{C}_1 \quad \Gamma_2, W_2 \vdash_a e_2 : \tau_2, \mathcal{C}_2}{\tau = \mathit{Fresh}(\sigma) \quad \tau' = \mathit{Fresh}(\sigma') \quad \beta \text{ is fresh}} \quad \mathcal{C} = \mathcal{C}_1 \cup \mathcal{C}_2 \cup \{\tau_1 \geq \tau \xrightarrow{\beta} \tau', \beta \geq 1, \tau_2 \geq \tau\}}{\Gamma_1 + \Gamma_2, W_1 + W_2 \vdash_b e_1^{\sigma \rightarrow \sigma'} e_2 : \tau', \mathcal{C}} \\
\\
\frac{\Gamma_1, W_1 \vdash_a e_1 : \tau_1, \mathcal{C}_1 \quad \Gamma_2, W_2 \vdash_a e_2 : \tau_2, \mathcal{C}_2}{\Gamma_1 + \Gamma_2, W_1 + W_2 \vdash_b e_1 \otimes e_2 : \tau_1 \otimes \tau_2, \mathcal{C}_1 \cup \mathcal{C}_2} \\
\\
\frac{\Gamma, W \vdash_a \pi_i(e) : \tau, \mathcal{C} \quad \tau_1 = \mathit{Fresh}(\sigma) \quad \tau_2 = \mathit{Fresh}(\sigma')}{\Gamma, W \vdash_b \pi_i(e^{\sigma_1 \otimes \sigma_2}) : \tau_i, \mathcal{C} \cup \{\tau \geq \tau_1 \otimes \tau_2\}} \\
\\
\frac{\Gamma_1, W_1 \vdash_a e_1 : \tau_1, \mathcal{C}_1 \quad \Gamma_2, W_2 \vdash_a e_2 : \tau_2, \mathcal{C}_2}{\Gamma_1 + \Gamma_2, W_1 + W_2 \vdash_b \mathit{ref} e_1 e_2^{\mathit{reg}(\rho)} : \mathit{ref}(\tau_1, \tau_2, \rho), \mathcal{C}_1 \cup \mathcal{C}_2} \\
\\
\frac{\frac{\Gamma_1, W_1 \vdash_a e_1 : \tau_1, \mathcal{C}_1 \quad \Gamma_2, W_2 \vdash_a e_2 : \tau_2, \mathcal{C}_2}{\Gamma_3, W_3 \vdash_a e_2 : \tau_3, \mathcal{C}_3} \quad \tau = \mathit{Fresh}(\sigma) \quad \tau' = \mathit{Fresh}(\sigma')}{\mathcal{C} = \mathcal{C}_1 \cup \mathcal{C}_2 \cup \{\tau_1 \geq \mathit{ref}(\tau, \tau', \rho), \tau_2 \geq \tau', \tau_3(\rho) \geq 1\}} \quad \Gamma = \Gamma_1 + \Gamma_2 + \Gamma_3 \quad W = W_1 + W_2 + W_3}{\Gamma, W \vdash_b \mathit{write} e_1^{\mathit{ref}(\sigma, \sigma', \rho)} e_2 e_3 : \tau_3, \mathcal{C}} \\
\\
\frac{\frac{\Gamma_1, W_1 \vdash_a e_1 : \tau_1, \mathcal{C}_1 \quad \Gamma_2, W_2 \vdash_a e_2 : \tau_2, \mathcal{C}_2}{\tau = \mathit{Fresh}(\sigma) \quad \tau' = \mathit{Fresh}(\sigma')} \quad \mathcal{C} = \mathcal{C}_1 \cup \mathcal{C}_2 \cup \{\tau_1 \geq \mathit{ref}(\tau, \tau', \rho), \tau_2(\rho) > 0\}}{\Gamma_1 + \Gamma_2, W_1 + W_2 \vdash_b \mathit{read} e_1^{\mathit{ref}(\sigma, \sigma', \rho)} e_2 : \tau \otimes \tau_2, \mathcal{C}} \\
\\
\frac{\Gamma, W \vdash_a e : \tau, \mathcal{C}}{\Gamma, W \vdash_b \mathit{letreg} x^{\mathit{reg}(\rho)} e : \tau, \mathcal{C} \cup \{W(\rho) \leq 1\}} \\
\\
\frac{\Gamma_1, W_1 \vdash_a e_1 : \tau_1, \mathcal{C}_1 \quad \Gamma_2, W_2 \vdash_a e_2 : \tau_2, \mathcal{C}_2}{\Gamma_1 + \Gamma_2, W_1 + W_2 \vdash_b \mathit{join} e_1 e_2 : \tau_1 + \tau_2, \mathcal{C}_1 \cup \mathcal{C}_2} \\
\\
\frac{\frac{\Gamma_1, W_1 \vdash_a e_1 : \tau_1, \mathcal{C}_1 \quad \Gamma_2, W_2 \vdash_a e_2 : \tau_2, \mathcal{C}_2}{x \notin \mathit{free}(e_1) \quad \mathcal{C} = \mathcal{C}_1 \cup \mathcal{C}_2 \cup \{\tau_1 \geq \Gamma_2(x)\}}}{\Gamma_1 + \Gamma_2, W_1 + W_2 \vdash_b \mathit{let} x = e_1 \mathit{in} e_2 : \tau_2, \mathcal{C}} \\
\\
\frac{\frac{\Gamma_1, W_1 \vdash_a e_1 : \tau_1, \mathcal{C}_1 \quad \Gamma_2, W_2 \vdash_a e_2 : \tau_2, \mathcal{C}_2}{x \in \mathit{free}(e_1)}}{\mathcal{C} = \mathcal{C}_1 \cup \mathcal{C}_2 \cup \{\tau_1 \geq \Gamma_1(x) \times \infty, \tau_1 \geq \Gamma_2(x)\}} \quad \Gamma_1 + \Gamma_2, W_1 + W_2 \vdash_b \mathit{let} x = e_1 \mathit{in} e_2 : \tau_2, \mathcal{C}
\end{array}$$

Figure 13: Type inference  $\vdash_b$ .

beit with regions, which is no harder than type variables. We omit the details of this phase. We may safely reject the program if the first phase fails. Otherwise we annotate each subexpression  $e$  by its inferred type  $\sigma: e^\sigma$ . In the second phase, we use the annotated program to generate the appropriate rational arithmetic constraints via bottom-up type-inference. Let  $e$  be an annotated program. Then the generated constraints for  $e$  is  $\mathcal{C}$  where  $\Gamma, W \vdash_a e : \tau, \mathcal{C}$  for some  $\Gamma, W$ , and  $\tau$ .

The second-phase type inference rules are separated into two kinds,  $\vdash_a$  (Figure 12) and  $\vdash_b$  (Figure 13), which must occur in strictly interleaving manner. The purpose of  $\vdash_a$  is to account for the type judgment rule **Source**, whereas  $\vdash_b$  accounts for all other rules.

We should note that, strictly speaking, types  $\tau$  appearing in the algorithm are different from the ones in the type judgment rules. That is, instead of rational numbers, the types  $\tau$  in the algorithm are qualified by *rational number variables*  $\alpha, \beta, \gamma$ , etc. Also the domain of a witness type  $W$  is not the entire **RegIDs** set but only some subset of it. In other words, a witness type  $W$  is a partial function from **RegIDs** to rational number variables. We re-define the addition of witness types as follows to reflect the change:

$$\begin{aligned}
W + W' = & \{ \rho \mapsto W(\rho) + W'(\rho) \mid \rho \in \mathit{dom}(W) \wedge \rho \in \mathit{dom}(W') \} \\
& \cup \{ \rho \mapsto W(\rho) \mid \rho \in \mathit{dom}(W) \wedge \rho \notin \mathit{dom}(W') \} \\
& \cup \{ \rho \mapsto W'(\rho) \mid \rho \in \mathit{dom}(W') \wedge \rho \notin \mathit{dom}(W) \}
\end{aligned}$$

We also re-define the addition of type environments:

$$\begin{aligned}
\Gamma + \Gamma' = & \{ x : \Gamma(x) + \Gamma'(x) \mid x \in \mathit{dom}(\Gamma) \wedge x \in \mathit{dom}(\Gamma') \} \\
& \cup \{ x : \Gamma(x) \mid x \in \mathit{dom}(\Gamma) \wedge x \notin \mathit{dom}(\Gamma') \} \\
& \cup \{ x : \Gamma'(x) \mid x \in \mathit{dom}(\Gamma') \wedge x \notin \mathit{dom}(\Gamma) \}
\end{aligned}$$

Note that we omit annotations when they are not used (i.e., we say  $e$  instead of  $e^\sigma$ , etc.). There are only two cases for  $\vdash_a$ . The first case is for expressions that were given a witness type  $I$  in the first phase. In this case, we add  $\mathit{Fresh}(I)$  to  $\tau$  and  $W$  to account for a possible application of **Source**.  $\mathit{Fresh}$  is defined in Figure 11. The second case is for expressions that were not given a witness type. In this case, we simply pass the result of the subderivation  $\vdash_b$  up.

We discuss a few representative  $\vdash_b$  rules. Note that  $\vdash_b$  rules are syntax directed. In the case of a variable  $x^\sigma$ , we create a fresh  $\tau$  from  $\sigma$  and pass  $\{x:\tau\}, \emptyset \vdash_b x^\sigma : \tau, \emptyset$  up to the parent derivation. (Recall our type inference is bottom-up.) The case for integers and dummy witnesses are trivial. In the case of an abstraction  $\lambda x.e$ , we multiply  $\Gamma$  and  $W$  passed from the subderivation by  $\beta$ . In the case of a function application  $e_1^{\sigma \rightarrow \sigma'} e_2$ , we add the constraints  $\{\tau_1 \geq \tau \xrightarrow{\beta} \tau', \beta \geq 1, \tau_2 \geq \tau\}$  to connect arguments and returns as well as requiring  $\beta$  to be at least 1. Note that the type rule **Sub** is implicitly incorporated in the constraints. In the case of **write**  $e_1^{\mathit{ref}(\sigma, \sigma', \rho)} e_2 e_3$ , we add the constraint  $\tau_2(\rho) \geq 1$  to match the type rule **Write**. Note that the first phase guarantees that  $\rho \in \mathit{dom}(\tau_3)$ . In the case **letreg**  $x^{\mathit{reg}(\rho)} e$ , the constraint  $W(\rho) \leq 1$  is effective only when  $\rho \in \mathit{dom}(W)$  as  $\rho \notin \mathit{dom}(W)$  implies that the region was not used at all. Note that there is no case corresponding to the type rule **LetA**. Prior to running the algorithm, we replace each occurrence of the expression **let**  $x = e$  **in**  $e'$  in the program by the expression  $e'[(\mathit{let} x = e \mathit{in} x)/x]$  whenever  $e \in V$  and  $x$  occurs more than once in  $e'$ .

As an example, consider the following program (a  $\lambda_{wit}^{reg}$  version of the last example from Section 2):

```

letreg r
  let x = ref 1 r in
    let w = write x 2 • in
      let f =  $\lambda y.$ read x w in
        let z = (f 0)  $\otimes$  (f 0) in
          let w = write x 3 join  $\pi_2(\pi_1(z)) \pi_2(\pi_2(z))$  in z

```

Suppose the first phase assigns  $r$  the type  $reg(\rho)$ . Assume each **let**-bound variable is treated monomorphically. The second phase generates the following constraints for the **let**-bound expressions (slightly simplified for readability):

$$\{r:reg(\rho)\}; \emptyset \vdash_a \mathbf{ref} \ 1 \ r : \mathbf{ref} \ (int, int, \rho), \emptyset$$

$$\Gamma; \{\rho \mapsto \gamma_1 + \gamma_2\} \vdash_a \mathbf{write} \ x \ 2 \ \bullet : \{\rho \mapsto \gamma_1 + \gamma_2\}, \{\gamma_1 \geq 1\}$$

where  $\Gamma = \{x:ref(int, int, \rho)\}$

$$\Gamma; W \vdash_a \lambda y. \mathbf{read} \ x \ w : int \xrightarrow{\beta_1} int \otimes \{\rho \mapsto \alpha_1 + \gamma_3\}, \mathcal{C}$$

where  $\Gamma = \{x:ref(int, int, \rho), w:\{\rho \mapsto \alpha_1 \times \beta_1\}\}$   
and  $W = \{\rho \mapsto \gamma_3 \times \beta_1\}$   
and  $\mathcal{C} = \{\alpha_1 + \gamma_3 > 0\}$

$$\Gamma; \emptyset \vdash_a (f \ 0) \otimes (f \ 0) : \tau, \mathcal{C}$$

where  $\Gamma = \{f:int \xrightarrow{\beta_2+\beta_3} int \otimes \{\rho \mapsto \alpha_2\}\}$   
and  $\tau = (int \otimes \{\rho \mapsto \alpha_2\}) \otimes (int \otimes \{\rho \mapsto \alpha_2\})$   
and  $\mathcal{C} = \{\beta_2 \geq 1, \beta_3 \geq 1\}$

$$\Gamma; \{\rho \mapsto \gamma_4\} \vdash_a \mathbf{write} \ x \ 3 \ \dots : \{\rho \mapsto \alpha_3 + \alpha_4 + \gamma_4\}, \mathcal{C}$$

where  $\Gamma = \{x:ref(int, int, \rho), z:\tau\}$   
and  $\tau = (int \otimes \{\rho \mapsto \alpha_3\}) \otimes (int \otimes \{\rho \mapsto \alpha_4\})$   
and  $\mathcal{C} = \{\alpha_3 + \alpha_4 + \gamma_4 \geq 1\}$

The final constraints, after some simplification, is as follows:

$$\begin{aligned} \gamma_1 \geq 1, 1 \geq \gamma_1 + \gamma_2 + \gamma_3 \times \beta_1 + \gamma_4, \beta_1 \geq \beta_2 + \beta_3, \\ \beta_2 \geq 1, \beta_3 \geq 1, \gamma_1 + \gamma_2 \geq \alpha_1 \times \beta_1, \alpha_1 + \gamma_3 > 0, \\ \alpha_1 + \gamma_3 \geq \alpha_2, \alpha_3 + \alpha_4 + \gamma_4 \geq 1, \alpha_2 \geq \alpha_3, \alpha_2 \geq \alpha_4 \end{aligned}$$

Note that the constraints are satisfiable, e.g., by the substitution

$$\begin{aligned} \beta_1 = 2 \quad \beta_2 = \beta_3 = \gamma_1 = 1 \\ \alpha_1 = \alpha_2 = \alpha_3 = \alpha_4 = 0.5 \\ \gamma_2 = \gamma_3 = \gamma_4 = 0 \end{aligned}$$

In general, a program  $e$  is well-typed if and only if the constraints  $\mathcal{C}$  generated by type inference are satisfiable. So it suffices to show that the satisfaction problem for any  $\mathcal{C}$  generated by the type inference algorithm can be solved.

To this end, we first note that because of the first phase, any constraint  $\tau \geq \tau' \in \mathcal{C}$  can be reduced to a set of rational arithmetic constraints of the form  $p \geq p'$  where  $p, p'$  are rational polynomials. The troublesome non-linearity comes from  $\Gamma \times \beta$  and  $W \times \beta$  in the  $\lambda x.e$  case. Let us focus our attention on the set  $\mathcal{B}$  of variables used in such multiplications. (Note that we have used  $\beta$  instead of  $\alpha$  just for this case in the pseudo-code to make it clear that these variables are special.) We can show that the following holds:

**THEOREM 7.** *Let  $p \triangleright p' \in \mathcal{C}$  where  $\triangleright \in \{\geq, >\}$ . If  $\beta \in \mathcal{B}$  occurs in the polynomial  $p$ , then it must be the case that  $\triangleright = \geq$ ,  $p = \beta$ , and that the polynomial  $p'$  consists only of symbols in the set  $\{+, \times, 1, \infty\} \cup \mathcal{B}$ .*

**Proof (Sketch):** Let  $\Gamma, W \vdash_a e : \tau, \mathcal{C}$ . For any  $\tau' \in \Gamma$ ,  $\times$  and  $\infty$  only appear at the top-level, i.e., not within argument and return types of a function type. Secondly, we can show by induction that within the type  $\tau$ ,  $\times$  only appears in negative positions. More precisely, for any  $p \in Pos(\tau)$ , the polynomial  $p$  contains no  $\times$

where  $Pos$  is defined as follows:

$$\begin{aligned} Pos(\tau \xrightarrow{p} \tau') &= Neg(\tau) \cup Pos(\tau') \cup \{p\} \\ Pos(\tau \otimes \tau') &= Pos(\tau) \cup Pos(\tau') \\ Pos(ref(\tau, \tau', \rho)) &= Pos(\tau) \cup Neg(\tau') \\ Pos(W) &= ran(W) \\ Neg(int) &= Pos(int) = Neg(reg(\rho)) = Pos(reg(\rho)) = \emptyset \\ Neg(\tau \xrightarrow{p} \tau') &= Pos(\tau) \cup Neg(\tau') \\ Neg(\tau \otimes \tau') &= Neg(\tau) \cup Neg(\tau') \\ Neg(ref(\tau, \tau', \rho)) &= Neg(\tau) \cup Pos(\tau') \\ Neg(W) &= \emptyset \end{aligned}$$

Third, for any  $\times$  that appears in a positive position of  $\tau$ , i.e. in some  $p \in Pos(\tau)$ , the polynomial  $p$  does not contain any  $\beta \in \mathcal{B}$ . Then the result follows from inspection of the subtyping rules.  $\square$

The theorem implies that we can compute all assignments to the variables in  $\mathcal{B}$  by computing the minimum satisfying assignment for  $\mathcal{C}' = \{\beta \geq p \mid \beta \in \mathcal{B}\} \subseteq \mathcal{C}$ . It is easy to see that such an assignment always exists. (Recall that the range is non-negative.) This problem can be solved in quadratic time by an iterative method in which all variables are initially set to 0, and at each iteration the new values for the variables are computed by taking the maximum of the right hand polynomials evaluated at the current values. It is possible to show that if the minimum satisfying assignment for a variable  $\beta$  is some  $q < \infty$ , then the iterative method finds  $q$  for  $\beta$  in  $2|\mathcal{C}'|$  iterations. Hence any variable changing after the  $2|\mathcal{C}'|$ th iteration can be safely set to  $\infty$ . All variables are then guaranteed to converge within  $3|\mathcal{C}'|$  iterations. Because each iteration examines every constraint, the overall time is at most quadratic in the size of  $\mathcal{C}'$ .

Substituting the computed assignments for  $\mathcal{B}$  in  $\mathcal{C}$  results in linear rational constraints, which can be solved efficiently by a linear programming algorithm.

## 5. Related Work

Adding side-effects to a functional language is an old problem with many proposed solutions. Here we compare our technique against two of the more prominent approaches: linear types [15, 6, 1] and monads [10, 11, 7, 9].

In linear types, there is an explicit *world* program value (or one world per region for languages with regions) that conceptually represents the current program state. By requiring each world have a linear type, the type system ensures that the world can be updated in place.

The linear type system can be expressed in our type system by restricting every flow to 1, every witness to contain only one flow, and designating one dummy witness to serve the role of the “starting” witness (or for regions, one dummy witness per region). Thus our approach is more expressive than such an approach. Note that this also implies that every function type can be restricted to have either 1 or  $\infty$  as the qualifier. It is easy to see that the program is well-typed under this restriction if and only if it is well-typed by the linear type system. The restriction limits programs to manipulate witnesses only in a linear fashion. In practice, this implies that there can be no parallel reads, no dead witnesses, no redundant witnesses, and no duplication of values containing witnesses.

Our approach can implement monadic primitives as follows (for concrete comparison, we use *state monads* [9]):

```

newVar =  $\lambda x.\lambda y.\text{let } z = (\text{ref } x \pi_2(y)) \text{ in } y \otimes z$ 
readVar =  $\lambda x.\lambda y.\text{let } z = (\text{read } x \pi_1(y)) \text{ in}$ 
            $(\pi_2(z) \otimes \pi_2(y)) \otimes \pi_1(z)$ 
writeVar =  $\lambda x.\lambda w.\lambda y.\text{let } z = (\text{write } x w \pi_1(y)) \text{ in}$ 
            $(z \otimes \pi_2(y)) \otimes w$ 
>>= =  $\lambda f.\lambda g.\lambda x.\text{let } y = (f x) \text{ in } g \pi_2(y) \pi_1(y)$ 
returnST =  $\lambda x.\lambda y.y \otimes x$ 
runST e = letreg x  $\pi_2(e (\bullet \otimes x))$ 

```

The idea behind these definitions is to implement each state monad of the type  $ST(\alpha, \tau)$  as a function that takes a witness and the region  $\alpha$  as arguments and returns a witness, the region  $\alpha$ , and a value of the type  $\tau$ . It is easy to see that if a state monad program is well-typed by the monadic type system, then it is also well-typed with our type system using the above definitions for the monadic primitives. Thus, our approach is more expressive than monads.

In practice, a monadic approach shares essentially the same limitations as linear types; for example, side-effects are restricted to a linear, sequential order. (In fact, it is not too hard to see that we can actually implement monadic primitives with the linear types restriction with only slightly longer code.) On the other hand, a monadic type system has an engineering advantage as it only requires Hindley-Milner type inference.

In addition to the above technical differences, our approach differs from previous approaches in its design motivation. That is, while our language feature, witnesses, is motivated by a pragmatic observation, understanding the motivation behind linear types and monads (i.e., not just knowing how to use them) arguably requires an appreciation of their underlying theory.

The technique used in our type system has some resemblance to *fractional permissions* [2] which can guarantee non-interference in imperative programs. Indeed, it may be possible to give an “interpretation” of witnesses as some relaxation of permissions or *capabilities* [3] (after promoting permissions or capabilities to first class values) in situations where the program is witness race free.

## 6. Conclusions

We have presented a new approach to adding side-effects in purely functional languages based on witnesses. We have stated a natural semantic correctness condition called witness race freedom and proposed a type-based approach for statically checking witness race freedom.

## References

[1] P. Achten, J. H. G. van Groningen, and R. Plasmeijer. High level specification of I/O in functional languages. In *Proceedings of the 1992 Glasgow Workshop on Functional Programming*, pages 1–17. Springer-Verlag, 1993.

[2] J. Boyland. Checking interference with fractional permissions. In R. Cousot, editor, *Static Analysis, Tenth International Symposium*, volume 2694 of *Lecture Notes in Computer Science*, pages 55–72, San Diego, CA, June 2003. Springer-Verlag.

[3] K. Cray, D. Walker, and G. Morrisett. Typed memory management in a calculus of capabilities. In *Proceedings of the 26th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, pages 262–275, San Antonio, Texas, Jan. 1999.

[4] M. Fähndrich, J. S. Foster, Z. Su, and A. Aiken. Partial online cycle elimination in inclusion constraint graphs. In *Proceedings of the 1998 ACM SIGPLAN Conference on Programming Language Design and Implementation*, pages 85–96, Montreal, Canada, June 1998.

[5] D. Grossman, G. Morrisett, T. Jim, M. Hicks, Y. Wang, and J. Cheney. Region-based memory management in Cyclone. In *Proceedings of the 2002 ACM SIGPLAN Conference on Programming Language Design and Implementation*, Berlin, Germany, June 2002.

[6] J. C. Guzman and P. Hudak. Single-threaded polymorphic lambda calculus. In *Proceedings, Fifth Annual IEEE Symposium on Logic in Computer Science*, pages 42–51, Philadelphia, Pennsylvania, June 1990.

[7] S. L. P. Jones and P. Wadler. Imperative functional programming. In *Proceedings of the 20th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, pages 71–84, Charleston, South Carolina, Jan. 1993.

[8] J. Launchbury. A natural semantics for lazy evaluation. In *Proceedings of the 20th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, pages 144–154, Charleston, South Carolina, Jan. 1993.

[9] J. Launchbury and A. Sabry. Monadic state: Axiomatization of type safety. In *Proceedings of the 1997 ACM SIGPLAN International Conference on Functional Programming*, pages 227–238, Amsterdam, The Netherlands, June 1997.

[10] E. Moggi. Notions of computation and monads. *Information and Computation*, 93(1):55–92, 1991.

[11] M. Odersky, D. Rabin, and P. Hudak. Call by name, assignment, and the lambda calculus. In *Proceedings of the 20th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, pages 43–56, Charleston, South Carolina, Jan. 1993.

[12] G. D. Plotkin. Call-by-name, call-by-value and the lambda-calculus. *Theoretical Computer Science*, 1:125–159, 1975.

[13] T. Terauchi and A. Aiken. Witnessing side-effects. Technical report, Computer Science Division, EECS Department, University of California, Berkeley.

[14] M. Tofte and J.-P. Talpin. Implementation of the typed call-by-value  $\lambda$ -calculus using a stack of regions. In *Proceedings of the 21st Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, pages 188–201, Portland, Oregon, Jan. 1994.

[15] P. Wadler. Linear types can change the world! In *Proceedings of the third ACM SIGPLAN International Conference on Functional Programming*, pages 63–74, Baltimore, Maryland, Sept. 1998.