

# Safe—A Semantic Technique for Transforming Programs in the Presence of Errors

ALEXANDER AIKEN

University of California, Berkeley

and

JOHN H. WILLIAMS and EDWARD L. WIMMERS

IBM Almaden Research Center

---

Language designers and implementors have avoided specifying and preserving the meaning of programs that produce errors. This is apparently because being forced to preserve error behavior severely limits the scope of program optimization, even for correct programs. However, preserving error behavior is desirable for debugging, and error behavior must be preserved in any language that permits user-generated errors (i.e., exceptions).

This paper presents a technique for expressing general program transformations for languages that possess a rich collection of distinguishable error values. This is accomplished by defining a higher-order function called “Safe”, which can be used to annotate those portions of a program that are guaranteed not to produce errors. It is shown that this facilitates the expression of very general program transformations, effectively giving program transformations in a language with many error values the same power and generality as program transformations in a language with only a single error value.

Using the semantic properties of `Safe`, it is possible to provide some useful sufficient conditions for establishing the correctness of transformations in the presence of errors. In particular, a Substitutability Theorem is proven, which can be used to justify “in-context” optimizations—transformations that alter the meanings of subexpressions without changing the meaning of the whole program. Finally, the effectiveness of the technique is demonstrated by some examples of its use in an optimizing compiler.

Categories and Subject Descriptors: D.3.4 [Processors]: Optimization

General Terms: Languages, Theory

Additional Key Words and Phrases: equational reasoning, exceptions, program optimization, program transformation

---

## 1. INTRODUCTION

A laudable trend of the past two decades has been the increased use of denotational semantics to guide the design and implementation of programming languages. Semantics-driven language design has produced cleaner and simpler languages and provided more precise standards for testing the correctness of language implementations.

An apparent exception to this trend is the treatment of error handling. All too

---

Address for first author: Computer Science Division, 773 Soda Hall, University of California, Berkeley, 94720-1776, email: aiken@cs.berkeley.edu. Address for second and third authors: IBM Almaden Research Center, 650 Harry Rd., San Jose, CA, 95120-6077, email: {williams,wimmers}@almaden.ibm.com. A preliminary version of this paper was presented at the 1990 POPL Conference [Aiken et al. 1990].

often, errors are considered to be outside the scope of the denotational semantics; if anything is specified about error behavior, it is usually through some *ad hoc* mechanism. Some language features—such as strong typing—reduce the negative impact of such a design but cannot avoid it completely; runtime errors exist in every language and must be handled in an implementation.

Many languages simply omit errors altogether from their formal semantic specification. For example, early Fortran implementations were free to report errors as the implementor saw fit, and transformations to improve performance could change the behavior of error-producing programs. This approach is still taken with some modern languages, e.g., Haskell [Hudak et al. 1988] and the languages of Turner [Turner 1985], where strong typing and lazy evaluation lessen the need for error recovery and exception handling.

Other languages include errors in the domain of values and provide some mechanism for computing with (or recovering from) errors, but the formal specification allows considerable variation in the behavior of implementations. For example, in [Steele 1984] Steele writes:

“The definition of Common Lisp . . . explicitly requires the interpreter and compiler to impose identical semantics on *correct* programs *so far as possible*.” (emphasis added)

Indeed, in one Common Lisp implementation, taking the `car` of an atom produces a run time error when interpreted but returns the current package when compiled!

One study of program transformations carefully accounts for the fact that many common optimizations do not preserve error behavior by giving a precise, denotational treatment of such transformations [Cartwright and Felleisen 1989]. In this approach, errors are considered to approximate other values, and program transformations are lifting operations that can (possibly) make programs more defined. Another approach [Hennessy 1981]

“. . . examines the optimization difficulties imposed by common exception handling facilities [and] proposes restrictions on these mechanisms that make the optimization of programs possible.”

Why have language designers and implementors avoided specifying and preserving the meaning of errors? The answer appears to be that not preserving error behavior increases the power and effectiveness of transforming correct programs, i.e., the “important” programs. For instance, substituting `s2` (which selects the second element of a sequence) for `s1 o t1` (which selects the first element of the tail of a sequence) may improve a program’s running time, but a programmer who hadn’t used the primitive `s2` and was unaware of its existence would be confused by the run-time error message “`s2` incorrectly applied to a non-sequence argument.”

If the language designer opts for language clarity and ease of debugging by making a semantic distinction between `t1` errors and `s2` errors, then the general transformation becomes invalid, and some weaker version must be substituted, perhaps in the form of a set of rules identifying particular contexts in which the replacement is valid. This can be a significant loss, since identifying such contexts in general requires knowledge of the entire program. Thus, including errors as semantic objects in order to make a language easier to use appears to weaken the generality

and power of the language’s program transformations.

This paper presents a technique for preserving the power of general program transformations in the presence of a rich collection of distinguishable error values. This is accomplished by introducing an annotation, “**Safe**”, to mark occurrences of functions that cannot produce errors. Succinct and general algebraic laws can be expressed using **Safe**, thereby giving program transformations in a language with many error values the same power and generality as program transformations in a language with only a single error value (such as FP [Backus 1978]). In fact, the **Safe** mechanism accomplishes much more. It actually strengthens equational reasoning by providing a sufficient condition on a program context  $\mathbf{E}(\cdot)$  and functions  $\mathbf{f}$  and  $\mathbf{g}$ , such that  $\mathbf{E}(\mathbf{f}) \equiv \mathbf{E}(\mathbf{g})$  even if  $\mathbf{f} \neq \mathbf{g}$ .

The **Safe** mechanism is presented in the context of the functional language FL [Backus et al. 1986], but it should be applicable in other source-to-source program transformation systems. In fact, the only requirement for using these techniques in other contexts is that a suitable definition of the function **Safe** can be found for the programming language under consideration. Section 2 describes enough of FL to illustrate the technique and prove its soundness. Section 3 introduces the **Safe** mechanism and gives a simple example illustrating that having just two distinguishable errors causes as much loss of algebraic generality as having arbitrarily many different kinds of errors. This shows that it is not possible through careful language design to make a gradual trade-off between the expressiveness of error reporting and algebraic generality. Section 3 also contains the Substitutability Theorem, which provides a criterion for proving the soundness of transformations involving **Safe**. Some examples of optimization using **Safe** are given in Section 4. Section 5 discusses some pragmatic concerns, while Section 6 reports on experience with the use of these techniques in an optimizing compiler for FL. Section 7 concludes with a discussion of related work and suggestions for further work.

## 2. AN OVERVIEW OF FL

FL [Backus et al. 1986] is the result of an effort to design a practical functional language based on FP [Backus 1978]. FL has many of the features expected in modern functional languages, including higher-order functions, recursive definitions, user-defined functions, and datatypes. FL also differs from other functional languages in important ways. In particular, FL is combinator-based, has no static type system, and has first-class exceptions. Only the last feature is important to the results of this paper.

Figure 1 gives the subset of FL needed to understand the laws and examples that follow. Some features of the language are ignored altogether; in particular, function definitions, input/output functions, syntactic sugar are omitted. The evaluation order of FL expressions is leftmost-innermost; thus, in  $[\mathbf{f} : \mathbf{x}, \mathbf{g} : \mathbf{x}]$ , first  $\mathbf{f} : \mathbf{x}$  is evaluated and then  $\mathbf{g} : \mathbf{x}$  is evaluated.

In Figure 1, a definition of each function is given only for some arguments; for all other arguments, the function  $\mathbf{f}$  returns an error value  $\mathbf{ferr}$  (e.g.,  $\mathbf{s1:0} \equiv \mathbf{s1err}$ ). FL functions produce more informative errors than just the name of the function, but this countable set of errors is sufficient for the purposes of this paper. In fact, the methods presented in subsequent sections are actually independent of the descriptiveness of the error values and work even with, e.g., error values that encode

$$\begin{aligned}
& \mathbf{f}: \mathbf{x} && \text{denotes function application} \\
\langle \mathbf{x}_1, \dots, \mathbf{x}_n \rangle &&& \text{denotes sequence construction, } n \geq 0 \\
(\mathbf{f} \circ \mathbf{g}): \mathbf{x} &= \mathbf{f}: (\mathbf{g}: \mathbf{x}) \\
[\mathbf{f}_1, \dots, \mathbf{f}_n]: \mathbf{x} &= \langle \mathbf{f}_1: \mathbf{x}, \dots, \mathbf{f}_n: \mathbf{x} \rangle \\
(\mathbf{p} \rightarrow \mathbf{q}; \mathbf{r}): \mathbf{x} &= \begin{cases} \mathbf{p}: \mathbf{x} & \text{if } \mathbf{p}: \mathbf{x} \in \mathcal{E}_{\text{FL}} \\ \mathbf{r}: \mathbf{x} & \text{if } \mathbf{p}: \mathbf{x} = \mathbf{false} \\ \mathbf{q}: \mathbf{x} & \text{otherwise} \end{cases} \\
\sim \mathbf{x}: \mathbf{y} &= \mathbf{x} \\
\alpha: \mathbf{f}: \langle \mathbf{x}_1, \dots, \mathbf{x}_n \rangle &= \langle \mathbf{f}: \mathbf{x}_1, \dots, \mathbf{f}: \mathbf{x}_n \rangle \\
+ : \langle \mathbf{x}_1, \dots, \mathbf{x}_n \rangle &= \mathbf{x}_1 + \dots + \mathbf{x}_n \\
\mathbf{si}: \langle \mathbf{x}_1, \dots, \mathbf{x}_n \rangle &= \mathbf{x}_i \\
\mathbf{tl}: \langle \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \rangle &= \langle \mathbf{x}_2, \dots, \mathbf{x}_n \rangle \\
\mathbf{rev}: \langle \mathbf{x}_1, \dots, \mathbf{x}_n \rangle &= \langle \mathbf{x}_n, \dots, \mathbf{x}_1 \rangle \\
\mathbf{al}: \langle \mathbf{x}, \langle \mathbf{y}_1, \dots, \mathbf{y}_n \rangle \rangle &= \langle \mathbf{x}, \mathbf{y}_1, \dots, \mathbf{y}_n \rangle \\
\mathbf{distl}: \langle \mathbf{x}, \langle \mathbf{y}_1, \dots, \mathbf{y}_n \rangle \rangle &= \langle \langle \mathbf{x}, \mathbf{y}_1 \rangle, \dots, \langle \mathbf{x}, \mathbf{y}_n \rangle \rangle \\
\mathbf{id}: \mathbf{x} &= \mathbf{x} \\
\mathbf{dom}: \mathbf{f} &= \mathbf{s1} \circ [\mathbf{id}, \mathbf{f}] \\
\mathbf{catch}: \langle \mathbf{f}, \mathbf{g} \rangle: \mathbf{x} &= \begin{cases} \mathbf{g}: \langle \mathbf{x}, \mathbf{y} \rangle & \text{if } \mathbf{f}: \mathbf{x} = \mathbf{y}_{\text{err}} \\ \mathbf{f}: \mathbf{x} & \text{otherwise} \end{cases}
\end{aligned}$$

Fig. 1. A subset of FL.

the position in the program where the error arose.

In designing FL, it was recognized that one of the deficiencies of FP is that it has a single error message  $\perp$  (or *Wrong!*) for all exceptional circumstances. Error messages and exception handling are an integral part of FL; as in the current version of Standard ML [Harper et al. 1989], errors are first class values rather than the special results of functions that fail to produce values. Semantically, error values in FL are treated differently than ordinary values. All functions are strict with respect to errors, so that  $\mathbf{f}: \mathbf{x}_{\text{err}} \equiv \mathbf{x}_{\text{err}}$  for any function  $\mathbf{f}$  and error value  $\mathbf{x}_{\text{err}}$ . Sequence construction is also strict with respect to errors; a sequence collapses to the leftmost error it contains. This behavior is justified by the intended use of errors in FL: errors represent a situation in which something extraordinary has happened, and therefore an error should persist until caught or until it escapes from (and becomes the result of) the program. Some of the semantic treatment of errors can be seen in the recursive domain equations for FL:

$$\begin{aligned}
\mathcal{D}_{\text{FL}} &= \mathcal{D}_{\text{FL}}^+ \cup \mathcal{E}_{\text{FL}} \\
\mathcal{D}_{\text{FL}}^+ &= A \cup \text{Seqs}(\mathcal{D}_{\text{FL}}^+) \cup (\mathcal{D}_{\text{FL}}^+ \rightarrow \mathcal{D}_{\text{FL}}) \\
&\quad \text{(the ordinary values)} \\
\mathcal{E}_{\text{FL}} &= \text{Err}(\mathcal{D}_{\text{FL}}^+) \cup \{\perp\} \\
&\quad \text{(the error values)}
\end{aligned}$$

In these equations,  $A$  is the set of atoms,  $\text{Seqs}$  is sequence construction, and  $\text{Err}$  is error construction. The ordering on  $\mathcal{D}_{\text{FL}}^+$  is the standard one; in  $\mathcal{E}_{\text{FL}}$ ,  $\mathbf{x}_{\text{err}} \leq \mathbf{y}_{\text{err}} \Leftrightarrow \mathbf{x} \leq \mathbf{y}$  and  $\perp \leq \mathbf{x}$  for all  $\mathbf{x}$ . Note that because other kinds of errors are given distinct

values, the value  $\perp$  is used only to denote non-terminating computations.

### 3. THE SAFE MECHANISM

One of the principles underlying FL is that a programming language should have a rich algebra useful for reasoning about and optimizing programs. Errors have a great impact on the algebra; for example, if two expressions can produce distinct errors, then the order of evaluation of the expressions usually cannot be changed without changing the error produced. Even with errors, however, there are many general identities that hold between FL programs; as usual,  $\mathbf{f} \equiv \mathbf{g}$  means that  $\mathbf{f}$  and  $\mathbf{g}$  denote the same semantic value.

$$\mathbf{f} \circ \text{id} \equiv \mathbf{f} \quad (1)$$

$$\text{id} \circ \mathbf{f} \equiv \mathbf{f} \quad (2)$$

$$[\mathbf{f}_1, \dots, \mathbf{f}_n] \circ \mathbf{g} \equiv [\mathbf{f}_1 \circ \mathbf{g}, \dots, \mathbf{f}_n \circ \mathbf{g}] \quad (3)$$

$$\mathbf{f} \circ (\mathbf{p} \rightarrow \mathbf{q}; \mathbf{r}) \equiv \mathbf{p} \rightarrow \mathbf{f} \circ \mathbf{q}; \mathbf{f} \circ \mathbf{r} \quad (4)$$

$$(\mathbf{p} \rightarrow \mathbf{q}; \mathbf{r}) \circ \mathbf{f} \equiv \mathbf{p} \circ \mathbf{f} \rightarrow \mathbf{q} \circ \mathbf{f}; \mathbf{r} \circ \mathbf{f} \quad (5)$$

These laws hold because the order of application of the component functions is unchanged.

However,

$$[\mathbf{g}, \mathbf{f}] \equiv \text{rev} \circ [\mathbf{f}, \mathbf{g}] \quad (6)$$

which is a law in FP, is not valid for all FL functions  $\mathbf{f}$  and  $\mathbf{g}$ ; if  $\mathbf{f}$  produces  $\perp$  and  $\mathbf{g}$  produces  $\text{t1err}$ , then  $[\mathbf{g}, \mathbf{f}]$  produces  $\text{t1err}$ , but  $\text{rev} \circ [\mathbf{f}, \mathbf{g}]$  produces  $\perp$ . This example shows there cannot be a gradual trade-off of expressiveness of error reporting for generality of program transformations. Having two errors is as limiting as having arbitrarily many, since the existence of just one error distinguishable from  $\perp$  is sufficient to invalidate any “law” that does not preserve the order of evaluation of its constituent functions.

#### 3.1 An Informal Treatment of Safe

Although (6) is not always true, it is the case that

$$[\mathbf{g}, \mathbf{f}] \equiv \text{rev} \circ [\mathbf{f}, \mathbf{g}] \text{ if neither side makes an error.} \quad (7)$$

That is, there are *contexts* in which  $[\mathbf{g}, \mathbf{f}]$  can be substituted for  $\text{rev} \circ [\mathbf{f}, \mathbf{g}]$ . There are many other examples of rewrite rules that are correct provided neither side produces an error, and including them greatly enhances the power of a program transformation system. The following informal examples illustrate the notion of “rewriting in context”, i.e., using rules whose validity depends on the context in which they are applied. In these examples, an occurrence of a function  $\mathbf{f}$  is annotated as being “safe” (written  $\mathbf{S}(\mathbf{f})$ ), if it is known that that occurrence is guaranteed not to produce an error when applied to a non-error value. The notation  $\mathbf{f} \mapsto \mathbf{g}$  simply indicates that  $\mathbf{f}$  is rewritten to  $\mathbf{g}$ . (N.B. For now,  $\mathbf{S}(\mathbf{f})$  is an extra-linguistic notion; the phrase “annotating  $\mathbf{f}$  with  $\mathbf{S}(\mathbf{f})$ ” has no more semantic content than the phrase “painting  $\mathbf{f}$  green”.)

Consider the program  $\mathbf{rev} \circ [\sim 0, \sim 1]$ . Since the two constant functions cannot produce an error unless applied to an error, neither can the construction of the two constant functions. Thus the program can be rewritten:

$$\begin{aligned} \mathbf{rev} \circ [\sim 0, \sim 1] &\mapsto \\ \mathbf{rev} \circ [\mathbf{S}(\sim 0), \mathbf{S}(\sim 1)] &\mapsto \\ \mathbf{rev} \circ \mathbf{S}([\sim 0, \sim 1]) & \end{aligned}$$

Now note that because the argument to reverse is safe, the order of evaluation of the elements of the sequence must be irrelevant, so the application of  $\mathbf{rev}$  can be eliminated. As a last step, the annotation  $\mathbf{S}$  can be dropped:

$$\begin{aligned} \mathbf{rev} \circ \mathbf{S}([\sim 0, \sim 1]) &\mapsto \\ \mathbf{S}([\sim 1, \sim 0]) &\mapsto \\ [\sim 1, \sim 0] & \end{aligned}$$

Intuitively, each of the above steps preserved the meaning of the program, so this appears to be a proper ‘optimization’. Indeed, in this case the final program is equivalent to the original program; i.e.,  $\mathbf{rev} \circ [\sim 0, \sim 1] \equiv [\sim 1, \sim 0]$ .

So far, however, the safe mechanism is informal, and it is easy to make mistakes. Consider the program  $\mathbf{rev} \circ [\mathbf{dist1}, \mathbf{t1} \circ \mathbf{a1}]$ . Because  $\mathbf{a1}$  follows  $\mathbf{dist1}$  in the order of evaluation, and because  $\mathbf{a1}$  produces an error for exactly the same arguments as  $\mathbf{dist1}$ ,  $\mathbf{a1}$  can be marked safe.

$$\mathbf{rev} \circ [\mathbf{dist1}, \mathbf{t1} \circ \mathbf{a1}] \mapsto \mathbf{rev} \circ [\mathbf{dist1}, \mathbf{t1} \circ \mathbf{S}(\mathbf{a1})]$$

Next notice that if  $\mathbf{a1}$  is error-free, then  $\mathbf{t1}$  always succeeds and returns the second component of the original argument. Therefore  $\mathbf{S}(\mathbf{s2})$  can be substituted for  $\mathbf{t1} \circ \mathbf{S}(\mathbf{a1})$ .

$$\mathbf{rev} \circ [\mathbf{dist1}, \mathbf{t1} \circ \mathbf{S}(\mathbf{a1})] \mapsto \mathbf{rev} \circ [\mathbf{dist1}, \mathbf{S}(\mathbf{s2})]$$

At this point, one might suppose that  $\mathbf{rev}$  can be eliminated as in the previous example, since only one of the functions in the sequence to be reversed can produce an error and therefore the evaluation order of the elements of the sequence is irrelevant:

$$\mathbf{rev} \circ [\mathbf{dist1}, \mathbf{S}(\mathbf{s2})] \stackrel{?}{\mapsto} [\mathbf{S}(\mathbf{s2}), \mathbf{dist1}]$$

Now, however, something has gone wrong:  $\mathbf{rev} \circ [\mathbf{dist1}, \mathbf{t1} \circ \mathbf{a1}] \not\equiv [\mathbf{s2}, \mathbf{dist1}]$ , because  $\mathbf{rev} \circ [\mathbf{dist1}, \mathbf{t1} \circ \mathbf{a1}] : 3$  produces  $\mathbf{dist1\_err}$ , whereas  $\mathbf{rev} \circ [\mathbf{s2}, \mathbf{dist1}] : 3$  produces  $\mathbf{s2\_err}$  (which isn’t even mentioned in the original program). The fact that intuition can fail on such a small and simple example is strong motivation to provide a precise formalism for stating and verifying these transformations.

### 3.2 The Formal Development of Safe

A first step towards formalizing the rewrite rules is to express qualified laws such as (7) equationally. Recall the intuition that a function  $\mathbf{f}$  is safe if  $\mathbf{f}$  cannot produce errors in the context in which it appears. In other words, the fact that  $\mathbf{f} : \mathbf{x}$  is an error is irrelevant if the context of  $\mathbf{f}$  guarantees that  $\mathbf{f}$  cannot be applied to  $\mathbf{x}$ . One way to make this intuition explicit is to map all error values to a single “don’t care”

value. In the following definition of  $\text{Safe}$ ,  $\perp$  is chosen as the “don’t care” value in addition to its role representing non-termination.<sup>1</sup>

DEFINITION 3.1. For every FL function  $f$ ,  $\text{Safe}:f$  denotes the function:

$$\text{Safe}:f:x = \begin{cases} x & \text{if } x \in \mathcal{E}_{\text{FL}} \\ \perp & \text{if } x \notin \mathcal{E}_{\text{FL}} \text{ and } f:x \in \mathcal{E}_{\text{FL}} \\ f:x & \text{if } x \notin \mathcal{E}_{\text{FL}} \text{ and } f:x \notin \mathcal{E}_{\text{FL}} \end{cases}$$

For convenience,  $\text{Safe}:f$  is often abbreviated  $S:f$ . With this definition of  $\text{Safe}$ , Law (7) can be expressed as:

$$S:[g, f] \equiv S:(\text{rev} \circ [f, g]) \quad (8)$$

Moreover, many other useful laws are expressible:

$$f \equiv S:f \circ \text{dom}:f \quad (9)$$

$$f \circ S:(\text{dom}:f) \equiv S:f \quad (10)$$

$$S:f \circ S:g \equiv S:(f \circ g) \quad (11)$$

$$S:p \rightarrow S:q; S:r \equiv S:(p \rightarrow q; r) \quad (12)$$

$$[f, g] \equiv [f, g \circ S:(\text{dom}:f)] \quad (13)$$

$$[S:f_1, \dots, S:f_n] \equiv S:[f_1, \dots, f_n] \quad (14)$$

$$\text{rev} \circ S:[f_1, \dots, f_n] \equiv S:[f_n, \dots, f_1] \quad (15)$$

$$S:(s_1 \circ t_1) \equiv S:s_2 \quad (16)$$

$$\text{catch}: \langle S:f, g \rangle \equiv S:f \quad (17)$$

$$\text{catch}: \langle f, g \rangle \circ S:(\text{dom}:h) \equiv \text{catch}: \langle f \circ S:(\text{dom}:h), g \rangle \quad (18)$$

Note that some of the informal  $\mapsto$  steps have been captured as equivalences; e.g., Laws (14) and (15). Unfortunately, others cannot be expressed as equivalences. For example, the rule  $t_1 \circ S:a_1 \mapsto S:s_2$  cannot be written as  $t_1 \circ S:a_1 \equiv s_2$ , since  $(t_1 \circ S:a_1):\langle 1, 2, 3 \rangle \equiv \perp$  whereas  $s_2:\langle 1, 2, 3 \rangle \equiv 2$ ; nor could it be  $t_1 \circ S:a_1 \equiv S:s_2$ , for the same reason.

The difficulty is that  $\equiv$  is a symmetric relation, whereas the desired property is inherently asymmetric:  $t_1 \circ S:a_1$  can be rewritten to  $S:s_2$ , because in any program in which  $t_1 \circ S:a_1$  could appear, the function  $S:s_2$  produces the same result. Thus, at least some of the desired rewrite rules  $f \mapsto g$  are valid only when  $f$  appears in a “good” context; i.e., in one which enforces the condition that  $S:f$  cannot produce  $\perp$ . The following definitions develop the relation  $\triangleright$  (read “rewrites in context to”), which both captures this asymmetry and defines a set of contexts in which such rewrite rules can be applied.

DEFINITION 3.2.

- (1) The set of *simple expressions* is the smallest set of FL functions such that:
  - $s_i$ ,  $t_1$ ,  $\text{rev}$ ,  $a_1$ ,  $\text{dist}t_1$ ,  $\text{id}$ , and  $\sim a$  for all atoms  $a$  are simple expressions.
  - If  $e_1, \dots, e_n$  are simple expressions, then  $e_1 \circ e_2$ ,  $[e_1, \dots, e_n]$ ,  $(e_1 \rightarrow e_2; e_3)$ ,  $\text{dom}:e_1$ ,  $\text{catch}: \langle e_1, e_2 \rangle$ ,  $\text{Safe}:e_1$ , and  $\alpha:e_1$  are simple expressions.

<sup>1</sup>The term “don’t care” is meant to suggest it doesn’t matter what this value is, because the context guarantees that it cannot arise. However, the choice of  $\perp$  is not arbitrary—see Section 3.3.

- (2)  $\mathbf{E}$  is a *simple context* iff  $\mathbf{E}(\mathbf{f})$  is a simple expression for every simple expression  $\mathbf{f}$ .

The simple expressions are first-order functions that can be built from the subset of FL given in Figure 1. Because the simple expressions are first-order, a great deal can be shown about their termination behavior; it is this property that will be exploited below.

It can be assumed that all expressions written by users have no occurrences of **Safe**, because only the language processor introduces and manipulates **Safe** expressions. Because all transformations preserve the meaning of expressions, the transformations are required to be correct only for expressions that also can be written without any occurrences of **Safe**.

DEFINITION 3.3. A simple expression  $\mathbf{f}$  is *user-definable* iff there exists a simple expression  $\mathbf{u}$  with no occurrences of **Safe** such that  $\mathbf{f} \equiv \mathbf{u}$ .

DEFINITION 3.4. Let  $\mathbf{f}, \mathbf{g}$  be simple expressions.  $\mathbf{f} \triangleright \mathbf{g}$  iff for every simple context  $\mathbf{E}$  such that  $\mathbf{E}(\mathbf{f})$  is user-definable,  $\mathbf{E}(\mathbf{f}) \equiv \mathbf{E}(\mathbf{g})$ .

This definition of  $\triangleright$  allows the expression of a large number of “in-context” transformations. Note that (22) is a correct version of the incorrect rewriting step discussed above.

$$\mathbf{t1} \circ \mathbf{S}: \mathbf{a1} \triangleright \mathbf{s2} \quad (19)$$

$$\mathbf{S}: \mathbf{f} \triangleright \mathbf{f} \quad (20)$$

$$\mathbf{dom}: (\mathbf{S}: \mathbf{f}) \triangleright \mathbf{id} \quad (21)$$

$$\mathbf{rev} \circ [\mathbf{S}: \mathbf{f}, \mathbf{g}] \triangleright [\mathbf{g}, \mathbf{S}: \mathbf{f}] \quad (22)$$

$$\alpha: \mathbf{f} \circ \alpha: (\mathbf{S}: \mathbf{g}) \triangleright \alpha: (\mathbf{f} \circ \mathbf{S}: \mathbf{g}) \quad (23)$$

$$\mathbf{dom}: (\mathbf{S}: (\alpha: \mathbf{f})) \triangleright \alpha: (\mathbf{S}: (\mathbf{dom}: \mathbf{f})) \quad (24)$$

There is one nagging problem. Definition 3.4 provides little assistance in establishing that  $\mathbf{f} \triangleright \mathbf{g}$ , because it requires reasoning about all possible contexts. The purpose of the Substitutability Theorem (given below) is to provide a sufficient condition that is easier to check. The following definition gives this condition; the idea is that  $\mathbf{f}$  should rewrite to  $\mathbf{g}$  if  $\mathbf{f}$  and  $\mathbf{g}$  agree wherever  $\mathbf{f}$  does not return the “don’t care” value.

DEFINITION 3.5.  $\mathbf{f} \leq_{\mathbf{s}} \mathbf{g}$  iff  $\mathbf{f}: \mathbf{x} \equiv \mathbf{g}: \mathbf{x}$  whenever  $\mathbf{f}: \mathbf{x} \neq \perp$ .

The following two lemmas precisely capture the properties of simple expressions that are needed to make the technique of “in-context substitutions” work.

LEMMA 3.6. If  $\mathbf{E}$  is a simple context and  $\mathbf{f} \leq_{\mathbf{s}} \mathbf{g}$ , then  $\mathbf{E}(\mathbf{f}) \leq_{\mathbf{s}} \mathbf{E}(\mathbf{g})$ .

PROOF. First note that if  $\mathbf{E}(\mathbf{f})$  does not contain  $\mathbf{f}$ , then  $\mathbf{E}(\mathbf{f}) = \mathbf{E}(\mathbf{g})$ . The rest of the proof is by induction on the structure of  $\mathbf{E}$ . For the base case, if  $\mathbf{E}(\mathbf{f}) = \mathbf{f}$  and  $\mathbf{E}(\mathbf{g}) = \mathbf{g}$  then the result is immediate since  $\mathbf{f} \leq_{\mathbf{s}} \mathbf{g}$ . In other base cases  $\mathbf{E}$  is a single primitive function not containing  $\mathbf{f}$  (Definition 3.2). The inductive step is proven only for function composition; the other cases are similar. Let  $\mathbf{E}(\mathbf{f}) = \mathbf{E}_1(\mathbf{f}) \circ \mathbf{E}_2(\mathbf{f})$  and assume that  $\mathbf{E}_1(\mathbf{f}): (\mathbf{E}_2(\mathbf{f}): \mathbf{x}) \neq \perp$ . Then  $\mathbf{E}_2(\mathbf{f}): \mathbf{x} \neq \perp$  so  $\mathbf{E}_2(\mathbf{f}): \mathbf{x} \equiv \mathbf{E}_2(\mathbf{g}): \mathbf{x}$  by induction, and therefore  $\mathbf{E}_1(\mathbf{f}): (\mathbf{E}_2(\mathbf{f}): \mathbf{x}) \equiv \mathbf{E}_1(\mathbf{g}): (\mathbf{E}_2(\mathbf{g}): \mathbf{x})$ , also by induction.  $\square$



Note that Lemma 3.6 is false in general for non-simple contexts. For example, let  $\mathbf{f}$  and  $\mathbf{g}$  be functions such that  $\mathbf{f} \leq_{\mathbf{s}} \mathbf{g}$  and  $\mathbf{f} \neq \mathbf{g}$ . Then  $\sim \mathbf{f} \leq_{\mathbf{s}} \sim \mathbf{g}$  does not hold.

LEMMA 3.7. If  $\mathbf{e}$  is a user-definable simple expression, then  $\mathbf{e}: \mathbf{x} \neq \perp$  whenever  $\mathbf{x} \neq \perp$ .

PROOF. Let  $\mathbf{e} \equiv \mathbf{u}$  where  $\mathbf{u}$  contains no occurrence of **Safe**. The proof is a simple induction on the structure of  $\mathbf{u}$  using part (1) of Definition 3.2.  $\square$

Lemma 3.7 provides the reason for working with simple expressions. Informally, user-definable simple expressions do not return  $\perp$  unless applied to  $\perp$ . Now, program transformation may introduce a function  $\mathbf{S}: \mathbf{f}$  into a user-definable simple expression  $\mathbf{E}(\mathbf{S}: \mathbf{f})$ , and it may be the case that  $(\mathbf{S}: \mathbf{f}): \mathbf{x} = \perp$ . Because program transformation must preserve the meaning of the simple expression and because expressions are strict in  $\perp$ , we know that the function  $\mathbf{S}: \mathbf{f}$  can never be applied to  $\mathbf{x}$  in the context  $\mathbf{E}(\mathbf{S}: \mathbf{f})$ . In short, a use of **Safe**:  $\mathbf{f}$  within a user-definable simple expression indicates that  $\mathbf{f}$  can never produce an exception in that context. This intuitive discussion is formalized by the Substitutability Theorem, which reduces the problem of verifying that  $\mathbf{f} \triangleright \mathbf{g}$  to the easier problem of checking that  $\mathbf{f} \leq_{\mathbf{s}} \mathbf{g}$  in the case that  $\mathbf{f}$  appears in a simple context.

THEOREM 3.8. If  $\mathbf{f}$  and  $\mathbf{g}$  are simple expressions and  $\mathbf{f} \leq_{\mathbf{s}} \mathbf{g}$ , then  $\mathbf{f} \triangleright \mathbf{g}$ .

PROOF. Let  $\mathbf{E}$  be a simple context and assume that  $\mathbf{E}(\mathbf{f})$  is user-definable. By Lemma 3.6,  $\mathbf{E}(\mathbf{f}) \leq_{\mathbf{s}} \mathbf{E}(\mathbf{g})$ . By Lemma 3.7,  $\mathbf{E}(\mathbf{f}): \mathbf{x} \neq \perp$  if  $\mathbf{x} \neq \perp$ . Together these facts imply that  $\mathbf{E}(\mathbf{f}): \mathbf{x} \equiv \mathbf{E}(\mathbf{g}): \mathbf{x}$  if  $\mathbf{x} \neq \perp$ . By strictness,  $\mathbf{E}(\mathbf{f}): \perp \equiv \perp \equiv \mathbf{E}(\mathbf{g}): \perp$ . Therefore,  $\mathbf{E}(\mathbf{f}) \equiv \mathbf{E}(\mathbf{g})$ .  $\square$

Using the Substitutability Theorem, transformations (19)-(24) are easily verified. Note that the proof of the Substitutability Theorem depends only on Lemmas 3.6 and 3.7; therefore, this approach works with any extension of the simple expressions that preserves these two Lemmas. Also note that pure identities, such as transformations (1)-(5) and (8)-(17), apply to *all* expressions, not merely the simple expressions.

### 3.3 Beyond Simple Expressions

At this point, it is natural to ask why the set of expressions to which the Substitutability Theorem can be applied must be restricted at all. We conjecture that no restriction is needed but, in fact, achieving an extension that works for all expressions appears to be a difficult technical challenge. This section explains the difficulty.

For the theory presented in Section 3.2 to work, it is necessary to distinguish between  $\perp$ 's role as the “don't care” value and  $\perp$ 's role representing non-termination. This is accomplished by restricting attention to the simple expressions, which are guaranteed to terminate. Thus, the straightforward approach to generalizing the Safe mechanism to all expressions is to use something other than  $\perp$  as the value to which **Safe** maps all errors. Even better, this value should be a new element of the domain, not available to programmers. It would then be easy to distinguish “don't care” values from values that a programmer can define.

Let  $\dagger$  be a new value and let **Safe** be redefined as follows:

DEFINITION 3.9. For every FL function  $\mathbf{f}$ ,  $\mathbf{Safe}:\mathbf{f}$  denotes the function:

$$\mathbf{Safe}:\mathbf{f}:\mathbf{x} = \begin{cases} \mathbf{x} & \text{if } \mathbf{x} \in \mathcal{E}_{\text{FL}} \\ \dagger & \text{if } \mathbf{x} \notin \mathcal{E}_{\text{FL}} \text{ and } \mathbf{f}:\mathbf{x} \in \mathcal{E}_{\text{FL}} \\ \mathbf{f}:\mathbf{x} & \text{if } \mathbf{x} \notin \mathcal{E}_{\text{FL}} \text{ and } \mathbf{f}:\mathbf{x} \notin \mathcal{E}_{\text{FL}} \end{cases}$$

Clearly,  $\dagger$  must be included somewhere in the domain ordering. Herein lies the problem—it appears that the only possibility is that  $\dagger = \perp$ ! The reasoning goes as follows:

The proof of Lemma 3.6 and Law 10 (among others) depend on the fact that  $\dagger$ , once produced, cannot be “changed” by the surrounding context. Thus, all expressions must be strict in  $\dagger$ .

Let  $\mathbf{f}$  be any function and  $\mathbf{x}$  any normal value such that  $\mathbf{f}:\mathbf{x}$  is an error value. Then  $(\mathbf{S}:\mathbf{f}):\mathbf{x} = \dagger$  by Definition 3.9. Now,  $(\mathbf{S}:\mathbf{f}):\perp = \perp$  by strictness, which is also taken into account in Definition 3.9. Clearly,  $\perp \leq \mathbf{x}$ , so by monotonicity of application  $\perp \leq \dagger$ .

Let  $\mathbf{K}(\mathbf{y})$  be the FL function such that  $\mathbf{K}(\mathbf{y}):\mathbf{x} = \mathbf{y}$  for all normal values  $\mathbf{x}$ . Let  $\mathbf{x}$  be any normal value. Then  $\dagger \leq \mathbf{x}$  since  $\mathbf{K}(\dagger) = \mathbf{S}:\mathbf{K}(\perp) \leq \mathbf{S}:\mathbf{K}(\mathbf{x}) = \mathbf{K}(\mathbf{x})$ .

Finally, note that for any normal value  $\mathbf{x}$ , we have  $\mathbf{K}(\perp):\mathbf{x} = \perp$  and  $\mathbf{K}(\perp):\dagger = \dagger$  by strictness of  $\dagger$ . Since  $\dagger \leq \mathbf{x}$ , it follows by monotonicity that  $\dagger \leq \perp$ . But  $\perp \leq \dagger$  as well. Therefore  $\perp = \dagger$ .

The argument presented above is not a complete proof that a distinct  $\dagger$  cannot be included in the domain, because it rests on several assumptions that could be modified. Unfortunately, there is very little room for modifications that might provide a solution within the framework of standard domain theory. For example, we have used the fact that FL functions are strict in  $\perp$ . However, the argument doesn’t use the fact that all functions are strict, just that there is at least one strict function in the domain. Thus, the problem is not strictness vs. laziness—the argument applies to non-strict languages as well.

After considerable experimentation, we have not discovered any way to extend the **Safe** mechanism to all user-definable expressions. We conjecture that it cannot be done within standard domain theory; some quite different construction appears to be required. It is surprising (at least to us) that the simple concept of a result that “can’t happen” during evaluation is difficult to formalize.

#### 4. USING SAFE

This section shows the usefulness of **Safe** with a few short examples illustrating the elimination of function calls, the use of **Safe** in code generation, the optimization of exception handling, and the use of in-context laws. Recall that  $\mathbf{dom}:\mathbf{f}:\mathbf{x}$  is  $\mathbf{f}:\mathbf{x}$  if  $\mathbf{f}:\mathbf{x}$  is an error value and  $\mathbf{x}$  otherwise (see Figure 1). The examples use the following laws involving  $\mathbf{dom}$ :

$$\mathbf{dom}:\mathbf{id} \equiv \mathbf{id} \tag{25}$$

$$\mathbf{dom}:[\mathbf{f}_1, \dots, \mathbf{f}_n] \equiv \mathbf{dom}:\mathbf{f}_n \circ \dots \circ \mathbf{dom}:\mathbf{f}_1 \tag{26}$$

$$\mathbf{dom}:\mathbf{al} \circ \mathbf{dom}:\mathbf{dist1} \equiv \mathbf{dom}:\mathbf{dist1} \tag{27}$$

$$\mathbf{S}:(\mathbf{dom}:\mathbf{t1}) \equiv \mathbf{S}:(\mathbf{dom}:\mathbf{s1}) \tag{28}$$

In the first example, nothing is known about  $f$ , but the fact that  $id$  is a total function allows the construction of the two functions to be reversed.

$$\begin{aligned}
& \mathbf{rev} \circ [f, id] \\
\equiv & \mathbf{rev} \circ S:[f, id] \circ \mathbf{dom}:[f, id] && \text{by 9} \\
\equiv & S:[id, f] \circ \mathbf{dom}:[f, id] && \text{by 15} \\
\equiv & S:[id, f] \circ \mathbf{dom}:id \circ \mathbf{dom}:f && \text{by 26} \\
\equiv & S:[id, f] \circ id \circ \mathbf{dom}:f && \text{by 25} \\
\equiv & S:[id, f] \circ \mathbf{dom}:f && \text{by 2} \\
\equiv & S:[id, f] \circ \mathbf{dom}:f \circ id && \text{by 1} \\
\equiv & S:[id, f] \circ \mathbf{dom}:f \circ \mathbf{dom}:id && \text{by 25} \\
\equiv & S:[id, f] \circ \mathbf{dom}:[id, f] && \text{by 26} \\
\equiv & [id, f] && \text{by 9}
\end{aligned}$$

This transformation is an optimization, because the end result eliminates the application of  $\mathbf{rev}$ . Note, however, that the intermediate steps are not necessarily optimizations, because they involve computing some values twice; in particular,  $f$  could be arbitrarily expensive to compute. The law  $f \equiv S:f \circ \mathbf{dom}:f$  is very useful for introducing safe functions, but if the added  $\mathbf{dom}$  cannot be discharged, then the resulting program could be less efficient than the original program. In the worst case, using Law 9 could result in a program that computes  $f$  once very slowly to preserve errors and then once again to produce the result! Section 6 shows how this is avoided in practice in the optimizing compiler for FL.

The second example illustrates that there are additional advantages to marking functions safe.

$$\begin{aligned}
& \mathbf{rev} \circ [dist1, al] \\
\equiv & \mathbf{rev} \circ S:[dist1, al] \circ \mathbf{dom}:[dist1, al] && \text{by 9} \\
\equiv & S:[al, dist1] \circ \mathbf{dom}:[dist1, al] && \text{by 15} \\
\equiv & S:[al, dist1] \circ \mathbf{dom}:al \circ \mathbf{dom}:dist1 && \text{by 26} \\
\equiv & S:[al, dist1] \circ \mathbf{dom}:dist1 && \text{by 27} \\
\equiv & [S:al, S:dist1] \circ \mathbf{dom}:dist1 && \text{by 14}
\end{aligned}$$

In this case, the end result is an optimization not only because the application of  $\mathbf{rev}$  is eliminated, but also because  $\mathbf{dom}:dist1$  permits the rest of the program to be executed without checking the arguments of any of the functions. The use of **Safe** makes it easy for a code generator to take advantage of this fact. When a primitive is marked as safe, a code generator can produce a version of the primitive that does not check its argument; in this example, both  $al$  and  $dist1$  can run unchecked.

As noted above, it is in general undesirable to use Law 9 because it duplicates computation. In this particular example, it is worth factoring  $\mathbf{dom}:dist1$  out of the computation, because  $\mathbf{dom}:dist1$  can be replaced by a function that merely checks whether the argument is a pair of which the second component is a sequence. This single check is more efficient than the two checks performed by the primitive functions  $dist1$  and  $al$  in the original expression.

The third example presents a more substantial optimization (similar to loop-jamming optimizations for imperative languages [Aho et al. 1986]) and illustrates the use of in-context laws:

$$\begin{aligned}
& [\alpha: \mathbf{s1}, \alpha: + \circ \alpha: \mathbf{t1}] \\
\equiv & [\alpha: \mathbf{s1}, \alpha: + \circ \alpha: \mathbf{t1} \circ \mathbf{S}: (\mathbf{dom}: (\alpha: \mathbf{s1}))] \text{ by 13} \\
\equiv & [\alpha: \mathbf{s1}, \alpha: + \circ \alpha: \mathbf{t1} \circ \alpha: (\mathbf{S}: (\mathbf{dom}: \mathbf{s1}))] \text{ by 24} \\
\equiv & [\alpha: \mathbf{s1}, \alpha: + \circ \alpha: (\mathbf{t1} \circ \mathbf{S}: (\mathbf{dom}: \mathbf{s1}))] \text{ by 23} \\
\equiv & [\alpha: \mathbf{s1}, \alpha: + \circ \alpha: (\mathbf{t1} \circ \mathbf{S}: (\mathbf{dom}: \mathbf{t1}))] \text{ by 28} \\
\equiv & [\alpha: \mathbf{s1}, \alpha: + \circ \alpha: (\mathbf{S}: \mathbf{t1})] \text{ by 10} \\
\equiv & [\alpha: \mathbf{s1}, \alpha: (+ \circ \mathbf{S}: \mathbf{t1})] \text{ by 23}
\end{aligned}$$

Even though the second, third, and last steps use in-context transformations, these steps are actually equivalences, because they occur in simple contexts. Since the first and last lines are equivalent they can be substituted freely one for the other in any program. Note that this shows that it is not necessary that the entire program be simple for an in-context law to apply—it is sufficient that in-context laws be used within simple sub-expressions.

The final example illustrates how **Safe** is used to optimize exception handling. Suppose a programmer defines a function `newt1` that returns the empty sequence whenever `t1` would return an error. A simple definition of `newt1` is `catch: <t1, []>`. If `newt1` appears in a context where it always gets an argument in the proper domain of `t1`, `newt1` can be transformed to `S:t1` as follows:

$$\begin{aligned}
& \mathbf{newt1} \circ \mathbf{S}: (\mathbf{dom}: \mathbf{t1}) \\
\equiv & \mathbf{catch}: \langle \mathbf{t1}, [] \rangle \circ \mathbf{S}: (\mathbf{dom}: \mathbf{t1}) \text{ by def. of newt1} \\
\equiv & \mathbf{catch}: \langle \mathbf{t1} \circ \mathbf{S}: (\mathbf{dom}: \mathbf{t1}), [] \rangle \text{ by 18} \\
\equiv & \mathbf{catch}: \langle \mathbf{S}: \mathbf{t1}, [] \rangle \text{ by 10} \\
\equiv & \mathbf{S}: \mathbf{t1} \text{ by 17}
\end{aligned}$$

## 5. ON SAFE AND DOM

This paper takes an algebraic approach to program optimization, expressing source-to-source program transformations as algebraic laws. There are practical limitations to this approach. This section discusses some of these limitations, particularly the difficulties that arise from the relationship between **Safe** and **dom**.

Because the function **Safe** is not available to programmers, a mechanism is required for introducing **Safe** annotations into a program. Law 9 is a general rule that introduces **Safe**. Recall that the law is:

$$\mathbf{f} \equiv \mathbf{S}: \mathbf{f} \circ \mathbf{dom}: \mathbf{f}$$

Law 9 expresses a fundamental principle of the algebra of **Safe** and, in principle, it can be used to introduce **Safe** wherever needed. From a pragmatic point of view, however, using Law 9 is a disaster. To understand why, consider how a program optimization system based on source-to-source program transformation works. Very briefly, the compiler repeatedly selects program transformations to apply to the program. Ideally, the compiler is able to judge the change in program cost (say, execution time) that results from a transformation, and the compiler selects transformations that decrease program cost. This high-level description ignores several practical problems, of which one of the most significant is that it

is hard to estimate the true value of transformation at the time it is applied. In particular, it is very difficult to include in the estimate of the value of a particular transformation  $T$  the value of additional transformations that  $T$  enables.

Law 9 aggravates this problem, since it duplicates the computation of an arbitrarily expensive function. Thus, application of Law 9 is purely speculative—at the time it is applied, it is not known whether the program is ultimately improved. In our opinion, it is desirable to guarantee that an optimizing compiler does not, at the very least, make programs worse. A simple way to achieve this guarantee is to ensure that each individual transformation proceeds “downhill” in the direction of improvement. Introducing arbitrary “uphill” steps in the transformation process removes this guarantee both in theory and, very often in our experience, in practice.

Section 6 presents a program analysis that conservatively infers where **Safe** can be inserted in a program without also introducing **dom**. Thus, this analysis avoids the “uphill” step of applying Law 9.

There is another aspect to the relationship between **Safe** and **dom**. By introducing uses of **dom**, it is possible to replace almost all in-context laws by equations. For example, Law 19 can be expressed equationally as:

$$t1 \circ S:a1 \equiv s2 \circ \text{dom}:(S:a1)$$

While this equation has the advantage of avoiding the additional concept of in-context transformations, it has the disadvantage that it is not clearly an optimization, whereas Law 19 clearly is. However, Law 19 can be recovered by further transformations:

$$\begin{aligned} & s2 \circ \text{dom}:(S:a1) \\ \triangleright & s2 \circ \text{id} && \text{by (21)} \\ \equiv & s2 && \text{by (1)} \end{aligned}$$

Note that this reasoning requires the use of an in-context law to remove the application of **dom**. Thus, it is not possible to eliminate in-context transformations, although in principle their use could be restricted to applications of Law 21. Formulating the laws in this way only makes sequences of transformations longer without reducing the number of rules; for this reason, we prefer the more direct versions of in-context laws given in Section 3.

There are instances where the introduction of **dom** is useful. As shown above, a function of the form  $\text{dom}:(S:f)$  or (equivalently)  $S:(\text{dom}:f)$  has little or no cost, since it always can be eliminated in favor of **id**. It is sometimes useful to introduce such a function to transfer context information between portions of a program (e.g., Law 13). Finally, there are even limited instances where an unsafe **dom** is introduced, which may or may not be eliminated by subsequent program transformation; an example is given on page 11 in Section 4. However, unlike that example, such introductions are tightly controlled by our compiler and limited to special cases where the transformation can still be shown to be “downhill.”

## 6. INTRODUCING SAFE

This section describes a program analysis method that conservatively identifies functions that can be annotated with **Safe**. That is, the analysis proves theorems of the form  $E(f) \equiv E(S:f)$ . The program analysis is a general type inference system,

but the details of type inference are not important to understand how **Safe** is introduced. This section states properties satisfied by the type inference system and shows how these properties are sufficient to introduce **Safe**. The interested reader is referred to [Aiken and Wimmers 1993; Aiken et al. 1994] for more information on the type inference algorithm.

*Types* are certain subsets of the domain  $\mathcal{D}_{\text{FL}}$ . The full definition of type requires more development [MacQueen et al. 1984], but the only property required here is that a type is a set of values. A *typed function* has the form  $\mathbf{f} :: \mathbf{A} \rightarrow \mathbf{B}$  where  $\mathbf{A}$  and  $\mathbf{B}$  are types and  $\mathbf{A} \subseteq \mathcal{D}_{\text{FL}}^+$  and  $\mathbf{B} \subseteq \mathcal{D}_{\text{FL}}$ . The notation  $\mathbf{f} :: \mathbf{A} \rightarrow \mathbf{B}$  is read “ $\mathbf{f}$  has type  $\mathbf{A} \rightarrow \mathbf{B}$ .” A formal semantics for the phrase  $\mathbf{f} :: \mathbf{A} \rightarrow \mathbf{B}$  is given below. For the moment, the intuition is that  $\mathbf{f} :: \mathbf{A} \rightarrow \mathbf{B}$  is an assertion saying two things:

- (1) For any  $\mathbf{x} \in \mathbf{A}$ , it is the case that  $\mathbf{f}:\mathbf{x} \in \mathbf{B}$ .
- (2) The function  $\mathbf{f}$  appears in a context where the only normal values it is applied only to are in  $\mathbf{A}$ .

For example, consider the following:

$$\begin{aligned} \text{NonEmptySeq} &= \{\langle x_1, \dots, x_n \rangle \mid n \geq 1, x_i \in \mathcal{D}_{\text{FL}}^+\} \\ \text{NotNonEmptySeq} &= \mathcal{D}_{\text{FL}}^+ - \text{NonEmptySeq} \end{aligned}$$

Then  $\mathbf{s1} :: \text{NonEmptySeq} \rightarrow \mathcal{D}_{\text{FL}}^+$  is interpreted as saying that  $\mathbf{s1}$  appears in a context where it is guaranteed to be applied to a non-empty sequence and (therefore) to produce a non-error value. Similarly,  $\mathbf{s1} :: \text{NotNonEmptySeq} \rightarrow \mathcal{E}_{\text{FL}}$  is interpreted as saying that  $\mathbf{s1}$  appears in a context where it is guaranteed to be applied to something that is not a non-empty sequence and (therefore) is guaranteed to produce an error. Finally,  $\mathbf{s1} :: \mathcal{D}_{\text{FL}}^+ \rightarrow \mathcal{D}_{\text{FL}}$  says that nothing is known about the context in which  $\mathbf{s1}$  appears.

The informal part of the description of  $\mathbf{f} :: \mathbf{A} \rightarrow \mathbf{B}$  above is the idea that  $\mathbf{f}$  appears in a context where it is applied to elements of  $\mathbf{A}$ . Projections are a well-known technique for making precise the idea of a “context” [Wadler and Hughes 1987].

DEFINITION 6.1. A *projection* is an idempotent function less than the identity.

The function  $\Delta:\mathbf{A}$  is a projection for any type  $\mathbf{A}$ :

$$\Delta:\mathbf{A}:\mathbf{x} = \begin{cases} \mathbf{x} & \text{if } \mathbf{x} \in \mathbf{A} \\ \perp & \text{otherwise} \end{cases}$$

Because  $\Delta:\mathbf{A}:\mathbf{x}$  is either  $\mathbf{x}$  or  $\perp$  for all  $\mathbf{x}$ , it is easy to see that  $\Delta:\mathbf{A} \leq \text{id}$  and that  $\Delta:\mathbf{A}$  is idempotent. Therefore,  $\Delta:\mathbf{A}$  is in fact a projection.<sup>2</sup> The formal meaning of  $\mathbf{f} :: \mathbf{A} \rightarrow \mathbf{B}$  is given in terms of projections.

DEFINITION 6.2. The expression  $\mathbf{f} :: \mathbf{A} \rightarrow \mathbf{B}$  is the function  $\text{Propagate}:(\Delta:\mathbf{B} \circ \mathbf{f} \circ \Delta:\mathbf{A})$ , where  $\text{Propagate}:\mathbf{f}:\mathbf{x}_{\text{err}} = \mathbf{x}_{\text{err}}$  and  $\text{Propagate}:\mathbf{f}:\mathbf{x} = \mathbf{f}:\mathbf{x}$  otherwise.

<sup>2</sup>A technical note: The function  $\Delta:\mathbf{A}$  is not necessarily continuous if  $\mathbf{A}$  is an arbitrary set of values. For any type  $\mathbf{A}$ , however,  $\Delta:\mathbf{A}$  is continuous. Note also that we abuse notation by allowing  $\mathbf{A}$  to contain errors even though  $\Delta:\mathbf{A}$  is not strict in the error values of  $\mathbf{A}$ .

This definition formalizes the intuition given above by defining  $\mathbf{f} :: \mathbf{A} \rightarrow \mathbf{B}$  to be a function that can produce a non-bottom result only when it is applied to some  $\mathbf{x} \in \mathbf{A}$  such that  $\mathbf{f}:\mathbf{x} \in \mathbf{B}$ . The functional **Propagate** simply ensures that if the argument to the typed function is an exception that exception is propagated in keeping with the normal semantics of FL. A *typed expression* is an FL expression that may include typed functions as subexpressions. The following definition explains what it means for the types on subexpressions to be correct.

**DEFINITION 6.3.** Let  $\mathbf{E}$  be a typed expression and let  $\mathbf{E}'$  be  $\mathbf{E}$  with all types erased (i.e.,  $\mathbf{f} :: \mathbf{A} \rightarrow \mathbf{B}$  in  $\mathbf{E}$  is replaced by  $\mathbf{f}$  in  $\mathbf{E}'$ ). Then  $\mathbf{E}$  is *well-typed* iff  $\mathbf{E} \equiv \mathbf{E}'$ .

The FL type inference system takes an FL expression and produces a well-typed expression. For example, consider the program  $\mathbf{s1} \circ [\sim 1, \sim 2]$ . The inference system produces the well-typed program

$$(\mathbf{s1} :: \mathbf{Pair}(\mathbf{Int}) \rightarrow \mathbf{Int}) \circ [\sim 1 :: \mathcal{D}_{\text{FL}}^+ \rightarrow \mathbf{Int}, \sim 2 :: \mathcal{D}_{\text{FL}}^+ \rightarrow \mathbf{Int}].$$

The type  $\mathbf{Int}$  is the set of integers and  $\mathbf{Pair}(\mathbf{X})$  is all sequences of length 2 with elements drawn from  $\mathbf{X}$ . The system also infers types for the higher-functions composition and sequence construction, but these types are omitted for readability. The following lemma states the main result of this section: if  $\mathbf{E}(\mathbf{f} :: \mathbf{A} \rightarrow \mathbf{B})$  is well-typed and the type  $\mathbf{B}$  contains no error values, then  $\mathbf{f}$  can be replaced by  $\mathbf{S}:\mathbf{f}$ .

**LEMMA 6.4.** Let  $\mathbf{E}(\mathbf{f} :: \mathbf{A} \rightarrow \mathbf{B})$  be well-typed and let  $\mathbf{B} \cap \mathcal{E}_{\text{FL}} \subseteq \{\perp\}$ . Then  $\mathbf{E}(\mathbf{f}) \equiv \mathbf{E}(\mathbf{S}:\mathbf{f})$ .

**PROOF.** We prove that  $\mathbf{S}:\mathbf{f} \geq \mathbf{f} :: \mathbf{A} \rightarrow \mathbf{B}$ . The result then follows because  $\mathbf{E}(\mathbf{f}) \geq \mathbf{E}(\mathbf{S}:\mathbf{f}) \geq \mathbf{E}(\mathbf{f} :: \mathbf{A} \rightarrow \mathbf{B})$  and  $\mathbf{E}(\mathbf{f}) \equiv \mathbf{E}(\mathbf{f} :: \mathbf{A} \rightarrow \mathbf{B})$  since  $\mathbf{E}(\mathbf{f} :: \mathbf{A} \rightarrow \mathbf{B})$  is well-typed.

If  $(\mathbf{f} :: \mathbf{A} \rightarrow \mathbf{B}):\mathbf{x} \equiv \perp$  there is nothing to prove. Assume  $(\mathbf{f} :: \mathbf{A} \rightarrow \mathbf{B}):\mathbf{x} \not\equiv \perp$ . If  $\mathbf{x} \in \mathcal{E}_{\text{FL}}$  then  $(\mathbf{f} :: \mathbf{A} \rightarrow \mathbf{B}):\mathbf{x} \equiv \mathbf{x}$  and  $\mathbf{S}:\mathbf{f}:\mathbf{x} \equiv \mathbf{x}$ . Now assume that  $\mathbf{x} \notin \mathcal{E}_{\text{FL}}$ . Then  $(\mathbf{f} :: \mathbf{A} \rightarrow \mathbf{B}):\mathbf{x} = (\Delta:\mathbf{B} \circ \mathbf{f} \circ \Delta):\mathbf{x} \not\equiv \perp$ . Therefore,  $\Delta:\mathbf{A}:\mathbf{x} \equiv \mathbf{x}$  and  $\mathbf{f}:\mathbf{x} \not\equiv \perp$  and  $\Delta:\mathbf{B}:(\mathbf{f}:\mathbf{x}) \equiv \mathbf{f}:\mathbf{x}$ . Since  $\mathbf{B} \cap \mathcal{E}_{\text{FL}} \subseteq \{\perp\}$  and  $\mathbf{f}:\mathbf{x} \in \mathbf{B} - \{\perp\}$ , it follows that  $\mathbf{f}:\mathbf{x} \notin \mathcal{E}_{\text{FL}}$ . Therefore  $\mathbf{S}:(\mathbf{f}:\mathbf{x}) = \mathbf{f}:\mathbf{x} = (\mathbf{f} :: \mathbf{A} \rightarrow \mathbf{B}):\mathbf{x}$ .  $\square$

Returning to the example above, Lemma 6.4 shows that

$$(\mathbf{s1} :: \mathbf{Pair}(\mathbf{Int}) \rightarrow \mathbf{Int}) \circ [\sim 1 :: \mathcal{D}_{\text{FL}}^+ \rightarrow \mathbf{Int}, \sim 2 :: \mathcal{D}_{\text{FL}}^+ \rightarrow \mathbf{Int}] \equiv \mathbf{S}:\mathbf{s1} \circ [\mathbf{S}:\sim 1, \mathbf{S}:\sim 2].$$

Type inference must be very accurate to be useful for proving that functions are safe. For example, in the example with  $\mathbf{s1}$  above, it is necessary to prove not just that  $\mathbf{s1}$  is applied to a sequence but that  $\mathbf{s1}$  is applied to a non-empty sequence. The standard Hindley/Milner type system [Damas and Milner 1982] cannot prove such properties. A type inference system that can infer types accurate enough to be useful in program optimization has been implemented for FL and is described in [Aiken and Wimmers 1993; Aiken et al. 1994].

## 7. SAFE IN AN OPTIMIZING COMPILER FOR FL

The **Safe** technique is used heavily in an optimizing compiler for FL developed at IBM Almaden. The FL compiler makes use of literally hundreds of laws, most of which involve safety information. This section gives a brief qualitative summary of

several years experience, both good and bad, with using **Safe** in the optimization of FL programs.

The overall structure of the FL compiler is simple. After parsing and desugaring, the type inference algorithm described in Section 6 annotates a program with types. Based on the types assigned, functions are then marked as **Safe** using Lemma 6.4. At this point, the system enters a “match-transform” loop in which subexpressions are matched against the left-hand sides of laws. When a match is found, the matched subexpression is replaced by the right-hand side of the law and the entire process repeats. The strategy used to search for matches is not fixed (there are several such strategies in the compiler) but normally the program is scanned repeatedly until no laws match. A few optimization phases that have well-defined algorithms are implemented directly outside of the match-transform paradigm.

The laws of the FL compiler are organized into groups according to purpose. Some laws are aimed at a particular kind of optimization (e.g., simplifying nested conditionals) or are used merely as an adjunct to compiler phases that have a large component implemented in a way other than via laws (e.g., compile-time evaluation). Although the focus of these phases is something other than safety information, almost all compiler phases make some use of safety information.

The primary pass of the compiler that relies on safety information is a set of approximately 200 laws that are best described as “peephole” optimizations. These laws all have the property that (1) they are improvements to the program (“down-hill” transformations) and (2) they reduce the size of the program. The laws are applied repeatedly using a greedy strategy until no more laws apply. Most of the rules in this paper are included in this phase of the compiler. Property (1) guarantees that the result is an optimization of the original program; property (2) guarantees that eventually no more rules apply and the process terminates.

There is another interesting property, namely the Church-Rosser property, that these laws do not enjoy. In our experience, this property is neither crucial nor really practical to obtain for large collections of program transformations. There are several reasons for this conclusion. First, experience suggests it is better to include an optimization rather than exclude it on the grounds that it makes a set of laws non-Church-Rosser. After all, the empty set of laws is Church-Rosser. Second, Church-Rosser does not imply termination of a set of laws and in fact these two properties are sometimes in conflict (e.g., adding a law that makes the system Church-Rosser may destroy termination). We regard the paramount concerns to be optimization and termination (properties (1) and (2) above). Finally, it is an enormous task to prove Church-Rosser properties for large sets of transformations. Even if there were some chance that the laws were Church-Rosser (they are not) it would be prohibitively difficult to prove this fact.

The overall experience with **Safe** has been very good. It is easy to write and add very general laws to the compiler. Without **Safe** many of these laws could not be expressed at all or could be expressed only as an enormous number of special cases. The denotational basis for **Safe** has also proven to be crucial. The example in Section 3.1 of how informal reasoning leads to incorrect transformations is not contrived; on several occasions similar incorrect transformations were added by the authors. Having a denotational basis for reasoning about **Safe** proved indispensable in these cases for identifying the incorrect transformations from among the hundreds



of laws in the compiler.

On an engineering level, **Safe** has also been successful. A key concern in the design of the compiler was the speed of the match-transform step. Because the FL compiler relies heavily on matching and transforming expressions, it is very important that matching common constructs be fast. The function **Safe** was engineered to be fast to match; it is implemented simply as a one-bit flag attached to functions and adds a negligible overhead to the matching process.

Using **Safe** does have two drawbacks, however, both of which are related to the use of in-context transformations. The first problem is that when writing an in-context law it is easy to write a law that is less accurate than desired. To see this, recall that  $S: f \triangleright f$  (Law 20), so it is always permissible to drop **Safe** from a function. Thus, for example,  $S: (s1 \circ t1) \equiv S: s2$  also can be written  $S: (s1 \circ t1) \triangleright s2$ . The disadvantage of the second form is that safety information is lost, so that additional laws are less likely to match, which degrades the quality of program optimization. To see how easy it is to write laws with “safety leaks,” note that in Law 19 the  $s2$  on the right-hand side is in fact **Safe**. Even the authors failed to notice this through many revisions of this paper.

Safety leaks are hard to detect. The only hint of a problem is that a program that should optimize well does not, and lost safety information is not the only potential cause. Losing safety information turned out to be a real but not insurmountable problem in the FL compiler, although considerable time was spent inspecting for and removing safety leaks from laws.

The second problem with in-context laws is that they are not invertible. This is, of course, inherent. The whole point of a law  $f \triangleright g$  is that  $g$  can replace  $f$  but not vice-versa. In-context laws are one-way tickets: if the compiler makes a poor choice in applying an in-context law, there is no way to undo it. This should not be taken as an argument against in-context laws. Such laws are very useful and even necessary for good optimization. The problem is that it may be difficult to decide when to apply in-context laws because other (perhaps better) avenues for optimization may be lost.

The FL experience is that, while this latter problem does occur, the practical effect is small. There are programs for which some optimization is missed because an in-context law is applied that cannot be undone in a later phase. However, these examples are relatively few and the cost in lost optimization is usually small. It was judged not to be worthwhile to make any special effort to selectively apply in-context laws. The FL compiler currently makes no distinction between laws that are equivalences and laws that apply only in-context.

## 8. RELATED WORK, CONCLUSIONS, AND FUTURE WORK

The first mention of qualified or “in context” laws for FP appeared in Backus’ original paper [Backus 1978]. Backus proposed *qualified equational laws*, which consist of an equation between functions  $f \equiv g$  and a predicate  $p$ . The meaning of qualified equational law is that  $f: x = g: x$  whenever  $p: x$  is true. In the context of this paper, the combinator **Safe** replaces the predicate so that equations can be expressed without side conditions.

Safety analysis is related to *projection analysis* [Wadler and Hughes 1987]. The idea of projection analysis is to express formally properties of programs using pro-

jections (see Section 7). The following lemma shows that projections can be used to characterize safety analysis:

LEMMA 8.1.  $S:f \equiv f \circ p$  for some projection  $p$ .

PROOF. The projection  $p$  is  $\text{dom}(S:f)$ . Since  $\text{dom}(S:f):x \equiv x$  or  $\text{dom}(S:f):x \equiv \perp$  for all  $x$ , it follows that  $\text{dom}(S:f) \leq \text{id}$  and that  $\text{dom}(S:f)$  is idempotent. Hence,  $\text{dom}(S:f)$  is in fact a projection.

It is easy to check that  $\text{dom}(S:f) \equiv \text{dom}f \circ \text{dom}(S:f)$ . Then,

$$\begin{aligned}
 S:f & \\
 \equiv S:(S:f) \circ \text{dom}(S:f) & \quad \text{by Law 9} \\
 \equiv S:f \circ \text{dom}(S:f) & \quad S \text{ is idempotent} \\
 \equiv S:f \circ \text{dom}f \circ \text{dom}(S:f) & \\
 \equiv f \circ \text{dom}(S:f) & \quad \text{by Law 9}
 \end{aligned}$$

□

Beyond the connection with the theory of projections, there are other similarities in the use of the two techniques. Both techniques deal with manipulating “annotations”. For projection analysis, these annotations are projections; for safety analysis, the annotation is the function **Safe**. Many of the techniques for manipulating the annotations are also similar.

However, safety analysis and projection analysis are addressed at two different problems. Projection analysis is primarily concerned with determining (in a lazy system) whether a function is strict in its arguments and gives a nice way of addressing that problem. Safety analysis (as presented here for a strict language) is concerned not only with determining when a function is “safe” but also with trying to use that fact to facilitate program transformations. For example, to the best of our knowledge the counterpart of in-context laws for projection analysis has not been developed, although it certainly could be using the techniques presented here.

There is another, perhaps deeper, difference between safety analysis and projection analysis. In the basic formulation of projection analysis, a lifted domain is used with a new element  $\dagger$  below  $\perp$  to which projections map elements outside the domain of interest [Wadler and Hughes 1987]. As discussed in Section 3.3, it appears that no standard domain with **Safe** allows  $\dagger$  to be separated from  $\perp$ . The crux of the difference between the two techniques is that projection analysis is concerned with program analysis only and thus can avoid giving semantics to programs that contain projections. However, **Safe** is intended for use in a program transformation system where programs contain occurrences of **Safe**, and thus it is crucial that programs with **Safe** be given semantics directly.

It would strengthen the theoretical treatment if the restriction to simple expressions of in-context laws could be removed. As discussed in Section 3.3, such a generalization appears to present difficult semantic problems. On a practical level, most of the common uses of **Safe** are covered even with the restriction to simple expressions. However, a model that worked for all expressions would help clarify the semantic role of **Safe** and add some power to the algebra. We leave as future work whether techniques that extend projection analysis to higher-order programs (e.g., [Hunt and Sands 1991]) can be used to generalize safety analysis.

The **Safe** mechanism resolves the tension between the desire to make functional programs run fast through optimization and the desire to have a language in which it is easy to write and debug programs. This tension is perhaps at a maximum in FL, because no distinction is made between user-generated exceptions and system-generated errors—both are legitimate error values. Thus, it would be disastrous for an FL compiler to fail to preserve the error behavior of a program; on the other hand, preserving errors creates problems for optimization. **Safe** solves this dilemma by providing a way to express the program transformations of a language with a single error value in a language with many error values.

#### ACKNOWLEDGMENTS

The congenial atmosphere of the FL group (John Backus, Thom Linden, Peter Lucas, and Paul Tucker) contributed significantly to this work. It is also a pleasure to thank Luca Cardelli, Robert Cartwright, David Chase, Jim Donahue, Joe Halpern, Paul Hudak, Matthias Felleisen, Phil Wadler, Jennifer Widom, and the members of IFIP WG 2.8 for their helpful comments and suggestions. In particular, Joe Halpern's persistent critiques greatly improved the presentation of Section 3. Finally, thanks go to the anonymous referees for many useful suggestions.

#### REFERENCES

- AHO, A. V., SETHI, R., AND ULLMAN, J. D. 1986. *Compilers: Principles, Techniques, and Tools*. Addison-Wesley.
- AIKEN, A., WILLIAMS, J. H., AND WIMMERS, E. L. 1990. Program transformation in the presence of errors. In *Seventeenth Annual ACM Symposium on Principles of Programming Languages*, pp. 210–217.
- AIKEN, A. AND WIMMERS, E. 1993. Type inclusion constraints and type inference. In *Proceedings of the 1993 Conference on Functional Programming Languages and Computer Architecture*, Copenhagen, Denmark, pp. 31–41.
- AIKEN, A., WIMMERS, E., AND LAKSHMAN, T. 1994. Soft typing with conditional types. In *Twenty-First Annual ACM Symposium on Principles of Programming Languages*, Portland, Oregon, pp. 163–173.
- BACKUS, J. 1978. Can programming be liberated from the von Neumann style? A functional style and its algebra of programs. *Commun. ACM* 21, 8 (Aug.), 613–641.
- BACKUS, J., WILLIAMS, J. H., AND WIMMERS, E. L. 1986. The FL language manual. Tech. Rep. RJ 5339 (54809), IBM.
- CARTWRIGHT, R. AND FELLEISEN, M. 1989. The semantics of program dependence. In *Proceedings of the 1989 Conference on Programming Language Design and Implementation*.
- DAMAS, L. AND MILNER, R. 1982. Principle type-schemes for functional programs. In *Ninth Annual ACM Symposium on Principles of Programming Languages*, pp. 207–212.
- HARPER, R., MILNER, R., AND TOFTE, M. 1989. The definition of standard ML—version 3. Tech. Rep. ECFS-LFCS-89-81, Laboratory for Foundations of Computer Science, University of Edinburgh.
- HENNESSY, J. 1981. Program optimization and exception handling. In *Proceedings of the 1981 Symposium on Principles of Programming Languages*.
- HUDAK, P., WADLER, P., ARVIND, BOUTEL, B., FAIRBAIRN, J., FASEL, J., HUGHES, J., JOHNSON, T., KIEBURTZ, D., JONES, S. P., NIKHIL, R., REEVE, M., WISE, D., AND YOUNG, J. 1988. Report on the functional programming language Haskell. Tech. Rep. DCS/RR-666 (Dec.), Yale University.
- HUNT, S. AND SANDS, D. 1991. Binding time analysis: A new PERSpective. In *Symposium on Partial Evaluation and Semantics-Based Program Manipulation*, pp. 154–165.

- MACQUEEN, D., PLOTKIN, G., AND SETHI, R. 1984. An ideal model for recursive polymorphic types. In *Eleventh Annual ACM Symposium on Principles of Programming Languages*, pp. 165–174.
- STEELE, G. L. 1984. *Common Lisp: The Language*. Digital Press.
- TURNER, D. A. 1985. Miranda: A non-strict functional language with polymorphic types. In *Proceedings of the IFIP International Conference on Functional Programming and Computer Architecture*. Springer Verlag Lecture Notes in Computer Science no. 201.
- WADLER, P. AND HUGHES, R. J. M. 1987. Projections for strictness analysis. In *Proceedings of the Symposium on Functional Programming Languages and Computer Architecture*, pp. 385–407. Springer Verlag Lecture Notes in Computer Science no. 274.