

Checkpoints and Continuations instead of Nested Transactions

Maurice Herlihy Eric Koskinen

Computer Science Department, Brown University

{mph,ejk}@cs.brown.edu

Abstract

We present a mechanism for partially aborting transactions through the use of data structure checkpoints and control-flow continuations. In particular, we show that boosted transactions [9] already have built-in restoration points and afford a simple, efficient implementation. Our mechanism is far simpler than previous work, which relied on complex nesting schemes to establish checkpoints. We quantify the overhead of checkpoints and explore several examples, illustrating the utility of partially aborting transactions.

As part of our evaluation, we implemented abstract locks as a novel, queue-based spin-lock which allows threads to timeout and differ in priority. Unlike the known lock due to Craig [3], this new lock is more efficient for priority schemes of few levels.

1. Introduction

Software Transactional Memory (STM) has emerged as an alternative to traditional mutual exclusion primitives such as monitors and locks, which scale poorly and do not compose cleanly. In an STM system, programmers organize activities as *transactions*, which are executed atomically: steps of two different transactions do not appear to be interleaved. A transaction may *commit*, making its effects appear to take place atomically, or it may *abort*, making its effects appear not to have taken place at all.

Being able to partially abort a transaction is useful for both performance and semantic reasons. Aborts are necessary to resolve conflicts, but often the conflict can be resolved by only rolling back some of a transaction's operations. Thus, partial aborts may yield higher performance since non-conflicting operations need not be reverted. Partial roll-back can also support semantic constructs such as conditional and or-else synchronization, where a transaction decides to roll back part of a computation to pursue another path.

A partial abort requires the notion of a *checkpoint*: a program location within a transaction to which control may jump during a partial abort. For example, after a transaction removes an element from list *A*, it may set a checkpoint before trying to insert the element into list *B*. Thus, if there is

a conflict adding to *B* the transaction can partially abort and try to insert the element again into list *B* or, alternatively, try to insert the element into list *C*. To accommodate partial aborts, each checkpoint must save and restore its *continuation*: the control state of the program location including the program counter, stack variables, and any allocated heap objects. One candidate checkpoint is the program location prior to each write operation. In practice this is too fine grained: the cost of storing a continuation per write operation is high, and most saved continuations would never be used.

In the literature thus far, checkpoints have been emulated as nested transactions [12, 13, 7]. The beginning of a nested transaction is a user-defined checkpoint. Further, each checkpoint already has a continuation available: the continuation which is captured as part of the transactional infrastructure. Nested transactions additionally solve a second *modularity* issue: during a transaction, a subroutine may be invoked which itself initiates a transaction. Modularity is essential, but not tied to "Nested Transactions" (though the name may lead one to believe otherwise); modularity can be solved syntactically, as we will later discuss.

Existing nesting approaches are based on transactional memory systems which use read/write sets. They therefore involve infrastructure to separately maintain levels of read/write sets so that recovery can occur on multiple nesting levels, a complexity which we argue is unnecessary.

In this paper, we show that transactions can be partially aborted without using nested transactions. We use the far simpler approach of data structure checkpoints and control-flow continuations to be able to revert transactions to intermediate program locations. For efficiency reasons checkpoints must be used sparingly and are best suited to operationally meaningful program locations rather than after each write operation; in much the same way, initiating a nested transaction after each write operation is not advisable.

In recent work, we have proposed Transactional Boosting [9] which solves synchronization and recovery through data structure semantics rather than read/write sets. In this paper, we additionally show that boosted transactions already establish operationally meaningful program locations that are well-suited checkpoints. Each logical operation can be a checkpoint. For example, the logical operation of "adding a node to a list" may involve several write oper-

ations, but only the first operation is a natural restoration point. Boosting a transaction both removes the burden on the user to chose checkpoints, and also facilitates an efficient runtime storage of checkpoints.

Specifically, we make the following contributions:

- We show that nesting is unneeded to partially abort a transaction. Instead, transactions can be reverted to previous control-flow locations by checkpointing data structures and capturing continuations. This is discussed in Section 3.
- In Section 3.1, we show that a boosted transaction has natural pre-defined checkpoints: each logical operation. This both eliminates the burden of choosing checkpoints and also leads to an efficient implementation, with a sparing use of checkpoints. We further show that Transactional Boosting allows us to record checkpoints *operationally* rather than treating data structure manipulations as flat memory access. Consequently, the cost of storing checkpoints is reduced.
- We explore the utility of checkpoints through a number of examples of boosted transactions in Section 4, and discuss our evaluation in Section 5.
- Finally, we discuss our spin-lock for abstract locking in Section 5.3. Ours is a novel priority-based queue spin-lock with timeouts, and we contrast it to the lock due to Craig [3]. Our lock has theoretically higher performance for binary priority schemes.

We emphasize that the utility of checkpoints and continuations is not limited to transactional boosting. For example LogTM [12] could introduce a new keyword, and allow users to manually denote semantically meaningful control locations whose continuation could be captured during execution. Then, rather than using nested transactions, partial aborts can be realized by reverting write operations to restore the heap and then invoking stored continuations.

2. Background

2.1 Transactional Boosting

Most transactional memory approaches rely on a log of memory access to determine synchronization and recovery. As a transaction executes, each `read` and `write` operation is stored in a log along with the corresponding memory location. Two transactions *conflict* when they concurrently access the same memory location and at least one operation is a `write`. In the event of a conflict, one transaction must be aborted: all `write` operations must be reverted. The log is again used to either restore old values or discard commit-time effects.

In recent work [9] we introduced Transactional Boosting, which yields significant performance gains by shedding memory access logs and relying on data structure semantics to determine synchronization and recovery. A more through

discussion can be found in [9], but for completeness we summarize transactional boosting here.

Transactional Boosting is built upon linearizable base objects, which have an abstract state and a corresponding concrete implementation. The semantics of object methods are known, and often specified in terms of pre- and post-conditions describing the state of the object before and after the method invocation. Two method invocations are said to *commute* if applying them in either order causes the object to transition to the same state and respond with the same return values. Additionally, an *inverse* of a method is a second method which returns the object to the previous state.

Transactional Boosting relies on commutativity to define and detect conflicts. Users define *abstract locks* associated with invocations of methods on the boosted object. These locks conflict whenever a pair of methods do not commute. Before invoking a method, transactions must acquire all abstract locks associated with the invocation. In this manner, non-commutative operations never occur concurrently since one of the two operations will be delayed.

Inverses allow recovery to be performed at the granularity of method calls. While a transaction is executing, a log of operations is maintained. If the transaction aborts, this log is played in reverse order before abstract locks are released; a committed transaction can simply discard the log.

Here is an example. A Set may have methods `add(x)`, `remove(y)`, and `contains(z)` with the usual semantics. Pairs of method calls such as `add(x)` and `add(y)` commute for $x \neq y$. Indeed, any pair of methods with distinct arguments commute. Moreover, inverse operations are well defined: the inverse of `add(x)` is simply `remove(x)`.

2.2 Nested Transactions

Several nesting approaches have been proposed [12, 13, 7] for software transactional memory. The purpose of nested transactions¹ is twofold:

Checkpoints. A nested transaction establishes a new program location where control may be returned in the event of an abort. Specifically, the beginning of each nested transaction is a checkpoint. Consider the following example:

```

global int list [10];
atomic {
  int x = list [0];
  atomic {
    list [2] = x + 1;
  }
}

```

If access to `list [2]` causes a conflict, the inner transaction may be aborted rather than the entire outer transaction. Thus, each nested `atomic` statement demarcates a checkpoint.

Modularity. Often code that transactionally mutates data structures may be encapsulated in libraries. Nested transactions enables the sound interaction between transactional user code and invocations of transactional library methods.

¹By “nested transactions” we mean closed nested transactions that only commit if the top-level parent commits.

For example, before a transaction removes an element from a queue, it may use a library method to check that the queue is nonempty:

```

global Queue sharedQueue;

bool EmptyQueue(Queue q) {
    atomic { ... }
}

atomic {
    if (EmptyQueue(sharedQueue))
        throw EmptyQueueException;
    return Dequeue(sharedQueue);
}

```

In the next section we describe how partial aborts can be realized through first-class checkpoints and continuations. We will then revisit nested transactions, comparing the two approaches.

3. Checkpoints and Continuations

The key idea of this paper is that partial aborts can be realized without using nested transactions. In this section, we define partial aborts, in Section 4 we explore some applications of partial aborts, and in Section 4.1 we quantify the efficacy of partial aborts for one particularly useful application.

As a running example, consider a pool of threads transactionally accessing shared hash tables. One thread’s access pattern might be as follows:

```

global HashTable ht1, ht2;

atomic {
    Object oldVal = HashTable_Remove(ht1,7);
    Object newVal = foo(oldVal);
    HashTable_Insert(ht2,7,newVal);
}

```

In this example, the thread removes an element from the first hash table, performs some computation on `oldVal` in `foo()`, and then inserts `newVal` into a second hash table.

Now, consider what would happen if contention at the second hash table triggered an abort of this thread. In the simplest case, transactional implementations would undo all of the thread’s effects, returning the element to the first hash table.

Yet there may not be a conflict at `ht1`, so it is unfortunate that the transaction must be aborted entirely. A more satisfying solution would be to *partially* abort the transaction. If control is returned to the point before the insertion, then more progress is preserved. This is of course a trivial example with a low penalty for aborting entirely, but it is easy to imagine cases where the penalty is more severe.

The core challenge involved in partially aborting a transaction is to be able to soundly return to an earlier program location. In programming language theory, this facility is known as a *continuation* [14, 1]. There are two correctness criteria involved in implementing a continuation:

1. **Stack Restoration:** All stack variables in scope must be restored to their values at the time when control first passed through the program location. Formally, the continuation *captures the environment* of the program location.
2. **Heap Restoration:** All heap data structures scope must be restored to their state at the time when control first passed through the program location. Formally, the continuation *captures the store* of the program location.

Purely functional languages are heap-less and so implementations must only be concerned with restoring the stack. Imperative languages add the complication of restoring the heap. Finally, concurrent shared-memory imperative programs must be able to restore the heap in such a way that consistency is maintained with respect to other concurrently executing threads. Nested transactions accomplish this final issue with added complexity, whereas addressing the issue directly is the novelty of our work.

In the remainder of this paper, we focus on checkpoints and continuations as they apply to boosted transactions. As mentioned in Section 1, the principles we discuss are not limited to the realm of transactional boosting. Nonetheless, a boosted transaction has convenient properties (discussed below) which we exploit to gain an efficient implementation.

3.1 Boosted Checkpoints

Transactions manipulate a boosted object by executing a series of operations, each causing the object to transition from one abstract state to another. Continuing with the hash table example, one such operation is “inserting a key/value pair.” Although the single invocation may result in many read and write operations, it implements a well-defined abstract state transition.

We define a *boosted checkpoint* at the program location that precedes each operation on a boosted object. In the hash table example, explicit check points are defined before the removal and insertion and implicitly at the beginning of the transaction:

```

atomic {
    Object oldVal = CHKPT(HashTable_Remove(ht1,7)); // CP1
    Object newVal = foo(oldVal);
    CHKPT(HashTable_Insert(ht2,7,newVal)); // CP2
}

```

Boosted checkpoints are *not* snapshots of the object, but rather semantically rich locations in the program’s control-flow graph. Between each subsequent checkpoint, there is a well-defined transition, which is given by the semantics of the interleaved operation. For example, the transition between the checkpoints denoted CP1 and CP2 is `HashTable_Remove(ht1,7)`, in other words: the removal of the value at key 7 from hash table `ht1`. The transition from checkpoint CP2 and the end of the transaction is similarly defined.

Since boosted object state transitions are well-defined between checkpoints, heap recovery can be realized at *any* checkpoint without maintaining shadow copies. If, for example, this transaction is aborted after completing the insertion into `ht2`, it can be partially aborted to CP2 rather than to the beginning of the transaction, by simply executing the inverse of `HashTable_Insert(ht2,7,newVal)`. As discussed in Section 2.1, all boosted objects have known inverse operations.

Boosted checkpoints are fundamentally different from nested transactions. Most nested implementations maintain extra copies of objects. In this sense, they are agnostic of data structure semantics and treat memory as a flat array of locations and values. However, as we demonstrated in [9], this is often expensive.

3.2 Boosted Continuations

We now turn to the issue of returning the program control to an earlier checkpoint. In the above hash table transaction, we may for example want to return control to the CP1 checkpoint if an abort occurs before the end of the transaction. Such jumps in control are a special case of the general concept of non-local control flow, which is accomplished by *continuations* from programming language theory. Informally, continuations represent the current values of stack variables in scope (the *environment*) and the program location (a point on the control-flow graph). Our non-local control flow is accomplished with two steps.

1. *Continuation storage.* Each time program control reaches a checkpoint, the current continuation is captured. Indeed, the checkpoint itself is augmented with the storage of the corresponding captured continuation.
2. *Continuation invocation.* When a transaction is aborted, inverses are played backwards until the desired checkpoint is reached, discarding saved continuations along the way. Upon reaching the desired checkpoint, the corresponding continuation is invoked.

Nested Transactions rely on *implicit* continuations to return to the beginning of an inner transaction. Transactions by definition already have a built-in notion of returning control, so the nesting approach introduces new transactions at each program location. We suggest that conventional closed nested transaction schemes and data structures can be discarded in favor of this simple continuation scheme.

3.3 Implementation Notes

We discuss implementation details thoroughly in Section 5.1, but continuations deserve special treatment here. Unfortunately, availability of continuations differs from one language to the next. High-level languages such as Scheme and Standard ML of New Jersey support first class continuations, while most imperative languages offer dynamic variants of `goto` such as `setjmp/longjmp`.

```

global bucket_t data[1000];

atomic {
  key_list [] = decide();
  for(int key :: key_list []) {
    CHECKPOINT();
    if (QLOCK(key) == TIMEOUT)
      ABORT();
    Manipulate(data,key);
    LOG_INVERSE(ManipulateInverse,data,key);
  }
}

```

Figure 1. A transaction accessing a shared array.

Our prototype implementation is built on Transactional Locking 2 (TL2) [5] which is written in C. As such, the checkpoint continuations were captured with `getcontext` and invoked with `setcontext`. These standard C library functions allow us to capture the continuation of multiple checkpoints, whereas `setjmp/longjmp` capture only a single point in the control flow.

Unfortunately, these library functions do not capture the entire program environment. To compensate, our prototype allocates additional storage for stack variable values alongside the captured context. When control passes through a checkpoint, the macro `SAVE_ENV(var)` adds the current value of `var` to a data structure, which is later restored via `RESTORE_ENV(var)` when control returns to the context after a partial abort. Ideally, of course, this work could be done by a compiler.

4. Applications and Examples

In this section we explore some examples that illustrate the utility of partial aborts in Transactional Boosting.

4.1 Priority

In a priority scheme, it is desirable for high priority threads to have greater throughput than low priority threads. In the context of transactional memory, high-priority threads can force low priority threads to abort. Unfortunately, aborting a low priority thread drastically harms low priority throughput. Checkpoints and continuations allow us to partially abort low priority threads, rolling them back just far enough to allow high priority threads to commit.

Consider threads accessing a boosted shared array, as shown in Figure 1. For the moment, assume `QLOCK` is a table from keys to simple spin locks. If a lock is acquired then it is added to the transaction's log to be released when the transaction commits or aborts. Let us also assume that `Manipulate(data,key)` is some method which modifies the specified index and that `ManipulateInverse(data,key)` is the method which performs the corresponding *inverse* operation (one might increment and the other decrement, for example).

A correct semantics for this array states that two invocations *commute* if they access distinct indices. If threads ac-

cess multiple indices, then they will acquire an abstract lock for each before committing. For example, the `key_list` for two threads might include the following array indices:

Thread *A*: 28, 56, 41, 10
 Thread *B*: 21, 19, 41, 67, 95

If thread *A* acquires the abstract lock for 41 first, then thread *B* will be delayed until *A* commits.

Now let's add priority to the scenario: perhaps *B* has higher priority than *A*. We would therefore like to allow *B* to proceed and commit at the cost of *A*'s throughput. Indeed this is possible without checkpoints: we can simply *abort A*. However, we then lose the progress that *A* has already made. There is no need for *A* to undo the work which required the acquisition of locks 28 and 56 since *B* will never try to acquire either lock.

The main idea of this paper addresses this issue. If we capture the continuation at each checkpoint, then we can semantically undo the changes to the data structure made by *A* to the point before 41 was acquired by executing inverse operations.

4.1.1 CLH Queue Lock with Priority

Realizing partial aborts in the context of priority requires that our spin lock support both priority and timeouts. The queue-based spinlock due to Craig, Landin and Hagersten [4, 11] (CLH Lock) already supports timeouts, so we build on the CLH Lock and present a novel spin lock which accommodates thread priority in a distinctly different way than the priority queue lock due to Craig [3]. We qualitatively describe our lock here and a full design discussion is given in Section 5.3.

Threads trying to acquire a queue-based spin lock enqueue themselves into a list of acquirers. The front-most thread has acquired the lock and executes its critical section. When the lock is released, the front-most thread passes the lock to the next thread in line, and removes itself from the queue. In this way, cache coherency issues are minimized because threads spin on distinct memory locations.

We augment the spin lock by defining two queues: a queue of low priority threads and another of high priority threads. When a low priority thread releases the lock it passes availability either to the high priority queue (if any high priority threads exist) or to the next low priority thread. Analogously, a high priority thread either releases it to the next high priority thread or turns the lock over to the low priority queue.

4.1.2 Aborting Low Priority Threads

When a high priority thread needs to acquire a lock which is currently held by a low priority thread, we must abort the low priority thread. Our transactional model has no notion of preemption, but we can signal low priority threads to abort themselves. When a high priority thread tries to acquire a

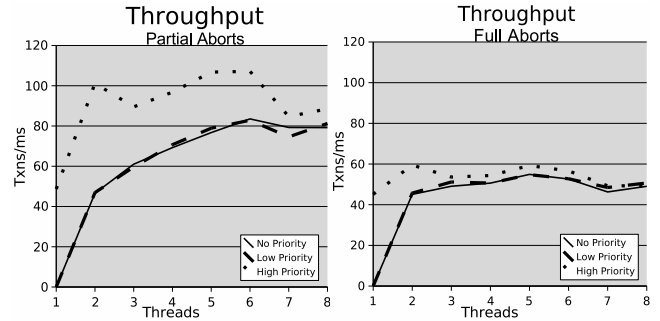


Figure 2. Throughput of high, low, and priority-free threads when they are partially aborted (to the left) or completely aborted (to the right).

lock held by a low priority thread, it can send a brief message (by writing to thread-local storage) to the low priority thread indicating which key it wishes to acquire. Low priority threads listen for messages while they are spinning on other locks.

When a low priority thread receives a message from a high priority thread, it initiates its own partial abort. The message contains the key desired by the high priority thread. The low priority thread can therefore partially abort itself to the checkpoint which immediately precedes the contested lock. When the lock is released the high priority thread, having earlier enqueued itself, immediately obtains it and continues with its transaction.

4.1.3 Performance

We extended our implementation of transactional boosting [9] (which is based on TL2) with the CLH-style queue lock discussed above. The benchmarks were run on a multiprocessor with four 2.0 GHz Xeon processors, each one two-way hyperthreaded for a total of eight threads.

We evaluated our priority example by confirming that high priority threads have higher throughput than (a) lower-priority threads and (b) the throughput they would have had in the absence of priority schemes. Figure 2 shows the throughput of high and low priority threads, compared with threads in a priority-free scheme. To the left, low priority threads are partially aborted, whereas they are completely aborted on the right. The two graphs illustrate that low, high, and priority-free threads have greater throughput when they are partially aborted rather than completely aborted. Performance is improved even though realizing partial aborts involves the overhead of capturing and storing continuations.

4.2 Conditional Synchronization

Checkpoints allow us to implement two forms of transactional conditional synchronization. STMHaskell introduced two constructs called `retry` and `orElse`. Both provide general mechanisms for incorporating conditional synchronization into transactional computation. Supporting them has

```
List inbound, outbound;
AbstractLock inLock, outLock;
Object in = NULL;
```

```
atomic {
  OP(inLock, 0, Prepare(inbound));
  do {
    in = OP(inLock, 0, Dequeue(inbound));
    if (!in.isSpecial)
      retry;
  }
  OP(outLock, 0, Enqueue(outbound,in));
}
```

Figure 3. Example using `retry` in a boosted transaction.

```
HashTable htA, htB, htC;
AbstractLock alA, alB, alC;
```

```
atomic {
  result_t result = OP(alA, key, Remove(htA,key));

  { OP(alB, key, Add(htB,key, result)) }
  orElse
  { OP(alC, key, Add(htC,key, result)) }
}
```

Figure 4. Example using `orElse` in a boosted transaction.

a pervasive effect on any STM implementation, while a boosted implementation of transactional objects does not. In this subsection, we use the notation:

$$\text{result} = \text{OP}(\text{aLock}, \text{lockKey}, \text{dataStructureMethod})$$

to represent the combined task of (a) declaring a checkpoint (b) acquiring abstract lock `aLock` for key `lockKey` (c) executing `dataStructureMethod` and (d) logging the corresponding inverse.

Figure 3 depicts an example of the `retry` construct. First the transaction prepares the inbound queue, and then begins dequeuing elements. Each dequeued element is examined, and if the thread has not found the desired element, it uses the `retry` construct to partially abort the transaction. The element is returned to the queue, and the transaction tries to dequeue again, hoping that some other thread has enqueued an element of interest. When such an element of interest arrives, it is finally enqueued onto the outbound queue.

The example in Figure 4 illustrates the utility of the `orElse` construct. In this example, the transaction removes an element from `htA` and then tries to add it to another hash table. While adding it to `htB`, it may experience a conflict with another transaction concurrently accessing `htB`. In this case, we can partially abort the transaction and, following the directions of the `orElse` statement, the transaction can attempt to add the element to `htC`. This example has applications towards various forms of the work queue paradigm.

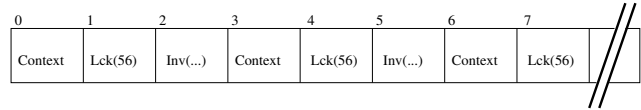


Figure 5. Diagram of the runtime *computation log*, which stores captured continuations (“Context”), abstract locks (“Lck”), and inverse methods (“Inv”).

4.3 Contention Management

Contention Managers are used in transactional memory to resolve conflicts with runtime heuristics to improve throughput. Boosted transactions empower contention managers with a more fine-grained abort mechanism. Rather than entirely aborting transactions, contention managers can partially abort transactions. Moreover, by examining the operation log, contention managers can revert transactions to just before the (first) conflicting operation. See [15] for a detailed discussion of contention management.

4.4 Modularity

Modularity is a desirable attribute of transactional programs. For example, during a transaction a user may want to invoke library methods which themselves initiate transactions. Consider the following code:

```
Object swap(HashTable ht, int key, Object newVal) {
  atomic {
    Object r = HashTable_Remove(ht,key);
    HashTable_Insert(ht, key, newVal);
    return r;
  }
}

atomic {
  Object x = swap(htA, key, y);
  HashTable_Insert(htB, key, x);
}
```

Here the transaction within `swap()` is “nested” within the outer transaction by the program’s control flow.

Rather than initiating an entirely new inner transaction, a boosted implementation with checkpoints can instead simply establish a new checkpoint at the beginning of each nested transaction. Finding nested transactions is a simple runtime task, keeping a flag per thread which indicates whether a transaction has been initiated. Moreover, this task can often be further simplified by statically detecting syntactically nested transactions and replacing them with checkpoints.

5. Evaluation

5.1 Checkpoints and Continuations

We implemented Checkpoints on top of our Transactional Boosting implementation in TL2 [9]. Much like abstract locks and inverses, checkpoints are stored in a runtime *computation log*. A diagram of the computation log is given in Figure 5. At each checkpoint, the continuation is captured

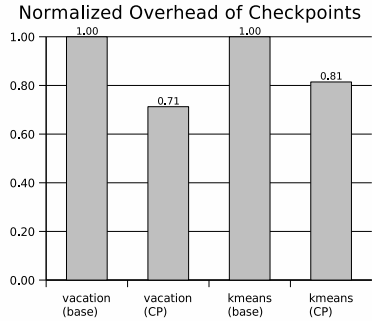


Figure 6. Normalized throughput of boosting with and without checkpoints for two *STAMP* [2] benchmarks.

and appended to the log. Checkpoints (“Context”) precede operations on shared data structures, which involves acquiring an abstract lock (“Lck”), performing the operation, and then recording the inverse operation (“Inv”) in the log.

When a transaction is aborted, as long as the log is traversed in the reverse direction, invoking any context produces correct behavior. As a context is passed it is deallocated, inverses are invoked to revert data structure operations, and abstract locks are released allowing other threads to acquire them. In Section 4 we discussed application-specific strategies for deciding how far to traverse.

5.2 Overhead

There is an overhead associated with capturing contexts and storing them in the computation log. The amount of overhead depends on the workload. When a transaction performs simple operations such as incrementing a counter, then the overhead is significant. More realistic workloads (indeed the type of workloads for which it is desirable to partially abort transactions in the first place) have a more manageable overhead.

We modified two of the Stanford *STAMP* benchmarks [2] (written in C) to use boosting with and without checkpointing. The benchmarks were run on a multiprocessor with four 2.0 GHz Xeon processors, each one two-way hyperthreaded for a total of eight threads. Figure 6 we compare the throughput of transactions with and without the capture of continuations for the vacation and kmeans benchmarks².

5.3 SpinLock Design and Implementation

The priority example in Section 4.1 requires a spin lock in the queue form which is due to Craig, Landin and Hagersten [4, 11], with the added need for timeouts and priority schemes.

As in [4, 11], threads attempting to acquire a lock append new “lock nodes” to the tail of a queue, and examine the previous tail element. The threads append themselves by swinging the tail pointer with a compare-and-swap (CAS) opera-

²For vacation we used the flags `-n4 -q90 -u80 -r65536 -t409` and kmeans we used `-m40 -n40 -t0.05 -i inputs/random1000_12`

tion to ensure a serial ordering. Threads then spin, waiting until their predecessor lock node is marked “unlocked,” at which point the thread has logically acquired the lock. After completing its critical section, the thread marks its own node as “unlocked,” allowing the next thread to acquire the lock. Threads may wish to timeout, in which case they can mark a node as “abandoned.” The predecessor thread spinning on the newly abandoned node deallocates the node and follows a pointer to spin on the subsequent node in the queue.

We augment the queue lock for priority by defining one queue per priority level. In the example of Section 4.1 only two queues are necessary: a high priority queue and a low priority queue. Threads trying to acquire the lock, add nodes to the queue corresponding to their priority level. When a low priority thread holding the lock completes its critical section, it passes the lock to the high priority queue if it is non-empty, or else passes it to the next node of the low priority queue. Similarly, when a high priority thread completes its critical section, it either passes the lock to the next high priority thread or else passes it back to the low priority queue. This generalizes to any number of priority levels.

We also introduce preemption to the queue lock. When a high priority thread tries to acquire a lock currently held by a low priority queue, the high priority thread can request that the low priority thread release the lock. The high priority thread writes the name of the lock into the holder’s thread-local field. The holder checks the field whenever it is spinning on other locks, and if it notices that it is non-null, then the low priority thread partially aborts itself to the checkpoint which precedes the contested lock. To avoid deadlock, we also allow low priority threads to abort each other, using the transaction start time as a tie-breaker.

Pseudo-code for our implementation is given in Figure 7. This implementation is different from the priority queue lock due to Craig [3]. The lock of Craig stores priority as a field within the queue. This has the advantage of compactness, at the cost of a queue traversal before each acquisition. By contrast, our implementation maintains separate queues for each priority level. Thus, acquisition is done in constant time. We also note the following implementation issues:

- The first enqueue case in Figure 7 is an important fast path. When it is empty, high priority threads can simply add themselves to the low priority queue.
- The second enqueue case occurs when the fast path does not apply and the thread has high priority. In addition to adding itself to the high priority queue, the thread must also enqueue in the low priority queue so that the lock can be passed from the low queue to the high queue.
- In the priority example, threads may try to acquire the same lock repeatedly during a single transaction. We added a counter to track how many times the lock had already been acquired, performing the actual `QLockKey_Lock()`.
- In order for high priority threads to find the low priority owner of a contested lock, threads “announce” them-

```

void QLockKey_Lock(QLockKey_t* ql) {
  // Enqueue
  if (! HighPriority || QLockKey_Empty(ql->low_queue))
    QLockKey_Enqueue(ql->low_queue);
  else {
    QLockKey_Enqueue(ql->high_queue);
    QLockKey_Enqueue(ql->low_queue);
    if (QLockKey_Acquired(ql->low_queue))
      QLockKey_Dequeue(ql->low_queue);
    else
      QLockKey_Abandon(ql->low_queue);
  }
  // Spin
  while(high || pattience--) {
    if (SpinNode->abandoned)
      QLockKey_Advance(SpinNode);
    else if (! SpinNode->locked)
      // Acquired!
      else if (ReceivedMessage(OtherLock))
        PartiallyAbortTo (OtherLock);
      else if (NeedToPreempt())
        SendMessage(SpinNode->owner());
      else
        // keep spinning!
  }
  // Timeout
  QLockKey_Abandon(ql->low_queue);
}

void QLockKey_Unlock(QLockKey_t *ql) {
  QLockKey_Dequeue(ql);
}

```

Figure 7. Pseudo-code for `QLock_Lock()` and `QLock_Unlock()`.

selves sometime after they acquire the lock, and “unannounce” before releasing the lock.

6. Formal Model

Theorem 6.1. *A transactional boosting implementation that obeys the correctness rules specified in [8] and permits partial aborts as described informally above yields histories whose committed transactions are strictly serializable.*

Proof. Trivial. Follows from the *Main Theorem* of [8], with an added “Partial Abort” event. \square

Informally, the system is correct because partial aborts “annihilate” a subsequence of a thread’s history, which is later regenerated. When the thread commits (if it does commit), commutativity ensures that the history defines the same state as if the thread’s operations occurred instantaneously.

7. Related Work

Our work stands in contrast to the related work on nesting approaches for transactional memory. We argue that data checkpoints and control-flow continuations alleviate the need for complex nesting strategies.

Closed nested transactions have been proposed by Moravan *et al.* [12] for the log-based transactional memory im-

plementation LogTM. LogTM implements STM by tracking read/write sets in a runtime log. Nested transactions are implemented by maintaining layers of read/write, analogous to activation records. Abstract nested transactions [7] are closed nested transactions with an added optimization. In this paper we argue that the need for nested transactions is artificial and suggest that such complex nesting strategies could be discarded in favor of simple data structure checkpoints and continuations.

More recently, open nested transactions were proposed by Ni *et al.* [13] and expose the layers of memory access to developers through a simple API. The goal is to allow programmers to manually avoid aborts which arise from false conflicts. Additionally, examples of how to implement collection classes with open nested transactions was given by Carlstrom *et al.* [6]. We again argue that nesting is unneeded. Moreover, false conflicts can be eliminated with the transactional boosting methodology [9] rather than exposing an API that is prone to developer mistakes.

Hardware support for nested transactions has also been proposed by Lev and Maessen [10], who show how nesting can be implemented in hardware by subdividing a (software) nested transaction into multiple hardware transactions.

8. Conclusion

We have shown that transactions on boosted objects can be partially aborted without the need for nesting. We store continuations which capture the context before each method invocation, and roll back transactions by operationally invoking inverse object operations. We explored the utility of partial aborts through several examples.

The added expressiveness of partial aborts comes with the cost of storing continuations with each checkpoint. However, since continuations are stored sparingly (per logical operation rather than per write operation) boosted checkpoints are more feasible than conventional read/write checkpoints. Nonetheless, it may be worth exploring applying the techniques presented here toward conventional transactional memory implementations such as LogTM [12].

Acknowledgments

This work was funded by NSF grant 0410042, and grants from Sun Microsystems, Microsoft, and Intel. We would like to thank Yossi Lev and the anonymous reviewers for their feedback.

References

- [1] APPEL, A. W. *Modern Compiler Implementation in ML*. Cambridge University Press, New York, NY, USA, 2004.
- [2] CAO MINH, C., TRAUTMANN, M., CHUNG, J., MCDONALD, A., BRONSON, N., CASPER, J., KOZYRAKIS, C., AND OLUKOTUN, K. An effective hybrid transactional memory system with strong isolation guarantees. In *Proceedings of the*

34th Annual International Symposium on Computer Architecture (ISCA 07). Jun 2007.

- [3] CRAIG, T. Queuing spin lock algorithms to support timing predictability. *Real-Time Systems Symposium, 1993., Proceedings*. (1-3 Dec 1993), 148–157.
- [4] CRAIG, T. Building fifo and priority-queueing spin locks from atomic swap. Tech. Rep. 93-02-02, Department of CSE, University of Washington, 1993.
- [5] DICE, D., SHALEV, O., AND SHAVIT, N. Transactional locking II. In *Proceedings of the 20th International Symposium on Distributed Computing (DISC '06)* (September 2006).
- [6] HAMMOND, L., WONG, V., CHEN, M., CARLSTROM, B. D., DAVIS, J. D., HERTZBERG, B., PRABHU, M. K., WIJAYA, H., KOZYRAKIS, C., AND OLUKOTUN, K. Transactional memory coherence and consistency. *SIGARCH Computer Architecture News* 32, 2 (2004), 102.
- [7] HARRIS, T., AND STIPIĆ, S. Abstract nested transactions.
- [8] HERLIHY, M., AND KOSKINEN, E. Transactional boosting: A methodology for highly concurrent transactional objects. Tech. Rep. CS-07-08, Brown University, 2007.
- [9] HERLIHY, M., AND KOSKINEN, E. Transactional boosting: A methodology for highly concurrent transactional objects. In *Proceedings of the 13th ACM SIGPLAN symposium on Principles and practice of parallel programming (PPoPP '08)* (2008).
- [10] LEV, Y., AND MAESSEN, J.-W. Split hardware transaction: True nesting of transactions using best-effort hardware transactional memory. In *Proceedings of the 13th ACM SIGPLAN symposium on Principles and practice of parallel programming (PPoPP '08)* (2008).
- [11] MAGNUSSON, P. S., LANDIN, A., AND HAGERSTEN, E. Queue locks on cache coherent multiprocessors. In *Proceedings of the 8th International Symposium on Parallel Processing* (Washington, DC, USA, 1994), IEEE Computer Society, pp. 165–171.
- [12] MORAVAN, M. J., BOBBA, J., MOORE, K. E., YEN, L., HILL, M. D., LIBLIT, B., SWIFT, M. M., AND WOOD, D. A. Supporting nested transactional memory in logtm. In *Proceedings of the 12th international conference on Architectural support for programming languages and operating systems (ASPLOS-XII)* (New York, NY, USA, 2006), ACM Press, pp. 359–370.
- [13] NI, Y., MENON, V. S., ADL-TABATABAI, A.-R., HOSKING, A. L., HUDSON, R. L., MOSS, J. E. B., SAHA, B., AND SHPEISMAN, T. Open nesting in software transactional memory. In *Proceedings of the 12th ACM SIGPLAN symposium on Principles and practice of parallel programming (PPoPP '07)* (New York, NY, USA, 2007), ACM Press, pp. 68–78.
- [14] STRACHEY, C., AND WADSWORTH, C. P. Continuations: A mathematical semantics for handling fulljumps. *Higher Order Symbol. Comput.* 13, 1-2 (2000), 135–152.
- [15] WILLIAM N. SCHERER, I., AND SCOTT, M. L. Advanced contention management for dynamic software transactional memory. In *PODC '05: Proceedings of the 24th annual ACM symposium on Principles of distributed computing* (New York, NY, USA, 2005), ACM, pp. 240–248.