

Hierarchical Memory Management for Mutable State

Extended Technical Appendix

Adrien Guatto
Carnegie Mellon University
adrien@guatto.org

Sam Westrick
Carnegie Mellon University
swestric@cs.cmu.edu

Ram Raghunathan
Carnegie Mellon University
ram.r@cs.cmu.edu

Umut Acar
Carnegie Mellon University
umut@cs.cmu.edu

Matthew Fluet
Rochester Institute of Technology
mtf@cs.rit.edu

Abstract

It is well known that modern functional programming languages are naturally amenable to parallel programming. Achieving efficient parallelism using functional languages, however, remains difficult. Perhaps the most important reason for this is their lack of support for efficient in-place updates, i.e., mutation, which is important for the implementation of both parallel algorithms and the run-time system services (e.g., schedulers and synchronization primitives) used to execute them.

In this paper, we propose techniques for efficient mutation in parallel functional languages. To this end, we couple the memory manager with the thread scheduler to make reading and updating data allocated by nested threads efficient. We describe the key algorithms behind our technique, implement them in the MLton Standard ML compiler, and present an empirical evaluation. Our experiments show that the approach performs well, significantly improving efficiency over existing functional language implementations.

CCS Concepts • **Software and its engineering** → **Garbage collection; Parallel programming languages; Functional languages;**

Keywords parallel functional language implementation, garbage collection, hierarchical heaps, mutation, promotion

1 Introduction

With the proliferation of parallel hardware, functional programming languages, such as Haskell and the ML family (OCaml, Standard ML), have received much attention from academia and industry. Even non-functional languages today such as C++, Python, and Swift support certain features of functional languages, including higher-order functions and, sometimes, rich type systems. An important virtue of

strongly typed functional languages is their ability to distinguish between pure and impure code. This aids in writing correct parallel programs by making it easier to avoid race conditions, which can become a formidable challenge in languages whose type systems don't distinguish between mutable and immutable data.

In the sequential realm, functional languages compete well with other garbage-collected languages such as Java and Go, often running within a factor of 2 or 3 and sometimes even faster. In many cases, functional languages even compete well with the C family, where memory is managed by the programmer [34].

In the parallel realm, however, the gap between functional and imperative languages is significantly larger. One reason for this is the (poor) support for mutation in parallel functional languages. A reality of modern hardware is that imperative algorithms can perform significantly better than pure functional algorithms by using constant-time random accesses and updates. Even when a parallel algorithm has a pure functional interface (immutable inputs and immutable outputs), it can be more efficient to use mutation internally. For example, the efficiency of a pure functional parallel merge-sort can be significantly improved by reverting to a sequential imperative quick-sort for small inputs. Committing to pure functional algorithms only does not completely avoid mutation: a language run-time system uses mutation to implement crucial facilities such as (thread) schedulers and synchronization primitives, which require communication between processors via shared memory.

Even though mutation is crucial for efficiency, it remains poorly supported in parallel functional languages and remains as an active area of research (see Section 6). For example, the Manticore project has developed rich extensions to the ML language to support parallelism but has focused on purely functional code where the programmer cannot use mutation [7, 16]. Other ML dialects such as OCaml and SML# continue to remain primarily sequential languages, though there is ongoing work in extending them to support modern parallelism features. Like Manticore, Haskell has a relatively rich set of parallelism features and its runtime must support efficient mutation [27]. Writing efficient parallel programs

PPoPP '18, February 24–28, 2018, Vienna, Austria

© 2018 Copyright held by the owner/author(s). Publication rights licensed to the Association for Computing Machinery.

This is the author's version of the work. It is posted here for your personal use. Not for redistribution. The definitive Version of Record was published in *PPoPP '18: Principles and Practice of Parallel Programming, February 24–28, 2018, Vienna, Austria*, <https://doi.org/10.1145/3178487.3178494>.

in Haskell, however, remains difficult in part because of lazy evaluation [29].

This state of the art raises the question of whether functional programming can be extended to support mutable data and parallelism efficiently. At the highest level of abstraction, this is a challenging problem because its parts—parallelism and efficient memory management—are individually challenging. The problem is further complicated by the fact that functional languages allocate and reclaim memory at a much faster rate than other languages [5–7, 13, 14, 18, 19, 27].

In this paper, we propose techniques for supporting mutable data efficiently within the run-time system of nested-parallel functional languages, focusing on strict languages in the style of the ML family. Our approach builds upon that of *hierarchical heaps*, a memory management technique developed in prior work [31]. The basic idea is to organize memory so that it mirrors the structure of the parallel computation. Specifically, each thread is assigned its own heap in a hierarchy (tree) of heaps which grows and shrinks as threads are forked and joined. Threads allocate data in their own heaps, and can read and update objects in ancestor heaps (including their own). A key invariant is that data in non-ancestor heaps remains unreachable to a thread. To enforce this invariant, we propose a *promotion* technique for copying data upwards in the hierarchy as necessary.

Our approach has several important benefits. First, threads can allocate, read, and update mutable objects in their heap, without synchronization or copying. This allows local mutable objects to be used efficiently. Second, because heaps are associated with threads rather than processors, a thread can be migrated between processors without copying data. These two properties contrast with the predominant approach to memory management in parallel functional languages with local heaps, where both mutation and thread migration require copying [7, 13, 14, 27]. Third, our techniques introduce no overhead for reads of immutable objects, which are pervasive in functional languages. Finally, any thread can collect its heap independently and, more generally, any subtree of heaps in the hierarchy could be collected independently.

The contributions of this paper include techniques and algorithms for handling mutable data in hierarchical heaps (Section 3), an implementation extending the MLton whole-program optimizing compiler for Standard ML [28], and an empirical evaluation considering a number of both pure and imperative benchmarks (Section 4). Our results show that these techniques can be implemented efficiently and can perform well in practice.

2 Overview

We present a brief overview of our techniques, using a simple example to illustrate both the programming model and details of memory management. In the process, we introduce terminology that will be used throughout the paper.

```

val GRAIN = ...
fun inplaceQSort s = ...
fun msort s =
  if Seq.length s <= GRAIN
  then let val a = Seq.toArray s
        val () = inplaceQSort a
        in Seq.fromArray a end
  else let val (l, r) = Seq.splitMid s
        val (l', r') = par (msort l, msort r)
        in Seq.merge (l', r') end

```

Figure 1. Code for parallel imperative merge sort.

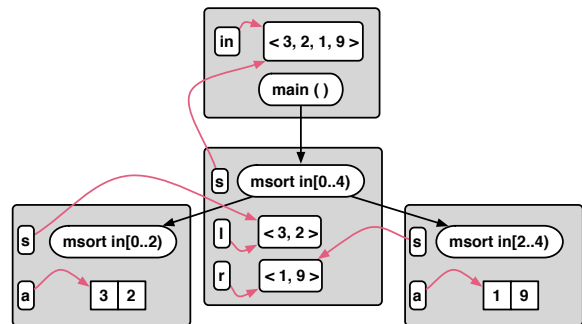


Figure 2. Hierarchical heap example for msort.

Consider the parallel merge sort in Figure 1. The implementation uses an immutable sequence data structure provided by a module `Seq`, whose details we omit. To sort an input sequence, the function `msort` first checks its length. If the length is less than some constant `GRAIN`, then the function uses an imperative in-place sequential quicksort to sort the input. Otherwise, the input is split into two halves and two recursive calls are performed in parallel.

Parallelism is exposed by the programmer through the `par` construct, which creates new *tasks*. Initially, there is only one task (corresponding to the execution of the entire program); thus all user code implicitly runs within the context of some task. With `par`, a task may spawn two new tasks. This establishes a parent/child nesting relationship where a parent task must wait until both its children complete before continuing its execution. Tasks are managed by a *scheduler*, which strives to minimize the completion time of the program by mapping tasks to processors.

As in Raghunathan et al. [31], to support parallel automatic memory management, each task is assigned its own *heap* in which it allocates new data. Heaps are organized into a *hierarchy* (tree) with the same parent/child relationships as their associated tasks. When both children of a task complete, their heaps are logically merged with the parent heap, allowing the parent to continue with child data stored locally as though the children had never existed.

Figure 2 illustrates an example where `msort` is called with the input sequence `(3, 2, 1, 9)` and `GRAIN = 2`. Tasks are drawn

as ellipses connected by straight, black arrows pointing from parent to child. The grey boxes drawn around each task delimit its heap, and red curved arrows show pointers in memory. At the root of the hierarchy is the initial task `main` which allocated the input sequence and then called `msort`. The middle task corresponds to this initial call of `msort`, which split the input into two sequences `l` and `r`, allocated locally. The two leaf tasks are the parallel recursive calls of `msort`, which have copied their inputs to local arrays `a`.

Consider the following definition. We say that the hierarchy is *disentangled* if, for any pointer from an object x in heap h_x to another object y in heap h_y , h_y is either equal to or an ancestor of h_x . In other words, in order to be disentangled, the hierarchy must not contain *down-pointers* that point from ancestor to descendant, nor may it contain *cross-pointers* that point between unrelated heaps (such as between two siblings). In a disentangled hierarchy, the lack of pointers into leaf heaps allows us to garbage collect the leaves independently, without synchronizing with other tasks, and in parallel with other leaf-heap collections. More generally, any two disjoint sub-trees of heaps may be collected independently in parallel (although the tasks within a sub-tree must cooperate). The execution of `msort` in Figure 2 is disentangled, thus the two leaves of the hierarchy could both be independently collected in parallel; for example, collecting the array `a` after computing `Seq.fromArray a` but before joining with the parent task. It can be shown that all purely functional programs naturally exhibit disentanglement [31].

In the presence of arbitrary mutation, disentanglement is not guaranteed. Consider for example a mutable reference `r` which is allocated at a parent task and then passed to two child tasks. One child could update `r` to point to a locally allocated object, creating a down-pointer. The second child could then read from `r` to create a cross-pointer.

To enforce disentanglement, we propose a *promotion* technique. The basic idea behind promotion is to detect when a down-pointer would be created in the hierarchy, and first promote (copy) the lower object upwards so that it lies in the same heap as the mutable object.

There are a number of challenges associated with promotion. In particular, promotion duplicates objects, which complicates the identity of mutable objects. We solve this issue by distinguishing one copy as the *master copy*, to which all accesses are redirected through *forwarding pointers*. Additionally, when data being promoted contains pointers to other objects, all transitively reachable data might need to be promoted. This introduces concurrency into the system even when none exists in the user code, since a task might access an object which is in scope of an in-progress promotion.

In our implementation, we prioritize the efficiency of updates to local objects. This facilitates an important idiom of practical parallelism where the overhead of parallelism is amortized by switching to a fast sequential algorithm on small inputs. As the fastest sequential algorithms on modern

```

1 type task
2 type data
3 type objptr
4 type field
5 type thunk = unit -> objptr
6 function forkjoin: thunk * thunk -> objptr * objptr
7 function alloc: field list -> objptr
8 function readImmutable: objptr * field -> data
9 function readMutable: objptr * field -> data
10 function writeNonPtr: objptr * field * data -> unit
11 function writePtr: objptr * field * objptr -> unit

```

Figure 3. High-Level Operations.

hardware are often imperative, local updates are thus crucial for efficiency. The `msort` program exemplifies this idiom, utilizing a fast imperative quicksort on small inputs. Indeed, in our results (Section 4.4), we see that `msort` can be up to twice as fast as a purely functional alternative.

3 Hierarchical Heaps for Mutable Data

3.1 Programming Model

Our proposal deals with the full ML language, including mutable data, extended with nested task parallelism. We extended an existing ML compiler to reduce this programming model to a small set of high-level operations that we implemented in the runtime system. In what follows, we describe the interface of these high-level operations in abstract terms and explain the challenges we face when implementing them in the context of hierarchical heaps.

High-Level Types. Figure 3 describes the types and operations we use to implement nested-parallel ML programs. While such operations are typically programmed in C or some other low-level language, here for the sake of simplicity we specify them in ML-like pseudocode.

Tasks compute on *data*. The exact definition of data does not matter at this level of description; we can assume that it consists in machine integers, floating-point numbers, pointers, etc. Amongst general data, we distinguish the type `objptr` of *object pointers*, that is, of pointers to allocated objects, or simply *objects*. Objects correspond to ML data types allocated during the course of execution, such as cons cells. For our purposes, an object consists of a finite list of *fields* storing its content. Objects are used for communication between tasks. The body of a task is a *thunk*, that is, a function expecting no argument and returning its result as an object pointer.

High-Level Operations. The types described above are used through six high-level operations. One of them deals with task management, the others with objects.

The compiler elaborates the `par` keyword into calls to the `forkjoin` operation. This operation takes two thunks,

```

1 function currentTask: unit -> task
2 function getField: objptr -> field -> ptr
3 function fields: objptr -> field list
4 function ptrFields: objptr -> field list
5 function nonptrFields: objptr -> field list
6 function fwdPtr: objptr -> ptr
7 function hasFwdPtr: objptr -> bool
8 type heap
9 function heapOfTask: task -> heap
10 function newChildHeap: heap -> heap
11 function joinHeap: heap * heap -> unit
12 function depth: heap -> int
13 function freshObj: heap * field list -> objptr
14 function heapOf: objptr -> heap
15 function lock: heap -> { READ, WRITE } -> unit
16 function unlock: heap -> unit

```

Figure 4. Low-level Primitives.

creates one task for each, and runs both tasks to completion, in parallel. This establishes a nesting relationship: we say that the task calling `forkJoin` is the *parent task*, while the tasks created by the operation are the *children tasks*. The parent task is suspended until both children return, at which point it receives their results and resumes.

Allocation arises from a variety of ML features: constructors of algebraic data types, explicit initialization of references and arrays, closures, etc. The `alloc` operation allocates a new object. The caller describes the list of fields of the new object and receives a pointer to it as a result.

The remaining primitives deal with reading from and writing to objects. Since the type system of ML distinguishes between mutable and immutable data, we have several reading and writing operations, depending on the type of the data being read or written. Later in this section, we will exploit these distinctions for efficiency purposes.

First, we distinguish between reading immutable data with the `readImmutable` operation and reading mutable data with the `readMutable` operation. Both operations take a pair of an object pointer and a field descriptor, and read the data held in the corresponding field of the object. If the field does not exist, the result is undefined.

Second, when writing to mutable data, we distinguish between writing non-pointer data (e.g., machine integers) and writing object pointers. The former takes the object and field to be written to, as well as the data to write. The latter takes the same argument, except that the data to write must be an object pointer. Note that, in the type signature of `writeNonptr`, we specify the general data type, which includes pointers; however, the behavior of the function is undefined if it is actually passed an object pointer.

3.2 Low-Level Primitives.

Our algorithms implement the interface of Figure 3 using hierarchical heaps. To do so, they rely on a number of low-level operations (Figure 4) provided by the runtime system.

Task-Related Operations. Calling `currentTask` retrieves the currently running task. Since user code is always implicitly running within some task, this never fails.

Memory-Related Operations. The `getField` operation returns a pointer to the specified field in an object pointer. This pointer does not point to (the beginning of) an object, but rather inside it. For this reason, the result type of `getField` is the abstract type `ptr`, rather than `objptr`.

Calling `fields(o)` operation returns all the fields of the memory object. Its variants `ptrFields` and `nonptrFields` return respectively the fields which hold object pointers and those which do not.

Finally, each object comes equipped with a special field for storing a *forwarding pointer*. Forwarding pointers are a classic ingredient used in copying collectors [22] to perform bookkeeping during collection. An object may or may not have a valid forwarding pointer in its dedicated field. This can be tested using the `hasFwdPtr` operation. The address of the field can be obtained by calling `fwdPtr`.

Heap-Related Operations. The first group of operations on heaps (abstract type `heap`) manages their relationships with tasks as well as with each other. Calling `heapOfTask(t)` returns the heap associated with task `t`. Like tasks, heaps are organized into a hierarchy which grows and shrinks through the `newChildHeap` and `joinHeap` primitives. Calling `depth(h)` returns the depth of a heap `h` in the hierarchy: the root is at depth zero, its children at depth one, etc.

The operation `freshObj` allocates a new object in the specified heap and with the specified fields. Correspondingly, the heap in which an object was allocated can be retrieved by the `heapOf` operation.

Finally, in order to deal with concurrency issues, every heap comes equipped with a readers-writers lock [20]. Such locks can be held in reading mode by several threads simultaneously, but by a single thread in writing mode (excluding any other readers or writers). The lock associated with a heap `h` can be acquired by calling `lock(h, m)`, with `m` being the `READ` or `WRITE` mode, and then released by calling `unlock(h)`.

3.3 Implementation of the High-Level Primitives

We can now explain our algorithms as implementations of the high-level operations in terms of the low-level ones, starting with an overview of the challenges involved.

Challenges. Our goal is to ensure that disentanglement of the hierarchy holds in the presence of calls to `writePtr(obj, field, ptr)`. If implemented naively

```

1 function forkjoin (f, g) =
2   heap ← heapOfTask(currentTask ())
3   heap_f ← newChildHeap(heap)
4   heap_g ← newChildHeap(heap)
5   (r_f, r_g) ← run t ∈ {f,g} in heapt and wait
6   joinHeap(heap, heap_f); joinHeap(heap, heap_g)
7   return (r_f, r_g)
    
```

Figure 5. Fork/join.

as $*getField(obj, field) \leftarrow ptr$, such writes would create down-pointers when $heapOf(obj)$ is an ancestor of $heapOf(ptr)$. In order to enforce disentanglement in this situation, our implementation of `writePtr` promotes (copies) the object at `ptr`, as well as all objects reachable from it, into $heapOf(obj)$. The address of the copy can then be written into $*getField(obj, field)$.

In the presence of repeated writes of `ptr` to objects held in heaps of decreasing depth, the object at `ptr` might be promoted several times. Ultimately, all copies but one should be eliminated, and pointers to them updated to point to the remaining one. Thus, we need a way to link an object with its copies. We use forwarding pointers to do so, and organize all the existing copies of an object into a singly-linked list.

All the copies of an object are equivalent as far as immutable fields are concerned, since by definition their content cannot change. In contrast, reads and writes of mutable fields cannot treat all copies as equivalent, for this could lead to lost updates (if an updated copy is then eliminated) or inconsistent reads (if a task updates one copy and then reads another one). Thus, we have to perform all mutable accesses on a unique, authoritative copy of an object, that we call its *master copy*. Given that all copies are arranged into a linked-list, we choose to take the last element of that list, that is the one in the shallowest heap, as the master copy.

A further difficulty is that several tasks might be trying to copy the same data, and thus update the relevant forwarding pointers, simultaneously. A task might also be updating forwarding pointers while another one is traversing them to find the master copy. Our algorithms avoid synchronization issues by acquiring the locks of the heaps they traverse during mutable accesses, from deepest to shallowest. Additionally, we propose several fast-paths avoiding locking in common cases.

Fork-join. Figure 5 shows a naïve implementation of the fork/join operation. First, create a heap attached to the heap of the running task for each child task. Then, run in parallel each t within its heap $heap_t$, for $t \in \{f, g\}$, and wait for them to complete (the exact realization of these operations depending on the scheduler at hand). Finally, join both child heap with their parent, and pass the results returned by each task to the caller. Joining heaps can be done without physically copying data.

```

1 function alloc (fields) =
2   return freshObj(heapOfTask(currentTask()), fields)
3 function readImmutable (obj, field) =
4   return *getField(obj, field)
5 function findMaster (obj) =
6   while true:
7     while hasFwdPtr(obj): obj ← *fwdPtr(obj)
8     lock(heapOf(obj), READ)
9     if not hasFwdPtr(obj): return obj
10    unlock(heapOf(obj))
11 function readMutable (obj, field) =
12   res ← *getField(obj, field)
13   if not hasFwdPtr(obj): return res
14   obj ← findMaster(obj)
15   res ← *getField(obj, field)
16   unlock(heapOf(obj))
17   return res
18 function writeNonptr (obj, field, val) =
19   *getField(obj, field) ← val
20   if not hasFwdPtr(obj): return
21   obj ← findMaster(obj)
22   *getField(obj, field) ← val
23   unlock(heapOf(obj))
    
```

Figure 6. Allocation, reads, non-pointer writes.

Allocation. Our implementation of `alloc` allocates the new object in the heap of the currently running task (l. 2).

Reading Immutable Data. ML programs read immutable data when destructuring values such as tuples or lists, e.g. through projections or pattern matching. Since these are very common operations, it is important to support them efficiently. Fortunately, all the potential copies of the same object hold the same value in their immutable fields, by definition. Thus, `readImmutable(obj, field)` does not care about the forwarding pointer slot of `obj` and can always access the contents of `field` without any indirection (l. 4).

Finding Master Copies. After several promotions occur, objects may exist in multiple copies linked in a chain by their forwarding pointers. This chain has to be taken into account when accessing mutable fields: intuitively, only the last copy, called the *master copy*, holds up to date information.

The function `findMaster(obj)` returns a pointer to the master copy of `obj`. Intuitively, it simply has to walk the chain of forwarding pointers, starting from `obj`. However, while doing so it might encounter a forwarding pointer installed by a promotion that is still ongoing. In that case, we should wait for the promotion to complete. We do so by acquiring the lock of the heap to which the copy belongs. Since we acquire the lock in shared mode, we do not block concurrent calls to `findMaster`. In contrast, promotion always locks in

```

1 function writePtr (obj, field, ptr) =
2   if heapOf(obj) = heapOfTask(currentTask())
3     and not hasFwdPtr(obj):
4     *getField(obj, field) ← ptr
5     return
6   obj ← findMaster(obj)
7   if depth(heapOf(obj)) ≥ depth(heapOf(ptr)):
8     *getField(obj, field) ← ptr
9     unlock(heap(obj))
10    return
11  unlock(heap(obj))
12  writePromote(obj, field, ptr)
13 function writePromote (obj, field, ptr) =
14   assert (depth(heapOf(obj)) < depth(heapOf(ptr)))
15   prev ← ptr
16   lock(heapOf(prev), WRITE)
17   while true:
18     for h from heapOf(prev) excluded
19       up to heapOf(obj) included: lock(h, WRITE)
20     if not hasFwdPtr(obj): break
21     else:
22       prev ← obj
23       obj ← *fwdPtr(obj)
24   promotedPtr ← promote(heapOf(obj), ptr)
25   *getField(obj, field) ← promotedPtr
26   for h from heapOf(obj) included
27     down to heapOf(ptr) included: unlock(h)
28 function promote (heap, obj) =
29   if depth(heapOf(obj)) ≤ depth(heap): return obj
30   if hasFwdPtr(obj):
31     return promote(heap, *fwdPtr(obj))
32   newObj ← freshObj(heap, sizeof(obj))
33   *fwdPtr(obj) = newObj
34   for field in nonptrFields(obj):
35     *getField(newObj, field) ←
36       *getField(obj, field)
37   for field in ptrFields(obj):
38     *getField(newObj, field) ←
39       promote(heap, *getField(obj, field))
40   return newObj

```

Figure 7. Pointer writes and promotion.

exclusive mode the heaps where it installs new forwarding pointers, ensuring mutual exclusion.

Our implementation of `findMaster` uses the classic double-checked locking pattern to reduce the cost of locking. As long as we observe forwarding pointers, we move up, without locking (l. 7). Once we see an object that is a candidate for being the master copy, we acquire the lock of its heap, and check whether the object has acquired a forwarding pointer in the meantime; if not, it is definitely the master copy, and can be returned (l. 8-9). Otherwise, we unlock the heap and

start walking the chain again (l. 10). Note that it is the caller’s responsibility to unlock the heap.

Reading Mutable Data. Mutators can read a mutable field in an object by calling `readMutable(obj, field)`. It would be correct to simply acquire the master copy, read the field, and release the lock (l. 14-17). We add a fast path: read the mutable field optimistically, then check for the absence of a forwarding pointer (l. 12-13). This way, accessing mutable fields in objects without copies only takes a couple of machine instructions.

Writing Non-pointer Data. Writing plain data such as integers or floating-point numbers cannot involve promotion, and thus is always relatively cheap. The implementation of `writeNonptr` mimicks that of `readMutable`. In the fast path, we optimistically write the value `val` in `field`, and then check whether `obj` was the master copy (l. 19-20). Otherwise, we find the master copy, write the value, and unlock (l. 21-23).

Writing Pointer Data. Writing pointer data is the most difficult case, as it might trigger a promotion. The code for `writePtr(obj, field, ptr)`, which attempts to write `ptr` to `field` in `obj`, is given in Figure 7. We may have to promote `ptr` to the heap of `obj` if writing it directly would result in entanglement. The algorithm can be decomposed into three cases: a fast path, non-promoting writes, and promoting writes. Let us describe each of these in turn.

The fast path of `writePtr` (l. 2-5) writes `ptr` into `obj` only if the latter has no forwarding pointer and is in the heap of the currently running task. Since this heap is necessarily a leaf in the hierarchy, promotion is never needed in this case. This qualifies as a fast path since testing whether an object pointer belongs to the currently running task can be implemented much more efficiently than computing the depth of an arbitrary heap.

On the slow path, we acquire the master copy of `obj` and obtain the depth of its heap. If it is deeper in the hierarchy than that of `ptr`, we are not creating entanglement, and can simply perform the write and unlock (l. 6-10). Otherwise, we have to promote, and thus unlock and call the dedicated function `writePromote(obj, field, ptr)` (l. 11-12).

Promoting Writes. We proceed in three phases. First, we lock in exclusive mode all the heaps on the path from `ptr` to the master copy of `obj` from the bottom up (l. 15-23). Then, we promote `obj` and write the address of the promoted copy to `field` (l. 24-25). Finally, we unlock the path from top to bottom (l. 26-27).

Acquiring the locks serve several purposes. Let us call h_1 the heap that contains `obj` and h_2 the heap containing the master copy of `obj`. By acquiring the locks on the path from h_1 to h_2 excluded, we take ownership of the forwarding pointer of any object we might need to promote. By acquiring the lock of h_2 , we ensure that no concurrent call to `findMaster` will return before we have finished promoting.

	Read		Write			
	Immutable	Mutable	Non-pointer	Non-promoting	Promoting	
Local	✓	✓✓	✓✓	✓✓		✓ single instruction
Distant	✓	✓✓	✓✓	~	≈	✓✓ few instructions
Promoted	✓	~	~	~	≈	~ single-heap locking
						≈ path locking + copying

Figure 8. Costs of memory operations.

The recursive function `promote` returns a pointer to a copy of `obj` held in `heap` or one of its parents. If `obj` already resides in `heap` or above, it can simply be returned (l. 29); if `obj` has a forwarding pointer, we follow it (l. 30-31). Taken together, these two tests ensure that objects with a chain of forwarding pointers leading above `heap` are not copied. If both fail, we have to introduce a new copy. We allocate a copy `newObj` of `obj` in `heap` (l. 38), set the forwarding pointer slot of `obj` to point to `newObj` (l. 32-33), and copy non-pointer fields (l. 34-36). We recursively promote pointer fields, since they might point strictly below `heap` (l. 37-39). At this point we can return `newObj` (l. 40) since it belongs to `heap`, as do every object reachable from it.

Note that while we have expressed `promote` as a recursive function for simplicity, it can be implemented using a work list. In particular, we were careful to set the forwarding pointer of `obj` before the recursive calls.

Cost Summary. Figure 8 summarizes the costs of each memory operation in various situations. Columns correspond to distinct memory operations, distinguishing between non-promoting and promoting pointer writes. Rows classify objects being read from or written to: *promoted* objects are those with forwarding pointers; *local* and *distant* objects have no forwarding pointers and belong either to the heap of the task performing the memory operation (local objects) or to one of its ancestors (distant objects).

3.4 Promotion-Aware Copy Collection

Promotion introduces redundant copies of objects. However, it is not difficult to eliminate these copies by piggybacking on the classic semispace garbage collection algorithm. We briefly explain our proposal; precise details are available in Appendix A.

Assume that we are collecting a sub-tree starting at some heap h ; each heap below h (included) thus acquires a to-space. When examining an object in from-space, our collector traverses its forwarding pointer chain, considering several possibilities in turn.

1. If the chain leads into a to-space, it points to a copy introduced during collection.
2. If the chain leads into a from-space that is strictly above h in the hierarchy, it leads to a copy introduced during promotion.

3. If the chain ends at an object that has no forwarding pointer, then this object is in a from-space below h .

In the first two cases, the address of the copy can be reused directly, whereas in the last case we have to introduce a new copy. The second case corresponds to the elimination of duplicates introduced during promotion. Note that since we do not attempt to access copies outside of the collection zone, no additional locking is needed. In the third case, we copy the last element of the forwarding chain into the to-space and update its forwarding pointer to point to this new copy. By doing so, we effectively make sure that all pointers to the object or its promoted copies will point to the new copy created in the to-space.

4 Implementation and Experiments

We implemented our techniques by building upon the parallel MLton compiler developed in prior work [31], from which we inherit the hierarchical heaps infrastructure, garbage collection policy, and scheduler. We extended this compiler with support for general mutation by closely following the algorithms described in Section 3. Further details on the implementation can be found in Appendix B.

We evaluate our techniques by considering a number of benchmarks compiled with several compilers for dialects of parallel ML. Thanks to the shared ML language, our benchmarks remain mostly identical across compilers except for minor compatibility edits. Our benchmark suite builds on previous suites from parallel functional languages [7, 16], and extends them, also to include imperative programs which use mutable data. Our benchmarks use several standard implementations of data types such as sequences and graphs. Unless stated otherwise, the elements of the sequences are 64-bit numeric types (integers or floating point) generated randomly with a hash function.

Experimental Setup. For the measurements, we use a 72-core (4 x 18 core Intel E7-8867 v4) Dell PowerEdge Server with 1 Terabyte of memory. For the sequential baselines, we use the whole-program optimizing MLton compiler [28], which we label `mlton`. We compare all of our benchmarks to the work of Blelloch, Spoonhower, and Harper [33], labeled `mlton-spoonhower`, which extends the MLton compiler to support nested fork-join parallelism and parallel allocation, but utilizes sequential, stop-the-world collection. For purely functional benchmarks, we also compare with the

Manticore compiler [7, 16], labeled `manticore`, which provides for parallel functional programming by using syntax similar to ours. We refer to our hierarchical-heaps compiler as `mlton-parmem`.

When taking timing measurements, we exclude initialization times. All reported timings include GC times and are reported as the median of five runs.

4.1 Pure Benchmarks

These benchmarks are purely functional, meaning that their source code does not use mutation.

fib. This benchmark computes the 42nd Fibonacci number via the naïve recursive formula $F(n) = F(n - 1) + F(n - 2)$, with a sequential threshold of $n = 25$.

tabulate, map, reduce, filter. These benchmarks each begin by generating an input sequence of size 10^8 . The `tabulate` benchmark completes once the input sequence is built. The `map` benchmark constructs a second sequence by applying a simple function to each element. The `reduce` benchmark sums the elements of the input sequence. The `filter` benchmark constructs a second sequence containing only the elements which satisfy a given predicate. They are each implemented with straightforward divide-and-conquer approaches, with a sequential threshold of 10^4 elements.

msort-pure. This benchmark first tabulates a sequence of size 10^7 . It then sorts the sequence with a function similar to the one shown in Figure 1, except that it uses a purely functional quick-sort below a sequential threshold of 10^4 instead of the imperative one.

dmm, smvm. These benchmarks operate on square matrices of size $n \times n$. Each matrix is represented by a sequence of rows (or columns). The `dmm` benchmark multiplies two dense matrices with the naïve $O(n^3)$ algorithm, where each of the n rows is implemented as another sequence of size n . The `smvm` benchmark multiplies a sparse matrix by a dense vector, where each row of the sparse matrix contains only the non-zero entries represented as index-value pairs. In `dmm`, $n = 600$. In `smvm`, $n = 20,000$ and each row has approximately 2,000 non-zero entries. The sequential threshold is one matrix row.

strassen. This benchmark multiplies two dense square matrices of size $n \times n$ using Strassen’s algorithm. The matrices are represented by quadrees with leaves of vectors of elements of size 64×64 . In our experiments, $n = 1024$ with a sequential threshold of 64 (that is, the leaves of the quadree are each processed sequentially).

raytracer. This benchmark is adapted from the raytracer benchmark written for the Manticore language [7], which was adapted from an `ld` program [30]. It renders a $600\text{px} \times 600\text{px}$ scene in parallel by tabulating a sequence of pixels with a sequential granularity of 300 pixels.

4.2 Imperative benchmarks.

These benchmarks are designed to exercise various different forms of mutation, as summarized in Figure 9. Due to mutation, they are not implementable in Manticore.

msort, dedup. These benchmarks begin by tabulating a sequence of size 10^7 before sorting it with a technique similar to that shown in Figure 1. The `msort` benchmark uses imperative quick-sort below the sequential threshold of 10^4 elements. The `dedup` benchmark is similar to `msort` but removes duplicate keys. Below the sequential threshold, this is accomplished by imperatively inserting elements into a hash set before sorting with the in-place quick-sort. For `dedup`, we guarantee the sequence has approximately 10^6 unique keys.

tourney. This benchmark tabulates a sequence of 10^8 contestants, and then computes a tournament tree. Each contestant is represented by an integer which measures their fitness. In the tournament tree, each contestant c has an associated parent pointer which points to the contestant that eliminated c from the tournament. This benchmark computes the tournament tree with a simple divide-and-conquer approach, using mutation at each join point in order to set a parent pointer.

usp, usp-tree, multi-usp-tree, reachability. These benchmarks consider variants of parallel breadth-first search (BFS) on directed, unweighted graphs. BFS visits vertices in rounds. At round r , BFS visits (in parallel) every vertex which has not previously been visited and is reachable by r hops from the source vertex. When a vertex is visited, a piece of mutable data associated with it is updated.

The BFS variants differ in the types of per-vertex mutable data. They also differ in the number of times a vertex is visited. Except in the `reachability` benchmark, we guarantee that each vertex is visited exactly once by marking vertices as visited with an atomic “compare-and-swap” operation. In the `reachability` benchmark, we check and update the visited status of vertices simply by reading and writing to a shared flag. This creates a data race and potentially causes some vertices to be visited multiple times (up to at most P visitations, $P = \text{number of processors}$). In practice, this variant of BFS often performs better on modern hardware because (a) atomic operations such as compare-and-swap are expensive, and (b) observing the data race within a particular execution is rare.

The specifics of these benchmarks are described below.

- `reachability` identifies which vertices are reachable from the source.
- `usp` computes the unweighted single-source shortest path length of all vertices. Every time a vertex is visited, the algorithm records the current round number as the distance to that vertex.
- `usp-tree` computes all unweighted single-source shortest paths. It is implemented with an array A of ancestor

Benchmark	Representative Operation
pure benchmarks	immutable reads
msort	local non-pointer writes
dedup	local non-pointer writes
tourney	local non-promoting writes
reachability	distant non-pointer writes
usp	distant non-pointer writes
usp-tree	distant promoting writes
multi-usp-tree	distant promoting writes

Figure 9. Representative operations of all benchmarks.

lists. When a vertex v is visited along an edge (u, v) , the ancestors of v are recorded as $A[v] := u :: A[u]$.

- **multi-usp-tree** runs 36 copies of `usp-tree` in parallel.

The input graph is the *orkut* social network graph [1], which has approximately 3 million vertices, 117 million edges, and a diameter of 9. Each benchmark begins by converting the input graph into a compact adjacency-sequence format suitable for parallel BFS.

4.3 Representative Operations

In Figure 9, using the terminology from Figure 8, we characterize each benchmark by a *representative memory operation*. Representative operations summarize which type of operation is most likely to be a dominant cost in execution time.¹ These in turn help understand and predict performance. For example, if a benchmark exhibits mostly “local, non-pointer” writes, then that benchmark likely has low overhead; in contrast, if a benchmark is characterized by many “distant, promoting” writes, then it might have high overhead and scale poorly.

4.4 Results

We collect the following statistics for each benchmark.

- T_s is the sequential execution time.
- T_1 and T_{72} are execution times on 1 and 72 processors.
- The *overhead* is T_1/T_s .
- The *speedup* is T_s/T_{72} . In general, the speedup on P processors is given by T_s/T_P .
- GC_s is the percent of time spent in GC during a sequential run.
- GC_{72} is the percent of processor time spent in GC during a 72-processor run.
- M_s is the memory consumption of the sequential run.
- I_1 and I_{72} are memory inflations on 1 and 72 processors.

The *memory consumption* statistic is an upper bound on the amount of physical memory required to store heap-allocated objects; it is computed by tracking the maximum heap occupancy within one execution, and includes fragmentation due

¹Note that immutable reads are pervasive in all benchmarks.

to parallel allocations. *Memory inflation* gives the memory consumption as a factor relative to the sequential memory consumption, M_s .

For `mlton-spoonhower`, GC_{72} includes processor time spent blocked during a stop-the-world collection. We do not report GC statistics for Manticore, because it is not able to collect statistics only for a specific region of code, which we need to have a meaningful comparison. We are also unable to report results for Manticore on `msort-pure` due to a compiler bug.

Our results are summarized in Figures 10, 11, 12, and 13. Figures 10 and 11 list the execution times, overheads, speedups, and GC percentages of pure and imperative benchmarks, respectively. Figure 12 shows the speedup of `mlton-parmem` on various benchmarks for processor counts between 1 and 72. Finally, Figure 13 lists the memory consumptions and inflations of all benchmarks.

Overheads. Inspecting Figures 10 and 11, we can make several observations and conclusions. First, for both pure and imperative benchmarks, the overheads of `mlton-parmem` are generally comparable to those of `mlton-spoonhower`, which serves as a good baseline because it does not support parallel memory management. This shows that our techniques for maintaining a dynamic memory structure based on hierarchical heaps can be implemented efficiently. Second, we observe that for pure benchmarks (Figure 10), our overheads are within a factor 2 of the sequential baseline and are consistently smaller than those of `manticore`. Third, for imperative benchmarks (Figure 11), our overheads are higher than those of the pure benchmarks but still remain within a factor of approximately 2.6 in comparison to the sequential baseline. The increase in overheads is due in part to the memory operations which are no longer plain loads and stores (see Figures 8 and 9).

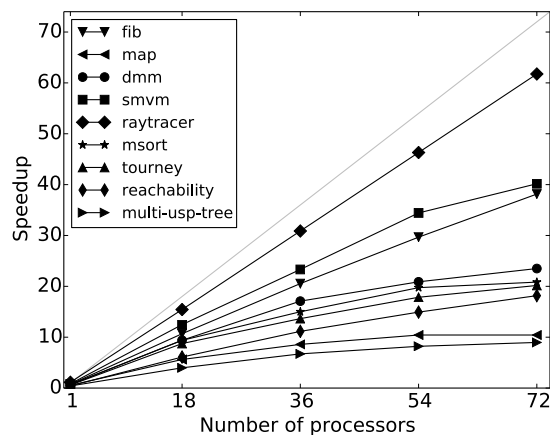
Speedups. Inspecting Figure 10, we observe that for pure benchmarks with 72 cores, `mlton-parmem` achieves speedups ranging between 10 and 62. Compared to `mlton-spoonhower`, our speedups are significantly higher, which is expected because `mlton-spoonhower` suffers from sequential, stop-the-world garbage collections. Our speedups are also significantly better than those of `manticore`, which is sometimes as low as 3x. This seems surprising because `manticore` is designed for purely functional programs. The reason is that `manticore` relies on imperative updates within the run-time to execute computations in parallel (e.g., to communicate the result of a remotely executed task to another processor) and employs a promotion technique to preserve local heap invariants [7]. To verify this, we measured that on the `map` benchmark with 72 cores, `manticore` promoted nearly 340MB of data in total, whereas `mlton-parmem` performed no promotions.

Inspecting the imperative benchmarks (Figure 11), we observe that `mlton-parmem` achieves good speedups on 72 processors (between 14 for high-overhead benchmarks and

	mlton		mlton-spoonhower					manticore				mlton-parmem (our compiler)				
	T_s	GC_s	T_1	Over head	T_{72}	Speed up	GC_{72}	T_1	Over head	T_{72}	Speed up	T_1	Over head	T_{72}	Speed up	GC_{72}
fib	2.67	0.0%	3.72	1.39	0.12	22.25	3.4%	5.19	1.94	0.12	22.25	3.63	1.36	0.07	38.14	0.0%
tabulate	1.11	39.4%	0.89	0.8	0.42	2.64	89.5%	1.92	1.73	0.12	9.25	1.62	1.46	0.07	15.86	15.4%
map	1.46	29.8%	1.31	0.9	0.48	3.04	82.2%	4.02	2.75	0.49	2.98	2.75	1.88	0.14	10.43	16.2%
reduce	1.13	40.5%	0.9	0.8	0.46	2.46	84.7%	1.93	1.71	0.13	8.69	1.43	1.27	0.09	12.56	11.7%
filter	3.62	11.9%	4.84	1.34	0.54	6.7	76.6%	6.15	1.7	0.49	7.39	5.75	1.59	0.18	20.11	12.0%
msort-pure	7.02	12.3%	8.52	1.21	1.79	3.92	75.8%	-	-	-	-	6.91	0.98	0.36	19.5	15.7%
dmm	3.76	15.2%	7.02	1.87	0.92	4.09	78.3%	8.3	2.21	0.19	19.79	5.83	1.55	0.16	23.5	6.8%
smvm	7.23	0.0%	9.93	1.37	0.2	36.15	0.0%	12.68	1.75	0.32	22.59	8.69	1.2	0.18	40.17	0.0%
strassen	2.54	1.6%	2.89	1.14	0.16	15.88	40.6%	4.36	1.72	0.12	21.17	2.94	1.16	0.12	21.17	7.9%
raytracer	7.41	1.3%	7.0	0.94	0.3	24.7	29.8%	6.97	0.94	0.17	43.59	6.52	0.88	0.12	61.75	0.5%

Figure 10. Execution times (in seconds), overheads, and speedups of purely functional benchmarks.

	mlton		mlton-spoonhower					mlton-parmem (our compiler)				
	T_s	GC_s	T_1	Over head	T_{72}	Speed up	GC_{72}	T_1	Over head	T_{72}	Speed up	GC_{72}
msort	3.75	3.4%	4.66	1.24	0.36	10.42	58.3%	5.33	1.42	0.18	20.83	7.7%
dedup	3.72	2.6%	4.05	1.09	0.32	11.63	52.1%	4.61	1.24	0.16	23.25	6.3%
tourney	4.64	6.3%	8.17	1.76	0.92	5.04	76.2%	7.86	1.69	0.23	20.17	7.1%
reachability	8.36	0.0%	21.59	2.58	0.52	16.08	0.0%	19.59	2.34	0.46	18.17	0.0%
usp	8.34	0.0%	23.38	2.8	0.61	13.67	0.0%	21.85	2.62	0.58	14.38	0.0%
usp-tree	8.63	0.0%	23.79	2.76	0.63	13.7	0.0%	22.3	2.58	7.93	1.09	0.0%
multi-usp-tree	100.25	3.0%	209.59	2.09	19.07	5.26	34.0%	245.6	2.45	11.18	8.97	8.2%

Figure 11. Execution times (in seconds), overheads, and speedups of imperative benchmarks.**Figure 12.** Speedups of mlton-parmem.

23 for those with low overhead) with a couple exceptions: the highly concurrent `usp-tree` and `multi-usp-tree` benchmarks. Poor performance on `usp-tree` and `multi-usp-tree` is expected, because these benchmarks exhibit close to pessimal cases for our techniques with frequent concurrent updates of shared pointer data. For example, every time a vertex is visited in `usp-tree`, one cell of a distant array (located at the root) is updated with a new list, triggering promotion from a

	mlton	mlton-spoonhower		mlton-parmem (our compiler)	
	M_s	I_1	I_{72}	I_1	I_{72}
fib	0.0	+0.0	+0.31	+0.0	+0.4
tabulate	3.48	0.23	0.34	0.97	1.25
map	1.6	1.01	1.42	3.07	5.07
reduce	0.8	1.01	2.02	3.15	4.95
filter	2.8	0.8	1.41	1.7	2.51
msort-pure	1.43	1.01	1.44	1.49	4.9
dmm	0.18	1.78	3.39	0.94	7.39
smvm	1.92	1.33	1.47	2.67	2.78
strassen	0.22	0.95	4.41	1.86	17.68
raytracer	0.13	1.15	3.77	1.46	5.38
msort	1.09	0.78	1.5	1.6	4.04
dedup	0.56	1.18	2.61	1.63	6.48
tourney	0.8	6.24	8.34	6.91	8.66
reachability	3.87	1.01	2.29	1.5	2.15
usp	3.87	1.01	2.35	1.5	2.19
usp-tree	3.97	1.01	2.25	1.55	2.21
multi-usp-tree	19.76	1.2	2.85	1.64	1.79

Figure 13. Memory consumption (in GB) and inflations.

leaf heap to the root heap. Since promotions require locking entire heaps, they can sequentialize otherwise parallel visitations, leading effectively to complete serialization of the

entire computation. However, when multiple `usp-tree` computations are performed in parallel in the `multi-usp-tree` benchmark, some promotions remain independent and execute in parallel because the updated array is not always at the root heap. We indeed see a 9-fold speedup in this benchmark.

In Figure 12, we observe that as the number of processors increases, the speedup of all benchmarks continues to increase. That is, there are no inversions. For multiple benchmarks the speedup improves nearly linearly, suggesting the possibility of further scalability to higher core counts.

Garbage Collection. Inspecting Figures 10 and 11, we observe that `mlton-parmem` only loses at most approximately 16% of its time to garbage collection on runs with 72 cores. As expected, `mlton-spoonhower` performs poorly due to its sequential GC. Some benchmarks (`smvm`, `reachability`, `usp`, `usp-tree`) spend no time in garbage collection, regardless of when run sequentially or in parallel. This is due to the fact that these benchmarks allocate memory in an already large heap, which was grown to accommodate the input. (In the case of `smvm`, the benchmark does not include input generation; for the graph algorithms, the benchmark does not include the time it takes to read the graph from disk into a large string in the heap.)

Memory Consumption. Inspecting Figure 13, with only a few exceptions, our compiler on 72 cores consumes at most 7x more memory than the sequential baseline. Note that in general, any P -processor execution scheduled via work-stealing can expect to see inflation of up to a factor P [10, 11]. Thus, using `mlton-spoonhower` as an alternative baseline helps determine how much inflation is due simply to parallel execution, versus how much is due to our techniques. Indeed, our inflation with respect to `mlton-spoonhower` is generally lower, staying consistently within a factor of approximately 4. Our implementation introduces additional inflation in part through (a) the need for a separate forwarding pointer on every object, and (b) greater fragmentation of allocation to distinguish heaps within the hierarchy.

5 Discussion and Future Work

The promotion techniques presented in this paper rely on coarse-grained locks to manage concurrent manipulations of overlapping data. This approach prevents certain promotions from proceeding in parallel, even when those promotions would otherwise be independent. For example, in the `usp-tree` benchmark, every visitation of a vertex triggers a promotion to the root of hierarchy, causing a serialization of visitations. However none of these promotions overlap, so they ought to be able to proceed in parallel. In future work, we intend to design a more fine-grained promotion strategy that would permit parallel promotions to the same heap.

With respect to garbage collection, our current implementation has two limitations: first, it can only collect leaf heaps;

second, each such collection is sequential. Parallelism is thus achieved by collecting many leaves independently. Our results suggest that this simple approach can perform well for highly parallel applications. However, if a collection takes place at the root heap (when there is no parallelism), such a collection would be sequential and effectively stop-the-world. More generally, when there is little parallelism, large collections can take place sequentially. In future work, we plan to complete the implementation by adding support for parallel collection of individual heaps, and in general parallel collection of sub-trees of heaps (not just leaves).

6 Related Work

With the proliferation of shared memory parallel computers using modern multicore processors, there has been significant work on the design and implementation of high-level programming languages for writing parallel programs [9, 12, 16, 17, 21, 23, 25, 26, 32]. For implicit memory allocation and reclamation with garbage collection, there are numerous techniques for incorporating parallelism, concurrency, and real-time features. Jones et al. [22] provides an excellent survey.

We contrast our work with a number of systems [4, 7, 13–15, 27, 31, 32] that use processor- or thread-local heaps which service (most) allocations of the executing computation and can be collected independently combined with a shared global heap that must be collected cooperatively.

The Doligez-Leroy-Gonthier (DLG) parallel collector [13, 14] employs this design, with the invariant that there are no pointers from the shared global heap into any processor-local heap and no pointers from one processor local-heap into another processor-local heap. To maintain this invariant, all mutable objects are allocated in the shared global heap and (transitively reachable) data is promoted (copied) from a processor-local heap to the shared global heap when updating a mutable object, which increases the cost of all mutable allocations and updates. In our approach, mutable objects are allocated in the thread-local heap and updates to such never need to promote data, making this common case significantly less expensive than in DLG. Moreover, in DLG, scheduling and communication actions, such as migrating a language-level thread from one processor to another or returning the result from a child task to a parent task, typically employ mutable objects and require promotions. With hierarchical heaps, the local heap is associated with the task, rather than the processor, and returning the result of a child task to the parent task is accomplished without copying.

Anderson [4] describes TCG, a variant of the DLG collector for a language with implicit parallelism serviced by a fixed number of worker threads pinned to processors. TCG allows mutable objects to be allocated in the processor-local heap and a processor-local collection copies live data to the global shared heap. When updating a mutable object

in the shared global heap with a processor-local pointer, a processor-local garbage collection is triggered, which copies the to-be-written object (and all other processor-local live data) to the shared global heap; updating a mutable object in the processor-local heap proceeds without a collection. Using collection to over-approximate promotion would not work well with our hierarchical heaps, because the collection would need to be triggered for *all* descendent heaps of the mutable object being written and could not be performed by the writing processor independently.

The Manticore garbage collector [7] is another variant of the DLG design, where the Appel semi-generational collector [5] is used for collection of the processor-local heaps. Although the high-level language is mutation-free, the implementation uses mutation to realize various parallel constructs and employs promotion to preserve the heap invariants. Recent work [24] has considered extending the Manticore language with mutable state via software transactional memory, but notes that promotions make a chronologically-ordered read set implemented as a mutable doubly-linked list inefficient. In contrast, such a data structure would be efficient in our hierarchical heaps, since a transaction’s read set is necessarily local to the thread executing the transaction.

The current Glasgow Haskell Compiler garbage collector [27] combines elements of the DLG and Domani et al. [15] collectors. Although Haskell is a pure language, there is significant mutation due to lazy evaluation. The collector allows mutable objects to be allocated in a dedicated portion of the processor-local heaps, which use a non-moving collector. The collector also allows pointers from the global heap to the processor-local heaps, mediated by proxy objects. When another processor accesses a proxy, it communicates with the owning processor to request that the object be promoted to the global shared heap. Proxy objects fit well with Haskell’s lazy evaluation, where all pointer accesses have a read barrier to check for unevaluated computations and where proxy objects can be incorporated into unevaluated computations without requiring promotion. Standard ML employs strict evaluation and eager promotion is a better fit.

The MultiMLton project [32] forked from `mlton-spoonhower` and shifted the domain to message-passing concurrency. While one could encode shared-state fork-join parallelism with their message-passing operations, the resulting overheads (emulating references with threads; eager thread creation vs lazy work-stealing) would not lead to a meaningful performance comparison. Their garbage collection strategy is tuned to their message-passing concurrency bias — starting from the DLG design and invariants, they avoid promotion through procrastination, blocking the writing thread (and executing another of the abundant concurrent threads) until a GC can be performed, which promotes the object and fixes references. They are concerned that promotion, leaving forwarding pointers that overwrite object data, would have unacceptable read-barrier

overhead, which also motivates their dynamic cleanliness analysis; in contrast, we have introduced a dedicated forwarding-pointer metadata field, which only requires a read barrier for mutable data.

Raghunathan et al. introduced hierarchical heaps [31] to mirror the hierarchy of tasks in a strict pure functional language with nested parallelism. They prove that the language enforces disentanglement, formulate a hierarchical garbage collection technique that allows independent heaps to be collected concurrently, and report the performance of an implementation in MLton. We extend this work to accommodate the mutable references and arrays of a strict impure functional language with nested parallelism. The hierarchical-heaps design is partly motivated by a desire to take advantage of the natural data locality of computation [2], which, with some care, could be preserved by thread schedulers (e.g., [3, 8]).

7 Conclusion

The high-level nature of functional programming languages makes them a good fit for parallel programming, but they require sophisticated memory managers which are challenging to get right in the joint presence of parallelism and mutation. In this paper, we showed how to provide efficient support for uses of mutation common in parallel programs by exploiting the hierarchical structure of functional computations. Our experiments suggest that these results could be an important step towards making functional programming a serious contender for performant parallel computing.

Acknowledgments

We thank Rohan Yadav for his assistance in the implementation of the benchmarks used in this paper. This material is based upon work supported by the National Science Foundation under Grant No. 1408940, Grant No. 1408981, and Grant No. 1629444. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author and do not necessarily reflect the views of the National Science Foundation. The first author was also partially supported by the German Research Council (DFG) under Grant No. ME14271/6-2.

References

- [1] [n. d.]. Stanford Large Network Dataset Collection. <http://snap.stanford.edu/>. ([n. d.]).
- [2] Umut A. Acar, Guy Blelloch, Matthew Fluet, and Stefan K. Mullerand Ram Raghunathan. 2015. Coupling Memory and Computation for Locality Management. In *Summit on Advances in Programming Languages (SNAPL)*.
- [3] Umut A. Acar, Guy E. Blelloch, and Robert D. Blumofe. 2002. The Data Locality of Work Stealing. *Theory of Computing Systems* 35, 3 (2002), 321–347.
- [4] Todd A. Anderson. 2010. Optimizations in a Private Nursery-Based Garbage Collector. In *9th International Symposium on Memory Management*, Jan Vitek and Doug Lea (Eds.). ACM Press, Toronto, Canada, 21–30.
- [5] Andrew W. Appel. 1989. Simple Generational Garbage Collection and Fast Allocation. *Software: Practice and Experience* 19, 2 (1989), 171–183.
- [6] Andrew W. Appel and Zhong Shao. 1996. Empirical and analytic study of stack versus heap cost for languages with closures. *Journal of Functional Programming* 6, 1 (Jan. 1996), 47–74. <https://doi.org/10.1017/S095679680000157X>
- [7] Sven Auhagen, Lars Bergstrom, Matthew Fluet, and John H. Reppy. 2011. Garbage collection for multicore NUMA machines. In *Proceedings of the 2011 ACM SIGPLAN workshop on Memory Systems Performance and Correctness (MSPC)*. 51–57.
- [8] Guy E. Blelloch, Jeremy T. Fineman, Phillip B. Gibbons, and Harsha Vardhan Simhadri. 2011. Scheduling irregular parallel computations on hierarchical caches. In *Proc. ACM Symposium on Parallelism in Algorithms and Architectures (SPAA)*. 355–366.
- [9] Guy E. Blelloch, Jonathan C. Hardwick, Jay Sipelstein, Marco Zagha, and Siddhartha Chatterjee. 1994. Implementation of a Portable Nested Data-Parallel Language. *J. Parallel Distrib. Comput.* 21, 1 (1994), 4–14.
- [10] Robert D. Blumofe and Charles E. Leiserson. 1998. Space-Efficient Scheduling of Multithreaded Computations. *SIAM J. Comput.* 27, 1 (1998), 202–229.
- [11] Robert D. Blumofe and Charles E. Leiserson. 1999. Scheduling multithreaded computations by work stealing. *J. ACM* 46 (Sept. 1999), 720–748. Issue 5.
- [12] Philippe Charles, Christian Grothoff, Vijay Saraswat, Christopher Donawa, Allan Kielstra, Kemal Ebcioglu, Christoph von Praun, and Vivek Sarkar. 2005. X10: an object-oriented approach to non-uniform cluster computing. In *Proceedings of the 20th annual ACM SIGPLAN conference on Object-oriented programming, systems, languages, and applications (OOPSLA '05)*. ACM, 519–538.
- [13] Damien Doligez and Georges Gonthier. 1994. Portable, Unobtrusive Garbage Collection for Multiprocessor Systems. In *21st Annual ACM Symposium on Principles of Programming Languages*. ACM Press, Portland, OR, 70–83. <https://doi.org/10.1145/174675.174673>
- [14] Damien Doligez and Xavier Leroy. 1993. A Concurrent Generational Garbage Collector for a Multi-Threaded Implementation of ML. In *20th Annual ACM Symposium on Principles of Programming Languages*. ACM Press, Charleston, SC, 113–123.
- [15] Tamar Domani, Elliot K. Kolodner, Ethan Lewis, Erez Petrank, and Dafna Sheinwald. 2002. Thread-Local Heaps for Java. In *3rd International Symposium on Memory Management (ACM SIGPLAN Notices 38(2 supplement))*, Hans-J. Boehm and David Detlefs (Eds.). ACM Press, Berlin, Germany, 76–87. <https://doi.org/10.1145/512429.512439>
- [16] Matthew Fluet, Mike Rainey, John Reppy, and Adam Shaw. 2011. Implicitly threaded parallelism in Manticore. *Journal of Functional Programming* 20, 5-6 (2011), 1–40.
- [17] golang [n. d.]. The Go Programming Language. <https://golang.org/>. ([n. d.]). Accessed: 2017.
- [18] Marcelo J. R. Gonçalves. 1995. *Cache Performance of Programs with Intensive Heap Allocation and Generational Garbage Collection*. Ph.D. Dissertation. Department of Computer Science, Princeton University.
- [19] Marcelo J. R. Gonçalves and Andrew W. Appel. 1995. Cache Performance of Fast-Allocating Programs. In *Conference on Functional Programming and Computer Architecture*. ACM Press, La Jolla, CA, 293–305. <https://doi.org/10.1145/224164.224219>
- [20] Maurice Herlihy and Nir Shavit. 2011. *The Art of Multiprocessor Programming*. Morgan Kaufmann.
- [21] Shams Mahmood Imam and Vivek Sarkar. 2014. Habanero-Java library: a Java 8 framework for multicore programming. In *2014 International Conference on Principles and Practices of Programming on the Java Platform Virtual Machines, Languages and Tools, PPPJ '14*. 75–86.
- [22] Richard Jones, Antony Hosking, and Eliot Moss. 2012. *The Garbage Collection Handbook: The Art of Automatic Memory Management*. Chapman & Hall.
- [23] Gabriele Keller, Manuel M.T. Chakravarty, Roman Leshchinskiy, Simon Peyton Jones, and Ben Lippmeier. 2010. Regular, shape-polymorphic, parallel arrays in Haskell. In *Proceedings of the 15th ACM SIGPLAN international conference on Functional programming (ICFP '10)*. 261–272.
- [24] Matthew Le and Matthew Fluet. 2015. Partial Aborts for Transactions via First-class Continuations. In *Proceedings of the 20th ACM SIGPLAN International Conference on Functional Programming (ICFP 2015)*. 230–242.
- [25] Doug Lea. 2000. A Java fork/join framework. In *Proceedings of the ACM 2000 conference on Java Grande (JAVA '00)*. 36–43.
- [26] Daan Leijen, Wolfram Schulte, and Sebastian Burckhardt. 2009. The design of a task parallel library. In *Proceedings of the 24th ACM SIGPLAN conference on Object Oriented Programming Systems Languages and Applications (OOPSLA '09)*. 227–242.
- [27] Simon Marlow and Simon L. Peyton Jones. 2011. Multicore Garbage Collection with Local Heaps. In *10th International Symposium on Memory Management*, Hans Boehm and David Bacon (Eds.). ACM Press, San Jose, CA, 21–32.
- [28] MLton [n. d.]. MLton web site. <http://www.mlton.org/>. ([n. d.]).
- [29] Ryan Newton. 2017. (2017). Indiana University. Personal Communication.
- [30] Rishiyur S. Nikhil. 1991. *Id Language Reference Manual*. (1991).
- [31] Ram Raghunathan, Stefan K. Muller, Umut A. Acar, and Guy Blelloch. 2016. Hierarchical Memory Management for Parallel Programs. In *ICFP 2016*. ACM Press.
- [32] K. C. Sivaramakrishnan, Lukasz Ziarek, and Suresh Jagannathan. 2014. MultiMLton: A multicore-aware runtime for standard ML. *Journal of Functional Programming* FirstView (6 2014), 1–62.
- [33] Daniel Spoonhower, Guy E. Blelloch, Robert Harper, and Phillip B. Gibbons. 2010. Space profiling for parallel functional programs. *Journal of Functional Programming* 20 (2010), 417–461. Issue Special Issue 5-6.
- [34] Philip Wadler. 1998. Why No One Uses Functional Languages. *SIGPLAN Notices* 33, 8 (1998), 23–27. <https://doi.org/10.1145/286385.286387>

```

1 function toSpaceOf: heap -> heap
2 function isToSpace: heap -> bool
3 function switchSemispaces: heap -> unit

1 function collect (topHeap) =
2   for r in current roots:
3     *r ← cheneyCopy(heap, *r)
4   for h below topHeap included:
5     switchSemispaces(h)
6 function cheneyCopy (topHeap, obj) =
7   heap ← heapOf(obj)
8   if depth(heap) < depth(topHeap): return obj
9   if isToSpace(heap): return obj
10  if hasFwdPtr(obj):
11    return cheneyCopy(topHeap, *fwdPtr(obj))
12  newObj ← freshObj(toSpaceOf(heap), sizeof(obj))
13  *fwdPtr(obj) = newObj
14  for field in nonptrFields(obj):
15    *getField(newObj, field) ←
16      *getField(obj, field)
17  for field in ptrFields(obj):
18    *getField(newObj, field) ←
19      cheneyCopy(topHeap, *getField(obj, field))
20  return newObj

```

Figure 14. Promotion-aware copy collection.

A Algorithms

A.1 Promotion-Aware Copy Collection

Promotion introduces redundant copies of memory objects. We now present a way to eliminate these copies by piggy-backing on the classic semispace garbage collection algorithm. Our proposal is given in Figure 14, including additional primitives specific to semispace collection.

We now assume that every hierarchical heap accessed by the mutator is paired with another heap used only during collection. Following standard terminology, we call the former a *from-space* and the latter a *to-space*. Our collection algorithm manipulates semispaces using three primitives: `toSpaceOf` returns the to-space associated with a given from-space; `isToSpace(heap)` returns true iff heap is a to-space; `switchSemispaces` swaps to-space and from-space.

The function `collect(topHeap)` is called to collect the subtree starting at `topHeap`. We assume that every task associated with a leaf heap below `topHeap` has been suspended by the runtime system. Thanks to disentanglement, this is sufficient to collect the entire subtree independently from other mutators and collectors. It copies every object reachable from a root to the to-spaces using `cheneyCopy` (l. 2-3) and then swaps the semispaces of every heap in the subtree (l. 4-5).

The function `cheneyCopy` takes a from-heap `topHeap` and an object pointer `obj`, and returns a copy of `obj`. This copy is guaranteed to be either in a to-heap below `topHeap`, or

in a from-heap strictly above `topHeap`. The latter case corresponds to the elimination of copies introduced during promotion: copy collection replaces a pointer to an old copy with a pointer to the a more recent one lying outside of the collection zone. In addition, since we do not follow forwarding pointers that belong to objects outside of the collection zone, we do not have to lock heaps during collection.

Like `promote`, we specify `cheneyCopy` as a recursive function. Let us call `heap` the heap where `obj` resides. If `heap` is strictly above `topHeap` (l. 9-10), or is a to-space (l. 11-12), `obj` can be returned. If `obj` has a forwarding pointer, `cheneyCopy` follows it (l. 13-14). Otherwise, as in `promote`, we create a new copy, setting up the forwarding pointer of `obj` to point to it, recursively call `cheneyCopy` on its pointer fields, and return it (l. 17-23).

B Implementation

Scheduler. Any implementation of a fork/join programming model requires a scheduler to coordinate work between worker threads. Tasks are evaluated within a “user-level thread” that is scheduled onto a “worker thread”. In addition to the `forkjoin` function exposed to the mutator, the scheduler also exposes a `schedule` function to idle worker threads to find waiting work.

A naive implementation of a work-stealing scheduler will create tasks for both thunks passed to it before evaluating them. However, task creation is expensive and its value is only realized upon a steal. Steals are far less frequent than calls to `forkjoin`, so our implementation ensures that calling `forkjoin` is cheap and expensive task creation is deferred to the steal. In addition, one of the thunks is evaluated immediately in the calling user-level thread, while only the other thunk is exposed to other worker threads. This reduces the number of user-level threads and thread switches by allowing a user-level thread to evaluate a path of tasks.

In our implementation, we use a work-stealing scheduler that has been annotated with heap management operations at the appropriate places. Worker threads are implemented as OS-level pthreads and user-level threads are implemented as the native user-level thread in MLton. As the scheduler operates outside of the computation, the objects it allocates, particularly in the `schedule` function, do not belong to any heap in the hierarchy. Our implementation has a separate “global heap” that is used to store scheduler data.

Superheaps. The goal of making calls to `forkjoin` cheap and deferring task creation to steals directly informs our design and implementation of hierarchical heaps. Our implementation uses a structure called a “superheap”, which is associated with a user-level thread. Superheaps contain a linked-list of heaps, each annotated with its depth. This set of heaps corresponds to the heaps of the path of tasks evaluated by the superheap’s associated user-level thread.

On a call to `forkjoin`, the current superheap's depth is incremented and future allocations take place in the new heap created for that depth. Once both thunks passed to `forkjoin` are completed, the depth is decremented. On depth decrement, the new heap will be joined to its parent heap as per the algorithm. As the set of heaps in the superheap is maintained as a linked list, and heaps are added and joined in LIFO order, the depth increment and decrement operations are very cheap.

If no steal occurs, both thunks passed to `forkjoin` will be evaluated in the newly created heap. If an idle worker thread steals one of the thunks, it creates a new user-level thread and associated superheap for that thunk. This superheap is then attached as a child to the parent superheap in order for the runtime to be aware of the complete hierarchy of heaps. When the stolen thunk is complete, it will reactivate its parent task which can then merge the superheap and extract the return value. Merging a superheap is just merging its set of heaps, which is a simple linked-list operation.

Heaps. We implement a heap as a linked-list of variable-sized memory regions called “chunks”. This formulation enables efficient implementations of key heap operations. Increasing heap size and joining heaps (`joinHeap`) are constant-time linked-list operations that do not require objects to be copied. Finding the heap of an arbitrary pointer (`heapOf`) is implemented by looking up the chunk metadata using address masking, which then contains a pointer to the heap associated with that chunk.