

Program Transformation in the Presence of Errors

Alexander Aiken

John H. Williams

Edward L. Wimmers

IBM Almaden Research Center

650 Harry Rd.

San Jose, CA 95120

email: *lastname*@ibm.com

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 often, errors are considered to be outside the pale 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 [HWA*88] and the languages of Turner [Tur85], 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 [Ste84] 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 [CF89]. 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 [Hen81]

“...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 [Bac78]). In fact, the **Safe** mechanism accomplishes much more. It actually strengthens equational reasoning by providing a sufficient condition on a program context $E(\cdot)$ and functions f and g , such that $E(f) \equiv E(g)$ even if $f \neq g$.

The **Safe** mechanism is presented in the context of the functional language FL [BWW86], but it should be applicable to any language whose program transformations can be expressed equationally. 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 6 concludes with suggestions for further work.

$f:x$	denotes function application
$\langle x_1, \dots, x_n \rangle$	denotes sequence construction
$(f \circ g):x$	$= f:(g:x)$
$[f_1, \dots, f_n]:x$	$= \langle f_1:x, \dots, f_n:x \rangle$
$(p \rightarrow q; r):x$	$= \begin{cases} p:x & \text{if } p:x \in \mathcal{E}_{\text{FL}} \\ r:x & \text{if } p:x = \mathbf{false} \\ q:x & \text{otherwise} \end{cases}$
$\tilde{x}:y$	$= x$
$\alpha:f:\langle x_1, \dots, x_n \rangle$	$= \langle f:x_1, \dots, f:x_n \rangle$
$+: \langle x_1, \dots, x_n \rangle$	$= x_1 + \dots + x_n$
$si:\langle x_1, \dots, x_n \rangle$	$= x_i$
$tl:\langle x_1, x_2, \dots, x_n \rangle$	$= \langle x_2, \dots, x_n \rangle$
$rev:\langle x_1, \dots, x_n \rangle$	$= \langle x_n, \dots, x_1 \rangle$
$al:\langle x, \langle y_1, \dots, y_n \rangle \rangle$	$= \langle x, y_1, \dots, y_n \rangle$
$distl:\langle x, \langle y_1, \dots, y_n \rangle \rangle$	$= \langle \langle x, y_1 \rangle, \dots, \langle x, y_n \rangle \rangle$
$id:x$	$= x$
$dom:f$	$= s1 \circ [id, f]$
$mkerr:"string"$	$= string_{err}$
$catch:\langle f, g \rangle:x$	$= \begin{cases} g:\langle x, y \rangle & \text{if } f:x = Err(y) \\ f:x & \text{otherwise} \end{cases}$

Figure 1: A subset of FL.

2 An Overview of FL

FL [BWW86] is the result of an effort to design a practical functional language based on FP [Bac78]. 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, input/output functions and syntactic sugar are omitted. A definition of each function is given only for some arguments; for all other arguments, the function f returns an error value f_{err} (e.g., $s1:0 \equiv s1_{err}$).¹ The evaluation order is leftmost-innermost; thus, in $[f, g]:x$, $f:x$ is evaluated and then $g:x$ is evaluated.

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 [HMT89], 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; that is, $f:x_{err} \equiv x_{err}$ for any function f and error value x_{err} . Sequence construction is also strict with respect to errors; a sequence collapses to the leftmost error it contains. This behavior is justified by

¹In fact, FL functions produce more informative errors than just the name of the function, but this countable set of errors is sufficient to illustrate the technique.

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, Seqs is sequence construction, and Err is error construction. The ordering on $\mathcal{D}_{\text{FL}}^+$ is the standard one; in \mathcal{E}_{FL} , $\text{Err}(x) \leq \text{Err}(y) \Leftrightarrow x \leq y$ and $\perp \leq x$ for all x .

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.

$$\mathbf{f} \circ \text{id} \equiv \mathbf{f} \tag{1}$$

$$\text{id} \circ \mathbf{f} \equiv \mathbf{f} \tag{2}$$

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

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

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

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

$$[\mathbf{g}, \mathbf{f}] \equiv \text{rev} \circ [\mathbf{f}, \mathbf{g}] \tag{6}$$

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 tl_{err} , then $[\mathbf{g}, \mathbf{f}]$ produces tl_{err} , 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. However, it is true that

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

That is, there are *contexts* in which $[g, f]$ can be substituted for $\text{rev} \circ [f, g]$. 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 f is annotated as being “safe” (written $S(f)$), if it is known that that occurrence is guaranteed not to produce an error when applied to a non-error value; the notation $f \mapsto g$ simply indicates that f is rewritten to g . (N.B. For now, $S(f)$ is an extra-linguistic notion; the phrase “annotating f with $S(f)$ ” has no more semantic content than the phrase “painting f green”.)

Consider the program $\text{rev} \circ [\sim 0, \sim 1]$. 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:

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

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 rev can be eliminated. As a last step, the annotation S can be dropped:

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

This rewriting sequence is correct in the sense that the final program is equivalent to the original program; i.e., $\text{rev} \circ [\sim 1, \sim 0] \equiv [\sim 1, \sim 0]$. At this point, however, the safe mechanism is informal, and it is easy to make mistakes. Consider the program $\text{rev} \circ [\text{distl}, \text{tl} \circ \text{al}]$. Because al follows distl in the order of evaluation, and because al produces an error for exactly the same arguments as distl , al can be marked safe.

$$\text{rev} \circ [\text{distl}, \text{tl} \circ \text{al}] \mapsto \text{rev} \circ [\text{distl}, \text{tl} \circ S(\text{al})]$$

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

$$\text{rev} \circ [\text{distl}, \text{tl} \circ S(\text{al})] \mapsto \text{rev} \circ [\text{distl}, S(\text{s2})]$$

At this point, one might suppose that 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:

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

Now, however, something has gone wrong: $\text{rev} \circ [\text{distl}, \text{tl} \circ \text{al}] \neq [\text{s2}, \text{distl}]$, because $\text{rev} \circ [\text{distl}, \text{tl} \circ \text{al}] : 3$ produces $\text{distl}_{\text{err}}$, whereas $\text{rev} \circ [\text{s2}, \text{distl}] : 3$ produces 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.

A first step towards formalizing the rewrite rules is to express qualified laws such as (7) equationally. One way to do this is to force both sides to produce some identical “don't care” value for all arguments not in the domain of interest. In the following definition of **Safe**, \perp plays the role of the “don't care” value.

Definition 3.1 (Safe) For every FL function f , $\text{Safe}:f$ denotes the function:

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

For convenience, $\text{Safe}:f$ is often abbreviated $S : f$. With this definition of **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:

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

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

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

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

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

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

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

$$S : (s1 \circ t1) \equiv S : s2 \quad (16)$$

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

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

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 $\text{tl} \circ S(\text{al}) \mapsto S(\text{s2})$ cannot be written as $\text{tl} \circ S : \text{al} \equiv \text{s2}$, since $(\text{tl} \circ S : \text{al}) : \langle 1, 2, 3 \rangle \equiv \text{al}_{\text{err}}$ whereas $\text{s2} : \langle 1, 2, 3 \rangle \equiv 2$; nor could it be $\text{tl} \circ S : \text{al} \equiv S : \text{s2}$, for the same reason.

The difficulty is that \equiv is a symmetric relation, whereas the desired property is inherently asymmetric: $\text{tl} \circ S : \text{al}$ can be rewritten to $S : \text{s2}$, because in any program in which $\text{tl} \circ S : \text{al}$ could appear, the function $S : \text{s2}$ 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:
 - `si`, `tl`, `rev`, `al`, `distl`, `id`, $\sim a$ for all atoms `a`, and `mkerr` are simple expressions.
 - If e_1, \dots, e_n are simple expressions, then $e_1 \circ e_2$, $[e_1, \dots, e_n]$, $(e_1 \rightarrow e_2; e_3)$, `dom`: e_1 , `catch`: $\langle e_1, e_2 \rangle$, `Safe`: e_1 , and α : e_1 are simple expressions.
2. E is a *simple context* iff $E(f)$ is a simple expression for every simple expression f .

Since only the language processor introduces and manipulates `Safe` expressions, and since transformations preserve the meaning of expressions, the transformations need to work only for expressions that can be written without `Safe`.

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

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

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.

$$tl \circ S:al \triangleright s2 \tag{19}$$

$$S:f \triangleright f \tag{20}$$

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

$$rev \circ [S:f, g] \triangleright [g, S:f] \tag{22}$$

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

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

There is one nagging problem. Definition 3.4 provides little assistance in establishing that $f \triangleright g$, 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 f should rewrite to g if f and g agree wherever f does not return the “don’t care” value.

Definition 3.5 $f \leq_s g$ iff $f:x \equiv g:x$ whenever $f:x \neq \perp$.

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 $\mathbf{f} \leq_s \mathbf{g}$, then $\mathbf{E}(\mathbf{f}) \leq_s \mathbf{E}(\mathbf{g})$.

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

The Substitutability Theorem reduces the problem of verifying that $\mathbf{f} \triangleright \mathbf{g}$ to the easier problem of checking that $\mathbf{f} \leq_s \mathbf{g}$ in the case where $\mathbf{E}(\mathbf{f})$ is a user-definable simple expression.

Theorem 3.8 (Substitutability Theorem)

If \mathbf{f} and \mathbf{g} are simple expressions and $\mathbf{f} \leq_s \mathbf{g}$, 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, $\mathbf{E}(\mathbf{f}) \leq_s \mathbf{E}(\mathbf{g})$. 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)-(17), apply to *all* expressions, not merely the simple expressions.

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 $\text{dom}:\mathbf{f}:\mathbf{x}$ 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 dom :

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

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

$$\text{dom}:\text{al} \circ \text{dom}:\text{distl} \equiv \text{dom}:\text{distl} \tag{27}$$

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

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

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

This transformation is an optimization, because the end result eliminates the application of `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 $f \equiv S:f \circ \text{dom}:f$ is very useful for introducing safe functions, but if the added `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!

The next example illustrates that there are intermediate cases where `dom` cannot be completely discharged but the result is still a useful optimization. Section 6 contains further discussion of the problem of discharging `dom`.

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

In this case, the end result is an optimization not only because the application of `rev` is eliminated, but also because `dom:distl` 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 `distl` can run unchecked.

The following example presents a more substantial optimization (similar to loop-jamming optimizations for imperative languages [ASU86]) and illustrates the use of in-context laws:

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

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 `newtl` that returns the empty sequence whenever `tl` would return an error. A simple definition of `newtl` is `catch: (tl, [])`. Now, if `newtl` appears in a context where it always gets an argument in the proper domain of `tl`, `newtl` can be transformed to `S:tl` as follows:

$$\begin{aligned}
& \text{newtl} \circ \text{S}:(\text{dom}:\text{tl}) \\
\equiv & \text{catch}:(\text{tl}, []) \circ \text{S}:(\text{dom}:\text{tl}) && \text{by def. of newtl} \\
\equiv & \text{catch}:(\text{tl} \circ \text{S}:(\text{dom}:\text{tl}), []) && \text{by 18} \\
\equiv & \text{catch}:(\text{S}:\text{tl}, []) && \text{by 10} \\
\equiv & \text{S}:\text{tl} && \text{by 17}
\end{aligned}$$

5 Related Work

Safety analysis bears some resemblance to *projection analysis* [WH87]. 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 other source-level optimizations.

6 Conclusions and Future Work

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.

The transformations involving **Safe** are useful only if **Safe** functions can be introduced by a compiler. While it remains to be proven that **Safe** enables large scale optimization in practice, one practical problem is readily apparent from the examples in this paper. As noted in Section 4, it is desirable to introduce **Safe** without using the law $f \equiv \text{S}:\text{f} \circ \text{dom}:\text{f}$, which requires computing `f` twice. In this law, `dom:f` ensures that the *type* of `S:f` is such that `S:f` cannot produce an error. Thus,

discharging dom:f is a type inference problem—a type inference algorithm can check whether \mathbf{f} is safe before Law 9 is applied, and if \mathbf{f} is safe the dom need not be introduced.

It would strengthen the theoretical treatment if the restriction of in-context laws to simple expressions could be removed. While this is not critical from a practical point of view—the class of simple expressions covers many commonly occurring situations—the added generality may help clarify the semantic role of `Safe` and add some power to the algebra.

7 Acknowledgements

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.

References

- [ASU86] A. V. Aho, R. Sethi, and J. D. Ullman. *Compilers: Principles, Techniques, and Tools*. Addison-Wesley, 1986.
- [Bac78] J. Backus. Can programming be liberated from the von Neumann style? A functional style and its algebra of programs. *Communications of the ACM*, 21(8):613–641, August 1978.
- [BWW86] J. Backus, J. H. Williams, and E. L. Wimmers. *The FL Language Manual*. Technical Report RJ 5339 (54809), IBM, 1986.
- [CF89] R. Cartwright and M. Felleisen. The semantics of program dependence. In *Proceedings of the 1989 Conference on Programming Language Design and Implementation*, June 1989.
- [Hen81] J. Hennessy. Program optimization and exception handling. In *Proceedings of the 1981 Symposium on Principles of Programming Languages*, January 1981.
- [HMT89] R. Harper, R. Milner, and M. Tofte. *The Definition of Standard ML—Version 3*. Technical Report ECFS-LFCS-89-81, Laboratory for Foundations of Computer Science, University of Edinburgh, 1989.
- [HWA*88] P. Hudak, P. Wadler, Arvind, B. Boutel, J. Fairbairn, J. Fasel, J. Hughes, T. Johnsson, D. Kieburtz, S. P. Jones, R. Nikhil, M. Reeve, D. Wise, and J. Young. *Report on the Functional Programming Language Haskell*. Technical Report DCS/RR-666, Yale University, December 1988.

- [Ste84] G. L. Steele. *Common Lisp: The Language*. Digital Press, 1984.
- [Tur85] D. A. Turner. 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, 1985.
- [WH87] P. Wadler and R. J. M. Hughes. Projections for strictness analysis. In *Proceedings of the Symposium on Functional Programming Languages and Computer Architecture*, pages 385–407, Springer Verlag Lecture Notes in Computer Science no. 274, September 1987.