

# Shrinking Reductions in SML.NET

Nick Benton<sup>1</sup>, Andrew Kennedy<sup>1</sup>, Sam Lindley<sup>2</sup>, and Claudio Russo<sup>1</sup>

<sup>1</sup> Microsoft Research, Cambridge {nick, akenn, crusso}@microsoft.com

<sup>2</sup> LFCS, University of Edinburgh Sam.Lindley@ed.ac.uk

**Abstract.** One performance-critical phase in the SML.NET compiler involves rewriting intermediate terms to monadic normal form and performing non-duplicating  $\beta$ -reductions. We present an imperative algorithm for this simplification phase, working with a mutable, pointer-based term representation, which significantly outperforms our existing functional algorithm. This is the first implementation and evaluation of a linear-time rewriting algorithm proposed by Appel and Jim.

## 1 Introduction

SML.NET [3,4] is a compiler for Standard ML that targets the .NET Common Language Runtime [7]. Like most other compilers for functional languages (e.g. GHC [10]), SML.NET is structured as the composition of a number of transformation phases on an intermediate representation of the user program. As SML.NET is a whole program compiler, the intermediate terms are typically rather large and good performance of the transformations is critical for usability.

Like MLj [5], SML.NET uses a monadic intermediate language (MIL) [2] that is similar to Moggi’s computational metalanguage. Most of the phases in SML.NET perform specific transformations, such as closure conversion, arity raising or monomorphisation, and are run only once. In between several of these phases, however, is a general-purpose ‘clean-up’ pass called *simplify*. Running *simplify* puts the term into *monadic normal form* [6,8], which we have previously called *cc-normal form* and is essentially the same as *A normal form* or *administrative normal form* for CPS [8]. The *simplify* pass also performs *shrinking reductions*:  $\beta$ -reductions for functions, computations, products that always reduce the size of the term.

Appel and Jim [1] describe three algorithms for shrinking reductions. The first ‘naïve’ and second ‘improved’ algorithms both have quadratic worst-case time complexity, and the third ‘imperative’ algorithm is linear, but requires a mutable representation of terms. Appel and Jim did not implement the third algorithm, which does not integrate easily in a mainly-functional compiler. Both SML/NJ and SML.NET use the ‘improved’ algorithm, which is reasonably efficient in practice. Nevertheless, SML.NET spends a significant amount of time performing shrinking reductions. We have now implemented a variant of the imperative algorithm in SML.NET, and achieved significant speedups.

This paper makes several contributions. It gives the first implementation and benchmarks of the imperative algorithm in a real compiler. It extends the

imperative algorithm to a richer language than considered by Appel and Jim. It introduces a ‘one-pass’ traversal strategy, giving a weak form of compositionality. An extended version of this work appears in the third author’s PhD thesis [9].

## 2 Simplified MIL

For purposes of exposition we present a simplified version of MIL:

Atoms	$a, b ::= x \mid c$
Values	$v, w ::= a \mid \text{pair}(a, b) \mid \text{proj}_1(a) \mid \text{proj}_2(a) \mid \text{inj}_1(a) \mid \text{inj}_2(a)$
Computations	$m, n, p ::= \text{app}(a, b) \mid \text{letfun } f(x) \text{ be } m \text{ in } n$ $\mid \text{val}(v) \mid \text{let } x \text{ be } m \text{ in } n \mid \text{case } a \text{ of } (x_1)n_1 ; (x_2)n_2$

where variables are ranged over by  $f, g, x, y, z$ , and constants are ranged over by  $c$ . Note that the `letfun` construct binds a possibly recursive function.

We say that a reduction is a *shrinking* reduction if it always reduces the size of terms (counting the number of nodes). The most important reductions are given by the shrinking  $\beta$ -rules:

$(\rightarrow.\beta_0)$	$\text{letfun } f(x) \text{ be } n \text{ in } m \longrightarrow m,$	$f \notin \text{fv}(m)$
$(\rightarrow.\beta_1)$	$\text{letfun } f(x) \text{ be } m \text{ in } C[\text{app}(f, a)] \longrightarrow C[m[x := a]],$	$f \notin \text{fv}(C[\cdot], m, a)$
$(T.\beta_0)$	$\text{let } x \text{ be val}(v) \text{ in } m \longrightarrow m,$	$x \notin \text{fv}(m)$
$(T.\beta_a)$	$\text{let } x \text{ be val}(a) \text{ in } m \longrightarrow m[x := a]$	
$(\times.\beta)$	$\text{let } y \text{ be val}(\text{pair}(a_1, a_2)) \text{ in } C[\text{proj}_i(y)]$ $\longrightarrow \text{let } y \text{ be val}(\text{pair}(a_1, a_2)) \text{ in } C[a_i]$	
$(+.\beta)$	$\text{let } y \text{ be val}(\text{inj}_i(a))$ $\text{in } C[\text{case } y \text{ of } (x_1)n_1 ; (x_2)n_2]$ $\longrightarrow \text{let } y \text{ be val}(\text{inj}_i(a)) \text{ in } C[n_i[x_i := a]]$	

We write  $R_\beta$  for the one-step reduction relation defined by the  $\beta$ -rules. The *simplify* transformation also performs commuting conversions. These ensure that bindings are explicitly sequenced, which enables further rewriting.

$(T.CC)$	$\text{let } y \text{ be } (\text{let } x \text{ be } m \text{ in } n) \text{ in } p$ $\longrightarrow \text{let } x \text{ be } m \text{ in let } y \text{ be } n \text{ in } p$
$(\rightarrow.CC)$	$\text{let } y \text{ be } (\text{letfun } f(x) \text{ be } m \text{ in } n) \text{ in } p$ $\longrightarrow \text{letfun } f(x) \text{ be } m \text{ in let } y \text{ be } n \text{ in } p$
$(+.CC)$	$\text{let } y \text{ be } (\text{case } a \text{ of } (x_1)n_1 ; (x_2)n_1) \text{ in } m$ $\longrightarrow \text{letfun } f(y) \text{ be } m \text{ in case } a \text{ of } (x_1)\text{let } y_1 \text{ be } n_1 \text{ in app}(f, y_1)$ $\quad ; (x_2)\text{let } y_2 \text{ be } n_2 \text{ in app}(f, y_2)$

We write  $R_{CC}$  for the one-step reduction relation defined by the CC-rules, and  $R$  for  $R_\beta \cup R_{CC}$ . Unlike the  $\beta$  rules, the commuting conversions are not actually

shrinking reductions. However,  $T.CC$  and  $\rightarrow.CC$  do not change the size, whilst  $+.CC$  gives only a constant increase in the size.

An alternative to the  $+.CC$  rule is:

$$\begin{aligned} (+.CC') \quad & \text{let } y \text{ be case } a \text{ of } (x_1)n_1 ; (x_2)n_2 \text{ in } m \\ & \longrightarrow \text{case } a \text{ of } (x_1)\text{let } y_1 \text{ be } n_1 \text{ in } m_1 ; (x_2)\text{let } y_2 \text{ be } n_2 \text{ in } m_2 \end{aligned}$$

where  $y_1, y_2$  are fresh,  $m_i = m[y := y_i]$ . This rule duplicates the term  $m$  and can exponentially increase the term's size. The  $+.CC$  rule instead creates a single new abstraction, shared across both branches of the **case**, though this inhibits some further rewriting. We write  $R'_{CC}$  for the one-step relation defined by the CC-rules where  $(+.CC)$  is replaced by  $(+.CC')$ , and  $R'$  for  $R_\beta \cup R'_{CC}$ .

**Proposition 1.**  *$R'$  is strongly-normalising.*

*Proof.* First, note that  $R_\beta$  is strongly-normalising as  $R_\beta$ -reduction strictly decreases the size of terms. We define two measures  $|\cdot|_\beta$  and  $|\cdot|_{cc}$  on terms:

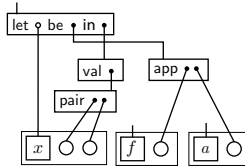
$$\begin{aligned} |a|_\beta &= 1 & |\text{letfun } f(x) \text{ be } m \text{ in } n|_\beta &= |m|_\beta + |n|_\beta + 1 \\ |\text{proj}_i(a)|_\beta &= |\text{inj}_i(a)|_\beta = 2 & |\text{let } x \text{ be } m \text{ in } n|_\beta &= |m|_\beta + |n|_\beta + 1 \\ |\text{app}(a, b)|_\beta &= |\text{pair}(a, b)|_\beta = 3 & |\text{val}(v)|_\beta &= |v|_\beta + 1 \\ |\text{case } a \text{ of } (x_1)n_1 ; (x_2)n_2|_\beta &= \max(|n_1|_\beta, |n_2|_\beta) + 2 \\ |a|_{cc} &= 1 & |\text{letfun } f(x) \text{ be } m \text{ in } n|_{cc} &= |m|_{cc} + |n|_{cc} + 1 \\ |\text{proj}_i(a)|_{cc} &= |\text{inj}_i(a)|_{cc} = 2 & |\text{let } x \text{ be } m \text{ in } n|_{cc} &= |m|_{cc}^2 + |n|_{cc} + 1 \\ |\text{app}(a, b)|_{cc} &= |\text{pair}(a, b)|_{cc} = 3 & |\text{val}(v)|_{cc} &= |v|_{cc} + 1 \\ |\text{case } a \text{ of } (x_1)n_1 ; (x_2)n_2|_{cc} &= \max(|n_1|_{cc}, |n_2|_{cc}) + 2 \end{aligned}$$

The lexicographic ordering  $(|\cdot|_\beta, |\cdot|_{cc})$  is a measure for  $R'$ -reduction. Each shrinking  $\beta$ -reduction decreases  $|\cdot|_\beta$ , whilst each CC-reduction decreases  $|\cdot|_{cc}$  and leaves  $|\cdot|_\beta$  unchanged.  $\square$

**Proposition 2.**  *$R$  is strongly-normalising.*

The proof uses  $R'$ -reduction to simulate  $R$ -reduction. The full details are omitted, but the idea is that for any  $R$ -reduction a corresponding non-empty sequence of  $R'$ -reductions can be performed. Thus, given that all  $R'$ -reduction sequences are finite, all  $R$ -reduction sequences must also be finite. The proof is slightly complicated by the fact that no non-empty sequence of  $R'$ -reductions corresponds with the  $\beta$ -reduction of a function introduced by the  $+.CC$  rule. A simple way of dealing with this is to count a  $+.CC'$ -reduction as two reductions.

Note that  $R$ -reductions are not confluent. The failure of confluence is due to the  $(+.CC)$  rule. Replacing  $(+.CC)$  with  $(+.CC')$  does give a confluent system. Confluence can make reasoning about reductions easier, but we do not regard failure of confluence as a problem. In our case, preventing exponential growth in the size of terms is far more important.



**Fig. 1.** Pictorial representation of `let  $x$  be  $\text{app}(f, a)$  in  $\text{val}(\text{pair}(x, x))$`

### 3 Previous Work

Appel and Jim [1] considered a calculus which is equivalent to a sub-calculus of our simplified MIL. In our setting the reductions that their algorithms perform are equivalent to:  $\rightarrow .\beta_1-$ ,  $\times.\beta-$ ,  $T.\beta_0-$ , and a restriction of  $\rightarrow .\beta_0$ -reduction. Appel and Jim show that their calculus is confluent in the presence of these reductions, and other ‘ $\delta$ -rules’ satisfying certain criteria.

The reductions rely on knowing the number of occurrences of a particular variable. The quadratic algorithms store this information in a table *Count* mapping variable names to their number of occurrences. Appel and Jim’s naïve algorithm repeatedly (i) zeros the usage counts, (ii) performs a *census* pass over the whole term to update the usage counts and then (iii) traverses the term performing reductions on the basis of the information in *Count*, until there are no redexes remaining.

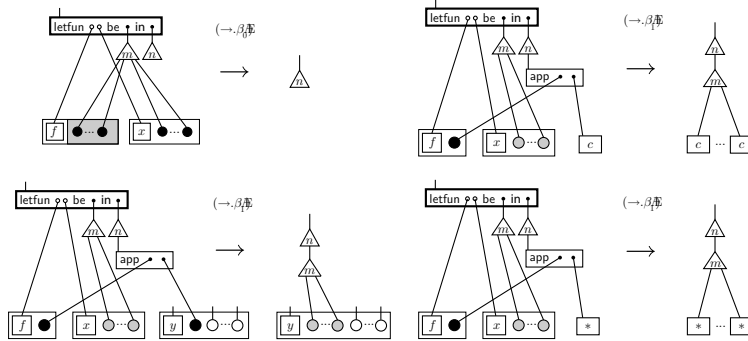
The improved algorithm, used in SML/NJ and SML.NET, dynamically updates the usage counts as reductions are performed. This allows more reductions to be performed on each pass, and only requires a full census to be performed once. The improved algorithm is better in practice, but both algorithms have worst-case time complexity  $\Theta(n^2)$  where  $n$  is the size of the input term.

Appel and Jim’s imperative algorithm runs in linear time and uses a pointer-based representation of terms which directly links all occurrences of a particular variable. This enables an efficient test to see if removing an occurrence will create any new redexes, and an efficient way of jumping to any such redexes. The algorithm first traverses the program tree collecting the set of all redexes. Then it repeatedly removes a redex from the set and reduces it in-place (possibly adding new redexes to the set), until none remain.

### 4 A Graph-based Representation

Our imperative algorithm works with a mutable graph representation comprising a doubly-linked expression tree and a list of pairs of circular doubly-linked lists collecting all the recursive (respectively non-recursive) uses of each variable. Such graphs can naturally be presented pictorially as shown by the example in Fig. 1.

Figure 2 shows the  $\beta$ -reductions for functions in this pictorial form. We find the pictorial representation intuitively very useful, but awkward to reason with



**Fig. 2.** Graph reductions

or use in presenting algorithms. Hence, like Appel and Jim, we will work with a more abstract structure comprising an expression tree and a collection of maps which capture the additional graphical structure between nodes of the tree.

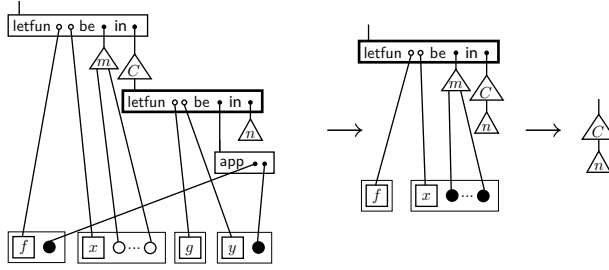
The structure of expression trees is determined by the abstract syntax of simplified MIL. In order to capture mutability we use ML-style references. Each node of the expression tree is a reference cell. We call the entities which reference cells contain *objects*. Given a reference cell  $l$ , we write  $!l$  to denote the object of  $l$ , and  $l := u$  to denote the assignment of the object  $u$  to  $l$ .

Atoms	$!a, !b ::= r \mid c$
Values	$!v, !w ::= a \mid \text{pair}(a, b) \mid \text{proj}_1(a) \mid \text{proj}_2(a) \mid \text{inj}_1(a) \mid \text{inj}_2(a)$
Computations	$!m, !n, !p ::= \text{app}(a, b) \mid \text{letfun } f(x) \text{ be } m \text{ in } n$ $\mid \text{val}(v) \mid \text{let } x \text{ be } m \text{ in } n \mid \text{case } a \text{ of } (x_1)n_1 ; (x_2)n_2$ $e ::= v \mid m \quad d ::= e \mid x \mid r$

where  $f, g, x, y, z$  range over defining occurrences, and  $r, s, t$  over uses. We write  $\text{parent}(e)$  for the parent of the node  $e$ . A distinguished sentinel node,  $\text{root}$ , marks the top of the expression tree. The object  $\text{dead}$  (omitted from the grammar) is used to indicate a *dead* node. If a node is dead then it has no parent. The  $\text{root}$  node is the parent of the proper expression tree and is always dead. We define  $\text{children}(e)$  of an expression node to be the set of nodes appearing in  $!e$ .

Initially both  $\text{parent}$  and  $\text{children}$  are entirely determined by the expression tree. However, in our algorithm we take advantage of the  $\text{parent}$  map in order to classify expression nodes as active or inactive. We ensure that the following invariant is maintained: for all expression nodes  $e$ , either

- $e$  is active:  $\text{parent}(d) = e$ , for all  $d \in \text{children}(e)$ ;
- $e$  is inactive:  $!(\text{parent}(d)) = \text{dead}$  for all  $d \in \text{children}(e)$ ; or
- $e$  is dead:  $!e = \text{dead}$ .



**Fig. 3.** Triggering non-local reductions

We define *splicing* as the operation which takes one subtree  $m$  and substitutes it in place of another subtree  $n$ . The subtree  $m$  is removed from the expression tree and then reintroduced in place of  $n$ . The parent map is adjusted accordingly for the children of  $m$ . We define *splicing a copy* as the corresponding operation which leaves the original copy of  $m$  in place. The operation  $[q]$  returns a new node containing  $q$ , with parent  $root$ . When embedded in an enclosing node  $e[[q]]$ , the parent of  $[q]$  is  $e$ . In patterns,  $[\cdot]$  matches against the contents of a node.

The *def-use* maps abstract the structures used for representing occurrences:

- $def(r)$  gives the defining occurrence of the use  $r$ .
- $non-rec-uses(x)$  is the set of non-recursive uses of the defining occurrence  $x$ .
- $rec-uses(x)$  is the set of recursive uses of the defining occurrence  $x$ .

In the real implementation occurrences are held in a pair of doubly-linked circular lists, such that each pair of lists intersects at a defining occurrence. We find it convenient to overload the maps to be defined over all occurrences and also define some additional maps:

$$\begin{aligned}
 non-rec-uses(r) &= non-rec-uses(def(r)) \\
 rec-uses(r) &= rec-uses(def(r)) & def(x) &= x \\
 occurrences(r) &= uses(r) \cup \{def(r)\} & uses(r) &= non-rec-uses(r) \cup rec-uses(r)
 \end{aligned}$$

None of these additional definitions affects the implementation.

The graph structure allows constant time movement up and down the expression tree in the normal way, but also allows constant time non-local movement via the occurrence lists. For example, consider the dead-function eliminations:

$$\begin{aligned}
 &letfun f(x) be m in C[letfun g(y) be app(f, y) in n] \\
 \longrightarrow_{(\rightarrow, \beta_0)} &letfun f(x) be m in C[n] \longrightarrow_{(\rightarrow, \beta_0)} C[n]
 \end{aligned}$$

where  $f, g \notin fv(C, n)$ , illustrated in Fig. 3. After one reduction,  $g$  is dead, so its definition can be deleted, removing the only use of  $f$ . Since this use is connected to its defining occurrence, we can detect that the definition of  $f$  is now dead. The defining occurrence is connected to its parent ( $root$ ) so the new dead-function redex can be reduced under the parent.

## 5 A One-pass Algorithm

In contrast to Appel and Jim's imperative algorithm, the algorithm we have implemented operates in one-pass. Essentially, the one-pass algorithm performs a depth-first traversal of the expression tree, reducing redexes on the way back up the tree. Of course, these reductions may trigger further reductions elsewhere in the tree. By carefully deactivating parts of the tree, we are able to control the reduction order and limit the testing required for new redexes. Here is an outline of our one-pass imperative algorithm:

```

contract(e) = reduceCCs(e)
              deactivate(e)
              apply contract to children of e
              reactivate(e)
              reduce(true, e)
reduce(initial, e) = if e is a redex then
                    reduce e in place
                    perform further reductions triggered by reducing e

```

The operation  $reduceCCs(e)$  performs commuting conversions on the way down the tree. The order of commuting conversions can have a significant effect on code quality, a poor choice leading to many jumps to jumps. We have found that the approach of doing them on the way down works well in practice (although the contract algorithm would still be valid without the call to  $reduceCCs$ ).

```

reduceCCs(e) = case !e of
  (let y be e' in p) =>
    if reduceCC(e, y, e', p) ≠ ∅ then reduceCCs(e) else skip
  (-) => skip
reduceCC(e, y, e', p) = case !e' of
  (letfun f(x) be m in n) =>
    splice [let y be n in p] in place of e'
    splice [letfun f(x) be m in e'] in place of e
    return {e'}
  (let x be m in n) =>
    splice [let y be n in p] in place of e'
    splice [let x be m in e'] in place of e
    return {e'}
  (case a of (x1)n1 ; (x2)n2) =>
    splice [let y1 be n1 in [app(f, y1)]] in place of n1
    splice [let y2 be n2 in [app(f, y2)]] in place of n2
    splice [letfun f(y) be p in [case a of (x1)n1 ; (x2)n2]]
    in place of e (where f is fresh)
    return {n1, n2}
  (-) => return ∅

```

Note that commuting conversions can also be triggered by other reductions. The return value for  $reduceCC$  will be used in the definition of  $reduce$  in order to catch reductions which are triggered by applying commuting conversions.

*deactivate*(*e*) deactivates *e*: *parent*(*d*) is set to **dead** for every  $d \in \text{children}(e)$ .  
*reactivate*(*e*) reactivates *e*: *parent*(*d*) is set to *e* for every  $d \in \text{children}(e)$ .

Deactivating nodes on the way down prevents reductions from being triggered above the current node in the tree. On the way back up the nodes are reactivated, allowing any new redexes to be reduced. Because subterms are known to be normalised, fewer tests are needed for new redexes. Consider, for example:

$$\text{let } y \text{ be } (\text{let } x \text{ be } m \text{ in } n) \text{ in } p \longrightarrow_{T.CC} \text{let } x \text{ be } m \text{ in let } y \text{ be } n \text{ in } p$$

Because we know that *let* *x* *be* *m* *in* *n* is in normal form, *m* cannot be of the form *let*(...), *letfun*(...), *case*(...) or *val*(...). Hence, it is not necessary to check whether *let* *x* *be* *m* *in* *let* *y* *be* *n* *in* *p* is a redex. (Of course, *let* *y* *be* *n* *in* *p* may still be a redex, and indeed exposing such redexes is one of the main purposes of performing CC-reduction.)

## 5.1 Reduction

The *reduce* function is the heart of the algorithm. Rather than maintaining a global ‘work-list’ of redexes, as Appel and Jim do, *reduce*(*initial*, *e*) reduces any new redexes created inside *e* (but none that are created above *e* in the expression tree). *initial* is boolean flag indicating whether this call to reduce originates from *contract* rather than some other recursive call. If *reduce*(*initial*, *e*) is invoked on an expression node which is not a redex, then no action is performed. The *reduce* function also returns a boolean to indicate whether a reduction took place. As we shall see, this is necessary in order to detect the triggering of new reductions. We now expand the definition of *reduce*.

```

reduce(initial, e) = case !e of
  (letfun f(x) be m in n) ⇒
    if non-rec-uses(f) = ∅ then
      splice n in place of e
      reduceOccs(cleanExp(m))
      return true
    else if rec-uses(f) = ∅ and non-rec-uses(f) = {f'} then
      let focus = parent(parent(f'))
      case !focus of
        (app(f', a) ⇒
          splice n in place of e
          splice m in place of focus
          let (occs, redexes) = substAtom(x, a)
          reduceOccs(occs ∪ cleanExp(a))
          reduceRedexes(redexes)
          return true
        (-) ⇒ return false
      else return false
  (let x be [val(v)] in n) ⇒

```



```

if  $uses(x) = \emptyset$  then
  splice  $n$  in place of  $e$ 
   $reduceOccs(cleanExp(parent(v)))$ 
  return true
else if  $v$  is an atom  $a$  then
  splice  $n$  in place of  $e$ 
  let  $(occs, redexes) = substAtom(x, a)$ 
   $reduceOccs(occs \cup cleanExp(parent(a)))$ 
   $reduceRedexes(redexes)$ 
  return true
else case ! $v$  of
(pair( $a, b$ ))  $\Rightarrow$ 
  if initial then
    let  $redexes = reduceProjections(e, x, a, b, uses(x))$ 
    if  $redexes = \emptyset$  then return false
    else
       $reduceRedexes(redexes)$ 
       $reduce(false, e)$ 
      return true
  else return false
(inj $i$ ( $a$ ))  $\Rightarrow$ 
  if initial then
    let  $(occs, redexes) = reduceCases(e, x, i, a, uses(x))$ 
    if  $redexes = \emptyset$  then return false
    else
       $reduceOccs(occs)$ 
       $reduceRedexes(redexes)$ 
       $reduce(false, e)$ 
      return true
  else return false
(-)  $\Rightarrow$  return false
(let  $y$  be  $e'$  in  $p$ )  $\Rightarrow$ 
  let  $redexes = reduceCC(e, y, e', p)$ 
  for  $e'' \in redexes$  do  $reduce(false, e'')$ 
  return true
(-)  $\Rightarrow$  return false

```

The first case covers  $\beta$ -reductions on functions, with two sub-cases:

- ( $\rightarrow \beta_0$ ) If the function is dead, its definition is removed, the continuation spliced in place of  $e$ , and any uses within the dead body deleted, possibly triggering new reductions.
- ( $\rightarrow \beta_1$ ) If the function has one occurrence, which is non-recursive, it is inlined. The continuation of  $e$  is spliced in place of  $e$ , the function body is inlined with the argument substituted for the parameter, and the argument deleted. Substitution may trigger further reductions.

The second case covers  $\beta$ -reductions on computations as well as some instances of  $\beta$ -reduction on products and sums. It is divided into four sub-cases.

- $(T.\beta_0)$  If a value is dead, then its definition can be removed. The continuation is spliced in place of  $e$ . Then the uses inside the dead function body are deleted, possibly triggering new reductions.
- $(T.\beta_a)$  If a value is atomic, then it can be inlined. First the continuation of  $e$  is spliced in place of  $e$ . Then the atom is substituted for the bound variable. Finally the atom is deleted.
- $(\times.\beta)$  If a pair is bound to a variable  $x$ , and this is the initial visit of  $e$ , then any projections of  $x$  are reduced. For efficiency, new projections will subsequently be reduced as and when they are created.
- $(+.\beta)$  This follows exactly the same pattern as  $\times.\beta$ -reduction. The only difference is that the reduction itself is more complex, so can trigger new reductions in different ways.

The third case deals with commuting conversions.

The algorithm ensures that the current reduction is complete before any new reductions are triggered. Potential new redexes created by the current reduction are encoded and executed after the current reduction has completed.

$reduceUp(e)$  reduces above  $e$  as far as possible:

$$reduceUp(e) = \text{if } reduce(\text{false}, e) \text{ then } reduceUp(\text{parent}(e)) \text{ else skip}$$

$reduceRedexes$  reduces a set of expression redexes, whilst  $reduceOccs$  reduces a set of occurrence redexes:

$$\begin{aligned} reduceRedexes(redexes) &= \text{for each } e \in redexes \text{ do } reduceUp(e) \\ reduceOccs(xs) &= \text{for each } r \in xs \text{ do} \\ &\quad \text{if } isSmall(r) \text{ then } reduceUp(\text{parent}(\text{def}(r))) \text{ else skip} \\ isSmall(r) &= r \notin rec\text{-uses}(r) \text{ and } |non\text{-rec}\text{-uses}(r)| \leq 1 \end{aligned}$$

$cleanExp(e)$  removes all occurrences and subexpressions inside  $e$  and returns a set of occurrence redexes.

$$\begin{aligned} cleanExp(e) &= \text{case !}e \text{ of} \\ (r) &\Rightarrow \\ &\quad e := \text{dead} \\ &\quad \text{return } deleteUse(r) \\ (\text{letfun } f(x) \text{ be } m \text{ in } n) &\Rightarrow \\ &\quad e, f, x := \text{dead} \\ &\quad \text{return } cleanExp(m) \cup cleanExp(n) \\ (\text{app}(a, b)) &\Rightarrow \\ &\quad e := \text{dead} \\ &\quad \text{return } cleanExp(a) \cup cleanExp(b) \\ \dots & \end{aligned}$$

*Remark* Marking nodes as dead ensures that unnecessary work is not done on dead redexes. A crucial difference between the imperative algorithms and the improved quadratic one is that reduction in the former immediately detects new redexes, whereas the improved quadratic algorithm only detects new (non-local) redexes on a subsequent traversal.

$deleteUse(r)$  removes  $r$  and returns a set of 0 or 1 occurrence redexes:

```

deleteUse(r) =
  if r is already dead then return  $\emptyset$ 
  let  $s = nextOcc(r)$ 
   $uses(s) := uses(s) - \{r\}$ 
  return  $\{s\}$ 

nextOcc(r) =
  let  $x = def(r)$ 
  if r is non-recursive then return  $s \in (non-rec-uses(x) \cup \{x\}) - \{r\}$ 
  else if r is recursive then return  $s \in (rec-uses(x) \cup \{x\}) - \{r\}$ 

```

$reduceProjections(e, x, a_1, a_2, xs)$  reduces projections indexed by  $xs$ .  $e$  is an expression node of the form  $let\ x\ be\ val(pair(a_1, a_2))$  in  $m$ , and  $xs$  is a subset of the uses of  $x$ .

```

reduceProjections(e, x, a_1, a_2, xs) =
  let  $redexes := \emptyset$ 
  for each  $s \in xs$  do
    let  $focus = parent(parent(s))$ 
    case ! $focus$  of
      ( $proj_i(s)$ )  $\Rightarrow$ 
        splice a copy of  $a_i$  in place of  $focus$ 
         $redexes := redexes \cup \{parent(focus)\}$ 
      (-)  $\Rightarrow$  skip
  return  $redexes$ 

```

All the projections in which a member of  $xs$  participates are reduced, and a set of expression redexes is constructed. Each projection can trigger the creation of a new  $T.\beta_a$ -redex. For instance, consider:

$$\begin{aligned}
 & \text{let } x \text{ be } val(pair(a, b)) \text{ in let } y \text{ be } val(proj_1(x)) \text{ in } m \\
 \longrightarrow_{x.\beta} & \text{let } x \text{ be } val(pair(a, b)) \text{ in let } y \text{ be } val(a) \text{ in } m \\
 \longrightarrow_{T.\beta_a} & \text{let } x \text{ be } val(pair(a, b)) \text{ in } m[y := a]
 \end{aligned}$$

$reduceCases(e, x, i, a, xs)$  reduces case-splits indexed by  $xs$ .  $e$  is an expression node of the form  $let\ x\ be\ val(inj_i(a))$  in  $m$ , and  $xs$  is a subset of the uses of  $x$ .

```

reduceCases(e, x, i, a, xs) =
  let  $occs := \emptyset$ 
  let  $redexes := \emptyset$ 

```

```

for each  $s \in xs$  do
  let  $focus = parent(parent(s))$ 
  case ! $focus$  of
  (case  $s$  of  $(x_1)n_1 ; (x_2)n_2$ )  $\Rightarrow$ 
     $occs := occs \cup cleanExp(n_{3-i})$ 
     $deleteUse(s)$ 
    splice  $n_i$  in place of  $focus$ 
    let  $(occs', redexes') = substAtom(x_i, a)$ 
     $occs := occs \cup occs'$ 
     $redexes := redexes \cup redexes' \cup \{parent(focus)\}$ 
     $x_1, x_2 := dead$ 
  (-)  $\Rightarrow$  skip
return  $(occs, redexes)$ 

```

The structure of *reduceCases* is similar to that of *reduceProjections*. However, it is slightly more complex because a single  $+. \beta$ -reduction inlines multiple atoms, splices one branch of a **case** and discards the other. Discarding the branch which is not taken gives a set of occurrence redexes as well as the expression redexes.

## 5.2 Substitution

$substAtom(x, a)$  substitutes the atom  $a$  for all the uses of the defining occurrence  $x$ . It returns a pair of a set of occurrence redexes and a set of expression redexes.

```

 $substAtom(x, a) =$  case ! $a$  of
  ( $r$ )  $\Rightarrow$   $substUse(x, r)$ 
  (-)  $\Rightarrow$ 
    for each  $r \in uses(x)$  do
      splice a copy of  $a$  in place of  $r$ 
       $x := dead$ 
    return  $(\emptyset, \emptyset)$ 

```

This is straightforward for non-variable atoms, as it cannot generate new redexes. In contrast, substituting a variable can trigger  $\times. \beta$ - and  $+. \beta$ -reductions.

$substUse(x, r)$  substitutes  $r$  for all the uses of the defining occurrence  $x$ .

```

 $substUse(x, r) =$ 
  let  $xs = uses(x)$ 
  if  $r \in rec-uses(r)$  then
     $rec-uses(r) := rec-uses(r) \cup xs$ 
  else if  $r \in non-rec-uses(r)$ 
     $non-rec-uses(r) := non-rec-uses(r) \cup xs$ 
   $x := dead$ 
  let  $e = parent(def(r))$ 
  case ! $e$  of
  (let  $y$  be val( $\lceil pair(a_1, a_2) \rceil$ ) in  $m$ )  $\Rightarrow$ 
    for each  $s \in xs$  do  $def(s) := def(r)$ 

```

```

    let redexes = reduceProjections(e, y, a1, a2, xs)
    return (∅, redexes)
  (let y be val([inji(ai)]) in m) ⇒
    for each s ∈ xs do def(s) := def(r)
    let (occs, redexes) = reduceCases(e, y, i, ai, xs)
    return (occs, redexes)
  (-) ⇒ return (∅, ∅)

```

Substitution is implemented by merging two sets together. Concretely, this amounts to the constant-time operation of inserting one doubly-linked circular list inside another. In addition, if  $x$  is bound to a pair, then projections are reduced, or if  $x$  is bound to an injection, then case-splits are reduced.

## 6 Analysis

There are two obvious operations mapping terms from the functional to the imperative representations, which we call *mutify* and *demutify*, respectively. We have a semi-formal argument for the following:

**Proposition 3.** *Let  $e$  be a term and  $e' = (\text{demutify} \circ \text{contract} \circ \text{mutify})(e)$ . Then  $e'$  is a normal form for  $e$ .*

The argument uses the invariants of Sect. 4, plus the invariant that the children of the current node are in normal form. When new redexes are created, this invariant is modified such that subterms may contain redexes, but only those stored in appropriate expression redex sets or occurrence redex sets. It is reasonably straightforward to verify that the operations which update the graph structure do in fact correspond to MIL reductions. When *contract* terminates, all the redex sets are empty and the term is in normal form.

### 6.1 Complexity without Commuting Conversions

Although our approach of performing CCs on the way down the tree works well in practice, the worst case time complexity is still quadratic in the size of the term. We define a version of our algorithm *contract*<sub>β</sub> which does not perform commuting conversions. This is obtained simply by removing the call to *reduceCCs* from *contract*, and the test for commuting conversions from *reduce*.

**Proposition 4.** *contract*<sub>β</sub>( $e$ ) is linear in the size of  $e$ .

The argument is very similar to that of Appel and Jim [1] for their imperative algorithm. Essentially most operations take constant time and shrink the size of the term. The only exception is substitution. In the case where a non-variable is substituted for a variable  $x$ , the operation is linear in the number of uses of  $x$ . But it is only possible to substitute a non-variable for a variable once, therefore the total time spent substituting atoms is linear. In the case where a variable  $y$  is substituted for a variable  $x$ , the operation is constant, providing  $y$  is not

bound to a pair or an injection. If  $y$  is bound to a pair or an injection, then the operation is linear in the number of uses of  $x$ . Again, once bound to a pair or an injection, a variable cannot be rebound, so the time remains linear.

Crucially, this argument relies on the fact that back pointers from uses back to defining occurrences are only maintained for pairs and injections. In our SML.NET implementation we found that maintaining back pointers from *all* uses back to defining occurrences does not incur any significant cost in practice. Even when bootstrapping the compiler ( $\sim 80,000$  lines of code) there was no discernible difference in compile time. Maintaining back pointers also allows us to perform various other rewrites including  $\eta$ -reductions. In the presence of all back pointers, optimising the union operation to always add the smaller list to the larger one guarantees  $O(n \log n)$  behaviour. Using an efficient *union-find* algorithm would restore essentially linear complexity.

## 6.2 Complexity with Commuting Conversions

Naively reducing commuting conversions can give quadratic behaviour. For instance, consider the following (innermost first) reductions:

$$\begin{aligned} & \text{let } x_k \text{ be (let } x_{k-1} \text{ be } \dots \text{ let } x_1 \text{ be } m_1 \text{ in } m_2 \text{ in } \dots m_k) \text{ in } n \\ & \longrightarrow^* (S(k-1) \text{ } T.CC\text{-reductions}) \\ & \text{let } x_k \text{ be (let } x_1 \text{ be } m_1 \text{ in } \dots \text{ let } x_{k-1} \text{ be } m_{k-1} \text{ in } m_k) \text{ in } n \\ & \longrightarrow^* (k-1 \text{ } T.CC\text{-reductions}) \\ & \text{let } x_1 \text{ be } m_1 \text{ in } \dots \text{ let } x_k \text{ be } m_k \text{ in } n \end{aligned}$$

The total number of reductions is given by the recurrence:  $S(1) = 0, S(k) = S(k-1) + k - 1$ . This has solution  $S(k) = k(k-1)/2$ . Assuming each of the  $m_i$ s and  $n$  have constant size, then  $k$  is linear in the size of the term. Hence the number of reductions is quadratic in the size of the term. If the *contract* function directly performed these reductions, then it would also be quadratic.

Another problem is that  $+CC$ -reductions can introduce ‘useless functions’:

$$\begin{aligned} & \text{let } z \text{ be (let } y \text{ be (case } a \text{ of } (x_1)n_1 ; (x_2)n_2) \text{ in } m) \text{ in } p \\ & \longrightarrow^* \text{letfun } f(y) \text{ be } m \\ & \quad \text{in let } z \text{ be case } a \text{ of } (x_1)\text{let } y_1 \text{ be } n_1 \text{ in app}(f, y_1) \\ & \quad \quad \quad ; (x_2)\text{let } y_2 \text{ be } n_2 \text{ in app}(f, y_2) \\ & \quad \text{in } p \\ & \longrightarrow^* \text{letfun } f(y) \text{ be } m \\ & \quad \text{in letfun } g(z) \text{ be } p \\ & \quad \quad \text{in case } a \text{ of } (x_1)\text{let } y_1 \text{ be } n_1 \text{ in let } z_1 \text{ be app}(f, y_1) \text{ in app}(g, z_1) \\ & \quad \quad \quad ; (x_2)\text{let } y_2 \text{ be } n_2 \text{ in let } z_2 \text{ be app}(f, y_2) \text{ in app}(g, z_2) \end{aligned}$$

The function  $g$  is useless in the sense that it is always applied to the result of applying  $f$  to an argument. One might hope that  $g$  be composed with  $f$ . If we

change the reduction order, such that the commuting conversions are performed outermost first, then it is:

$$\begin{aligned}
& \text{let } z \text{ be (let } y \text{ be (case } a \text{ of } (x_1)n_1 ; (x_2)n_2 \text{) in } m \text{) in } p \\
\longrightarrow^* & \text{let } y \text{ be (case } a \text{ of } (x_1)n_1 ; (x_2)n_2 \text{) in let } z \text{ be } m \text{ in } n \\
\longrightarrow^* & \text{letfun } f(y) \text{ be let } z \text{ be } m \text{ in } p \\
& \text{in case } a \text{ of } (x_1)\text{let } y_1 \text{ be } n_1 \text{ in app}(f, y_1) \\
& \quad ; (x_2)\text{let } y_2 \text{ be } n_2 \text{ in app}(f, y_2)
\end{aligned}$$

Fortunately, given the limited ways in which commuting conversions can trigger other reductions, the full imperative algorithm can get away with performing commuting conversions outermost first, with an initial call to *reduceCCse* before recursively contracting *e*'s children. The operation *reduceCCs(e)* repeatedly checks *e* to see if it is a CC-redex. If it is, then it performs the commuting conversion, and iterates. If not, then it returns.

The previous example of quadratic behaviour due to commuting conversions becomes linear with this reduction strategy. However, quadratic behaviour can still arise through inlining functions that trigger further commuting conversions:

$$\begin{aligned}
\text{letfun} & \quad f_k(x_k) \text{ be let } y_k \text{ be app}(g, x_k) \text{ in app}(g, y_k) \\
& \quad f_{k-1}(x_{k-1}) \text{ be let } y_{k-1} \text{ be app}(f_k, x_{k-1}) \text{ in app}(g, y_{k-1}) \\
& \quad \vdots \\
& \quad f_1(x_1) \text{ be let } y_1 \text{ be app}(f_2, x_1) \text{ in app}(g, y_1) \\
\text{in} & \quad \text{app}(f_1, a)
\end{aligned}$$

*contract* takes quadratic time to reduce this term. In order to get a linear number of reductions, one would have to inline all the functions first, before performing any commuting conversions.

## 7 Performance

We have extended our one-pass imperative algorithm *contract* to the whole of MIL and compared its performance with the current implementation of *simplify*. Replacing *simplify* with *contract* is not entirely straightforward, as all the other phases in the pipeline are written to work on a straightforward immutable tree datatype for terms, which is incompatible with the representation used in *contract*. We therefore make use of *mutify* and *demutify* to change representation before and after *contract*. Since both *mutify* and *demutify* completely rebuild the term, they are very expensive – calling *mutify* and *demutify* generally takes longer than *contract* itself. Ideally, of course, all the phases would use the same representation. However, using two representations allowed us to compare the running times of *simplify* and *contract* on real programs.

Table 1 compares the total compile times (*tsimplify* vs *tcontract*) of several benchmark programs for the existing compiler, using *simplify*, and for the modified one, using *demutify*  $\circ$  *contract*  $\circ$  *mutify*. Table 2 compares the time *simp*

**Table 1.** Total compile time (in seconds)

Benchmark	Lines of code	SML/NJ		MLton	
		<i>tsimplify</i>	<i>tcontract</i>	<i>tsimplify</i>	<i>tcontract</i>
<b>sort</b>	70	2.11	3.47	0.46	0.52
<b>xq</b>	1,300	13.1	14.4	2.46	1.76
<b>mllex</b>	1,400	11.6	16.0	2.39	2.03
<b>raytrace</b>	2,500	18.1	24.0	4.30	3.03
<b>mlyacc</b>	6,200	57.3	43.8	10.0	6.04
<b>hamlet</b>	20,000	219	156	43.7	26.2
<b>bootstrap</b>	80,000	1310	1190	289	221

**Table 2.** Shrinking reduction time (in seconds) under SML/NJ and MLton

Benchmark	Under SML/NJ					Under MLton				
	Total		Breakdown			Total		Breakdown		
	<i>simp</i>	<i>mcd</i>	<i>m</i>	<i>c</i>	<i>d</i>	<i>simp</i>	<i>mcd</i>	<i>m</i>	<i>c</i>	<i>d</i>
<b>sort</b>	1.00	2.00	0.87	0.70	0.43	0.22	0.11	0.02	0.07	0.02
<b>xq</b>	5.86	5.98	1.90	3.61	0.47	1.46	0.54	0.35	0.15	0.06
<b>mllex</b>	6.09	7.49	3.31	3.16	1.02	1.21	0.57	0.27	0.23	0.07
<b>raytrace</b>	9.32	11.8	5.16	5.44	1.17	2.13	0.65	0.37	0.19	0.09
<b>mlyacc</b>	33.2	20.0	9.42	8.60	1.94	5.63	1.26	0.68	0.37	0.21
<b>hamlet</b>	84.5	56.4	26.2	21.5	8.59	23.3	5.54	1.85	2.77	0.92
<b>bootstrap</b>	439	282	130	100	53.0	107	36.6	11.8	18.4	6.38

spent in *simplify* with the times *m*, *c*, *d* spent in each of *mutify*, *contract* and *demutify* and their sum *mcd*. Each benchmark was run under two different versions of SML.NET. One was compiled under SML/NJ [12] and the other under MLton [13]. Benchmarks were run on a 1.4Ghz AMD Athlon PC equipped with 512MB of RAM and Microsoft Windows XP SP1.

The first five benchmarks are demos distributed with SML.NET. The **sort** benchmark applies quicksort to a list of integers; **xq** is an interpreter for an XQuery-like language for querying XML documents; **mllex** and **mlyacc** are ports of SML/NJ's ml-lex and ml-yacc utilities; **raytrace** is a port to SML of the winning entry from the Third Annual ICFP Programming Contest. The remaining benchmarks are much larger: **hamlet** is Andreas Rossberg's SML interpreter, whilst **bootstrap** is SML.NET compiling itself.

Figure 4 gives a graphical comparison of both tables. On small benchmarks, the current compiler is faster (*tcontract/tsimplify*). But for medium and large benchmarks, we were surprised to discover that *contract* is faster than *simplify*, even though much of the time is spent in useless representation changes. Under SML/NJ, *tcontract/tsimplify* shows a decrease of nearly 30% in the total compile time in some cases; under MLton, there is a decrease of up to 40%. This is a significant improvement, given that in the existing compiler only around 50% of compile time is spent performing shrinking reductions. Comparing the actual shrinking reduction times *c* and *simp*, *contract* is up to four times faster than



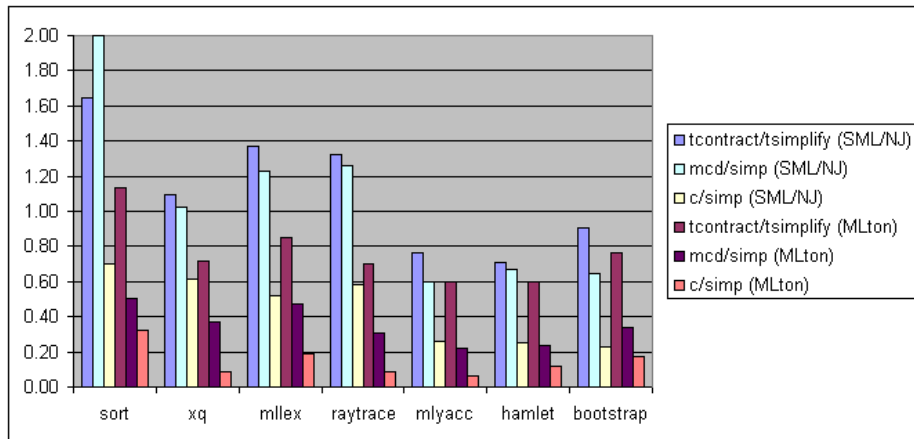


Fig. 4. Comparing *contract* with *simplify*

*simplify* under SML/NJ, and up to 15 times faster under MLton (on *mlyacc*). The level of improvement under MLton is striking. Our results suggest that MLton is considerably better than SML/NJ at compiling ML code which makes heavy use of references.

As an exercise, one of the other transformations *deunit*, which removes redundant unit values and types was translated to use the new representation. The *contract* function is called before and after *deunit*, so this enabled us to eliminate one call to *demutify* and one call to *mutify*. This translation was easy to do and did not change the performance of *deunit*. We believe that it should be reasonably straightforward, if somewhat tedious, to translate the rest of the transformations to work directly with the mutable representation.

## 8 Conclusions and Further Work

We have implemented and extended Appel and Jim’s imperative algorithm for shrinking reductions and shown that it can yield significant reductions in compile times relative to the algorithm currently used in SML/NJ and SML.NET. The improvements are such that, for large programs, it is even worth completely changing representations before and after *contract*, but this is clearly suboptimal. The results of this experiment indicate that it would be worth the effort of rewriting other phases of the compiler to use the graph-based representation.

Making more extensive use of the pointer-based representation would allow many transformations to be written in a different style, for example replacing explicit environments with extra information on binding nodes, though this does not interact well with the hash-consing currently used for types. We also believe that ‘code motion’ transformations can be more easily and efficiently expressed.

It is unfortunate that CCs and inlining conspire to produce quadratic complexity. Sabry and Wadler's study of CPS translations offers an interesting insight [11]. In their variant of Moggi's computational lambda calculus  $\lambda_{c^{**}}$ , terms are in CC-normal form by definition, and  $\beta$ -reduction of an application is combined with CC-normalisation of its enclosing let-expression: adopting this more refined notion of redex may allow us to achieve linear complexity.

More speculatively, we would like to investigate more principled mutable graph-based intermediate representations. There has been much theoretical work on graph-based representations of proofs and programs, yet these do not seem to have been exploited in compilers for higher-order languages (though of course, compilers for imperative languages have used a mutable flow-graph representations for decades). With a careful choice of representation, some of our transformations (such as *T.CC*) could simply be isomorphisms and we believe that a better treatment of shared continuations in the other commuting conversions would also be possible.

## References

1. Andrew W. Appel and Trevor Jim. Shrinking lambda expressions in linear time. *Journal of Functional Programming*, 7(5):515–540, 1997.
2. N. Benton and A. Kennedy. Monads, effects and transformations. In *3rd International Workshop on Higher Order Operational Techniques in Semantics (HOOTS), Paris*, volume 26 of *ENTCS*. Elsevier, September 1999.
3. N. Benton, A. Kennedy, and C. Russo. SML.NET. <http://www.cl.cam.ac.uk/Research/TSG/SMLNET/>, June 2002.
4. N. Benton, A. Kennedy, and C. Russo. Adventures in interoperability: The SML.NET experience. In *Proc. 6th ACM-SIGPLAN International Conference on Principles and Practice of Declarative Programming (PPDP)*, August 2004.
5. Nick Benton, Andrew Kennedy, and George Russell. Compiling Standard ML to Java bytecodes. In *Proc. ACM SIGPLAN International Conference on Functional Programming (ICFP '98)*, volume 34(1), pages 129–140, 1999.
6. O. Danvy. A new one-pass transformation into monadic normal form. In *Proc. 12th International Conference on Compiler Construction*, number 2622 in Lecture Notes in Computer Science, pages 77–89. Springer, 2003.
7. Ecma International. ECMA Common Language Infrastructure standard, December 2002. <http://www.ecma-international.org/publications/standards/Ecma-335.htm>.
8. J. Hatcliff and O. Danvy. A generic account of continuation-passing styles. In *Proc. 21st Annual Symposium on Principles of Programming Languages*. ACM, 1994.
9. Sam Lindley. *Normalisation by evaluation in the compilation of typed functional programming languages*. PhD thesis, The University of Edinburgh, 2005.
10. S. L. Peyton Jones and A. L.M. Santos. A transformation-based optimiser for Haskell. *Science of Computer Programming*, 1998.
11. Amr Sabry and Philip Wadler. A reflection on call-by-value. *ACM Transactions on Programming Languages and Systems*, 19(6):916–941, 1997.
12. Standard ML of New Jersey (SML/NJ) compiler: <http://smlnj.org/>.
13. Stephen Weeks, Matthew Fluet, Henry Cejtin, and Suresh Jagannathan. MLton whole-program optimizing compiler: <http://mlton.org/>.