

Selective Memoization*

Umut A. Acar

Guy E. Blelloch

Robert Harper

School of Computer Science
Carnegie Mellon University
Pittsburgh, PA 15213

{umut,blelloch,rwh}@cs.cmu.edu

Abstract

We present a framework for applying memoization selectively. The framework provides programmer control over equality, space usage, and identification of precise dependences so that memoization can be applied according to the needs of an application. Two key properties of the framework are that it is efficient and yields programs whose performance can be analyzed using standard techniques.

We describe the framework in the context of a functional language and an implementation as an SML library. The language is based on a modal type system and allows the programmer to express programs that reveal their true data dependences when executed. The SML implementation cannot support this modal type system statically, but instead employs run-time checks to ensure correct usage of primitives.

Categories and Subject Descriptors

D.3.0 [Programming Languages]: General; F.2.0 [Analysis of Algorithms and Problem Complexity]: [General]; D.3.1 [Programming Languages]: Formal Definitions and Theory; D.3.3 [Programming Languages]: Language Constructs and Features—Control Structures

General Terms

Languages, Performance, Algorithms

Keywords

Memoization, selective, programmer controlled, performance

*This research was supported in part by NSF grants CCR-9706572, CCR-0085982, and CCR-0122581.

1 Introduction

Memoization is a fundamental and powerful technique for result re-use. It dates back a half century [7, 21, 22] and has been used extensively in many areas such as dynamic programming [4, 9, 10, 19], incremental computation [11, 34, 12, 36, 16, 1, 37, 20, 14, 2], and many others [8, 23, 17, 25, 26, 20]. In fact, lazy evaluation provides a limited form of memoization [18].

Although memoization can dramatically improve performance and can require only small changes to the code, no language or library support for memoization has gained broad acceptance. Instead, many successful uses of memoization rely on application-specific support code. The underlying reason for this is one of control: since memoization is all about performance, the user must be able to control the performance of memoization. Many subtleties of memoization, including the cost of equality checking and the cache replacement policy for memo tables, can make the difference between exponential and linear running time.

To be general and widely applicable a memoization framework must provide control over these three areas: (1) the kind and cost of equality tests; (2) the identification of precise dependences between the input and the output of memoized code; and (3) space management. Control over equality tests is critical, because this is how re-usable results are identified. Control over identification of precise dependences is important to maximize result reuse. Being able to control when memo tables or individual their entries are purged is critical, because otherwise the user will not know whether or when results are re-used.

In this paper, we propose a framework for memoization that provides control over equality and identification of dependences, and some control over space management. We study the framework in the context of a small language called MFL and provide an implementation for it in the Standard ML language. We also prove the type safety and correctness of MFL—*i.e.*, that the semantics are preserved with respect to a non-memoized version. As an example, we show how to analyze the performance of a memoized version of Quicksort within our framework.

In the next section we describe background and related work. In Section 3 we introduce our framework via some examples. In Section 4 we formalize the MFL language and discuss its safety, correctness, and performance properties. In Section 5 we present a simple implementation of the framework as a Standard ML library. In Section 6 we discuss how the framework might be extended to allow for better control of space usage, and discuss the relationship of this work to our previous work on adaptive computation [2].

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.
POPL'03, January 15–17, 2003, New Orleans, Louisiana, USA.
Copyright 2003 ACM 1-58113-628-5/03/0001 ...\$5.00

2 Background and Related Work

A typical memoization scheme maintains a memo table mapping argument values to previously computed results. This table is consulted before each function call to determine if the particular argument is in the table. If so, the call is skipped and the result is returned; otherwise the call is performed and its result is added to the table. The semantics and implementation of the memo lookup are critical to performance. Here we review some key issues in implementing memoization efficiently.

Equality. Any memoization scheme needs to search a memo table for a match to the current arguments. Such a search will, at minimum, require a test for equality. Typically it will also require some form of hashing. In standard language implementations testing for equality on structures, for example, can require traversing the whole structure. The cost of such an equality test can negate the advantage of memoizing and may even change the asymptotic behavior of the function. A few approaches have been proposed to alleviate this problem. The first is based on the fact that for memoization equality need not be exact—it can return unequal when two arguments are actually equal. The implementation could therefore decide to skip the test if the equality is too expensive, or could use a conservative equality test, such as “location” equality. The problem with such approaches is that whether a match is found could depend on particulars of the implementation and will surely not be evident to the programmer.

Another approach for reducing the cost of equality tests is to ensure that there is only one copy of every value, via a technique known as “hash consing” [13, 5, 35]. If there is only one copy, then equality can be implemented by comparing locations. In fact, the location can also be used as a key to a hash table. In theory, the overhead of hash-consing is constant in the expected case (expectation is over internal randomization of hash functions). The reality, however, is rather different because of large memory demands of hash-consing and its interaction with garbage collection. In fact, several researchers have argued that hash-consing is too expensive for practical purposes [32, 33, 6, 24]. As an alternative to hash consing, Pugh proposed lazy structure sharing [32]. In lazy structure sharing whenever two equal values are compared, they are made to point to the same copy to speed up subsequent comparisons. As Pugh points out, the disadvantage of this approach is that the performance depends on the order comparisons and thus it is difficult to analyze.

We note that even with hash-consing, or any other method, it remains critical to define equality on all types including reals and functions. Claiming that functions are never equivalent, for example, is not satisfactory because the result of a call involving some function as a parameter will never be re-used.

Precise Dependences. To maximize result re-use, the result of a function call must be stored with respect to its true dependences. This issue arises when the function examines only parts or an approximation of its parameter. To enable “partial” equality checks, the unexamined parts of the parameter should be disregarded. To increase the likelihood of result re-use, one should be able to match on the approximation, rather than the parameter itself. As an example, consider the code

```
fun f(x,y,z) = if (x > 0) then fy(y) else fz(z)
```

The result of `f` depends on either (x,y) or (x,z) . Also, it depends on an approximation of x —whether or not it is positive—rather than

its exact value. Thus, the memo entry $(7,11,20)$ should match $(7,11,30)$ or $(4,11,50)$ since, when x is positive, the result depends only on y .

Several researchers have remarked that partial matching can be very important in some applications [28, 27, 1, 14]. Abadi, Lampson, Lévy [1], and Heydon, Levin, Yu [14] have suggested program analysis methods for tracking dependences for this purpose. Although their technique is likely effective in catching potential matches, it does not provide a programmer controlled mechanism for specifying what dependences should be tracked. Also, their program analysis technique can change the asymptotic performance of a program, making it difficult to assess the effects of memoization.

Space management. Another problem with memoization is its space requirement. As a program executes, its memo tables can become large limiting the utility of memoization. To alleviate this problem, memo tables or individual entries should be disposed of under programmer control.

In some application, such as in dynamic programming, most result re-use occurs among the recursive calls of some function. Thus, the memo table of such a function can be disposed of whenever it terminates. In other applications, where result re-use is less structured, individual memo table entries should be purged according to a replacement policy [15, 33]. The problem is to determine what exact replacement policy should be used and to analyze the performance effects of the chosen policy. One widely used approach is to replace the least recently used entry. Other, more sophisticated, policies have also been suggested [33]. In general the replacement policy must be application-specific, because, for any fixed policy, there are programs whose performance is made worse by that choice [33].

3 A Framework for Selective Memoization

We present an overview of our framework via some examples. The framework extends a purely functional language with several constructs to support selective memoization. In this section, we use an extension to an ML-like language for the discussion. We formalize the core of this language and study its safety, soundness, and performance properties in Section 4.

The framework enables the programmer to determine precisely the dependences between the input and the result of a function. The main idea is to deem the parameters of a function as *resources* and provide primitives to explore incrementally any value, including the underlying value of a resource. This incremental exploration process reveals the dependences between the parameter of the function and its result.

The incremental exploration process is guided by types. If a value has the modal type $! \tau$, then the underlying value of type τ can be bound to an ordinary, unrestricted variable by the `let!` construct; this will create a dependence between the underlying value and the result. If a value has a product type, then its two parts can be bound to two resources using the `let*` construct; this creates no dependences. If the value is a sum type, then it can be case analyzed using the `mcase` construct, which branches according to the outermost form of the value and assigns the inner value to a resource; `mcase` creates a dependence on the outer form of the value of the resource. The key aspect of the `let*` and `mcase` is that they bind resources rather than ordinary variables.

| Non-memoized | Memoized |
|---|--|
| <pre> fun fib (n:int)= if (n < 2) then n else fib(n-1) + fib(n-2) </pre> | <pre> mfun mfib (n':!int)= let !n = n' in return (if (n < 2) then n else mfib(!n-1) + mfib(!n-2)) end </pre> |
| <pre> fun f (x:int, y:int, z:int)= if (x > 0) then fy y else fz z </pre> | <pre> mfun mf (x':int, y':!int, z':!int)= mif (x' > 0) then let !y = y' in return (fy y) end else let !z = z' in return (fz z) end </pre> |

Figure 1. Fibonacci and expressing partial dependences.

Exploring the input to a function via `let!`, `mcas`, and `let*` builds a *branch* recording the dependences between the input and the result of the function. The `let!` adds to the branch the full value, the `mcas` adds the kind of the sum, and `let*` adds nothing. Consequently, a branch contains both data dependences (from `let!`'s) and control dependences (from `mcas`'s). When a `return` is encountered, the branch recording the revealed dependences is used to key the memo table. If the result is found in the memo table, then the stored value is returned, otherwise the body of the `return` is evaluated and the memo table is updated to map the branch to the result. The type system ensures that all dependences are made explicit by precluding the use of resources within `return`'s body.

As an example consider the Fibonacci function `fib` and its memoized counterpart `mfib` shown in Figure 1. The memoized version, `mfib`, exposes the underlying value of its parameter, a resource, before performing the two recursive calls as usual. Since the result depends on the full value of the parameter, it has a bang type. The memoized Fibonacci function runs in linear time as opposed to exponential time when not memoized.

Partial dependences between the input and the result of a function can be captured by using the incremental exploration technique. As an example consider the function `f` shown in Figure 1. The function checks whether `x` is positive or not and returns `fy(y)` or `fz(z)`. Thus the result of the function depends on an approximation of `x` (its sign) and on either `y` or `z`. The memoized version `mf` captures this by first checking if `x'` is positive or not and then exposing the underlying value of `y'` or `z'` accordingly. Consequently, the result will depend on the sign of `x'` and on either `y'` or `z'`. Thus if `mf` is called with parameters (1, 5, 7) first and then (2, 5, 3), the result will be found in the memo the second time, because when `x'` is positive the result depends only on `y'`. Note that `mif` construct used in this example is just a special case of the more general `mcas` construct.

A critical issue for efficient memoization is the implementation of memo tables along with lookup and update operations on them. In our framework we support expected constant time memo table lookup and update operations by representing memo tables using hashing. To do this, we require that the underlying type τ of a modal type $!\tau$ be an *indexable type*. An indexable type is associated with an injective function, called an *index function*, that maps each value of that type to a unique integer; this integer is called the *index* of the value. The uniqueness property of the indices for a given type ensures that two values are equal if and only if their indices are equal. In our framework, equality is only defined for

| Non-memoized | Memoized |
|---|---|
| <pre> type irl=(int*real) list fun ks (c:int,l:irl) = case l of nil => 0 (w,v)::t => if (c < w) then ks (c,t) else let v1 = ks(c,t) v2 = v+ks(c-w,t) in if (v1>v2) then v1 else v2 end end </pre> | <pre> type irl=(int*real) blist mfun mks (c':!int,l':!irl) let !c = c' in let !l = l' in return (case (unbox l) of NIL => 0 CONS((w,v),t) => if (c < w) then mks(!c,!t) else let v1 = mks(!c,!t) v2 = v+mks(!c-w,!t) in if (v1>v2) then v1 else v2 end end) end </pre> |

Figure 2. Memo tables for memoized Knapsack can be discarded at completion.

indexable types. This enables us to implement memo tables as hash tables keyed by branches consisting of indices.

We assume that each primitive type comes with an index function. For examples, for integers, the identity function can be chosen as the index function. Composite types such as lists or functions must be *boxed* to obtain an indexable type. A boxed value of type τ has type τbox . When a box is created, it is assigned a unique locations (or tag), and this location is used as the unique index of that boxed value. For example, we can define boxed lists as follows.

```

datatype  $\alpha$  blist' = NIL
                | CONS of  $\alpha * ((\alpha \text{ blist}') \text{ box})$ 

type  $\alpha$  blist = ( $\alpha$  blist') box

```

Based on boxes we implement hash-consing as a form of memoization. For example, hash-consing for boxed lists can be implemented as follows.

```

mfun hCons (h':! $\alpha$ , t':!( $\alpha$  blist')) =
  let !h = h' in let !t = t' in
    return (box (CONS(h,t)))
  end

```

The function takes an item and a boxed list and returns the boxed list formed by consing them. Since the function is memoized, if it is ever called with two values that are already hash-consed, then the same result will be returned. The advantage of being able to define hash-consing as a memoized function is that it can be applied selectively.

To control space usage of memo tables, our framework gives the programmer a way to dispose of memo tables by conventional scoping. In our framework, each memoized function is allocated its own memo table. Thus, when the function goes out of scope, its memo table can be garbage collected. For example, in many dynamic-programming algorithms result re-use occurs between recursive calls of the same function. In this case, the programmer can scope the memoized function inside an auxiliary function so that its memo table is discarded as soon as the auxiliary function returns. As an example, consider the standard algorithm for the Knapsack Problem `ks` and its memoized version `mks` Figure 2. Since result sharing mostly occurs among the recursive calls of `mks`, it can be scoped in some other function that calls `mks`; once `mks` returns its memo table will go out of scope and can be discarded.

We note that this technique gives only partial control over space usage. In particular it does not give control over when individual

| Non-memoized | Memoized |
|--|--|
| <pre> fun fil (g:int->bool, l:int list) = case l of nil => nil h::t => let tt = fil(g,t) in case (g h) of true => h::tt false => tt end end fun qs (l:int list) = case l of nil => nil cons(h,t) => let s = fil(fn x=>x<h,t) g = fil(fn x=>x>h,t) in (qs s)@(h::(qs g)) end end </pre> | <pre> empty = box NIL fun mfil (g:int->bool, l:int blist)= case (unbox l) of NIL => empty CONS(h,t) => let tt = mfil(g,t) in case (g h) of true => hCons(h,tt) false => tt end end mfun mqs (l':int blist) = let !l = l' in return (case (unbox l) of NIL => nil CONS(h,t) => let s = mfil(fn x=>x<h,t) g = mfil(fn x=>x>h,t) in (mqs s)@(h::(mqs g)) end end) end </pre> |

Figure 3. The Quicksort algorithm.

memo table entries are purged. In Section 6, we discuss how the framework might be extended so that each memo table is managed according to a programmer specified caching scheme. The basic idea is to require the programmer to supply a caching scheme as a parameter to the `mfun` and maintain the memo table according to the chosen caching scheme.

Memoized Quicksort. As a more sophisticated example, we consider Quicksort. Figure 3 show an implementation of the Quicksort algorithm and its memoized counterpart. The algorithm first divides its input into two lists containing the keys less than the pivot, and greater than the pivot by using the filter function `fil`. It then sorts the two sublists, and returns the concatenation of the results. The memoized filter function `mfil` uses hash-consing to ensure that there is only one copy of each result list. The memoized Quicksort algorithm `mqs` exposes the underlying value of its parameter and is otherwise similar to `qs`. Note that `mqs` does not build its result via hash-consing—it can output two copies of the same result. Since in this example the output of `mqs` is not consumed by any other function, there is no need to do so. Even if the result were consumed by some other function, one can choose not to use hash-consing because operations such as insertions to and deletions from the input list will surely change the result of Quicksort.

When the memoized Quicksort algorithm is called on “similar” inputs, one would expect that some of the results would be re-used. Indeed, we show that the memoized Quicksort algorithm computes its result in expected linear time when its input is obtained from a previous input by inserting a new key at the beginning. Here the expectation is over all permutations of the input list and also the internal randomization of the hash functions used to implement the memo tables. For the analysis, we assume, without loss of generality, that all keys in the list are unique.

Theorem 1

Let L be a list and let $L' = [a, L]$. Consider running memoized Quicksort on L and then on L' . The running time of Quicksort on the modified list L' is expected $O(n)$ where n is the length of L' .

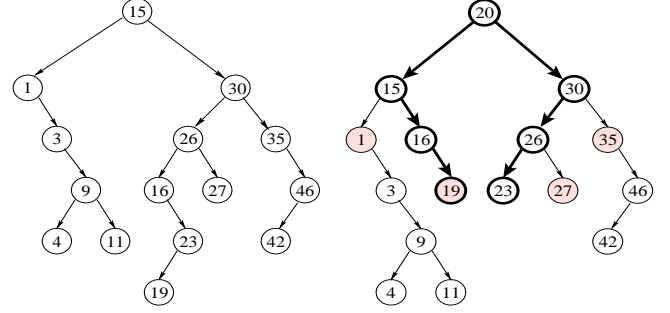


Figure 4. The recursion tree for Quicksort with inputs $L = [15, 30, 26, 1, 3, 16, 27, 9, 35, 4, 46, 23, 11, 42, 19]$ (left) and $L' = [20, L]$ (right).

Proof: Consider the recursion tree of Quicksort with input L , denoted $Q(L)$, and label each node with the pivot of the corresponding recursive call (see Figure 4 for an example). Consider any pivot (key) p from L and let L_p denote the keys that precede p in L . It is easy to see that a key k is in the subtree rooted at p if and only if the following two properties are satisfied for any key $k' \in L_p$.

1. If $k' < p$ then $k > k'$, and
2. if $k' > p$ then $k < k'$.

Of the keys that are in the subtree of p , those that are less than p are in its left subtree and those greater than p are in its right subtree.

Now consider the recursion tree $Q(L')$ for $L' = [a, L]$ and let p be any pivot in $Q(L')$. Suppose $p < a$ and let k be any key in the left subtree of p in $Q(L)$. Since $k < p$, by the two properties k is in the left subtree of p in $Q(L')$. Similarly if $p > a$ then any k in the right subtree of p in $Q(L)$ is also in the right subtree of p in $Q(L')$. Since filtering preserves the respective order of keys in the input list, for any $p, p < a$, the input to the recursive call corresponding to its left child will be the same. Similarly, for $p > a$, the input to the recursive call corresponding to its right child will be the same. Thus, when sorting L' these recursive calls will find their results in the memo. Therefore only recursive calls corresponding to the root, to the children of the nodes in the rightmost spine of the left subtree of the root, and the children of the nodes in the leftmost spine of the right subtree of the root may be executed (the two spines are shown with thick lines in Figure 4). Furthermore, the results for the calls adjacent to the spines will be found in the memo.

Consider the calls whose results are not found in the memo. In the worst case, these will be all the calls along the two spines. Consider the sizes of inputs for the nodes on a spine and define the random variables $X_1 \dots X_k$ such that X_i is the least number of recursive calls (nodes) performed for the input size to become $(\frac{3}{4})^i n$ or less after it first becomes $(\frac{3}{4})^{(i-1)} n$ or less. Since $k \leq \lceil \log_{4/3} n \rceil$, the total and the expected number of operations along a spine are

$$\begin{aligned}
 C(n) &\leq \sum_{i=1}^{\lceil \log_{4/3} n \rceil} X_i \left(\frac{3}{4}\right)^{i-1} n, \text{ and} \\
 E[C(n)] &\leq \sum_{i=1}^{\lceil \log_{4/3} n \rceil} E[X_i] \left(\frac{3}{4}\right)^{i-1} n.
 \end{aligned}$$

Since the probability that the pivot lies in the middle half of the list is $\frac{1}{2}$, $E[X_i] \leq 2$ for $i \geq 1$, and we have

$$E[C(n)] \leq \sum_{i=1}^{\lceil \log_{4/3} n \rceil} 2 \left(\frac{3}{4} \right)^{i-1} n.$$

Thus, $E[C(n)] = O(n)$. This bound holds for both spines; therefore the number of operations due to calls whose results are not found in the memo is $O(n)$. Since each operation, including hash-consing, takes expected constant time, the total time of the calls whose results are not in the memo is $O(n)$. Now, consider the calls whose results are found in the memo, each such call will be on a spine or adjacent to it, thus there are an expected $O(\log n)$ such calls. Since, the memo table lookup overhead is expected constant time the total cost for these is $O(\log n)$. We conclude that Quicksort will take expected $O(n)$ time for sorting the modified list L' . ■

It is easy to extend the theorem to show that the $O(n)$ bound holds for an insertion anywhere in the list. Although, this bound is better than a complete rerun, which would take $O(n \log n)$, we would like to achieve $O(\log n)$. In Section 6 we discuss how a combination of memoization and adaptivity [2] may be used to reduce the expected cost of a random insertion to $O(\log n)$.

4 MFL

In this section we study a small functional language, called MFL, that supports selective memoization. MFL distinguishes memoized from non-memoized code, and is equipped with a modality for tracking dependences on data structures within memoized code. This modality is central to our approach to selective memoization, and is the focus of our attention here. The main result is a soundness theorem stating that memoization does not affect the outcome of a computation compared to a standard, non-memoizing semantics. We also show that the memoization mechanism of MFL causes a constant factor slowdown compared to a standard, non-memoizing semantics.

4.1 Abstract Syntax

The abstract syntax of MFL is given in Figure 5. The meta-variables x and y range over a countable set of *variables*. The meta-variables a and b range over a countable set of *resources*. (The distinction will be made clear below.) The meta-variable l ranges over a countable set of *locations*. We assume that variables, resources, and locations are mutually disjoint. The binding and scope conventions for variables and resources are as would be expected from the syntactic forms. As usual we identify pieces of syntax that differ only in their choice of bound variable or resource names. A term or expression is *resource-free* if and only if it contains no free resources, and is *variable-free* if and only if it contains no free variables. A *closed* term or expression is both resource-free and variable-free; otherwise it is *open*.

The types of MFL include 1 (unit), `int`, products and sums, recursive data types $\mu u. \tau$, memoized function types, and bang types $! \eta$. MFL distinguishes *indexable types*, denoted η , as those that accept an injective function, called an *index function*, whose co-domain is integers. The underlying type of a bang type $! \eta$ is restricted to be an indexable type. For `int` type, identity serves as an index function; for 1 (unit) any constant function can be chosen as the index function. For non-primitive types an index can be supplied by boxing values of these types. Boxed values would be allocated

| | | | |
|---------------|--------|-------|--|
| <i>Types</i> | η | $::=$ | <code>1</code> <code>int</code> ... |
| <i>Types</i> | τ | $::=$ | η $! \eta$ $\tau_1 \times \tau_2$ $\tau_1 + \tau_2$ $\mu u. \tau$ $\tau_1 \rightarrow \tau_2$ |
| <i>Op's</i> | o | $::=$ | <code>+</code> <code>-</code> ... |
| <i>Expr's</i> | e | $::=$ | <code>return</code> (t) <code>let</code> $! x : \eta$ <code>be</code> t <code>in</code> e <code>end</code> <code>let</code> $a_1 : \tau_1 \times a_2 : \tau_2$ <code>be</code> t <code>in</code> e <code>end</code> <code>mc</code> <code>case</code> t <code>of</code> <code>inl</code> $(a_1 : \tau_1) \Rightarrow e_1$ <code>inr</code> $(a_2 : \tau_2) \Rightarrow e_2$ |
| <i>Terms</i> | t | $::=$ | v $o(t_1, \dots, t_n)$ $\langle t_1, t_2 \rangle$ <code>mf</code> <code>fun</code> $f(a : \tau_1) : \tau_2$ <code>is</code> e <code>end</code> $t_1 t_2$ $! t$ <code>inl</code> $\tau_1 + \tau_2$ t <code>inr</code> $\tau_1 + \tau_2$ t <code>roll</code> (t) <code>unroll</code> (t) |
| <i>Values</i> | v | $::=$ | x a \star n $! v$ $\langle v_1, v_2 \rangle$ <code>mf</code> <code>unl</code> $f(a : \tau_1) : \tau_2$ <code>is</code> e <code>end</code> |

Figure 5. The abstract syntax of MFL.

in a store and the unique location of a box would serve as an index for the underlying value. With this extension the indexable types would be defined as $\eta : : = 1 \mid \text{int} \mid \tau \text{ box}$. Although supporting boxed types is critical for practical purposes, we do not formalize this here to focus on the main ideas.

The syntax is structured into *terms* and *expressions*, in the terminology of Pfenning and Davies [30]. Roughly speaking, terms evaluate independently of their context, as in ordinary functional programming, whereas expressions are evaluated relative to a memo table. Thus, the body of a memoized function is an expression, whereas the function itself is a term. Note, however, that the application of a function is a *term*, not an *expression*; this corresponds to the encapsulation of memoization with the function, so that updating the memo table is benign. In a more complete language we would include case analysis and projection forms among the terms, but for the sake of simplicity we include these only as expressions. We would also include a plain function for which the body is a term. Note that every term is trivially an expression; the `return` expression is the inclusion.

4.2 Static Semantics

The type structure of MFL extends the framework of Pfenning and Davies [30] with a “necessitation” modality, $! \eta$, which is used to track data dependences for selective memoization. This modality does not correspond to a monadic interpretation of memoization effects ($\odot \tau$ in the notation of Pfenning and Davies), though one could imagine adding such a modality to the language. The introductory and eliminatory forms for necessity are standard, namely $! t$ for introduction, and `let` $! x : \eta$ `be` t `in` e `end` for elimination.

Our modality demands that we distinguish variables from resources. Variables in MFL correspond to the “validity”, or “unrestricted”, context in modal logic, whereas resources in MFL correspond to the “truth”, or “restricted” context. An analogy may also be made to the judgmental presentation of linear logic [29, 31]: variables correspond to the intuitionistic context, resources to the linear context.¹

¹Note, however, that we impose no linearity constraints in our type system!

| | |
|-------------|--|
| Var, Res | $\frac{(\Gamma(x) = \tau)}{\Gamma; \Delta \vdash x : \tau} \quad \frac{(\Delta(a) = \tau)}{\Gamma; \Delta \vdash a : \tau}$ |
| Unit & Nums | $\frac{}{\Gamma; \Delta \vdash n : \text{int}} \quad \frac{}{\Gamma; \Delta \vdash \star : 1}$ |
| Prim | $\frac{\Gamma; \Delta \vdash t_i : \tau_i \ (1 \leq i \leq n) \quad \vdash_o o : (\tau_1, \dots, \tau_n) \tau}{\Gamma; \Delta \vdash o(t_1, \dots, t_n) : \tau}$ |
| Pairs | $\frac{\Gamma; \Delta \vdash t_1 : \tau_1 \quad \Gamma; \Delta \vdash t_2 : \tau_2}{\Gamma; \Delta \vdash \langle t_1, t_2 \rangle : \tau_1 \times \tau_2}$ |
| Fun | $\frac{\Gamma, f : \tau_1 \rightarrow \tau_2; \Delta, a : \tau_1 \vdash e : \tau_2}{\Gamma; \Delta \vdash \text{mfun } f(a : \tau_1) : \tau_2 \text{ is } e \text{ end} : \tau_1 \rightarrow \tau_2}$ |
| FunVal | $\frac{\Gamma, f : \tau_1 \rightarrow \tau_2; \Delta, a : \tau_1 \vdash e : \tau_2}{\Gamma; \Delta \vdash \text{mfun}_l f(a : \tau_1) : \tau_2 \text{ is } e \text{ end} : \tau_1 \rightarrow \tau_2}$ |
| Apply | $\frac{\Gamma; \Delta \vdash t_1 : \tau_1 \rightarrow \tau_2 \quad \Gamma; \Delta \vdash t_2 : \tau_1}{\Gamma; \Delta \vdash t_1 t_2 : \tau_2}$ |
| Bang | $\frac{\Gamma; \emptyset \vdash t : \eta}{\Gamma; \Delta \vdash !t : !\eta}$ |
| Inl, Inr | $\frac{\Gamma; \Delta \vdash t : \tau_1}{\Gamma; \Delta \vdash \text{inl}_{\tau_1 + \tau_2} t : \tau_1 + \tau_2}$ $\frac{\Gamma; \Delta \vdash t : \tau_2}{\Gamma; \Delta \vdash \text{inr}_{\tau_1 + \tau_2} t : \tau_1 + \tau_2}$ |
| (Un)Roll | $\frac{\Gamma; \Delta \vdash t : [\mu u. \tau / u] \tau}{\Gamma; \Delta \vdash \text{roll}(t) : \mu u. \tau}$ $\frac{\Gamma; \Delta \vdash t : \mu u. \tau}{\Gamma; \Delta \vdash \text{unroll}(t) : [\mu u. \tau / u] \tau}$ |

Figure 6. Typing judgments for terms.

The inclusion, $\text{return}(t)$, of terms into expressions has no analogue in pure modal logic, but is specific to our interpretation of memoization as a computational effect. The typing rule for $\text{return}(t)$ requires that t be resource-free to ensure that any dependence on the argument to a memoized function is made explicit in the code before computing the return value of the function. In the first instance, resources arise as parameters to memoized functions, with further resources introduced by their incremental decomposition using $\text{let} \times$ and mcase . These additional resources track the usage of as-yet-unexplored parts of a data structure. Ultimately, the complete value of a resource may be accessed using the $\text{let}!$ construct, which binds its value to a variable, which may be used without restriction. In practice this means that those parts of an argument to a memoized function on whose value the function depends will be given modal type. However, it is not essential that all resources have modal type, nor that the computation depend upon every resource that does have modal type.

The static semantics of MFL consists of a set of rules for deriving typing judgments of the form $\Gamma; \Delta \vdash t : \tau$, for terms, and $\Gamma; \Delta \vdash e : \tau$, for expressions. In these judgments Γ is a *variable type assignment*, a finite function assigning types to variables, and Δ is a *resource type assignment*, a finite function assigning types to resources. The rules for deriving these judgments are given in Figures 6 and 7.

| | |
|--------------|--|
| Return | $\frac{\Gamma; \emptyset \vdash t : \tau}{\Gamma; \Delta \vdash \text{return}(t) : \tau}$ |
| Let! | $\frac{\Gamma; \Delta \vdash t : !\eta \quad \Gamma, x : \eta; \Delta \vdash e : \tau}{\Gamma; \Delta \vdash \text{let } !x : \eta \text{ be } t \text{ in } e \text{ end} : \tau}$ |
| Let \times | $\frac{\Gamma; \Delta \vdash t : \tau_1 \times \tau_2 \quad \Gamma; \Delta, a_1 : \tau_1, a_2 : \tau_2 \vdash e : \tau}{\Gamma; \Delta \vdash \text{let } a_1 : \tau_1 \times a_2 : \tau_2 \text{ be } t \text{ in } e \text{ end} : \tau}$ |
| Case | $\frac{\begin{array}{c} \Gamma; \Delta \vdash t : \tau_1 + \tau_2 \\ \Gamma; \Delta, a_1 : \tau_1 \vdash e_1 : \tau \\ \Gamma; \Delta, a_2 : \tau_2 \vdash e_2 : \tau \end{array}}{\Gamma; \Delta \vdash \text{mcase } t \text{ of } \text{inl}(a_1 : \tau_1) \Rightarrow e_1 : \tau \mid \text{inr}(a_2 : \tau_2) \Rightarrow e_2}$ |

Figure 7. Typing judgments for expressions.

4.3 Dynamic Semantics

The dynamic semantics of MFL formalizes selective memoization. Evaluation is parameterized by a store containing memo tables that track the behavior of functions in the program. Evaluation of a function expression causes an empty memo table to be allocated and associated with that function. Application of a memoized function is affected by, and may affect, its associated memo table. Should the function value become inaccessible, so also is its associated memo table, and hence the storage required for both can be reclaimed.

Unlike conventional memoization, however, the memo table is keyed by control flow information rather than by the values of arguments to memoized functions. This is the key to supporting selective memoization. Expression evaluation is essentially an exploration of the available resources culminating in a resource-free term that determines its value. Since the exploration is data-sensitive, only certain aspects of the resources may be relevant to a particular outcome. For example, a memoized function may take a pair of integers as argument, with the outcome determined independently of the second component in the case that the first is positive. By recording control-flow information during evaluation, we may use it to provide selective memoization.

For example, in the situation just described, all pairs of the form $\langle 0, v \rangle$ should map to the same result value, irrespective of the value v . In conventional memoization the memo table would be keyed by the pair, with the result that redundant computation is performed in the case that the function has not previously been called with v , even though the value of v is irrelevant to the result! In our framework we instead key the memo table by a “branch” that records sufficient control flow information to capture the general case. Whenever we encounter a return statement, we query the memo table with the current branch to determine whether this result has been computed before. If so, we return the stored value; if not, we evaluate the return statement, and associate that value with that branch in the memo table for future use. It is crucial that the returned term not contain any resources so that we are assured that its value does not change across calls to the function.

The dynamic semantics of MFL is given by a set of rules for deriving judgments of the form $\sigma, t \Downarrow^t v, \sigma'$ (for terms) and $\sigma, l : \beta, e \Downarrow^e v, \sigma'$ (for expressions). The rules for deriving these judgments are given in Figures 8 and 9. These rules make use of branches, memo tables, and stores, whose precise definitions are as follows.

A *simple branch* is a list of *simple events* corresponding to “choice points” in the evaluation of an expression.

$$\begin{array}{ll} \text{Simple Event } \varepsilon & ::= !v \mid \text{inl} \mid \text{inr} \\ \text{Simple Branch } \beta & ::= \bullet \mid \varepsilon \cdot \beta \end{array}$$

We write $\beta \hat{\varepsilon}$ to stand for the extension of β with the event ε at the end.

A *memo table*, θ , is a finite function mapping simple branches to values. We write $\theta[\beta \mapsto v]$, where $\beta \notin \text{dom}(\theta)$, to stand for the extension of θ with the given binding for β . We write $\theta(\beta) \uparrow$ to mean that $\beta \notin \text{dom}(\theta)$.

A *store*, σ , is a finite function mapping *locations*, l , to memo tables. We write $\sigma[l \mapsto \theta]$, where $l \notin \text{dom}(\sigma)$, to stand for the extension of σ with the given binding for l . When $l \in \text{dom}(\sigma)$, we write $\sigma[l \leftarrow \theta]$ for the store σ that maps l to θ and $l' \neq l$ to $\sigma(l')$.

Term evaluation is largely standard, except for the evaluation of (memoizing) functions and applications of these to arguments. Evaluation of a memoizing function term allocates a fresh memo table, which is then associated with the function’s value. Expression evaluation is initiated by an application of a memoizing function to an argument. The function value determines the memo table to be used for that call. Evaluation of the body is performed relative to that table, initiating with the null branch.

Expression evaluation is performed relative to a “current” memo table and branch. When a `return` statement is encountered, the current memo table is consulted to determine whether or not that branch has previously been taken. If so, the stored value is returned; otherwise, the argument term is evaluated, stored in the current memo table at that branch, and the value is returned. The `let!` and `mcas` expressions extend the current branch to reflect control flow. Since `let!` signals dependence on a complete value, that value is added to the branch. Case analysis, however, merely extends the branch with an indication of which case was taken. The `let×` construct does not extend the branch, because no additional information is gleaned by splitting a pair.

4.4 Soundness of MFL

We will prove the soundness of MFL relative to a non-memoizing semantics for the language. It is straightforward to give a purely functional semantics to the pure fragment of MFL by an inductive definition of the relations $t \Downarrow_p^t v$ and $e \Downarrow_p^e v$. Here t , e , and v are “pure” in the sense that they may not involve subscripted function values. The *underlying term*, t^- , of an MFL term, t , is obtained by erasing all location subscripts on function values occurring within t .

The soundness of MFL consists of showing that evaluation with memoization yields the same outcome as evaluation without memoization.

Theorem 2 (Soundness)

If $\emptyset, t \Downarrow^t v, \sigma$, where $\emptyset; \emptyset \vdash t : \tau$, then $t^- \Downarrow_p^t v^-$.

The full proof is given in [3]. The statement of the theorem must be strengthened considerably to account for both terms and expressions, and to take account of non-empty memoization contexts. The proof then proceeds by induction on evaluation.

It is easy to show that the non-memoizing semantics of MFL is type safe, using completely conventional techniques. It follows that the

| | |
|----------|--|
| Unit | $\sigma, \star \Downarrow^t \star, \sigma$ |
| Number | $\sigma, n \Downarrow^t n, \sigma$ |
| PrimOp | $\frac{\begin{array}{c} \sigma, t_1 \Downarrow^t v_1, \sigma_1 \\ \vdots \\ \sigma_{n-1}, t_n \Downarrow^t v_n, \sigma_n \\ v = \text{app}(o, (v_1, \dots, v_n)) \end{array}}{\sigma, o(t_1, \dots, t_n) \Downarrow^t v, \sigma_n}$ |
| Pair | $\frac{\begin{array}{c} \sigma, t_1 \Downarrow^t v_1, \sigma' \\ \sigma', t_2 \Downarrow^t v_2, \sigma'' \end{array}}{\sigma, \langle t_1, t_2 \rangle \Downarrow^t \langle v_1, v_2 \rangle, \sigma''}$ |
| Fun | $\frac{\begin{array}{c} (e = \text{mfun}_f(a : \tau_1) : \tau_2 \text{ is } e' \text{ end}) \\ (v = \text{mfun}_f(a : \tau_1) : \tau_2 \text{ is } e' \text{ end}) \\ (l \notin \text{dom}(\sigma), \quad \sigma' = \sigma[l \mapsto \emptyset]) \end{array}}{\sigma, e \Downarrow^t v, \sigma'}$ |
| FunVal | $\frac{(v = \text{mfun}_f(a : \tau_1) : \tau_2 \text{ is } e' \text{ end}, l \in \text{dom}(\sigma))}{\sigma, v \Downarrow^t v, \sigma}$ |
| Apply | $\frac{\begin{array}{c} \sigma, t_1 \Downarrow^t v_1, \sigma_1 \\ \sigma_1, t_2 \Downarrow^t v_2, \sigma_2 \\ \sigma_2, l : \bullet, [v_1, v_2 / f, a] e \Downarrow^e v, \sigma' \\ (v_1 = \text{mfun}_f(a : \tau_1) : \tau_2 \text{ is } e' \text{ end}) \end{array}}{\sigma, t_1 t_2 \Downarrow^t v, \sigma'}$ |
| Bang | $\frac{\sigma, t \Downarrow^t v, \sigma'}{\sigma, !t \Downarrow^t !v, \sigma'}$ |
| Inject | $\frac{\sigma, t \Downarrow^t v, \sigma'}{\sigma, \text{inl}_{\tau_1 + \tau_2} t \Downarrow^t \text{inl}_{\tau_1 + \tau_2} v, \sigma'}$ |
| | $\frac{\sigma, t \Downarrow^t v, \sigma'}{\sigma, \text{inr}_{\tau_1 + \tau_2} t \Downarrow^t \text{inr}_{\tau_1 + \tau_2} v, \sigma'}$ |
| (Un)Roll | $\frac{\sigma, t \Downarrow^t v, \sigma'}{\sigma, \text{roll}(t) \Downarrow^t \text{roll}(v), \sigma'}$ |
| | $\frac{\sigma, t \Downarrow^t \text{roll}(v), \sigma'}{\sigma, \text{unroll}(t) \Downarrow^t v, \sigma'}$ |

Figure 8. Evaluation of terms.

memoizing semantics is also type-safe, for if not, there would be a closed value of a type τ that is not canonical for that type. However, erasure preserves and reflects canonical forms, hence, by the Soundness Theorem, MFL must also be type safe.

4.5 Performance

We show that memoization slows down an MFL program by a constant factor (expected) with respect to a standard, non-memoizing semantics even when no results are re-used. The result relies on representing a branch as a sequence of integers and using this sequence to key memo tables, which are implemented as hash tables. To represent branches as integer sequences we use the property of MFL that the underlying type η of a bang type, $! \eta$, is an indexable

| | |
|------|--|
| | $\frac{\sigma(l)(\beta) = v}{\sigma, l: \beta, \text{return}(t) \Downarrow^e v, \sigma} \quad (\text{Found})$ |
| Ret | $\frac{\begin{array}{c} \sigma(l) = \theta \quad \theta(\beta) \uparrow \\ \sigma, t \Downarrow^t v, \sigma' \\ \sigma'(l) = \theta' \end{array}}{\sigma, l: \beta, \text{return}(t) \Downarrow^e v, \sigma' [l \leftarrow \theta'[\beta \mapsto v]]} \quad (\text{Not Found})$ |
| Let! | $\frac{\begin{array}{c} \sigma, t \Downarrow^t !v, \sigma' \\ \sigma', l: !v \cdot \beta, [v/x]e \Downarrow^t v', \sigma'' \end{array}}{\sigma, l: \beta, \text{let}! x: \eta \text{ be } t \text{ in } e \text{ end} \Downarrow^e v', \sigma''}$ |
| Let× | $\frac{\begin{array}{c} \sigma, t \Downarrow^t v_1 \times v_2, \sigma' \\ \sigma', l: \beta, [v_1/a_1, v_2/a_2]e \Downarrow^e v, \sigma'' \end{array}}{\sigma, l: \beta, \text{let } a_1 \times a_2 \text{ be } t \text{ in } e \text{ end} \Downarrow^t v, \sigma''}$ |
| Case | $\frac{\begin{array}{c} \sigma, t \Downarrow^t \text{inl}_{\tau_1 + \tau_2} v, \sigma' \\ \sigma', l: \text{inl} \cdot \beta, [v/a_1]e_1 \Downarrow^e v_1, \sigma'' \end{array}}{\sigma, l: \beta, \text{mcase } t \text{ of } \text{inl } (a_1: \tau_1) \Rightarrow e_1 \Downarrow^t v_1, \sigma'' \mid \text{inr } (a_2: \tau_2) \Rightarrow e_2}$ $\frac{\begin{array}{c} \sigma, t \Downarrow^t \text{inr}_{\tau_1 + \tau_2} v, \sigma' \\ \sigma', l: \text{inr} \cdot \beta, [v/a_2]e_2 \Downarrow^e v_2, \sigma'' \end{array}}{\sigma, l: \beta, \text{mcase } t \text{ of } \text{inl } (a_1: \tau_1) \Rightarrow e_1 \Downarrow^t v_2, \sigma'' \mid \text{inr } (a_2: \tau_2) \Rightarrow e_2}$ |

Figure 9. Evaluation of expressions.

type. Since any value of an indexable type has an integer index, we can represent a branch of dependencies as sequence of integers corresponding to the indices of `let!`ed values, and zero or one for `inl` and `inr`.

Consider a non-memoizing semantics, where the `return` rule always evaluates its body and neither looks up nor updates memo tables (stores). Consider an MFL program and let T denote the time it takes (the number of evaluation steps) to evaluate the program with respect to this non-memoizing semantics. Let T' denote the time it takes to evaluate the same program with respect to the memoizing semantics. In the worst case, no results are re-used, thus the difference between T and T' is due to memo-table lookups and updates done by the memoizing semantics. To bound the time for these, consider a memo table lookup or update with a branch β and let $|\beta|$ be the length of the branch. Since a branch is a sequence of integers, a lookup or update can be performed in expected $O(|\beta|)$ time using nested hash tables to represent memo tables. Now note that the non-memoizing semantics takes $|\beta|$ time to build the branch thus, the cost of a lookup or update can be charged to the evaluations that build the branch β , *i.e.*, evaluations of `let!` and `mcase`. Furthermore, each evaluation of `let!` and `mcase` can be charged by exactly one `return`. Thus, we conclude that $T' = O(T)$ in the expected case.

5 Implementation

We describe an implementation of our framework as a Standard ML library. The aspects of the MFL language that relies on the syntactic distinction between resources and variables cannot be enforced statically in Standard ML. Therefore, we use a separate type for resources and employ run-time checks to detect violations of correct usage.

```
signature MEMO =
sig
  (* Expressions *)
  type 'a expr
  val return: (unit -> 'a) -> 'a expr

  (* Resources *)
  type 'a res
  val expose: 'a res -> 'a

  (* Bangs *)
  type 'a bang
  val bang: ('a -> int) -> 'a -> 'a bang
  val letBang: ('a bang) -> ('a -> 'b expr) -> 'b expr

  (* Products *)
  type ('a, 'b) prod
  val pair: 'a -> 'b -> ('a, 'b) prod
  val letx: ('a, 'b) prod -> (('a res * 'b res) -> 'c expr) -> 'c expr
  val split: ('a, 'b) prod -> (('a * 'b) -> 'c) -> 'c

  (* Sums *)
  type ('a, 'b) sum
  val inl: 'a -> ('a, 'b) sum
  val inr: 'b -> ('a, 'b) sum
  val mcase: ('a, 'b) sum -> ('a res -> 'c expr) -> ('b res -> 'c expr) -> 'c expr
  val choose: ('a, 'b) sum -> ('a -> 'c) -> ('b -> 'c) -> 'c

  (* Memoized arrow *)
  type ('a, 'b) marrow
  val mfun: ('a res -> 'b expr) -> ('a, 'b) marrow
  val mfun_rec: (('a, 'b) marrow -> 'a res -> 'b expr) -> ('a, 'b) marrow
  val mapapply: ('a, 'b) marrow -> 'a -> 'b
end

signature BOX = sig
  type 'a box
  val init: unit -> unit
  val box: 'a -> 'a box
  val unbox: 'a box -> 'a
  val getKey: 'a box -> int
end
```

Figure 10. The signatures for the memo library and boxes.

The interface for the library (shown in Figure 10) provides types for expressions, resources, bangs, products, sums, memoized functions along with their introduction and elimination forms. All expressions have type `'a expr`, which is a monad with `return` as the inclusion and various forms of “bind” induced by the elimination forms `letBang`, `letx`, and `mcase`. A resource has type `'a res` and `expose` is its elimination form. Resources are only created by the library, thus no introduction form for resources is available to the user. The introduction and elimination form for bang types are `bang` and `letBang`. The introduction and elimination form for product types are `pair`, and `letx` and `split` respectively. The `letx` is a form of “bind” for the monad `expr`; `split` is the elimination form for the term context. The treatment of sums is similar to product types. The introduction forms are `inl` and `inr`, and the elimination forms are `mcase` and `choose`; `mcase` is a form of bind for the `expr` monad and `choose` is the elimination for the term context.

Memoized functions are introduced by `mfun` and `mfun_rec`; `mfun` takes a function of type `'a res -> 'b expr` and returns the memoized function of type `('a, 'b) marrow`; `mfun_rec` is similar to `mfun` but it also takes as a parameter its memoized version. Note that the result type does not contain the “effect” `expr`—we encapsulate memoization effects, which are benign, within the function. The elimination form for the `marrow` is the memoized apply function `mapapply`.


```

functor BuildMemo (structure Box: BOX
                    structure Memopad:MEMOPAD):MEMO =
struct

  type 'a expr = int list * (unit -> 'a)
  fun return f = (nil,f)

  type 'a res = 'a
  fun res v = v
  fun expose r = r

  type 'a bang = 'a * ('a -> int)
  fun bang h t = (t,h)
  fun letBang b f = let
    val (v,h) = b
    val (branch,susp) = f v
  in
    ((h v)::branch, susp)
  end

  type ('a,'b) prod = 'a * 'b
  fun pair x y = (x,y)
  fun split p f = f p
  fun letx p f = let
    val (x1,x2) = p
  in
    f (res x1, res x2)
  end

  datatype ('a,'b) sum = INL of 'a | INR of 'b
  fun inl v = INL(v)
  fun inr v = INR(v)
  fun mcase s f g = let
    val (lr,(branch,susp)) =
      case s of
        INL v => (0,f (res v))
      | INR v => (1,g (res v))
  in
    (lr::branch,susp)
  end
  fun choose s f g = case s of INL v => f v
    | INR v => g v

  type ('a,'b) marrow = 'a -> 'b
  fun mfun_rec f = let
    val mpad = Memopad.empty ()
    fun mf rf x = let
      val (branch,susp) = f rf (res x)
      val result =
        case Memopad.extend mpad branch of
          (NONE,SOME mpad') => let (* Not found *)
            val v = susp ()
            val _ = Memopad.add v mpad'
          in
            v
          end
        | (SOME v,NONE) => v (* Found *)
      in
        result
      end
    fun mf' x = mf mf' x
  in
    mf'
  end

  fun mfun f = ... (* Similar to mfun_rec *)

  fun mapply f v = f v
end

```

Figure 11. The implementation of the memoization library.

Figure 11 shows an implementation of the library without the run-time checks for correct usage. To incorporate the run-time checks, one needs a more sophisticated definition of resources in order to detect when a resource is exposed out of its context (*i.e.*, function instance). In addition, the interface must be updated so that the first parameter of `letBang`, `letx`, and `mcase`, occurs in suspended form. This allows us to update the state consisting of certain flags before forcing a term.

```

structure Examples =
struct
  type 'a box = 'a Box.box

  (** Some utilities **)
  fun iBang v = bang (fn i => i) v
  fun bBang b = bang (fn b => Box.key b) b

  (** Fibonacci **)
  fun mfib' f (n') =
    letBang (expose n') (fn n => return (fn()=>
      if n < 2 then
        n
      else
        mapply f (iBang(n-1))+
        mapply f (iBang(n-2))))
  fun mfib n = mapply (mfun_rec mfib') n

  (** Boxed lists **)
  datatype 'a blist' = NIL
    | CONS of ('a * (('a blist') box))
  type 'a blist = ('a blist') box

  (** Hash Cons **)
  fun hCons' (x') =
    letx (expose x') (fn (h',t') =>
      letBang (expose h') (fn h =>
        letBang (expose t') (fn t =>
          return (fn()=> box (CONS(h,t))))))
    val hCons = mfun hCons'

  (** Knapsack **)
  fun mks' mks (arg) =
    letx (expose arg) (fn (c',l') =>
      letBang (expose c') (fn c =>
        letBang (expose l') (fn l => return (fn () =>
          case (unbox l) of
            NIL => 0
          | CONS((w,v),t) =>
            if (c < w) then
              mapply mks (pair (iBang c) (bBang t))
            else
              let
                val arg1 = pair (iBang c) (bBang t)
                val v1 = mapply mks arg1
                val arg2 = pair (iBang (c-w)) (bBang t)
                val v2 = v + mapply mks arg2
              in
                if (v1 > v2) then v1
                else v2
              end))))
    fun mks x = mapply (mfun_rec mks') x

  (** Quicksort **)
  fun mqs () = let
    val empty = box NIL
    val hCons = mfun hCons'
    fun fil f l =
      case (unbox l) of
        NIL => empty
      | CONS(h,t) =>
        if (f h) then
          mapply hCons (pair (iBang h) (bBang (fil f t)))
        else
          fil f t
    fun qs' qs (l') =
      letBang (expose l') (fn l => return (fn () =>
        case (unbox l) of
          NIL => nil
        | CONS(h,t) =>
          let
            val ll = fil (fn x=>x<h) t
            val gg = fil (fn x=>x>=h) t
            val sll = mapply qs (bBang ll)
            val sgg = mapply qs (bBang gg)
          in
            sll@(h::sgg)
          end))
    in
      mfun_rec qs'
    end
  end
end

```

Figure 12. Examples from Section 3 in the SML library.

The implementation extends the operational semantics of the MFL language (Section 4.3) with boxes. The `bang` primitive takes a value and an injective function, called the index function, that maps the value to an integer, called the index. The index of a value is used to key memo tables. The restriction that the indices be unique, enables us to implement memo tables as a nested hash tables, which support update and lookup operations in expected constant time. The primitive `letBang` takes a value `b` of bang type and a body. It applies the body to the underlying value of `b`, and extends the branch with the index of `b`. The function `letx` takes a pair `p` and a body. It binds the parts of the pair to two resources and applies the body to the resources; as with the operational semantics, `letx` does not extend the branch. The function `mcase` takes value `s` of sum type and a body. It branches on the outer form of `s` and binds its inner value to a resource. It then applies the body to the resource and extends the branch with 0 or 1 depending on the outer form of `s`. The elimination forms of sums and products for the term context, `split` and `choose` are standard.

The `return` primitive finalizes the branch and returns its body as a suspension. The branch is used by `mfun_rec` or `mfun`, to key the memo table; if the result is found in the memo table, then the suspension is disregarded and the result is re-used; otherwise the suspension is forced and the result is stored in the memo table keyed by the branch. The `mfun_rec` primitive takes a recursive function `f` as a parameter and “memoizes” `f` by associating it with a memo pad. A subtle issue is that `f` must call its memoized version recursively. Therefore `f` must take its memoized version as a parameter. Note also that the memoized function internally converts its parameter to a resource before applying `f` to it.

The interface of the library provides no introduction form for resources. Indeed, all resources are created by the library inside the `letx`, `mcase`, `mfun_rec`, and `mfun`. The function `expose` is the elimination form for resources. If, for example, one would like to apply `letBang` to a resource, then he must first `expose` the resource, which “exposes” the underlying value.

Figure 12 show the examples from Section 3 written in the SML library. Note that the memoized Fibonacci function `mfib` creates a memo table every time it is called. When `mfib` finishes, this table can be garbage collected (the same applies to `mks`). For Quicksort, we provide a function `mqs` that returns an instance of memoized Quicksort when applied. Each such instance has its own memo table. Note also that `mqs` creates a local instance of the hash-cons function so that each instance of memoized Quicksort has its own memo table for hash-consing.

In the examples, we do not use the sum types provided by the library to represent boxed lists, because we do not need to. In general, one will use the provided sum types instead of their ML counterparts (for example if an `mcase` is needed). The examples in Figure 12 can be implemented using the following definition of boxed lists.

```
datatype 'a boxlist' =
  ROLL of (unit, (('a, 'a boxlist' box) prod)) sum
type 'a boxlist = ('a boxlist')
```

Changing the code in Figure 12 to work with this definition of boxed lists requires several straightforward modifications.

6 Discussion

Space and Cache Management. Our framework associates a separate memo table with each memoized function. This allows the programmer to control the life-span of memo tables by conventional scoping. This somewhat coarse degree of control is sufficient in certain applications such as in dynamic programming, but finer level of control may be desirable for applications where result reuse is less regular. Such an application can benefit from specifying a caching scheme for individual memo tables so as to determine the size of the memo table and the replacement policy. We discuss how the framework can be extended to associate a cache scheme with each memo table and maintain the memo table accordingly.

The caching scheme should be specified in the form of a parameter to the `mfun` construct. When evaluated, this construct will bind the caching scheme to the memo table and the memo table will be maintained accordingly. Changes to the operational semantics to accommodate this extension is small. The store σ will now map a label to a pair consisting of a memo table and its caching scheme. The handling of the `return` will be changed so that the stores do not merely expand but are updated according to the caching scheme before adding a new entry. The following shows the updated return rule. Here S denotes a caching scheme and θ denotes a memo table. The update function denotes a function that updates the memo table to accommodate a new entry by possibly purging an existing entry. The programmer must ensure that the caching scheme does not violate the integrity of the memo table by tampering with stored values.

$$\frac{\sigma(l) = (\theta, S) \quad \theta(\beta) = v}{\sigma, l:\beta, \text{return}(t) \Downarrow^e v, \sigma} \quad (\text{Found})$$

$$\frac{\sigma(l) = (\theta, S) \quad \theta(\beta) \uparrow \quad \sigma, t \Downarrow^t v, \sigma' \quad \theta' = \text{update}(\theta', S, (\beta, v))}{\sigma, l:\beta, \text{return}(t) \Downarrow^e v, \sigma' [l \leftarrow \theta']} \quad (\text{Not Found})$$

For example, we can specify that the memo table for the Fibonacci function, shown in Figure 1, can contain at most two entries and be managed using the least-recently-used replacement policy. This is sufficient to ensure that the memoized Fibonacci runs in linear time. This extension can also be incorporated into the type system described in Section 4. This would require that we associate types with memo stores and also require that we develop a type system for “safe” update functions if we are to enforce that the caching schemes are safe.

Local vs. Non-local Dependences. Our dependence tracking mechanism only captures “local” dependences between the input and the result of a function. A local dependence of a function `f` is one that is created inside the static scope of `f`. A non-local dependence of `f` is created when `f` passes its input to some other function `g`, which examines `f`’s input indirectly. In previous work, Abadi *et al.* [1] and Heydon *et al.* [14] showed a program analysis technique for tracking non-local dependences by propagating dependences of a function to its caller. They do not, however, make clear the performance implications of their technique.

Our framework can be extended to track non-local dependences by introducing an application form for memoized functions in the expression context. This extension would, for example, allow for dependences of non-constant length. We chose not to support non-local dependences because it is not clear if its utility exceeds its performance effects.

Memoization and Adaptivity. The work we present in this paper was motivated by our previous work on adaptive computation [2]. We briefly discuss the relationship between memoization and adaptivity and how they can be combined to obtain efficient dynamic or incremental algorithms.

An adaptive computation maintains a dynamic dependence graph representing data and control dependences. When the input is modified, a change propagation algorithm updates the output and the dependence graph. The adaptivity mechanism handles “deep” changes efficiently. We say that a change is *deep* if it affects calls that occur at leaves of the call tree for the computation. In contrast, a change is *shallow* if it affects by a calls that occur at the roots of the call tree.

As an example consider the Quicksort algorithm that picks the first key of its input as pivot. Inserting a new key at the end of the input list is a deep change because this change will affect the last recursive calls of some filter functions and will become pivot only at the end of some sequence of recursive calls to Quicksort. In contrast, inserting a new key at the beginning of the list is a shallow change for Quicksort, because the new key will be selected as a pivot immediately by the first call to Quicksort. The adaptivity mechanism based on dynamic dependence graphs handles an insertion at the end of the input, a deep change, in expected $O(\log n)$ time [2], whereas the insertion at the beginning of the list, a shallow change, will cause a complete rerun, which takes $O(n \log n)$ time. Using memoization, however, an insertion at the beginning of the list can be handled in $O(n)$ time as showed in Section 3.

Any change can be thought of a combination of shallow and deep changes. Since memoization and adaptivity complement each other in their handling of deep and shallow changes, we would expect that a combination of these two techniques would handle general changes efficiently. For example, in Quicksort, we expect that an insertion in a random position in the list would be handled in expected $O(\log n)$ time by a combination of these two techniques.

7 Conclusion

We presented a framework for selective memoization under programmer control. The framework makes explicit the performance effects of memoization and yields programs whose running times can be analyzed using standard techniques. A key aspect of the framework is that it can capture both control and data dependences between input and the result of a memoized function. The main contributions of the paper are the particular set of primitives we suggest and the semantics along with the proofs that it is sound. We gave a simple implementation of the framework in the Standard ML language. We expect that this framework can be implemented in any purely-functional language.

8 References

- [1] M. Abadi, B. W. Lampson, and J.-J. Levy. Analysis and caching of dependencies. In *International Conference on Functional Programming*, pages 83–91, 1996.
- [2] U. A. Acar, G. E. Blelloch, and R. Harper. Adaptive functional programming. In *Proceedings of the Twenty-ninth Annual ACM-SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, Jan. 2002.
- [3] U. A. Acar, G. E. Blelloch, and R. Harper. Selective memoization. Technical Report CMU-CS-02-194, Carnegie Mellon University, Computer Science Department, Dec. 2002.
- [4] A. V. Aho, J. E. Hopcroft, and J. D. Ullman. *The Design and Analysis of Computer Algorithms*. Addison-Wesley, 1974.
- [5] J. Allen. *Anatomy of LISP*. McGraw Hill, 1978.
- [6] A. W. Appel and M. J. R. Gonçalves. Hash-consing garbage collection. Technical Report CS-TR-412-93, Princeton University, Computer Science Department, 1993.
- [7] R. Bellman. *Dynamic Programming*. Princeton University Press, 1957.
- [8] R. S. Bird. Tabulation techniques for recursive programs. *ACM Computing Surveys*, 12(4):403–417, Dec. 2002.
- [9] N. H. Cohen. Eliminating redundant recursive calls. *ACM Transactions on Programming Languages and Systems*, 5(3):265–299, July 1983.
- [10] T. H. Cormen, C. E. Leiserson, and R. L. Rivest. *Introduction to Algorithms*. MIT Press/McGraw-Hill, 1990.
- [11] A. Demers, T. Reps, and T. Teitelbaum. Incremental evaluation of attribute grammars with application to syntax directed editors. In *Conference Record of the 8th Annual ACM Symposium on POPL*, pages 105–116, Jan. 1981.
- [12] J. Field and T. Teitelbaum. Incremental reduction in the lambda calculus. In *Proceedings of the ACM '90 Conference on LISP and Functional Programming*, pages 307–322, June 1990.
- [13] E. Goto and Y. Kanada. Hashing lemmas on time complexities with applications to formula manipulation. In *Proceedings of the 1976 ACM Symposium on Symbolic and Algebraic Computation*, pages 154–158, 1976.
- [14] A. Heydon, R. Levin, and Y. Yu. Caching function calls using precise dependencies. *ACM SIGPLAN Notices*, 35(5):311–320, 2000.
- [15] J. Hilden. Elimination of recursive calls using a small table of randomly selected function values. *BIT*, 16(1):60–73, 1976.
- [16] R. Hoover. Alphonse: incremental computation as a programming abstraction. In *Proceedings of the 5th ACM SIGPLAN conference on Programming language design and implementation*, pages 261–272. ACM Press, 1992.
- [17] R. J. M. Hughes. Lazy memo-functions. In *Proceedings 1985 Conference on Functional Programming Languages and Computer Architecture*, 1985.
- [18] S. P. Jones. *The Implementation of Functional Programming Languages*. Prentice-Hall, 1987.
- [19] Y. A. Liu and S. D. Stoller. Dynamic programming via static incrementalization. In *European Symposium on Programming*, pages 288–305, 1999.
- [20] Y. A. Liu, S. D. Stoller, and T. Teitelbaum. Static caching for incremental computation. *ACM Transactions on Programming Languages and Systems*, 20(3):546–585, 1 May 1998.
- [21] J. McCarthy. A Basis for a Mathematical Theory of Computation. In P. Braffort and D. Hirschberg, editors, *Computer Programming and Formal Systems*, pages 33–70. North-Holland, Amsterdam, 1963.
- [22] D. Michie. ‘memo’ functions and machine learning. *Nature*, 218:19–22, 1968.

- [23] J. Mostov and D. Cohen. Automating program speedup by deciding what to cache. In *Proceedings of the Ninth International Joint Conference on Artificial Intelligence*, pages 165–172, Aug. 1985.
- [24] T. Murphy, R. Harper, and K. Crary. The wizard of TILT: Efficient(?), convenient and abstract type representations. Technical Report CMU-CS-02-120, School of Computer Science, Carnegie Mellon University, Mar. 2002.
- [25] P. Norvig. Techniques for automatic memoization with applications to context-free parsing. *Computational Linguistics*, pages 91–98, 1991.
- [26] H. A. Partsch. *Specification and Transformation of Programs—A Formal Approach to Software Development*. Springer-Verlag, 1990.
- [27] M. Pennings. *Generating Incremental Attribute Evaluators*. PhD thesis, University of Utrecht, Nov. 1994.
- [28] M. Pennings, S. D. Swierstra, and H. Vogt. Using cached functions and constructors for incremental attribute evaluation. In *Seventh International Symposium on Programming Languages, Implementations, Logics and Programs*, pages 130–144, 1992.
- [29] F. Pfenning. Structural cut elimination. In D. Kozen, editor, *Proceedings of the Tenth Annual Symposium on Logic in Computer Science*, pages 156–166. Computer Society Press, 1995.
- [30] F. Pfenning and R. Davies. A judgmental reconstruction of modal logic. *Mathematical Structures in Computer Science*, 11:511–540, 2001. Notes to an invited talk at the *Workshop on Intuitionistic Modal Logics and Applications (IMLA'99)*, Trento, Italy, July 1999.
- [31] J. Polakow and F. Pfenning. Natural deduction for intuitionistic non-commutative linear logic. In J.-Y. Girard, editor, *Proceedings of the 4th International Conference on Typed Lambda Calculi and Applications (TLCA'99)*, pages 130–144. Springer-Verlag LNCS 1581, 1999.
- [32] W. Pugh. *Incremental computation via function caching*. PhD thesis, Department of Computer Science, Cornell University, 1987.
- [33] W. Pugh. An improved replacement strategy for function caching. In *Proceedings of the 1988 ACM conference on LISP and functional programming*, pages 269–276. ACM Press, 1988.
- [34] W. Pugh and T. Teitelbaum. Incremental computation via function caching. In *Conference Record of the 16th Annual Symposium on POPL*, pages 315–328, Jan. 1989.
- [35] J. M. Spitzen and K. N. Levitt. An example of hierarchical design and proof. *Communications of the ACM*, 21(12):1064–1075, 1978.
- [36] R. S. Sundaresh and P. Hudak. Incremental compilation via partial evaluation. In *Conference Record of the 18th Annual ACM Symposium on POPL*, pages 1–13, Jan. 1991.
- [37] Y. Zhang and Y. A. Liu. Automating derivation of incremental programs. In *Proceedings of the third ACM SIGPLAN international conference on Functional programming*, page 350. ACM Press, 1998.