

# **Selective Tail Call Elimination**

Yasuhiko Minamide

Institute of Information Sciences and Electronics  
University of Tsukuba

and  
PRESTO, JST

ISE-TR-03-192

## **Abstract**

Tail calls are expected not to consume stack space in most functional languages. However, there is no support for tail calls in some environments. Even in such environments, proper tail calls can be implemented with a technique called a trampoline. To reduce the overhead of trampolining while preserving stack space asymptotically we propose selective tail call elimination based on an effect system. The effect system infers the number of successive tail calls generated by the execution of an expression, and trampolines are introduced only when they are necessary.

# 1 Introduction

Tail calls are expected not to consume stack space in most functional languages. Implementation of proper tail calls requires some support from a target environment, but some environments including C and Java Virtual Machine (JVM) [10] do not provide such support. Even in such environments, proper tail calls can be implemented with a technique called a *trampoline*. However, the trampoline technique is based on a non-standard calling convention and is considered to introduce too much overhead. Thus, most compilers for such environments do not adopt trampolining and abandon proper tail calls [2, 17, 3].

To solve this problem we selectively introduce trampolines based on an effect system. We consider a typed call-by-value language as a source language of selective tail call elimination. Effect systems were originally proposed to infer side-effects of a program by Gifford and Lucassen [8, 11], and an extension was applied to estimate the execution time of a program [6, 15]. Our effect system infers the number of successive tail calls that can be generated by the execution of an expression. Based on effects, functions are divided into two kinds: those that lead to a finite number of successive tail calls and those that may lead to an infinite number of successive tail calls. Then, it is necessary to adopt trampolines only for the latter kind. In this manner we can reduce the overhead of trampolining while preserving stack space asymptotically.

Our effect system includes the rule of subtyping, and some applications may call both kinds of functions. To support subtyping on functions and enable selective elimination of tail calls, we introduce a transformation that translates a function into a record containing two functions supporting different calling conventions. We prove that the increase of stack space usage caused by this transformation is bounded by a factor determined by the effect system.

We have incorporated selective tail call elimination into the MLj compiler [2], which compiles Standard ML into Java bytecodes. Results for benchmark programs show that our effect system is strong enough to indicate that most functions are safe without tail call elimination. We show that there is little degradation of performance for most benchmark programs. We also measure impact of tail call elimination on stack space usage. It shows that general tail call elimination sometimes greatly reduces stack space usage.

This paper is organized as follows. In Section 2, we review the trampolining transformation and clarify its problem. In Section 3, we introduce our selective tail call elimination and prove its soundness. In Section 4, the effect system is extended with a wider class of effects. In Section 5 we outline effect inference for our effect system. Section 6 shows some examples where our effect system shows that a function is unsafe without tail call elimination. In Section 7 we describe our implementation and discuss the results of our experiments. Finally, we review related work and present our conclusions.

## 2 Tail Call Elimination with Trampolines

In functional languages, loops are often expressed with tail calls and thus they are expected not to consume stack space. Let us consider the following program written in Standard ML. The application `sum (x-1, a+x)` in this program is called a tail call because there is nothing to do in the function `sum` after the application.

```
let fun sum (x, a) = if x = 0 then a else sum (x-1, a+x)
in
  sum (y, 0)
end
```

If the tail call is properly implemented, the program above requires only constant stack space. On the other hand, if it is not properly implemented, it requires stack space proportional to `y`. Loops are often expressed in this way in functional languages and thus it is important to implement proper tail calls.

However, it is not straightforward to implement proper tail calls in environments without direct support of tail calls, such as in C and JVM. In such environments it is possible to implement proper tail calls with a non-standard calling convention called a trampoline [1]. We will explain the trampoline technique as a source-to-source transformation. The example above is transformed into the following program.

```
datatype 'a ret = Thunk (unit -> 'a ret) | Val 'a

fun loop (Val x) = x
  | loop (Thunk f) = loop (f ())

let fun sum (x, a) =
      if x=0 then Val a else Thunk (fn () => sum (x-1, a+x))
in
  loop (sum (y, 0))
end
```

The tail call in the program is translated into the creation of a closure. Then the closure is called from the loop.<sup>1</sup> This technique was used in a Standard ML to C compiler [18] and is also useful to implement various features of programming languages [7]. However, it is clear that this technique introduces a lot of overhead. Thus, most compilers into C and JVM do not adopt trampolining and abandon proper tail calls.

## 3 Selective Tail Call Elimination

To reduce the overhead introduced by implementation of proper tail calls by techniques such as trampolining, we propose selective elimination of tail calls that preserves the asymptotic complexity of the stack space.

---

<sup>1</sup>The function `loop` should be implemented so as not to consume stack space.

The basic idea is that if the number of successive tail calls generated by a function call is bounded by some constant, it is not necessary to adopt a trampoline for the function. Let us consider the following program.

```
let fun f x = x
    fun g x = f x
in
  g 0
end
```

There are two tail calls in this program: `f x` and `g 0`. However, it is not necessary to introduce a trampoline for this program. The execution of the function `f` leads to no tail call and thus `f x` generates only one successive tail call. The function `g` calls `f` at a tail-call position and thus the function call `g 0` generates two successive tail calls. Since the number of successive tail calls is bounded in this program, it is safe to execute this program without tail call elimination.

On the other hand, in the following program the number of successive tail calls that the application `h y` leads to cannot be bounded. Thus it is necessary to introduce a trampoline for the program.

```
let fun h x = if x > 0 then h (x - 1) else 0
in
  h y
end
```

If we can statically analyze the number of successive tail calls generated by execution of each function, we can avoid introducing trampolines for functions satisfying some safety condition and selectively eliminate tail calls with trampolines. In the next section we introduce an effect system to eliminate tail calls selectively.

### 3.1 Effect System

We introduce sized effects to check the number of successive tail calls generated by the execution of an expression. The following are the effects we consider.

$$\rho ::= \omega \mid 0 \mid 1 \mid \dots$$

We consider an order relation  $\rho \leq \rho'$  between effects: the usual order relation between natural numbers and  $i \leq \omega$  for any natural number  $i$ . For a finite effect  $i$ , the effect  $i^+$  is defined as  $i^+ = i + 1$ . Then we consider the following types where a function type is annotated with an effect.

$$\tau ::= \text{nat} \mid \tau \rightarrow^\rho \tau$$

We include `nat` as a base type for natural numbers and we also use other base types in examples. The functions we discussed above have the following types.

```
f : int  $\rightarrow^1$  int
g : int  $\rightarrow^2$  int
h : int  $\rightarrow^\omega$  int
```

$$\begin{array}{c}
\frac{x : \tau \in E}{E \vdash x : \tau} \quad E \vdash 0 : \text{nat} \quad \frac{E \vdash V : \text{nat}}{E \vdash \text{suc}(V) : \text{nat}} \quad \frac{E \vdash V : \tau}{E \vdash V : \tau!0} \\
\\
\frac{E, x : \tau_1 \rightarrow^{i^+} \tau_2, y : \tau_1 \vdash M : \tau_2!i}{E \vdash \text{fix } x. \lambda^{i^+} y. M : \tau_1 \rightarrow^{i^+} \tau_2} \quad \frac{E, x : \tau_1 \rightarrow^\omega \tau_2, y : \tau_1 \vdash M : \tau_2!\omega}{E \vdash \text{fix } x. \lambda^\omega y. M : \tau_1 \rightarrow^\omega \tau_2} \\
\\
\frac{E \vdash M_1 : \tau_1 \rightarrow^\rho \tau_2! \rho_1 \quad E \vdash M_2 : \tau_1! \rho_2}{E \vdash @^\rho M_1 M_2 : \tau_2! \rho} \quad \frac{E \vdash M : \tau'! \rho' \quad \rho' \leq \rho \quad \tau' \leq \tau}{E \vdash M : \tau! \rho} \\
\\
\frac{E \vdash V : \text{nat} \quad E \vdash M_1 : \tau! \rho \quad E, x : \text{nat} \vdash M_2 : \tau! \rho}{E \vdash \text{case } V \text{ of } 0 \Rightarrow M_1, \text{suc}(x) \Rightarrow M_2 : \tau! \rho}
\end{array}$$

Figure 1: Type system

Application of  $\mathbf{f}$  at a tail-call position leads to one tail call. Since the application of  $\mathbf{g}$  leads to a subsequent application of  $\mathbf{f}$ , the type of  $\mathbf{g}$  is annotated with 2. On the other hand, the number of successive tail calls generated by  $\mathbf{h}$  cannot be bounded. Thus its type is annotated with  $\omega$ .

We consider the subtyping relation defined by the following rules.

$$\tau \leq \tau \quad \frac{\tau_1 \leq \tau'_1 \quad \tau'_2 \leq \tau_2 \quad \rho' \leq \rho}{\tau'_1 \rightarrow^{\rho'} \tau'_2 \leq \tau_1 \rightarrow^\rho \tau_2}$$

We formalize our effect system for the following language, where abstractions and applications are annotated with effects.

$$\begin{array}{l}
V ::= x \mid 0 \mid \text{suc}(V) \mid \text{fix } x. \lambda^\rho y. M \\
M ::= V \mid @^\rho M M \mid \text{case } V \text{ of } 0 \Rightarrow M, \text{suc}(x) \Rightarrow M
\end{array}$$

We will discuss how to annotate the language without effect annotations in Section 5. The values 0 and  $\text{suc}(V)$  are of natural numbers. Judgments of the effect system have the following forms:

$$\begin{array}{c}
E \vdash V : \tau \\
E \vdash M : \tau! \rho
\end{array}$$

where  $\rho$  represents the maximum number of successive tail calls generated by evaluation of  $M$ . The rules of the effect system are defined in Figure 1. The rules of abstraction and application are explained as follows.

- If the body of a function leads to  $i$  successive tail calls, application of the function at a tail-call position leads to  $i^+$  successive tail calls.
- If the body of a function has effect  $\omega$ , the function has effect  $\omega$ . That means successive tail calls generated by application of the function cannot be bounded.

$$\begin{array}{c}
\vdash_T V \downarrow^0 V \quad \vdash_N V \downarrow^0 V \\
\hline
\vdash_N M_1 \downarrow^l \text{fix } x. \lambda^{\rho'} y. M_0 \quad \vdash_N M_2 \downarrow^m V_2 \quad \vdash_T M_0[\text{fix } x. \lambda^{\rho'} y. M_0/x][V_2/y] \downarrow^n V \\
\hline
\vdash_T @^\rho M_1 M_2 \downarrow^{\max(l,m,n)} V \\
\hline
\vdash_N M_1 \downarrow^l \text{fix } x. \lambda^{\rho'} y. M_0 \quad \vdash_N M_2 \downarrow^m V_2 \quad \vdash_T M_0[\text{fix } x. \lambda^{\rho'} y. M_0/x][V_2/y] \downarrow^n V \\
\hline
\vdash_N @^\rho M_1 M_2 \downarrow^{\max(l,m,n+1)} V \\
\hline
\vdash_\alpha M_1 \downarrow^n V \\
\hline
\vdash_\alpha \text{case } 0 \text{ of } 0 \Rightarrow M_1, \text{suc}(x) \Rightarrow M_2 \downarrow^n V \\
\hline
\vdash_\alpha M_2[V_0/x] \downarrow^n V \\
\hline
\vdash_\alpha \text{case } \text{suc}(V_0) \text{ of } 0 \Rightarrow M_1, \text{suc}(x) \Rightarrow M_2 \downarrow^n V
\end{array}$$

Figure 2: Operational semantics

- The effects of  $M_1$  and  $M_2$  are ignored in the rule of application because they correspond to evaluation at non-tail-call positions. Thus, the effect of the application  $@^\rho M_1 M_2$  is determined only by the effect annotation of the function type.

To discuss the soundness of selective tail elimination we introduce an operational semantics that profiles stack space and models evaluation with proper implementation of tail calls. We define a big-step operational semantics with the following judgments:

$$\vdash_T M \downarrow^i V$$

$$\vdash_N M \downarrow^i V$$

with the meanings that  $M$  is evaluated to  $V$  with  $i$  stack frames at a tail-call position or a non-tail-call position, respectively. A whole program is considered to be evaluated at a tail-call position. The rules are given in Figure 2 where  $\vdash_\alpha M \downarrow^n V$  means the rule holds both for  $\alpha = N$  and  $\alpha = T$ . At a tail-call position, the evaluation of the body of a function requires no new stack frame. Thus, the stack space required for evaluation of the application is  $\max(l, m, n)$ . On the other hand, at a non-tail-call position, it requires a new stack frame: the stack space is  $\max(l, m, n + 1)$ . This semantics models stack space usage when a program is executed after compilation. Correspondence to a semantics that models execution based on an interpreter is discussed in [12].

With respect to the operational semantics, the soundness of the type system in the usual sense is proved. However, the following lemma says nothing about effects inferred by the effect system.

- Lemma 1 (Soundness)**
1. If  $\emptyset \vdash M : \tau! \rho$  and  $\vdash_N M \downarrow^m V$  then  $\emptyset \vdash V : \tau$ .
  2. If  $\emptyset \vdash M : \tau! \rho$  and  $\vdash_T M \downarrow^m V$  then  $\emptyset \vdash V : \tau$ .

### 3.2 Transformation

We introduce a program transformation that selectively eliminates tail calls based on effect annotations. The idea is to eliminate tail calls of the form  $@^\omega MM$  with trampolining and to adopt the standard calling convention for  $@^i MM$ .<sup>2</sup>

However, implementation based on this idea is not so simple because  $\tau_1 \rightarrow^i \tau_2$  can be considered as  $\tau_1 \rightarrow^\omega \tau_2$  by subtyping. Let us consider the following program.

```

let fun f x = x
    fun g x = if ... then x else g (x - 1)
in
  (if ... then f else g) 0
end

```

The functions  $f$  and  $g$  have types  $\text{int} \rightarrow^1 \text{int}$  and  $\text{int} \rightarrow^\omega \text{int}$ , respectively. It is not possible to determine the kinds of functions that are called by the application in the body of the `let`-expression. Thus, it is not straightforward to compile the application in the body.

To solve this problem, we represent  $\tau_1 \rightarrow^i \tau_2$  as a record that contains two functions: one for a trampoline and one for the standard calling convention. Then subtyping on function types is translated into record (object) subtyping.

For the target language of the transformation we consider the following types.

$$\sigma ::= \text{nat} \mid \sigma \rightarrow \sigma \mid \sigma \rightarrow^t \sigma \mid \{\text{fun} : \sigma, \text{tfun} : \sigma\} \mid \{\text{tfun} : \sigma\}$$

There are two kinds of function types:  $\sigma \rightarrow \sigma$  uses the standard calling convention without proper tail calls, and  $\sigma \rightarrow^t \sigma$  uses the non-standard calling convention with tail call elimination. There is no subtyping relation between  $\sigma_1 \rightarrow \sigma_2$  and  $\sigma_1 \rightarrow^t \sigma_2$ . Two kinds of record types,  $\{\text{tfun} : \sigma\}$  and  $\{\text{fun} : \sigma, \text{tfun} : \sigma\}$ , are included to translate function types and we consider the following subtyping relation between them.

$$\{\text{fun} : \sigma_1, \text{tfun} : \sigma_2\} \leq \{\text{tfun} : \sigma_2\}$$

Then our transformation translates function types into record types so that the subtyping relation is preserved. The translation of types  $|\tau|$  is defined as follows:

$$|\tau_1 \rightarrow^\omega \tau_2| = \{\text{tfun} : |\tau_1| \rightarrow^t |\tau_2|\}$$

$$|\tau_1 \rightarrow^i \tau_2| = \{\text{fun} : |\tau_1| \rightarrow |\tau_2|, \text{tfun} : |\tau_1| \rightarrow^t |\tau_2|\}$$

We therefore have the following translation of subtyping.

$$|\tau_1 \rightarrow^i \tau_2| \leq |\tau_1 \rightarrow^\omega \tau_2|$$

This translation of subtyping is natural for compilation to Java bytecodes because JVM has subtyping on objects through inheritance of classes.

---

<sup>2</sup>We assume that the standard calling convention does not support tail call elimination.

The syntax of the target language is defined as follows. It includes two kinds of abstraction and application, and syntax for records and field selection.

$$\begin{aligned}
N & ::= W \mid @NN \mid @^t NN \mid N.\text{fun} \mid N.\text{tfun} \mid \text{case } W \text{ of } 0 \Rightarrow N, \text{ suc}(x) \Rightarrow N \\
W & ::= x \mid 0 \mid \text{suc}(W) \mid \lambda x.N \mid \lambda^t x.N \mid \text{fix } x.\{\text{fun} = \bar{W}, \text{tfun} = W\} \mid \\
& \quad \text{fix } x.\{\text{tfun} = W\}
\end{aligned}$$

The fields of a record expression are restricted to values: this restriction is sufficient for our transformation. The type system of the target language is standard and does not include effects. The definition of the type system is shown in Appendix A.

We define an operational semantics of the target language in the same manner as that of the source language: we define a big-step operational semantics with following judgments.

$$\begin{aligned}
& \vdash_T N \downarrow^i W \\
& \vdash_N N \downarrow^i W
\end{aligned}$$

The main rules are defined as follows. It should be noted that tail calls are not properly implemented for application  $@N_1N_2$  and thus the evaluation of the body of the function requires a new stack frame: the stack space required is not  $\max(l, m, n)$ , but  $\max(l, m, n+1)$ .

$$\begin{array}{c}
\frac{\vdash_N N_1 \downarrow^l \lambda^t x.N \quad \vdash_N N_2 \downarrow^m W_2 \quad \vdash_T N[W_2/x] \downarrow^n W}{\vdash_T @^t N_1 N_2 \downarrow^{\max(l,m,n)} W} \\
\frac{\vdash_N N_1 \downarrow^l \lambda x.N \quad \vdash_N N_2 \downarrow^m W_2 \quad \vdash_T N[W_2/x] \downarrow^n W}{\vdash_T @N_1 N_2 \downarrow^{\max(l,m,n+1)} W}
\end{array}$$

The other rules are shown in Appendix B.

The transformation of selective tail call elimination is defined as follows:

$$\begin{aligned}
[[x]] & = x \\
[[0]] & = 0 \\
[[\text{suc}(V)]] & = \text{suc}([[V]]) \\
[[\text{fix } x.\lambda^\omega y.M]] & = \text{fix } x.\{\text{tfun} = \lambda^t y. [[M]]\} \\
[[\text{fix } x.\lambda^i y.M]] & = \text{fix } x.\{\text{fun} = \lambda y. [[M]], \text{tfun} = \lambda^t y. [[M]]\} \\
[[@^i M_1 M_2]] & = @([[M_1].\text{fun}]) [[M_2]] \\
[[@^\omega M_1 M_2]] & = @^t([[M_1].\text{tfun}]) [[M_2]] \\
[[\text{case } V \text{ of } 0 \Rightarrow M_1, \text{ suc}(x) \Rightarrow M_2]] & = \text{case } [[V]] \text{ of } 0 \Rightarrow [[M_1]], \text{ suc}(x) \Rightarrow [[M_2]]
\end{aligned}$$

We extend the translation of types to type environments as  $|E|(x) = |E(x)|$ . Then the type correctness of this transformation is formulated as the following lemma and proved by induction on the derivation of  $E \vdash M : \tau! \rho$ .

**Lemma 2 (Type soundness)** *If  $E \vdash M : \tau! \rho$  then  $|E| \vdash [[M]] : |\tau|$ .*

To formalize the soundness of the transformation we introduce the following notation:  $\vdash_T M \downarrow^{\leq k} V$  if  $\vdash_T M \downarrow^{k'} V$  for some  $k' \leq k$ . The factor of increase of stack space usage by selective tail call elimination is determined by the maximum of the effect annotations in  $M$ , denoted by  $\max(M)$ .

**Theorem 1 (Soundness)** *Let  $C = \max(M) + 1$ .*

1. *If  $\emptyset \vdash M : \tau ! i$  and  $\vdash_T M \downarrow^k V$  then  $\vdash_T \llbracket M \rrbracket \downarrow^{\leq Ck+i} \llbracket V \rrbracket$ .*
2. *If  $\emptyset \vdash M : \tau ! \omega$  and  $\vdash_T M \downarrow^k V$  then  $\vdash_T \llbracket M \rrbracket \downarrow^{\leq Ck+C-1} \llbracket V \rrbracket$ .*
3. *If  $\emptyset \vdash M : \tau ! \rho$  and  $\vdash_N M \downarrow^k V$  then  $\vdash_N \llbracket M \rrbracket \downarrow^{\leq Ck} \llbracket V \rrbracket$ .*

This theorem ensures that the stack space usage of a program is preserved asymptotically.

For example,  $\text{@}^\omega(\text{fix } f.\lambda^\omega x.\text{@}^1(\text{fix } g.\lambda^1 y.y)x)0$  and its translation are evaluated as follows:

$$\begin{aligned} & \vdash_N \text{@}^\omega(\text{fix } f.\lambda^\omega x.\text{@}^1(\text{fix } g.\lambda^1 y.y)x)0 \downarrow^1 0 \\ & \vdash_N \text{@}^t(\text{fix } f.\{\text{tfun} = \lambda^t x.\text{@}(\text{fix } g.\{\text{fun} = \lambda y.y, \text{tfun} = \lambda^t y.y\}.\text{fun})x\}.\text{tfun})0 \downarrow^2 0 \end{aligned}$$

This example corresponds to the worst case:  $k = 1$  and  $C = 2$ . The proof of the theorem appears in Appendix C.

## 4 Extension of the Effect System

The effect system we have presented has one unnatural limitation:  $\omega$  must always be assigned to a function which calls a function with effect  $\omega$  at tail call position, even if the function is safe without tail call elimination. In this section, we extend our effect system to overcome this limitation by considering a wider class of effects.

We first show an example where the limitation of our effect system appears. Let us consider the following program.

```

fun f x = f x
fun g x = f x
fun h (0,x) = g x
  | h (n,x) = h (n-1,x)

```

The function  $g$  is safe without tail call elimination: the stack space usage is increased by 1 even if it is implemented with the standard calling convention. However, in our effect system the function is assigned the effect  $\omega$  to because it calls the function  $f$  of the effect  $\omega$  at a tail call position.

We solve this limitation by extending the effects in our type system into the following form.

$$\rho ::= \omega \cdot i + j$$

where  $i$  and  $j$  are natural numbers. The intuition is that the function with effect  $\omega \cdot i + j$  such that  $j > 0$  is safe without tail call elimination. We identifies  $\omega \cdot i + 0$  and  $\omega \cdot 0 + j$

with  $\omega \cdot i$  and  $j$ , respectively. The effect  $\rho^+$  and the subeffect relation  $\rho \leq \rho'$  are defined as follows:

$$(\omega \cdot i + j)^+ = \omega \cdot i + (j + 1)$$

$$\omega \cdot i + j \leq \omega \cdot i' + j' \quad \text{iff} \quad i < i', \text{ or } i = i' \text{ and } j \leq j'$$

The typing rules of abstraction in the effect system are extended as follows:

$$\frac{E, x : \tau_1 \rightarrow^{\rho^+} \tau_2, y : \tau_1 \vdash M : \tau_2 ! \rho}{E \vdash \text{fix } x. \lambda^{\rho^+} y. M : \tau_1 \rightarrow^{\rho^+} \tau_2} \quad \frac{E, x : \tau_1 \rightarrow^{\omega \cdot i} \tau_2, y : \tau_1 \vdash M : \tau_2 ! \omega \cdot i}{E \vdash \text{fix } x. \lambda^{\omega \cdot i} y. M : \tau_1 \rightarrow^{\omega \cdot i} \tau_2}$$

Then we can assign the following types and thus  $g$  can be safely implemented with the standard calling convention.

$$\begin{aligned} f &: \text{int} \rightarrow^{\omega} \text{int} \\ g &: \text{int} \rightarrow^{\omega+1} \text{int} \\ h &: \text{int} \times \text{int} \rightarrow^{\omega \cdot 2} \text{int} \end{aligned}$$

We also need to modify the transformation to implement selective tail call elimination. Since a function with effect  $\omega \cdot i$  can be considered to have effect  $\omega \cdot i + 1$  in this system, a function with effect  $\omega \cdot i$  must support both calling conventions. The transformation is modified as follows:

$$\begin{aligned} \llbracket \text{fix } x. \lambda^{\rho} y. M \rrbracket &= \text{fix } x. \{ \text{fun} = \lambda y. \llbracket M \rrbracket, \text{tfun} = \lambda^t y. \llbracket M \rrbracket \} \\ \llbracket @^{\omega \cdot i + j} M_1 M_2 \rrbracket &= @(\llbracket M_1 \rrbracket. \text{fun}) \llbracket M_2 \rrbracket \\ \llbracket @^{\omega \cdot i} M_1 M_2 \rrbracket &= @^t(\llbracket M_1 \rrbracket. \text{tfun}) \llbracket M_2 \rrbracket \end{aligned}$$

where  $j > 0$ . The intuitive meaning of extended effects can be explained with this transformation.

- A tail call with effect  $\omega$  generates successive tail calls of  $@^t$  and then successive tail calls of  $@$ . The successive tail calls of  $@$  may be generated by subeffect relation  $i \leq \omega$ .
- A tail call with effect  $\omega \cdot i$  may repeat  $i$  times the pattern of tail calls for  $\omega$
- A tail call with effect  $\omega \cdot i + j$  may generate  $j$  successive tail calls of  $@$  and then generates tail calls of the pattern for  $\omega \cdot i$ .

The soundness theorem is extended in the following form. We write  $\max^i(M)$  for the maximum  $j$  of  $\omega \cdot i + j$  appearing in  $M$ . The proof of the theorem appears in Appendix D.

**Theorem 2 (Soundness)** *Let  $C = \sum_{i=0}^{\infty} \max^i(M) + 1$  and  $D(j) = \sum_{i=0}^{j-1} \max^i(M)$ .*

1. *If  $\emptyset \vdash M : \tau ! (\omega \cdot i + j)$  and  $\vdash_T M \downarrow^k V$  then  $\vdash_T \llbracket M \rrbracket \downarrow^{\leq Ck + D(i) + j} \llbracket V \rrbracket$ .*
2. *If  $\emptyset \vdash M : \tau ! \rho$  and  $\vdash_N M \downarrow^k V$  then  $\vdash_N \llbracket M \rrbracket \downarrow^{\leq Ck} \llbracket V \rrbracket$ .*

## 5 Effect Inference

We show how to infer effects in this section. The effect inference can be formalized as a type system with constraints, where a constraint generated by the type system is solved with a simple graph-based algorithm. We assume that types are already inferred with the standard type inference and consider the following explicitly-typed language for effect inference.

$$\begin{aligned}\tau & ::= \text{nat} \mid \tau \rightarrow^\alpha \tau \\ V & ::= 0 \mid \text{succ}(V) \mid x \mid \text{fix } x : \tau. \lambda y. M \\ M & ::= V \mid @^\alpha M M \mid \text{case } V \text{ of } 0 \Rightarrow M, \text{succ}(x) \Rightarrow M\end{aligned}$$

where  $\alpha$  denotes an effect variable. The effect annotation of a lambda abstraction can be determined from the type annotation of the `fix`-expression. We assume effect variables appearing in a program are distinct.

Judgments of the effect system have the following forms:  $E; C \vdash V : \tau$  and  $E; C \vdash M : \tau! \alpha$  where  $C$  is a set of subeffect relations:  $\alpha < \alpha'$  and  $\alpha \leq \alpha'$ . A constraint  $\alpha < \alpha'$  holds if  $\alpha \leq \alpha'$  and  $\alpha \neq \alpha'$ , or  $\alpha = \alpha' = \omega \cdot i$  for some  $i$ . The main rules of the effect system are given as follows:

$$\frac{E; C \vdash V : \tau \quad \alpha \text{ is fresh}}{E; C \vdash V : \tau! \alpha} \quad \frac{E, x : \tau_1 \rightarrow^\alpha \tau_2, y : \tau_1; C \vdash M : \tau_2! \alpha'}{E; C \cup \{\alpha' < \alpha\} \vdash \text{fix } x : \tau_1 \rightarrow^\alpha \tau_2. \lambda y. M : \tau_1 \rightarrow^\alpha \tau_2}$$

$$\frac{E; C_1 \vdash M_1 : \tau_1 \rightarrow^{\alpha'} \tau_2! \alpha_1 \quad E; C_2 \vdash M_2 : \tau_1'! \alpha_2}{E; C_1 \cup C_2 \cup \{\alpha' \leq \alpha\} \cup C_{\leq}(\tau_1', \tau_1) \vdash @^\alpha M_1 M_2 : \tau_2! \alpha}$$

where  $C_{\leq}(\tau_1', \tau_1)$  is the constraint to obtain the subtyping  $\tau_1' \leq \tau_1$ .

$$\begin{aligned}C_{\leq}(\text{nat}, \text{nat}) &= \emptyset \\ C_{\leq}(\tau_1 \rightarrow^\alpha \tau_2, \tau_1' \rightarrow^{\alpha'} \tau_2') &= C_{\leq}(\tau_1', \tau_1) \cup C_{\leq}(\tau_2, \tau_2') \cup \{\alpha \leq \alpha'\}\end{aligned}$$

The constraint obtained by the rules above can be solved in the following manner. We consider the graph of the relations  $\alpha < \alpha'$  and  $\alpha \leq \alpha'$ , and compute the strongly connected components of the graph. If a strongly connected component contains a relation of the form  $\alpha < \alpha'$ , the effect of the form  $\omega \cdot i$  must be assigned to the effect variables appearing in the component. It is clear that an effect  $\omega \cdot i + j$  ( $j > 0$ ) can be assigned to an effect variable not belonging to such components.

## 6 Examples

There are several situations where an effect of the form  $\omega \cdot i$  is assigned to a function. Although some of them are actually unsafe without tail call elimination, our effect system sometimes assigns  $\omega \cdot i$  to functions safe without tail call elimination. In this section we show some examples of both the situations.

In our effect system,  $\omega \cdot i$  must be assigned to tail recursive functions in general. However, tail recursive calls in a single recursive function can be implemented as a loop and thus trampolining can be avoided for such tail calls. Our effect system can be extended to be consistent with this implementation and then such functions are not assigned  $\omega \cdot i$  to. Then there are two common examples where recursive functions are unsafe without tail call elimination: mutually tail recursive functions and higher order recursive functions.

In the following example, the functions `f` and `g` contains mutually recursive tail calls and thus must be assigned  $\omega$  to.

```
fun f 0 = 0
  | f n = g (n-1)
and g n = f (n-1)
```

However, it is possible to implement the tail calls as a loop if only one of `f` and `g` are used from the other part of a program or the functions are copied into two definitions.

The following is an example with a higher order function.

```
fun h 0 y = y
  | h x y = h (x-1) (x+y)
```

The function `h` has type `int → int → int`. The tail recursive call of `h (x-1) (x+y)` cannot be implemented as a loop. However, if the function is uncurried, the tail call in the function can be implemented as a loop and thus the function can be safely implemented without tail call elimination.

As the third example, we show a program that is safe without tail call elimination, but is assigned  $\omega$  to with our effect system.

```
let fun f (g:int -> int) = g 0
in
  f (fn x => f (fn y => y))
end
```

You can check this program is safe without tail call elimination. Let us infer the type of `f`. By assuming that `f` has type  $(\text{int} \rightarrow^{\alpha_1} \text{int}) \rightarrow^{\alpha_2} \text{int}$ , the constraints  $\alpha_1 < \alpha_2$  and  $\alpha_2 < \alpha_1$  must be satisfied where the first constraint is obtained from the definition of `f` and the second constraint is obtained from the body of the `let`-expression. Then the effects  $\alpha_1$  and  $\alpha_2$  must be  $\omega \cdot i$  for some  $i$ . This weakness of our effect system appears in the benchmark “logic” we will discuss in the next section: many functions similar to the function above appear in the program.

## 7 Implementation and Measurements

We have incorporated our selective tail call elimination into the MLj compiler [2], which translates Standard ML into Java bytecodes. In MLj, tail calls are compiled into Java method invocations except recursive calls which are implemented as loops. MLj uses

	total	(A)	(B)	(C)	(D)
barnes-hut	25	0	0	0	2
boyer	24	0	0	0	2
fft	28	0	0	0	3
knuth-bendix	83	5	5	5	3
lexgen	110	13	11	2	3
life	30	1	1	0	2
logic	58	51	36	36	2
mandelbrot	6	0	0	0	1
nucleic	38	0	0	0	3
ratio-regions	51	0	0	0	2
ray	82	2	2	0	3
simple	175	0	0	0	4
tsp	20	0	0	0	2
vliw	463	19	18	16	4

Table 1: Results of effect inference

an intermediate language based on monads to represent effects such as IO operations, exception, and non-termination. We extended the effects of the intermediate language with our effect. The effects for selective tail call elimination are inferred in a late stage of compilation and the transformation is merged into the code generation phase of the compiler.

The translation of a function presented in Section 3.2 and 4 has one problem: the body of  $\lambda^i x.M$  is copied into two functions and thus the translation may increase the code size. This problem can be solved by the following translation.

$$\lambda^i y.M = \text{let } y = \lambda x.[M] \text{ in } \{\text{fun} = y, \text{tfun} = \lambda^i x.@yx\}$$

However, we do not adopt this translation because it makes it difficult to compare stack space usage. The worst case increase of code size observed for the benchmark programs we will discuss later is about 35 %.<sup>3</sup>

We measured the effectiveness of our selective tail call elimination for the most benchmark programs obtained from the following URL.<sup>4</sup>

<ftp://ftp.research.bell-labs.com/dist/smlnj/benchmarks/>

<sup>3</sup>The pair representation is not used for the known function that are safe without tail call elimination.

<sup>4</sup>We excluded two programs count-graphs and mlyacc that are difficult to compile with MLj because of the limitation of MLj.

	MLj	TCE	STCE
knuth-bendix	3895	3373	3485
lexgen	1259	94	94
life	298	49	49
logic	3814	236	260

Table 2: Maximum stack size: (in number of frames)

Table 1 summarizes the results of our effect analysis. The column total shows the number of the functions generated for each benchmark program. The columns (A), (B) and (C) are the numbers of functions the analysis assigns an effect of the form  $\omega \cdot i$ : (A), (B) and (C) are the numbers for selective tail call elimination without extension, with extension and with extension and an extra phase of optimization, respectively. The column (D) shows  $\sum_{i=0}^{\infty} \max(P)$  for each program  $P$ , which determines the theoretical upper bound of increase of stack space usage.

- Eight programs out of 13 are shown safe without tail call elimination. Even for the other programs except for the program logic, the ratio of the function of effect  $\omega \cdot i$  is small.
- The most functions of the program logic have effect  $\omega \cdot i$  by the reason we described in Section 6. The extension reduces the number, but more than half of the functions still have  $\omega \cdot i$ .
- Since the maximum of the numbers in column (D) is 4, the theoretical upper bound of stack space increase for selective tail call elimination compared to tail call elimination is 5.
- The effectiveness of our selective tail call elimination depends on other phases of compilation. An extra phase of optimization including uncurrying decreased the number of functions of effect  $\omega \cdot i$  for three programs.

Table 2 shows the maximum stack size during execution measured by the number of frames. The table shows the results for the benchmark programs where tail call elimination has some impact on the results. For all the other program, the numbers are between 33 and 103. The results supports that tail call elimination is desirable: stack sizes are greatly reduced for several programs. Selective tail call elimination may increase stack size compared to tail call elimination. However, the increase is relatively small, compared to the theoretical upper bound.

Table 3 shows execution times. The columns TCE and STCE are the results for tail call elimination and selective tail call elimination, respectively. The numbers in the parenthesis are the ratios to those of MLj. Even for TCE, all the non-tail-calls are implemented with the standard calling convention based on the pair representation of

	Interpreted-mode			HotSpot Client VM		
	MLj	TCE	STCE	MLj	TCE	STCE
barnes-hut	5.77	6.30(109.2)	5.83(101.0)	1.22	1.21(99.2)	1.23(100.8)
boyer	1.41	1.65(117.0)	1.41(100.0)	0.81	0.58(71.6)	0.82(101.2)
fft	1.71	2.12(124.0)	1.71(100.0)	0.57	0.60(105.3)	0.64(112.3)
knuth-bendix	11.09	10.80(97.4)	8.19(73.9)	2.00	2.12(106.0)	1.61(80.5)
lexgen	3.50	3.40(97.1)	3.49(99.7)	0.61	0.75(123.0)	0.63(103.3)
life	1.24	1.24(100.0)	1.15(92.7)	0.36	0.37(102.8)	0.35(97.2)
logic	17.40	18.90(108.6)	16.49(94.8)	4.51	2.79(61.9)	2.70(59.9)
mandelbrot	4.18	6.45(154.3)	4.22(101.0)	0.52	1.49(286.5)	0.55(105.8)
nucleic	0.68	0.73(107.4)	0.67(98.5)	0.34	0.44(129.4)	0.34(100.0)
ratio-regions	212.93	216.38(101.6)	209.11(98.2)	33.62	42.34(125.9)	33.67(100.1)
ray	5.77	6.20(107.5)	5.68(98.4)	1.88	1.13(60.1)	2.12(112.8)
simple	6.61	7.35(111.2)	6.38(96.5)	1.49	1.63(109.4)	1.54(103.4)
tsp	4.78	5.13(107.3)	4.75(99.4)	0.88	0.96(109.1)	0.84(95.5)
vliw	6.26	7.32(116.9)	6.42(102.6)	1.34	2.37(176.9)	1.46(109.0)

Table 3: Execution time (in seconds) : 1.8GHz Pentium 4, Linux, JDK 1.4.0

functions. Measurements were done using Sun JDK 1.4.0, Java HotSpot Client VM on a Linux PC with 1.8GHz Pentium 4. We measured execution time on the interpreted-mode with the `-Xint` option, and on the mode where HotSpot compilation is enabled because it is sometimes difficult to interpret results on the HotSpot VM. Each benchmark was run five times and we chose the fastest run.

- TCE sometimes degrades the performance a lot. The worst case overhead is 54.3 % and 186.5 % for the interpreted-mode and the HotSpot VM, respectively. Compared to TCE, STCE causes little overhead: the worst case overhead is 2.6 % and 12.8 %, respectively.
- For benchmark programs where stack size is reduced by tail call elimination, execution times are sometimes reduced for both TCE and STCE: knuth-bendix and logic. This can be explained as a reduction of garbage collection (GC) time. For example, the GC times for logic are 2.43, 0.49 and 0.49 for MLj, TCE and STCE, respectively. The same phenomenon was observed by Schinz and Odersky in their tail call elimination for JVM [16].
- There are unexpected results on boyer and ray: the big improvement of execution time over MLj and STCE is observed for TCE. We checked the profiling data of executions and found that better decisions on compilation are made by the HotSpot VM for TCE and the programs compiled by MLj and STCE spent more time on interpreted methods.

	Interpreted-mode			HotSpot Client VM		
	MLj	TCE	STCE	MLj	TCE	STCE
barnes-hut	13.95	16.29(116.8)	14.05(100.7)	2.49	2.36(94.8)	2.49(100.0)
boyer	3.32	3.78(113.9)	3.34(100.6)	1.90	1.19(62.6)	1.88(98.9)
fft	3.52	4.37(124.1)	3.53(100.3)	0.98	1.24(126.5)	0.96(98.0)
knuth-bendix	20.80	25.00(120.2)	20.75(99.8)	3.22	3.85(119.6)	3.27(101.6)
lexgen	8.47	8.81(104.0)	8.68(102.5)	1.47	1.75(119.0)	1.48(100.7)
life	3.35	3.64(108.7)	3.41(101.8)	0.65	0.72(110.8)	0.63(96.9)
logic	39.84	51.31(128.8)	43.78(109.9)	9.42	7.24(76.9)	6.32(67.1)
mandelbrot	9.75	16.04(164.5)	9.82(100.7)	1.01	1.71(169.3)	1.02(101.0)
nucleic	1.56	1.68(107.7)	1.50(96.2)	0.69	0.77(111.6)	0.67(97.1)
ratio-regions	597.95	672.76(112.5)	580.94(97.2)	100.76	111.16(110.3)	99.99(99.2)
ray	13.60	15.32(112.6)	13.86(101.9)	4.76	2.24(47.1)	4.73(99.4)
simple	18.00	19.96(110.9)	17.84(99.1)	4.20	3.25(77.4)	4.29(102.1)
tsp	14.12	14.67(103.9)	14.13(100.1)	1.46	1.57(107.5)	1.41(96.6)
vliw	15.74	18.35(116.6)	15.83(100.6)	3.86	5.17(133.9)	4.14(107.3)

Table 4: Execution time (in seconds) : 650MHz UltraSPARC Iii, Solaris 9, JDK 1.4.1

We also measured execution times on a workstation with 650MHz Ultra SPARC Iii. The results are shown in Figure 4 and similar to the previous results.

## 8 Related Work

Dornic, Jouvelot and Gifford proposed an effect system to estimate execution time of a program [6], and their work was extended by Reistad and Gifford [15] with sized types [9]. By adapting the effect system extended with sized types we may obtain more information about the stack usage of a program. However, our simple effect system gives enough information for selective tail call elimination.

Implementation of proper tail calls and space safety are discussed by Clinger [5]. He considered that an implementation is safe with respect to space if it does not increase the asymptotic space complexity of programs. Our selective tail call elimination satisfies the criterion on stack space, but the factor of increase of stack space depends on the program and is determined by the effect system.

Schinz and Odersky proposed tail call elimination for the Java virtual machine [16] that preserves complexity on stack space. Their method is dynamic: it keeps track of the number of successive tail calls and execution is returned to a trampoline if some predefined limit is exceeded. We think that it is possible to reduce the overhead of their method with selective elimination of tail calls.

We have translated a function in the source language into a record with two func-

tions supporting different calling conventions. Similar translation was used to support multiple calling conventions in type-directed unboxing by Minamide and Garrigue [14], and the vectorized functions of Chakravarty and Keller [4].

## 9 Conclusion and Future Work

We have presented an effect system and a program transformation to eliminate tail calls selectively. The transformation translates a function into a record with two functions supporting different calling conventions. The transformation preserves stack space asymptotically.

Our effect system will be useful even for target environments that directly supports tail calls. Various program transformations sometimes translate tail calls into non-tail calls. With our effect system, it is possible to check if such translation is safe for each tail call.

We incorporated our effect system into the MLj compiler and measured the proportion of functions that are unsafe without tail call elimination. The results indicated that selective tail call elimination is very effective for most programs and most functions can be implemented with the standard calling convention. However, there is a limitation that our effect system is monovariant. This limitation may be solved if we extend our effect system with effect polymorphism or intersection types.

By selective tail call elimination, the asymptotic complexity of stack space is preserved. The factor of the increase is determined by effect analysis and depends on a program. However, it is also possible to guarantee the factor of stack space increase by translating applications with annotations greater than the predefined factor as applications with annotation  $\omega$ .

## Acknowledgments

This work is partially supported by Grant-in-Aid for Encouragement of Young Scientists, No. 13780193.

## References

- [1] H. Baker. Cons should not cons its arguments, part II: Cheney on the M.T.A. *SIGPLAN Notices*, 30(9):17–20, 1995.
- [2] N. Benton, A. Kennedy, and G. Russell. Compiling Standard ML to Java bytecodes. In *Proceedings of the Third ACM SIGPLAN International Conference on Functional Programming (ICFP '98)*, pages 129–140, 1998.
- [3] P. Bothner. Kawa - compiling dynamic languages to the Java VM. In *Proceedings of the USENIX 1998 Technical Conference*, 1998.

- [4] M. M. T. Chakravarty and G. Keller. More types for nested data parallel programming. In *Proceedings of the Fifth ACM SIGPLAN International Conference on Functional Programming*, pages 94–105, 1999.
- [5] W. D. Clinger. Proper tail recursion and space efficiency. In *Proceedings of the ACM SIGPLAN '98 Conference on Programming Language Design and Implementation*, pages 174–185. ACM Press, 1998.
- [6] V. Dornic, P. Jouvelot, and D. K. Gifford. Polymorphic time systems for estimating program complexity. *ACM Letters on Programming Languages and Systems (LOPLAS)*, 1(1):33–45, 1992.
- [7] S. D. Ganz, D. P. Friedman, and M. Wand. Trampolined style. In *Proceedings of the 4th ACM SIGPLAN International Conference on Functional Programming (ICFP '99)*, pages 18–22, 1999.
- [8] D. K. Gifford and J. M. Lucassen. Integrating functional and imperative programming. In *Proceedings of the ACM Conference on Lisp and Functional Programming*, pages 28–38, 1986.
- [9] J. Hughes, L. Pareto, and A. Sabry. Proving the correctness of reactive systems. In *Proceedings of the 23rd ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, pages 410–423, 1996.
- [10] T. Lindholm and F. Yellin. *The Java Virtual Machine Specification*. Addison Wesley, 1999.
- [11] J. M. Lucassen and D. K. Gifford. Polymorphic effect systems. In *Proceedings of the Fifteenth ACM SIGACT-SIGPLAN Symposium on Principles of Programming Languages*, pages 47–57, 1988.
- [12] Y. Minamide. A new criterion for safe program transformations. In *Proceedings of the Forth International Workshop on Higher Order Operational Techniques in Semantics (HOOTS)*, volume 41(3) of *ENTCS*, Montreal, 2000.
- [13] Y. Minamide. Selective tail call elimination. Technical Report ISE-TR-03-192, Institute of Information Sciences and Electronics, University of Tsukuba, 2003.
- [14] Y. Minamide and J. Garrigue. On the runtime complexity of type-directed unboxing. In *Proceedings of the Third ACM SIGPLAN International conference on Functional Programming*, pages 1–12, 1998.
- [15] B. Reistad and D. K. Gifford. Static dependent costs for estimating execution time. In *Proceedings of the 1994 ACM Conference on LISP and Functional Programming*, pages 65–78, 1994.
- [16] M. Schinz and M. Odersky. Tail call elimination on the Java virtual machine. In *Proceedings of the First International Workshop on Multi-Language Infrastructure and Interoperability (BABEL)*, volume 59(1) of *ENTCS*, 2001.

- [17] B. Serpette and M. Serrano. Compiling Scheme to JVM bytecode: a performance study. In *Proceedings of the Seventh ACM SIGPLAN International Conference on Functional Programming*, pages 259–270, 2002.
- [18] D. Tarditi, A. Acharya, and P. Lee. No assembly required: Compiling standard ML to C. *ACM Letters on Programming Languages and Systems (LOPLAS)*, 1(2):161–177, 1992.

## A Type System of the Target Language

The type system of the target language is defined as a deductive system with judgments of the form  $E \vdash N : \sigma$ . We do not show the rules for  $0$ ,  $\text{succ}(W)$ , and  $\text{case}$ -expressions that are the same as those for the source language.

$$\begin{array}{c}
\frac{x : \sigma \in E}{E \vdash x : \sigma} \quad \frac{E \vdash N : \sigma' \quad \sigma' \leq \sigma}{E \vdash N : \sigma} \\
\\
\frac{E, x : \sigma_1 \vdash N : \sigma_2}{E \vdash \lambda x. N : \sigma_1 \rightarrow \sigma_2} \quad \frac{E, x : \sigma_1 \vdash N : \sigma_2}{E \vdash \lambda^t x. N : \sigma_1 \rightarrow^t \sigma_2} \\
\\
\frac{E \vdash N_1 : \sigma_1 \rightarrow \sigma_2 \quad E \vdash N_2 : \sigma_1}{E \vdash @N_1 N_2 : \sigma_2} \quad \frac{E \vdash N_1 : \sigma_1 \rightarrow^t \sigma_2 \quad E \vdash N_2 : \sigma_1}{E \vdash @^t N_1 N_2 : \sigma_2} \\
\\
\frac{E, x : \{\text{fun} : \sigma_1, \text{tfun} : \sigma_2\} \vdash W_1 : \sigma_1 \quad E, x : \{\text{fun} : \sigma_1, \text{tfun} : \sigma_2\} \vdash W_2 : \sigma_2}{E \vdash \text{fix } x. \{\text{fun} = W_1, \text{tfun} = W_2\} : \{\text{fun} : \sigma_1, \text{tfun} : \sigma_2\}} \\
\\
\frac{E, x : \{\text{tfun} : \sigma\} \vdash W : \sigma}{E \vdash \text{fix } x. \{\text{tfun} = W\} : \{\text{tfun} : \sigma\}} \\
\\
\frac{E \vdash N : \{\text{fun} : \sigma_1, \text{tfun} : \sigma_2\}}{E \vdash N. \text{fun} : \sigma_1} \quad \frac{E \vdash N : \{\text{tfun} : \sigma\}}{E \vdash N. \text{tfun} : \sigma}
\end{array}$$

## B Operational Semantics of the Target Language

The following are the rules of the operational semantics of the target language. We write  $\vdash_\alpha N \downarrow^n W$  if the rule holds for both  $\vdash_N N \downarrow^n W$  and  $\vdash_T N \downarrow^n W$ .

$$\begin{array}{c}
\frac{\vdash_N N_1 \downarrow^l \lambda x. N \quad \vdash_N N_2 \downarrow^m W_2 \quad \vdash_\alpha N[W_2/x] \downarrow^n W}{\vdash_\alpha @N_1 N_2 \downarrow^{\max(l,m,n+1)} W} \\
\\
\vdash_\alpha W \downarrow^0 W \quad \frac{\vdash_N N \downarrow^n \text{fix } x. \{\text{tfun} = W\}}{\vdash_\alpha N. \text{tfun} \downarrow^n W[\text{fix } x. \{\text{tfun} = W\}/x]} \\
\\
\frac{\vdash_N N \downarrow^n \text{fix } x. \{\text{fun} = W_1, \text{tfun} = W_2\}}{\vdash_\alpha N. \text{tfun} \downarrow^n W_2[\text{fix } x. \{\text{fun} = W_1, \text{tfun} = W_2\}/x]} \\
\\
\frac{\vdash_N N \downarrow^n \text{fix } x. \{\text{fun} = W_1, \text{tfun} = W_2\}}{\vdash_\alpha N. \text{fun} \downarrow^n W_1[\text{fix } x. \{\text{fun} = W_1, \text{tfun} = W_2\}/x]}
\end{array}$$

## C Proof of Soundness

The following lemma is crucial to establish soundness of transformation.

**Lemma 3**  $\llbracket M[V/x] \rrbracket \equiv \llbracket M \rrbracket \llbracket [V]/x \rrbracket$ .

The following lemma is also used to prove the soundness. It is shown by case analysis.

**Lemma 4** *If  $\emptyset \vdash \text{fix } x.\lambda^{\rho'} y.M : \tau_1 \rightarrow^{\rho} \tau_2$ , then  $x : \tau_1 \rightarrow^{\rho'} \tau_2, y : \tau_1 \vdash M : \tau_2! \rho''$  for some  $\rho'' \leq \rho$ .*

We prove the main theorem in the following form to simplify case-analysis.

**Lemma 5** *Let  $C = \max(M) + 1$ .*

1. *If  $\emptyset \vdash M : \tau! \rho$  and  $\vdash_T M \downarrow^k V$  then  $\vdash_T \llbracket M \rrbracket \downarrow^{\leq Ck + D(\rho)} \llbracket V \rrbracket$ .*

2. *If  $\emptyset \vdash M : \tau! \rho$  and  $\vdash_N M \downarrow^k V$  then  $\vdash_N \llbracket M \rrbracket \downarrow^{\leq Ck} \llbracket V \rrbracket$ .*

where  $D(\rho)$  is a function such that  $D(i) = i$  and  $D(\omega) = \max(M)$ .

Proof. By mutual induction on the derivations of  $\vdash_T M \downarrow^k V$  and  $\vdash_N M \downarrow^k V$ .

Proof of Property 1.

Case:  $\vdash_T @^i M_1 M_2 \downarrow^k V$  is derived from  $\vdash_N M_1 \downarrow^l V_1$  and  $\vdash_N M_2 \downarrow^m V_2$  and  $\vdash_T M[V_1/x][V_2/y] \downarrow^n V$  where  $V_1 \equiv \text{fix } x.\lambda^{j^+} y.M$  and  $k = \max(l, m, n)$ . From the definition of the type system,  $j < i$ . From  $\emptyset \vdash @^i M_1 M_2 : \tau! \rho, i \leq \rho$ . We also have  $x : \tau' \rightarrow^{j^+} \tau, y : \tau' \vdash M : \tau! j$ .

By the induction hypothesis,  $\vdash_N \llbracket M_1 \rrbracket \downarrow^{\leq Cl} \llbracket V_1 \rrbracket$  and  $\vdash_N \llbracket M_2 \rrbracket \downarrow^{\leq Cm} \llbracket V_2 \rrbracket$ . From  $\emptyset \vdash M[V_1/x][V_2/y] : \tau! j$ , by the induction hypothesis,

$$\vdash_T \llbracket M \rrbracket \llbracket [V_1]/x \rrbracket \llbracket [V_2]/y \rrbracket \downarrow^{\leq Cn+j} \llbracket V \rrbracket$$

This case is proved by the following derivation.

$$\frac{\begin{array}{c} \vdash_N \llbracket M_1 \rrbracket.\text{fun} \downarrow^{\leq Cl} \lambda y. \llbracket M \rrbracket \llbracket [V_1]/x \rrbracket \quad \vdash_N \llbracket M_2 \rrbracket \downarrow^{\leq Cm} \llbracket V_2 \rrbracket \\ \vdash_T \llbracket M \rrbracket \llbracket [V_1]/x \rrbracket \llbracket [V_2]/y \rrbracket \downarrow^{\leq Cn+j} \llbracket V \rrbracket \end{array}}{\vdash_T @(\llbracket M_1 \rrbracket.\text{fun}) \llbracket M_2 \rrbracket \downarrow^{\leq \max(Cl, Cm, Cn+j+1)} \llbracket V \rrbracket}$$

where  $\max(Cl, Cm, Cn + j + 1) \leq Ck + D(\rho)$ .

Case:  $\vdash_T @^\omega M_1 M_2 \downarrow^k V$  is derived from  $\vdash_N M_1 \downarrow^l V_1$  and  $\vdash_N M_2 \downarrow^m V_2$  and  $\vdash_T M[V_1/x][V_2/y] \downarrow^n V$  where  $V_1 \equiv \text{fix } x.\lambda^{\rho'} y.M$  and  $k = \max(l, m, n)$ . From  $\emptyset \vdash @^\omega M_1 M_2 : \tau! \rho$ ,  $\rho$  must be  $\omega$ . By Lemma 4,  $x : \tau' \rightarrow^{\rho'} \tau, y : \tau' \vdash M : \tau! \rho''$  where  $\rho'' \leq \rho'$ . By the induction hypothesis,  $\vdash_N \llbracket M_1 \rrbracket \downarrow^{\leq Cl} \llbracket V_1 \rrbracket$  and  $\vdash_N \llbracket M_2 \rrbracket \downarrow^{\leq Cm} \llbracket V_2 \rrbracket$ . From  $\emptyset \vdash M[V_1/x][V_2/y] : \tau! \rho''$ , by the induction hypothesis

$$\vdash_T \llbracket M[V_1/x][V_2/y] \rrbracket \downarrow^{\leq Cn + D(\rho'')} \llbracket V \rrbracket$$

This case is proved by the following derivation.

$$\frac{\begin{array}{c} \vdash_N \llbracket M_1 \rrbracket . \text{tfun} \downarrow^{\leq Cl} \lambda^t y. \llbracket M \rrbracket \llbracket \llbracket V_1 \rrbracket \rrbracket \quad \vdash_N \llbracket M_2 \rrbracket \downarrow^{\leq Cm} \llbracket V_2 \rrbracket \\ \vdash_T \llbracket M \rrbracket \llbracket \llbracket V_1 \rrbracket / x \rrbracket \llbracket \llbracket V_2 \rrbracket / y \rrbracket \downarrow^{\leq Cn+D(\rho'')} \llbracket V \rrbracket \end{array}}{\vdash_T @^t(\llbracket M_1 \rrbracket . \text{tfun}) \llbracket M_2 \rrbracket \downarrow^{\leq \max(Cl, Cm, Cn+D(\rho''))} \llbracket V \rrbracket}$$

where  $\max(Cl, Cm, Cn + D(\rho'')) \leq Ck + C - 1 = Ck + D(\omega)$ .

Proof of Property 2.

Case:  $\vdash_N @^\omega M_1 M_2 \downarrow^k V$  is derived from  $\vdash_N M_1 \downarrow^l V_1$  and  $\vdash_N M_2 \downarrow^m V_2$  and  $\vdash_N M[V_1/x][V_2/y] \downarrow^n V$  where  $V_1 \equiv \text{fix } x. \lambda^{\rho'} y. M$  and  $k = \max(l, m, n + 1)$ . From  $\emptyset \vdash @^\omega M_1 M_2 : \tau! \rho$ ,  $\rho$  must be  $\omega$ . By Lemma 4, we have  $x : \tau' \rightarrow^{\rho'} \tau, y : \tau' \vdash M : \tau! \rho''$  where  $\rho'' \leq \rho'$ .

By the induction hypothesis,  $\vdash_N \llbracket M_1 \rrbracket \downarrow^{\leq Cl} \llbracket V_1 \rrbracket$  and  $\vdash_N \llbracket M_2 \rrbracket \downarrow^{\leq Cm} \llbracket V_2 \rrbracket$ . From  $\emptyset \vdash M[V_1/x][V_2/y] : \tau! \rho''$ , by the induction hypothesis,

$$\vdash_T \llbracket M \rrbracket \llbracket \llbracket V_1 \rrbracket / x \rrbracket \llbracket \llbracket V_2 \rrbracket / y \rrbracket \downarrow^{\leq Cn+D(\rho'')} \llbracket V \rrbracket$$

This case is proved by the following derivation.

$$\frac{\begin{array}{c} \vdash_N \llbracket M_1 \rrbracket . \text{tfun} \downarrow^{\leq Cl} \lambda^t y. \llbracket M \rrbracket \llbracket \llbracket V_1 \rrbracket / x \rrbracket \quad \vdash_N \llbracket M_2 \rrbracket \downarrow^{\leq Cm} \llbracket V_2 \rrbracket \\ \vdash_T \llbracket M \rrbracket \llbracket \llbracket V_1 \rrbracket / x \rrbracket \llbracket \llbracket V_2 \rrbracket / y \rrbracket \downarrow^{\leq Cn+D(\rho'')} \llbracket V \rrbracket \end{array}}{\vdash_N @^t(\llbracket M_1 \rrbracket . \text{tfun}) \llbracket M_2 \rrbracket \downarrow^{\leq \max(Cl, Cm, Cn+D(\rho'')+1)} \llbracket V \rrbracket}$$

where  $\max(Cl, Cm, Cn + D(\rho'') + 1) \leq Ck$ .

Case:  $\vdash_N @^i M_1 M_2 \downarrow^k V$  is derived from  $\vdash_N M_1 \downarrow^l V_1$  and  $\vdash_N M_2 \downarrow^m V_2$  and  $\vdash_N M[V_1/x][V_2/y] \downarrow^n V$  where  $V_1 \equiv \text{fix } x. \lambda^{j^+} y. M$  and  $k = \max(l, m, n + 1)$ . From the definition of the type system,  $j < i$ . From  $\emptyset \vdash @^i M_1 M_2 : \tau! \rho$ ,  $i \leq \rho$ . We also have  $x : \tau' \rightarrow^{j^+} \tau, y : \tau' \vdash M : \tau! j$ .

By the induction hypothesis,  $\vdash_N \llbracket M_1 \rrbracket \downarrow^{\leq Cl} \llbracket V_1 \rrbracket$  and  $\vdash_N \llbracket M_2 \rrbracket \downarrow^{\leq Cm} \llbracket V_2 \rrbracket$ . From  $\emptyset \vdash M[V_1/x][V_2/y] : \tau! j$ , by the induction hypothesis,

$$\vdash_T \llbracket M \rrbracket \llbracket \llbracket V_1 \rrbracket / x \rrbracket \llbracket \llbracket V_2 \rrbracket / y \rrbracket \downarrow^{Cn+j} \llbracket V \rrbracket$$

This case is proved by the following derivation.

$$\frac{\begin{array}{c} \vdash_N \llbracket M_1 \rrbracket . \text{fun} \downarrow^{\leq Cl} \lambda x. \llbracket M \rrbracket \quad \vdash_N \llbracket M_2 \rrbracket \downarrow^{\leq Cm} \llbracket V_2 \rrbracket \\ \vdash_N \llbracket M \rrbracket \llbracket \llbracket V_1 \rrbracket / x \rrbracket \llbracket \llbracket V_2 \rrbracket / y \rrbracket \downarrow^{\leq Cn+j} \llbracket V \rrbracket \end{array}}{\vdash_N @^t(\llbracket M_1 \rrbracket . \text{fun}) \llbracket M_2 \rrbracket \downarrow^{\leq \max(Cl, Cm, Cn+j+1)} \llbracket V \rrbracket}$$

where  $\max(Cl, Cm, Cn + j + 1) \leq Ck$ .  $\square$

## D Proof of Soundness for the Extended Effect System

We use the following lemma to simplify case-analysis.

- Lemma 6** • *If  $\emptyset \vdash \text{fix } x.\lambda^{\rho'} y.M : \tau_1 \rightarrow^{\rho} \tau_2$ , then  $x : \tau_1 \rightarrow^{\rho'} \tau_2, y : \tau_1 \vdash M : \tau_2! \rho''$  and  $\rho'' \leq \rho$  for some  $\rho''$ .*
- *If  $\emptyset \vdash \text{fix } x.\lambda^{\omega \cdot i + j} y.M : \tau_1 \rightarrow^{\rho} \tau_2$  for  $j > 0$ , then  $x : \tau_1 \rightarrow^{\omega \cdot i + j} \tau_2, y : \tau_1 \vdash M : \tau_2! \rho''$  and  $\rho''^+ \leq \omega \cdot i + j$  for some  $\rho''$ .*

We prove the soundness theorem in the following form. The structure of the proof is almost the same as that for the unextended system.

**Lemma 7** *Let  $C = \sum_{k=0}^{\infty} \max^k(M) + 1$  and  $D(\omega \cdot i + j) = \sum_{k=0}^{i-1} \max^k(M) + j$ .*

1. *If  $\emptyset \vdash M : \tau! \rho$  and  $\vdash_T M \downarrow^k V$  then  $\vdash_T \llbracket M \rrbracket \downarrow^{\leq Ck + D(\rho)} \llbracket V \rrbracket$ .*
2. *If  $\emptyset \vdash M : \tau! \rho$  and  $\vdash_N M \downarrow^k V$  then  $\vdash_N \llbracket M \rrbracket \downarrow^{\leq Ck} \llbracket V \rrbracket$ .*

Proof. By mutual induction on the derivations of  $\vdash_T M \downarrow^k V$  and  $\vdash_N M \downarrow^k V$ . We use the following property of  $D(\rho)$ .

$$\begin{aligned} \rho_1 \leq \rho_2 &\Rightarrow D(\rho_1) \leq D(\rho_2) \\ \rho_1^+ \leq \rho_2 &\Rightarrow D(\rho_1) + 1 \leq D(\rho_2) \end{aligned}$$

Proof of Property 1.

Case:  $\vdash_T @^{\omega \cdot i + j} M_1 M_2 \downarrow^k V$  is derived from  $\vdash_N M_1 \downarrow^l V_1$  and  $\vdash_N M_2 \downarrow^m V_2$  and  $\vdash_T M[V_1/x][V_2/y] \downarrow^n V$  where  $j > 0$  and  $V_1 \equiv \text{fix } x.\lambda^{\rho'} y.M_0$  and  $k = \max(l, m, n)$ . By definition of the type system,  $\omega \cdot i + j \leq \rho$ .

From Lemma 6 we have  $x : \tau' \rightarrow^{\rho'} \tau, y : \tau' \vdash M : \tau! \rho''$  and  $\rho''^+ \leq \omega \cdot i + j \leq \rho$  for some  $\rho''$ .

By the induction hypothesis,  $\vdash_N \llbracket M_1 \rrbracket \downarrow^{\leq Cl} \llbracket V_1 \rrbracket$  and  $\vdash_N \llbracket M_2 \rrbracket \downarrow^{\leq Cm} \llbracket V_2 \rrbracket$ . From  $\emptyset \vdash M[V_1/x][V_2/y] : \tau! \rho''$ , by the induction hypothesis,

$$\vdash_T \llbracket M_0 \rrbracket \llbracket \llbracket V_1 \rrbracket / x \rrbracket \llbracket \llbracket V_2 \rrbracket / y \rrbracket \downarrow^{\leq Cn + D(\rho'')} \llbracket V \rrbracket$$

This case is proved by the following derivation.

$$\frac{\begin{array}{c} \vdash_N \llbracket M_1 \rrbracket . \text{fun} \downarrow^{\leq Cl} \lambda y. \llbracket M_0 \rrbracket \llbracket \llbracket V_1 \rrbracket / x \rrbracket \quad \vdash_N \llbracket M_2 \rrbracket \downarrow^{\leq Cm} \llbracket V_2 \rrbracket \\ \vdash_T \llbracket M_0 \rrbracket \llbracket \llbracket V_1 \rrbracket / x \rrbracket \llbracket \llbracket V_2 \rrbracket / y \rrbracket \downarrow^{\leq Cn + D(\rho'')} \llbracket V \rrbracket \end{array}}{\vdash_T @(\llbracket M_1 \rrbracket . \text{fun}) \llbracket M_2 \rrbracket \downarrow^{\leq \max(Cl, Cm, Cn + D(\rho'') + 1)} \llbracket V \rrbracket}$$

where  $\max(Cl, Cm, Cn + D(\rho'') + 1) \leq Ck + D(\rho)$ .

Case:  $\vdash_T @^{\omega \cdot i} M_1 M_2 \downarrow^k V$  is derived from  $\vdash_N M_1 \downarrow^l V_1$  and  $\vdash_N M_2 \downarrow^m V_2$  and  $\vdash_T M[V_1/x][V_2/y] \downarrow^n V$  where  $V_1 \equiv \text{fix } x. \lambda^{\rho'} y. M$  and  $k = \max(l, m, n)$ . By definition of the type system,  $\omega \cdot i \leq \rho$ .

From Lemma 6,  $x : \tau' \rightarrow^{\rho'} \tau, y : \tau' \vdash M : \tau! \rho''$  where  $\rho'' \leq \omega \cdot i \leq \rho$ .

By the induction hypothesis,  $\vdash_N \llbracket M_1 \rrbracket \downarrow^{\leq Cl} \llbracket V_1 \rrbracket$  and  $\vdash_N \llbracket M_2 \rrbracket \downarrow^{\leq Cm} \llbracket V_2 \rrbracket$ .

From  $\emptyset \vdash M[V_1/x][V_2/y] : \tau! \rho'$ , by the induction hypothesis

$$\vdash_T \llbracket M \rrbracket \llbracket \llbracket V_1 \rrbracket / x \rrbracket \llbracket \llbracket V_2 \rrbracket / y \rrbracket \downarrow^{\leq Cn + D(\rho'')} \llbracket V \rrbracket$$

This case is proved by the following derivation.

$$\frac{\begin{array}{c} \vdash_N \llbracket M_1 \rrbracket. \text{tfun} \downarrow^{\leq Cl} \lambda^t y. \llbracket M \rrbracket \llbracket \llbracket V_1 \rrbracket \rrbracket \quad \vdash_N \llbracket M_2 \rrbracket \downarrow^{\leq Cm} \llbracket V_2 \rrbracket \\ \vdash_T \llbracket M \rrbracket \llbracket \llbracket V_1 \rrbracket / x \rrbracket \llbracket \llbracket V_2 \rrbracket / y \rrbracket \downarrow^{\leq Cn + D(\rho'')} \llbracket V \rrbracket \end{array}}{\vdash_T @^t (\llbracket M_1 \rrbracket. \text{tfun}) \llbracket M_2 \rrbracket \downarrow^{\leq \max(Cl, Cm, Cn + D(\rho''))} \llbracket V \rrbracket}$$

where  $\max(Cl, Cm, Cn + D(\rho'')) \leq Ck + D(\rho)$ .

Proof of Property 2.

Case:  $\vdash_N @^{\omega \cdot i} M_1 M_2 \downarrow^k V$  is derived from  $\vdash_N M_1 \downarrow^l V_1$  and  $\vdash_N M_2 \downarrow^m V_2$  and  $\vdash_N M[V_1/x][V_2/y] \downarrow^n V$  where  $V_1 \equiv \text{fix } x. \lambda^{\rho'} y. M$  and  $k = \max(l, m, n + 1)$ .

From Lemma 6,  $x : \tau' \rightarrow^{\rho'} \tau, y : \tau' \vdash M : \tau! \rho''$  and  $\rho'' \leq \omega \cdot i \leq \rho$  for some  $\rho''$ .

By the induction hypothesis,  $\vdash_N \llbracket M_1 \rrbracket \downarrow^{\leq Cl} \llbracket V_1 \rrbracket$  and  $\vdash_N \llbracket M_2 \rrbracket \downarrow^{\leq Cm} \llbracket V_2 \rrbracket$ .

From  $\emptyset \vdash M[V_1/x][V_2/y] : \tau! \rho''$ , by the induction hypothesis,

$$\vdash_T \llbracket M \rrbracket \llbracket \llbracket V_1 \rrbracket / x \rrbracket \llbracket \llbracket V_2 \rrbracket / y \rrbracket \downarrow^{\leq Cn + D(\rho'')} \llbracket V \rrbracket$$

This case is proved by the following derivation.

$$\frac{\begin{array}{c} \vdash_N \llbracket M_1 \rrbracket. \text{tfun} \downarrow^{\leq Cl} \lambda^t y. \llbracket M \rrbracket \llbracket \llbracket V_1 \rrbracket / x \rrbracket \quad \vdash_N \llbracket M_2 \rrbracket \downarrow^{\leq Cm} \llbracket V_2 \rrbracket \\ \vdash_T \llbracket M \rrbracket \llbracket \llbracket V_1 \rrbracket / x \rrbracket \llbracket \llbracket V_2 \rrbracket / y \rrbracket \downarrow^{\leq Cn + D(\rho'')} \llbracket V \rrbracket \end{array}}{\vdash_N @^t (\llbracket M_1 \rrbracket. \text{tfun}) \llbracket M_2 \rrbracket \downarrow^{\leq \max(Cl, Cm, Cn + D(\rho'') + 1)} \llbracket V \rrbracket}$$

where  $\max(Cl, Cm, Cn + D(\rho'') + 1) \leq Ck$  because  $C \geq D(\rho'') + 1$ .

Case:  $\vdash_N @^{\omega \cdot i + j} M_1 M_2 \downarrow^k V$  is derived from  $\vdash_N M_1 \downarrow^l V_1$  and  $\vdash_N M_2 \downarrow^m V_2$  and  $\vdash_N M[V_1/x][V_2/y] \downarrow^n V$  where  $j > 0$  and  $V_1 \equiv \text{fix } x. \lambda^{\rho'} y. M$  and  $k = \max(l, m, n + 1)$ . By definition of the type system,  $\omega \cdot i + j \leq \rho$ .

From Lemma 6 we have  $x : \tau' \rightarrow^{\rho'} \tau, y : \tau' \vdash M : \tau! \rho''$  and  $\rho'' \leq \omega \cdot i + j \leq \rho$  for some  $\rho''$ .

By the induction hypothesis,  $\vdash_N \llbracket M_1 \rrbracket \downarrow^{\leq Cl} \llbracket V_1 \rrbracket$  and  $\vdash_N \llbracket M_2 \rrbracket \downarrow^{\leq Cm} \llbracket V_2 \rrbracket$ .

From  $\emptyset \vdash M[V_1/x][V_2/y] : \tau! \rho''$ , by the induction hypothesis,

$$\vdash_T \llbracket M \rrbracket \llbracket \llbracket V_1 \rrbracket / x \rrbracket \llbracket \llbracket V_2 \rrbracket / y \rrbracket \downarrow^{\leq Cn + D(\rho'')} \llbracket V \rrbracket$$

This case is proved by the following derivation.

$$\frac{\begin{array}{c} \vdash_N \llbracket M_1 \rrbracket. \text{fun} \downarrow^{\leq Cl} \lambda y. \llbracket M \rrbracket \llbracket \llbracket V_1 \rrbracket / x \rrbracket \quad \vdash_N \llbracket M_2 \rrbracket \downarrow^{\leq Cm} \llbracket V_2 \rrbracket \\ \vdash_T \llbracket M \rrbracket \llbracket \llbracket V_1 \rrbracket / x \rrbracket \llbracket \llbracket V_2 \rrbracket / y \rrbracket \downarrow^{\leq Cn + D(\rho'')} \llbracket V \rrbracket \end{array}}{\vdash_N @ (\llbracket M_1 \rrbracket. \text{fun}) \llbracket M_2 \rrbracket \downarrow^{\leq \max(Cl, Cm, Cn + D(\rho'') + 1)} \llbracket V \rrbracket}$$

where  $\max(Cl, Cm, Cn + D(\rho'') + 1) \leq Ck$  because  $C \geq D(\rho'') + 1$ .  $\square$