

Stretching Transactional Memory

Aleksandar Dragojević Rachid Guerraoui Michał Kapałka
EPFL, Lausanne, Switzerland

Abstract

Transactional memory (TM) is an appealing abstraction for programming multi-core systems. Potential target applications for TM, such as business software and video games, are likely to involve complex data structures and large transactions, requiring specific software solutions (STM). So far, however, STMs have been mainly evaluated and optimized for smaller scale benchmarks.

We revisit the main STM design choices from the perspective of complex workloads and propose a new STM, which we call SwissTM. In short, SwissTM is lock- and word-based and uses (1) optimistic (commit-time) conflict detection for read/write conflicts and pessimistic (encounter-time) conflict detection for write/write conflicts, as well as (2) a new two-phase contention manager that ensures the progress of long transactions while inducing no overhead on short ones. SwissTM outperforms state-of-the-art STM implementations, namely RSTM, TL2, and TinySTM, in our experiments on STMBench7, STAMP, Lee-TM and red-black tree benchmarks.

Beyond SwissTM, we present the most complete evaluation to date of the individual impact of various STM design choices on the ability to support the mixed workloads of large applications.

1. Introduction

Transactional memory (TM) is an appealing abstraction for making concurrent programming accessible to a wide community of non-expert programmers while avoiding the pitfalls of critical sections. With a TM, application threads communicate by executing operations on shared data inside lightweight in-memory *transactions*. A transaction performs a number of actions and then either *commits*, in which case all the actions are applied to shared data *atomically*, or *aborts*, in which case the effects of those actions are rolled back and never visible to other transactions. From a programmer's perspective, the TM paradigm is very promising as it promotes program composition [17], in contrast to explicit locking, while still providing the illusion that all shared objects are protected by some global lock. Yet, it offers the possibility of performance comparable to hand-crafted, fine-grained locking.

A major target of TMs are large applications such as business software or video games: the size of these applications make them ideal candidates to benefit from emerging multi-core architectures. Such applications typically involve dynamic and non-uniform data structures consisting of many objects of various complexity. For example, a video gameplay simulation can use up to 10,000 active

interacting game objects, each having mutable state, being updated 30–60 times per second, and causing changes to 5–10 other objects on every update [31]. Unless a TM is used, making such code thread-safe and scalable on multi-cores is a daunting task [31]. The big size and complexity of such applications can, in turn, easily lead to large transactions, for these can naturally be composed [17]. Some TM interfaces [1], in fact, promote the encapsulation of entire applications within very few transactions.

The motivation of this work is to explore the ability of software mechanisms to effectively support mixed workloads consisting of small and large transactions, as well as possibly complex data structures. We believe this to be of practical relevance because even if hardware TM support becomes widely available in the future, it is likely that only smaller-scale transactional workloads will be fully executed in hardware, while software support will still be needed for transactions with large read and write sets. For example, the hybrid hardware/software scheme proposed in [22] switches from full hardware TM to full software TM when it encounters large transactions. The ability of STM systems to effectively deal with large transactions will be crucial in these settings as well.

Since the seminal paper on a *software* TM (STM) that supported *dynamic* data structures and *unbounded* transactions [19], all modern STMs are supposed to handle complex workloads [19, 23, 9, 24, 18, 2]. A wide variety of STM techniques, mainly inspired by database algorithms, have been explored. The big challenge facing STM researchers is to determine the right combination of strategies that suit the requirements of concurrent applications—requirements that are significantly different than of database applications. So far, however, most STM experiments have been performed using benchmarks characterized by small transactions, simple and uniform data structures, or regular data access patterns. While such experiments reveal performance differences between STM implementations, these are not fully representative of complex workloads that STMs are likely to get exposed to once used in real applications. Worse, they can mislead STM implementors by promoting certain strategies that may perform well in small-scale applications but are counter-productive with complex workloads. Examples of such strategies, which we discuss in more details later in the paper, include the following.

1. The *commit-time locking* scheme, used for instance in TL2 [9], is indeed effective for short transactions, but might waste significant work of longer transactions that eventually abort due to write/write conflicts. This is because write/write conflicts, which usually lead to transaction aborts, are detected too late.
2. The *encounter-time locking* scheme, used by most STMs, e.g., TinySTM [24] and McRT-STM [2], immediately aborts a transaction that tries to read a memory location locked by another transaction. Hence, read/write conflicts, which can often be handled without aborts, are detected very early and resolved by aborting readers. Long transactions that write memory locations commonly read by other transactions might thus end up block-

ing many other transactions, and for a long time, thus slowing down the system overall.

3. The *timid contention management* scheme, used by many STMs, especially word-based ones such as TL2 and TinySTM, and which aborts transactions immediately upon a conflict, favors short transactions. Contention managers such as Greedy [14] or Serializer [26] are more appropriate for large transactions, but are hardly ever used due to the overhead they impose on short transactions.

It is appealing but challenging to come up with strategies that account both for long transactions and complex workloads, as well as for short transactions and simple data structures: these might indeed typically co-exist in real applications. This paper is a first step towards taking that challenge up. We perform that step through SwissTM, a new lock- and word-based STM. The main distinctive features of SwissTM are:

- A conflict detection scheme that detects (a) write/write conflicts eagerly, in order to prevent transactions that are doomed to abort from running and wasting resources, and (b) read/write conflicts late, in order to optimistically allow more parallelism between transactions. In short, transactions eagerly acquire objects for writing, which helps detect write/write conflicts as soon as they appear. This also avoids wasting work of transactions that are already doomed to abort after a write/write conflict. By using invisible reads and allowing transactions to read objects acquired for writing, SwissTM detects read/write conflicts late, thus increasing inter-transaction parallelism. A time-based scheme [9, 25] is used to reduce the cost of transaction validation with invisible reads.
- A two-phase contention manager that incurs no overhead on read-only and short read-write transactions while favoring the progress of transactions that have performed a significant number of updates. Basically, transactions that are short or read-only use the simple but inexpensive timid contention management scheme, aborting on first encountered conflict. Transactions that are more complex switch dynamically to the Greedy mechanism that involves more overhead but favors these transactions, preventing starvation. Additionally, transactions that abort due to write/write conflicts back-off for a period proportional to the number of their successive aborts, hence reducing contention on memory hot spots.

We evaluate SwissTM with state-of-the-art STMs by using benchmarks that cover a large part of the complexity space. We start with STMBench7 [16], which involves (1) non-uniform data structures of significant size, and (2) a mix of operations of various length and data access patterns. Then, we move to Lee-TM [3]—a benchmark with large but regular transactions—and STAMP [7]—a collection of realistic medium-scale workloads. Finally, we evaluate SwissTM with a red-black tree microbenchmark that involves very short and simple transactions. SwissTM outperforms state-of-the-art STMs—RSTM [23], TL2 [9], and TinySTM [24]—in all the considered benchmarks. For example, in the read-dominated workload of STMBench7 (90% of read-only operations), SwissTM outperforms the other STMs by up to 65%, and in the write-dominated workload (10% of read-only operations)—by up to 10%. Also, SwissTM provides a better scalability than the other STMs, especially for read-dominated and read-write (60% of read-only operations) workloads of STMBench7.

We compare SwissTM to RSTM, TL2, and TinySTM for two reasons.

- They constitute the state-of-the-art performance-wise, among the publicly available library-based STMs. Furthermore, just like SwissTM, they can be used to manually instrument con-

<i>STM design choices</i>			
<i>Acquire</i>	<i>Reads</i>	<i>CM</i>	<i>Effectiveness</i>
lazy	invisible	any	+
eager	visible	any	+
eager	invisible	Polka	+
eager	invisible	timid or Greedy	++
mixed	invisible	timid or Greedy	+++
mixed	invisible	2-phase	++++

Table 1. A summary comparison of the effectiveness of selected combinations of STM design choices in mixed workloads.

current applications with transactional accesses. Indeed, our goal is to evaluate the performance of the core STM algorithm, not the efficiency of the higher layers such as STM compilers. We did not use for instance Intel’s STM [2], because it does not expose such a low-level API to a programmer. Evaluating STM-aware compilers (which naturally introduce additional overheads above the low-level STM interface [5]) is an orthogonal issue;

- They represent a wide spectrum of known TM design choices: *obstruction-free* vs. *lock-based* implementation, *eager* vs. *lazy* updates, *invisible* vs. *visible* reads, and *word-level* vs. *object-level* access granularity. They also allow for experiments with a variety of contention management strategies, from simply aborting a transaction on a conflict, through exponential back-off, up to advanced contention managers like Greedy [14], Serializer [26], or Polka [32].

We report on our SwissTM (trial-and-error) experience, which we believe is interesting in its own right. It is the first to date that evaluates the ability of software solutions to provide good performance to large transactions and complex objects without introducing significant overheads on short transactions and simple data structures. We evaluate the individual impact of various STM design choices on the ability to support mixed workloads. A summary of our observations, is presented in Table 1.

From an implementation perspective, we also evaluate the impact of the locking granularity. Word-based STM implementations used so far either word-level locking (e.g., TL2 and TinySTM) or cache-line level locking (e.g., Intel’s STM). Our sensitivity analysis shows that a lock granularity of four words (32 bit machine) outperforms both word level and cache-line level locking by 4% and 5% respectively across all benchmarks we considered.

To summarize, the main contributions of this paper are (1) the design and implementation of an STM that performs particularly well with large-scale complex transactional workloads while having good performance in small-scale ones, and (2) an extensive experimental evaluation of STM strategies and implementations from the perspective of complex applications with mixed workloads.

The rest of the paper is organized as follows. In Section 2, we give a short overview of STM design space and benchmarks. We then present SwissTM in Section 3. In Sections 4 and 5, we present the results of our experimental evaluation: first, we compare the performance of SwissTM to that of TL2, TinySTM, and RSTM, and, second, we evaluate the individual impact of the design choices underlying SwissTM.

2. Background

Transactional memory was first proposed in hardware (HTM) [20]. So far, HTMs support only limited-size transactions and often do not ensure transaction progress upon specific system events, e.g., interrupts, context switches, or function calls [8]. Hybrid approaches either execute short transactions in hardware and fall back

to software for longer ones (e.g., [22]), or accelerate certain operations of an STM in hardware. This work focuses on pure software solutions (STM) [29]. In this section, we survey some distinctive features of STMs and discuss the three representative STMs we focus on in our evaluation: RSTM [23], TL2 [9], and TinySTM [24] (see [21] for a full survey). We also give a short description of the benchmarks used in our experiments.

2.1 STM Design Space

The main task of an STM is to *detect* conflicts among concurrent transactions and *resolve* them. Deciding what to do when conflicts arise is performed by a (conceptually) separate component called a *contention manager* [19]. A concept closely related to conflict detection is that of *validation*. Validating a transaction consists of checking its read set (i.e., the set of locations¹ the transaction has already read) for consistency.

Two classes of STMs can be distinguished, *word*-based and *object*-based, depending on the granularity at which they perform logging. RSTM is object-based while TL2 and TinySTM are word-based. There are also two general classes of STM implementations: *lock*-based and *obstruction-free*. Lock-based STMs, first proposed in [11], implement some variant of the two-phase locking protocol [12]. Obstruction-free STMs [19] do not use any blocking mechanisms (such as locks), and guarantee progress even when some of the transactions are delayed. RSTM (version 3) is obstruction-free, while TL2 and TinySTM internally use locks.

Conflict detection. Most STMs employ the single-writer-multiple-readers strategy; accesses to the same location by concurrent transactions *conflict* when at least one of the accesses is a write (update). In order to commit, a transaction T must eventually *acquire* every location x that is updated by T . Acquisition can be *eager*, i.e., at the time of the first update operation of T on x , or *lazy*, i.e., at the commit time of T . A transaction T that reads x can be either *visible* or *invisible* [23] to other transactions accessing x . When T is invisible, T has the sole responsibility of detecting conflicts on x with transactions that write x concurrently, i.e., validating its read set. The time complexity of a basic validation algorithm is proportional to the size of the read set, but can be boosted with a global commit counter heuristic (RSTM), or a time-based scheme [9, 24] (TL2 and TinySTM).

A *mixed invalidation* conflict detection scheme (first proposed in [30]) eagerly detects write/write conflicts while lazily detecting read/write conflicts (it is a mix between pure lazy and pure eager schemes). A similar conflict detection scheme is provided by more general (but also more expensive) *multi-versioning* schemes used in LSA-STM [25] and JVSTM [6]. Mixed invalidation, which underlies SwissTM, has never been used with lock-based or word-based STMs, nor has it been evaluated with any large-scale workload.

RSTM supports lazy and eager acquisition, as well as visible and invisible reads (i.e., four algorithm variants). TL2 and TinySTM use, respectively, lazy and eager acquisition. Both TL2 and TinySTM employ invisible reads.

Contention management. The contention manager decides what a given transaction (*attacker*) should do in case of a conflict with another transaction (*victim*). Possible outcomes are: aborting the attacker, aborting the victim, or forcing the attacker to retry after some period.

The simplest scheme (which we call *timid*) is to always abort the attacker (possibly with a short back-off). This is the default scheme in TL2 and TinySTM. More involved contention managers were proposed in [32, 28, 14], and are provided with RSTM. They can

¹These are memory words in word-based STMs and objects in object-based STMs.

also be combined at run-time [13]. *Polka* [32] assigns every transaction a priority that is equal to the number of objects the transaction accessed so far. Whenever the attacker waits, its priority is temporarily increased by one. If the attacker has a lower priority than the victim, it will be forced to wait (using exponential back-off to calculate the wait interval), otherwise the victim gets aborted. *Greedy* assigns each transaction a unique, monotonically increasing timestamp on its start. The transaction with the lower timestamp always wins. An important property of Greedy is that, unlike other contention managers we mention, it avoids starvation of transactions. Polka has been shown to provide best performance in smaller-scale benchmarks previously [32], while our experiments show that Greedy performs better in large-scale workloads (Section 5). *Serializer* is very similar to Greedy except that it assigns a new timestamp to a transaction on every restart, and thus does not prevent starvation or even livelocks of transactions.

2.2 STM Benchmarks

In this section, we give an overview of the benchmarks we use in our experiments. These represent a large spectrum of workload types: from simple data structures with small transactions (the red-black tree microbenchmark) to complex applications with possibly long transactions (STMBench7).

STMBench7. STMBench7 [16] is a synthetic benchmark whose workloads aim at representing realistic, complex, object-oriented applications that are an important target for STMs. STMBench7 exhibits a large variety of operations (from very short, read-only operations to very long ones that modify large parts of the data structure) and workloads (from workloads consisting mostly of read-only transactions to write-dominated ones). The data structure used by STMBench7 is many orders of magnitude larger than in other typical STM benchmarks. Also, its transactions are longer and access larger numbers of objects.

STMBench7 is inherently object-based and its implementations also use standard language libraries. A thin wrapper, described in [10], is thus necessary to use STMBench7 with word-based STMs (TL2, TinySTM, and SwissTM).

STAMP. STAMP [7] is a TM benchmarking suite that consists of eight different transactional programs and ten workloads.² STAMP applications are representative of various real-world workloads, including bioinformatics, engineering, computer graphics, and machine learning. While STAMP covers a broad range of possible STM uses, it does not involve very long transactions, such as those that might be produced by average, non-expert programmers or generated automatically by a compiler along the lines of [1]. Furthermore, some STAMP algorithms (e.g., *bayes*) split logical operations into multiple transactions and use intricate programming techniques that might not be representative of average programmers' skill.

Lee-TM. Lee-TM [3] is a benchmark that offers large, realistic workloads and is based on Lee's circuit routing algorithm. The algorithm takes pairs of points (e.g., of an integrated circuit) as its input and produces non-intersecting routes between them. While transactions of Lee-TM are significant in size, they exhibit very regular access patterns—every transaction first reads a large number of locations (searching for suitable paths) and then updates a small number of them (setting up a path). Moreover, the benchmark uses very simple objects (each can be represented as a single integer variable). It is worth noting that STAMP contains an application (called *labyrinth*) that uses the same algorithm as Lee-TM. However, Lee-TM uses real-world input sets that make it more realistic than *labyrinth*.

²We used STAMP version 0.9.9.

Red-black tree. The prevailing way of measuring the performance of STMs has been through microbenchmarks. The widely used (first in [19]) red-black tree microbenchmark consists of short transactions that insert, lookup, and remove elements from a red black tree data structure. While short and simple transactions of microbenchmarks are good for testing mechanics of STM itself and comparing low-level details of various implementations, it is questionable whether they are representative of real-world systems.

So far, experimental evaluations of STM strategies [30, 32, 27] focused mostly on microbenchmarks. These gave little insight into the behavior of STMs in large-scale workloads. Recent STM papers, e.g., [2, 18, 5], focus on compiler optimizations in STM-aware compilers. Their results are complementary to ours: compiler optimizations and the strategies used by the base STM algorithm are orthogonal issues.

3. SwissTM

SwissTM is a lock-based STM that uses invisible reads and a counter based heuristics (the same as in TinySTM and TL2). It features eager write/write and lazy read/write conflict detection, as well as a two-phase contention manager with random linear back-off. The API of SwissTM is word-based, as it enables transactional access to arbitrary memory words. SwissTM uses a redo-logging scheme (partially to support its conflict detection scheme). Similarly to most other STM libraries, SwissTM guarantees opacity [15].

3.1 Algorithm

We give the pseudo-code of SwissTM in Algorithm 1. All transactions share a global commit counter *commit-ts* incremented by every non-read-only transaction upon commit. Each memory word *m* is mapped to a pair of locks in a global *lock table*: *r-lock* (read) and *w-lock* (write). Lock *w-lock* is acquired by a writer *T* of *m* (eagerly) to prevent other transactions from *writing* to *m*. Lock *r-lock* is acquired by *T* at commit time to prevent other transactions from *reading* word *m* and, as a result, observing inconsistent states of words written by *T*. Every transaction *T* has a *transaction descriptor tx* that contains (among other data): (1) the value of *commit-ts* read at the start or subsequent validation of *T*, and (2) read and write logs of *T*.

Transaction start. Every transaction *T*, upon its start, reads the global counter *commit-ts* and stores its value to *tx.valid-ts* (line 2).

Reading. When reading location *addr*, transaction *T* first reads the value of *w-lock* to detect possible read-after-write cases. If *T* is the owner of *w-lock*, then *T* can return the value from its write log immediately, which is the last value *T* has written to *addr* (line 6). Otherwise, i.e., when some other transaction owns *w-lock* or when *w-lock* is unlocked, *T* reads the value of *r-lock*, then the value of *addr*, and then again the value of *r-lock*. Transaction *T* repeats these three reads until (1) two values of *r-lock* are the same, meaning that *T* has read consistent values of *r-lock* and *addr*, and (2) *r-lock* is unlocked (lines 8–15). When *r-lock* is unlocked, it contains the current version *v* of *addr*. If *v* is lower or equal to the validation timestamp *tx.valid-ts* of *T* (which means that *addr* has not changed since *T*'s last validation or start), *T* returns the value at *addr* read at line 18. Otherwise, *T* revalidates its read set. If the re-validation does not succeed, *T* rolls back (line 17). If it succeeds, the read operation returns.

Writing. Whenever some transaction *T* writes to a memory location *addr*, *T* first checks if *T* is the owner of the lock *w-lock* corresponding to *addr*. If it is, *T* updates the value of *addr* in its write log and returns (lines 21–23). Otherwise, *T* tries to acquire *w-lock* by atomically replacing, using a compare-and-swap (CAS)

Algorithm 1: Pseudo-code representation of SwissTM.

```

1 function start(tx)
2   tx.valid-ts ← commit-ts;
3   cm-start(tx);

4 function read-word(tx, addr)
5   (r-lock, w-lock) ← map-addr-to-locks(addr);
6   if is-locked-by(w-lock, tx) then return get-value(w-lock, addr);
7   version ← read(r-lock);
8   while true do
9     if version = locked then
10      version ← read(r-lock);
11      continue;
12     value ← read(addr);
13     version2 ← read(r-lock);
14     if version = version2 then break;
15     version2 ← version;
16   add-to-read-log(tx, r-lock, version);
17   if version > tx.valid-ts and not extend(tx) then rollback(tx);
18   return value;

19 function write-word(tx, addr, value)
20   (r-lock, w-lock) ← map-addr-to-locks(addr);
21   if is-locked-by(w-lock, tx) then
22     update-log-entry(w-lock, addr, value);
23     return;
24   while true do
25     if is-locked(w-lock) then
26       if cm-should-abort(tx, w-lock) then rollback(tx);
27       else continue;
28     log-entry ← add-to-write-log(tx, w-lock, addr, value);
29     if compare&swap(w-lock, unlocked, log-entry) then
30       break;
31   if read(r-lock) > tx.valid-ts and not extend(tx) then
32     rollback(tx);
33   cm-on-write(tx);

34 function commit(tx)
35   if is-read-only(tx) then return;
36   for log-entry in tx.read-log do write(log-entry.r-lock, locked);
37   ts ← increment&get(commit-ts);
38   if ts > tx.valid-ts + 1 and not validate(tx) then
39     for log-entry in tx.read-log do
40       write(log-entry.r-lock, log-entry.version);
41     rollback(tx);
42   for log-entry in tx.write-log do
43     write(log-entry.r-lock, ts);
44     write(log-entry.w-lock, unlocked);

45 function rollback(tx)
46   for log-entry in tx.write-log do
47     write(log-entry.w-lock, unlocked);
48   cm-on-rollback(tx);

49 function validate(tx)
50   for log-entry in tx.read-log do
51     if log-entry.version ≠ read(log-entry.r-lock) and not
52     is-locked-by(log-entry.r-lock, tx) then return false;
53   return true;

53 function extend(tx)
54   ts ← read(commit-ts);
55   if validate(tx) then tx.valid-ts ← ts; return true;
56   return false;

```

operation, value `unlocked` with the pointer to the T 's write log entry that contains the new value of `addr` (line 29). If CAS does not succeed, T asks the contention manager whether to rollback and retry or wait for the current owner of the lock to finish (line 26). In order to guarantee opacity, T has to revalidate its read set if the current version of `addr` (contained in `r-lock`) is higher than its validity timestamp `tx.valid-ts` (lines 31–32).

Validation. To validate itself, T compares the versions of all memory locations read so far to their versions at the point they were initially read by T (lines 50–51). These versions are stored in T 's read log. If there is a mismatch between any version numbers, the validation fails (line 51). After a successful validation of its read set during read operation T (line 17) extends its validation timestamp `tx.valid-ts` to the current value of `commit-ts` (line 55).

Commit. A read-only transaction T can commit immediately, as its read log is guaranteed to be consistent (line 35). A transaction T that is not read-only first locks all read locks of memory locations T has written to (line 36). Then, T increments `commit-ts` (line 37) and re-validates its read log. If the validation does not succeed, T rollbacks and restarts (lines 38–41). Upon successful validation, T traverses its write set, updates values of all written memory locations, and releases the corresponding read and write locks (lines 42–44). When releasing read locks, T writes the new value of `commit-ts` to those locks.

Rollback. On rollback, transaction T releases all write locks it holds (lines 46–47), and then restarts itself.

Algorithm 2: Pseudo-code of the two-phase contention manager.

```

1 function cm-start(tx)
2   if not-restart(tx) then tx.cm-ts ← ∞;

3 function cm-on-write(tx)
4   if tx.cm-ts = ∞ and size(tx.write-log) =  $W_n$  then
     tx.cm-ts ← increment&get(greedy-ts);

5 function cm-should-abort(tx, w-lock)
6   lock-owner = owner(w-lock);
7   if lock-owner.cm-ts < tx.cm-ts then return true;
8   else abort(lock-owner); return false;

9 function cm-on-rollback(tx)
10  wait-random(tx.succ-abort-count);

```

Contention management. We give the pseudo-code of our two-phase contention manager in Algorithm 2. The contention manager gets invoked by the main algorithm (1) at transaction start (line 3), (2) on a write/write conflict (line 26), (3) after a successful write (line 33), and (4) after restart (line 48). Every transaction, upon executing its W_n^{th} write, increments global counter `greedy-ts` and stores its value to `tx.cm-ts` (line 4). Upon a conflict, the transaction with the higher value of `cm-ts` is aborted (lines 7–8). Conceptually, every transaction starts with an infinite value of `cm-ts`. After restarting, transactions are delayed using a randomized back-off scheme (line 10).

3.2 Implementation Highlights

We implemented SwissTM in C++ (g++ 4.0.1 compiler). We used the (fairly portable) `atomic_ops` library [4] for atomic operations implementation. Currently, SwissTM works on 32-bit x86 Linux 2.6.x and OS X 10.5 platforms (64-bit port is in progress).

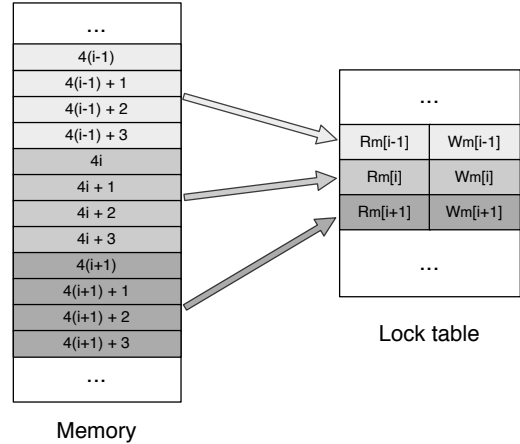


Figure 1. Mapping of memory words to lock table entries.

Lock table. To map memory word m to a lock table entry, we take the address a of m , shift it to the right by 4 (it would be 5 with 64 bit words). This makes each lock map to consecutive four memory words (we empirically selected this value, as explained in Section 5). Then, we set all high order bits to zero. As the lock table contains 2^{22} entries in our implementation, we just perform logical AND operation between shifted address and $2^{22} - 1$ to get the index into the table. Figure 1 depicts the mapping scheme.

Every lock is implemented as a single memory word w . The write lock is equal to 0 in its unlocked state and contains a pointer to the corresponding write log entry when locked. Acquiring a write lock is done using a compare-and-swap (CAS) operation. When releasing write lock w , transactions simply write 0 to w . The read lock has its least significant bit set to 0 when unlocked, while other bits store version number of memory locations corresponding to w . When locked, read lock is equal to 1. Both locking and unlocking of read locks is performed by simply writing a new value to w (CAS is not used), as only the transaction that already acquired a write lock can acquire the corresponding read lock.

4. Evaluating SwissTM

We compare the performance of SwissTM to that of TL2, TinySTM, and RSTM.

We performed all measurements on a 4 processor dual-core AMD Opteron 8216 2.4 GHz 1024 KB cache machine with 8 GB of RAM running Linux operating system. This provided us with 8 cores to experiment on. All results were averaged over multiple runs, where the length and the number of runs were chosen to reduce variations in collected data. We typically used 20 runs for STMBench7 and STAMP, 10 runs for LeeTM and 80 runs for red-black tree microbenchmark. We used the TL2 x86 implementation provided with the STAMP benchmark suite (version 0.9.5), and the RSTM (version 3) and TinySTM (version 0.9.5) implementations available from respective sites. Unless stated otherwise, we configured RSTM to use eager conflict detection, invisible reads with the commit counter heuristic, and the Polka contention manager. We used default configurations of TL2 (i.e., lazy conflict detection, GV4) and TinySTM (i.e., encounter time locking, timid contention manager).

STMBench7. Figure 2 shows the performance of SwissTM, TL2, TinySTM, and RSTM with STMBench7. We configured RSTM to use the Serializer contention manager, as this gave best performing RSTM configuration in STMBench7. SwissTM significantly out-

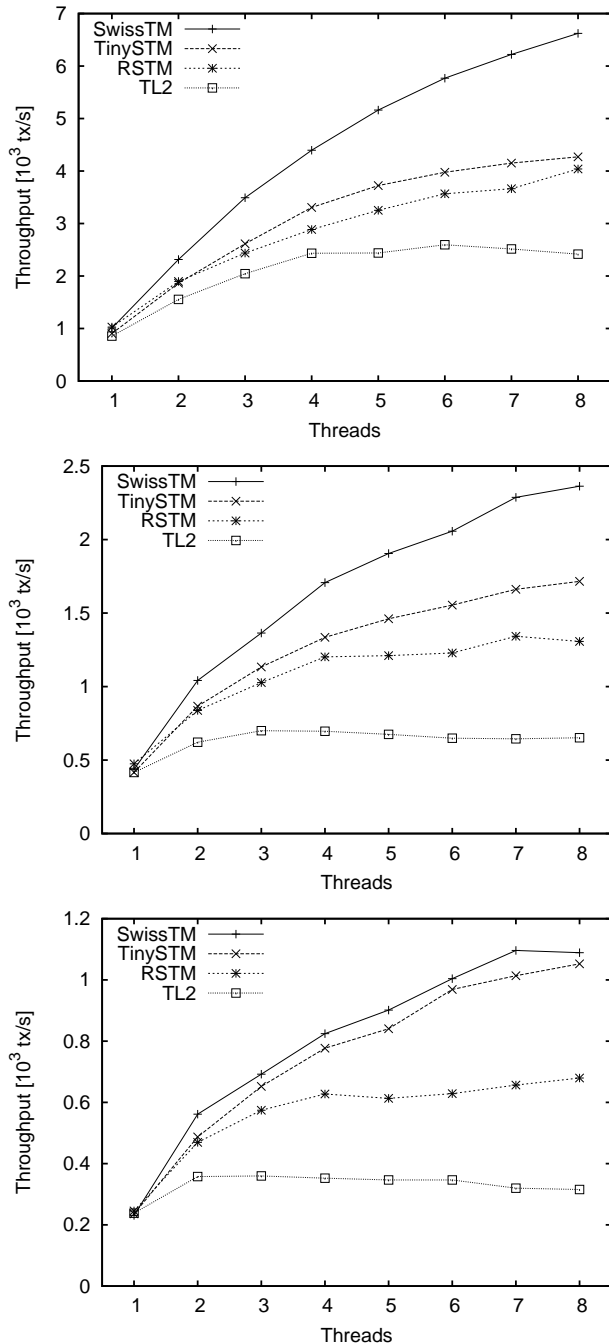


Figure 2. Throughput of SwissTM, RSTM, TL2, and TinySTM with STMBench7; top to bottom: read-dominated, read-write, and write-dominated workload

performs all other STMs for both read-dominated and read-write workloads, while also achieving superior scalability. SwissTM also outperforms other STMs in high-contention write-dominated workload, but it is only marginally faster than TinySTM.

The main reason for the good performance of SwissTM is (a) its optimism in detecting read/write conflicts when compared to RSTM and TinySTM, and (b) its conservatism in detecting write/write conflicts when compared to TL2. The contention man-

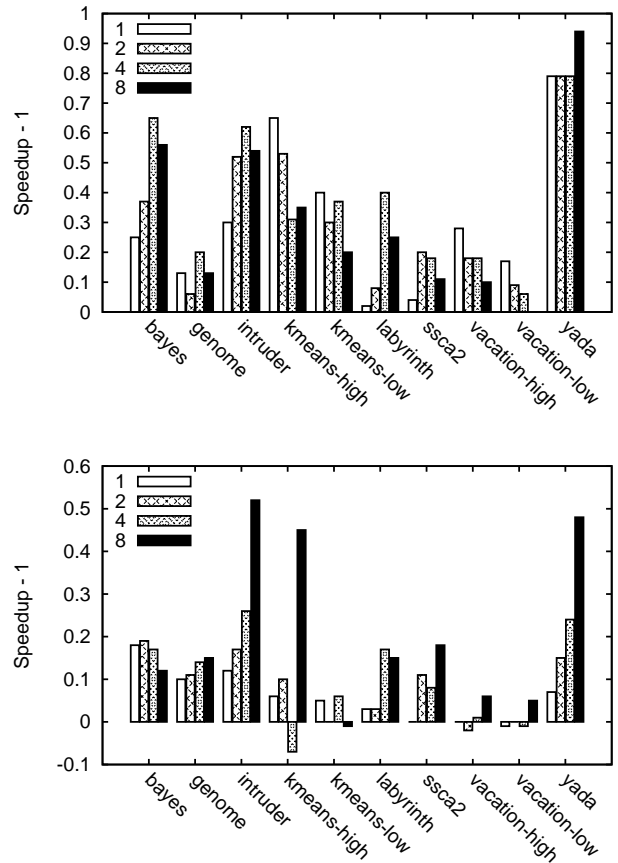


Figure 3. SwissTM vs. TL2 (top) and TinySTM (bottom) in STAMP. The figure conveys the speedup of SwissTM compared to TL2 and TinySTM (with 1 subtracted, i.e., positive numbers mean that SwissTM is faster and negative that it is slower) for 1, 2, 4, and 8 threads.

agement scheme used in SwissTM also helps boost performance, as we illustrate in Section 5.

TL2 performs poorly even in the read-dominated workload—it does not scale after 4 threads. Its performance gets even worse with higher contention. The main reason for this is the lazy conflict detection scheme of TL2, which wastes more work of transactions than the eager write/write conflict detection of other STMs.

STAMP. Figure 3 compares the performance of SwissTM, TL2, and TinySTM in the STAMP benchmark suite workloads.³ SwissTM outperforms TL2 in all STAMP workloads, for all thread counts, excluding the *vacation* benchmark under low contention where TL2 and SwissTM have the same performance. SwissTM outperforms TL2 by over 50% with eight threads for the *bayes*, *intruder*, and *yada* benchmarks (being almost twice as fast as TL2 in *yada*), and by over 20% in *kmeans* (both variants) and *labyrinth*, while being about 10% faster than TL2 in *genome*, *ssc2*, and *vacation* under high contention. SwissTM outperforms TinySTM in ten STAMP workloads with eight threads, except for the *kmeans* benchmark under low contention where TinySTM has a slightly better performance (1% of difference). SwissTM outperforms TinySTM by over 45% with eight threads in *intruder*, *kmeans* under high contention, and *yada*, and by over 12% in

³ There is no RSTM implementation of STAMP.

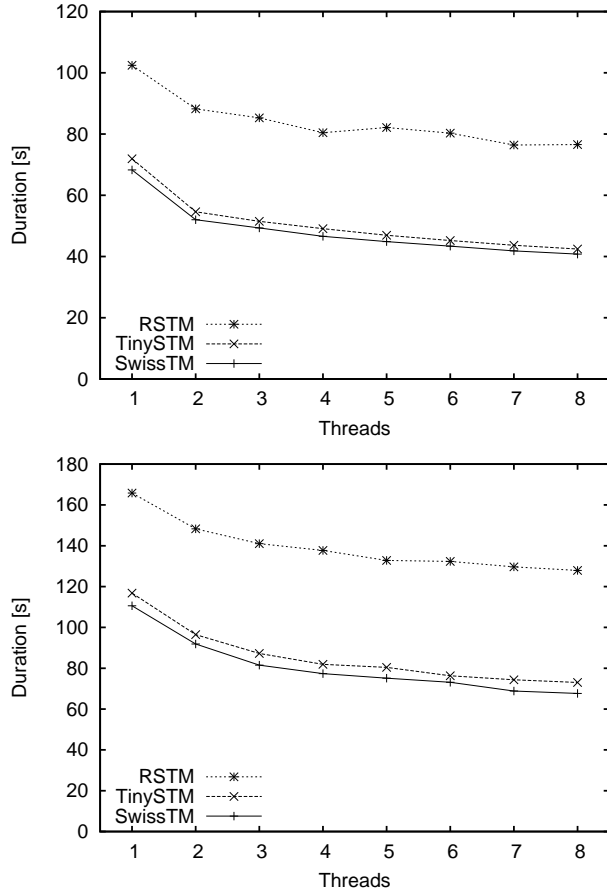


Figure 4. Execution time of SwissTM, TinySTM and RSTM in the Lee-TM benchmark; top is *memory*, bottom is *main board* input

bayes, *genome*, *labyrinth*, and *ssca2*, while being about 5% faster in *vacation* (both variants). SwissTM has good performance with lower thread counts, while scaling well as the number of concurrent threads increases. To summarize, SwissTM outperforms both TL2 and TinySTM in all STAMP benchmark workloads.

Lee-TM. Figure 4 compares the performance of SwissTM, RSTM, and TinySTM in the Lee-TM benchmark.⁴ RSTM has the lowest performance mainly because of higher single object access overheads (objects in Lee-TM are very simple—consisting of a single integer variable). SwissTM and TinySTM have very similar performance, although SwissTM is faster by a small margin for all thread counts.

Red-black tree. Finally, Figure 5 compares the performance of SwissTM to TL2, TinySTM, and RSTM in the commonly used red-black tree microbenchmark. RSTM and TL2 deliver significantly lower performance than both TinySTM and SwissTM due to their higher overheads on single memory location accesses. These low-level overheads have most significant impact in microbenchmarks like this one. SwissTM outperforms TinySTM after four threads (except for the small drop at six threads) by a small margin.⁵

⁴ When running Lee-TM with TL2, the program hangs even with a single thread.

⁵ It is worth noting that in this particular benchmark an older version of TinySTM outperformed SwissTM by a small and constant margin (< 5%).

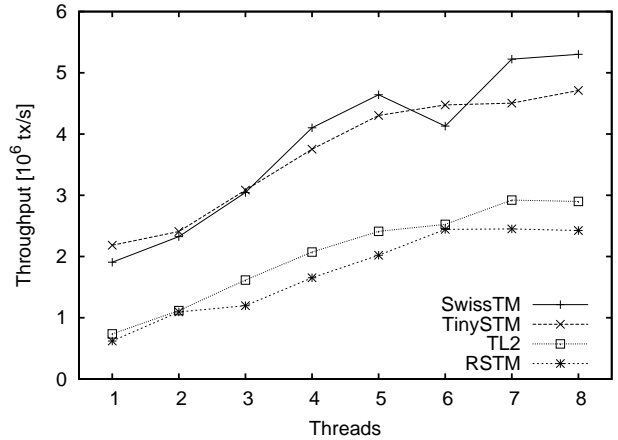


Figure 5. Throughput of SwissTM, TL2, TinySTM, and RSTM on red-black tree with range of 16384 and 20% of update operations

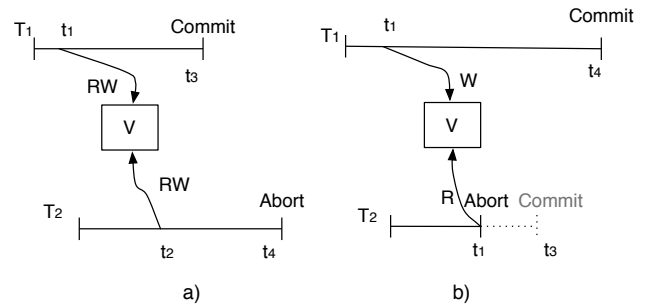


Figure 6. Disadvantages of lazy (left) and eager (right) conflict detection strategies

5. Dissecting the SwissTM Experience

We evaluate here individually the design choices underlying SwissTM: its conflict detection strategy, the two-phase contention manager, and the granularity of the lock table.

Conflict detection. Current state-of-the-art STMs typically detect both read/write and write/write conflicts in the same way—either as soon as conflicts occur (eagerly, e.g., TinySTM), or at commit time (lazily, e.g., TL2). Detecting conflicts eagerly helps avoid wasting work of transactions that are doomed to abort after a conflict. Lazy conflict detection, however, is more optimistic and gives transactions more possibilities to commit. For example, Figure 6a depicts an execution of an STM that uses lazy conflict detection. There, transaction T_2 spends time between t_3 (commit time of T_1) and t_4 (commit time of T_2) performing work that is doomed to be rolled-back. The period between t_3 and t_4 can be significant with long transactions. Figure 6b, gives an example execution of an STM that uses eager conflict detection. There, transaction T_2 has to wait until time t_4 before continuing, although it could commit already at time t_3 if lazy scheme was used. The waiting time of T_2 might be significant if T_1 is very long.

SwissTM takes the best of both strategies—it detects write/write conflicts eagerly and read/write conflicts lazily. This combined strategy is beneficial for complex workloads with long transactions because it (1) prevents transactions with write/write conflicts from running for a long time before detecting the conflict, and (2) allows

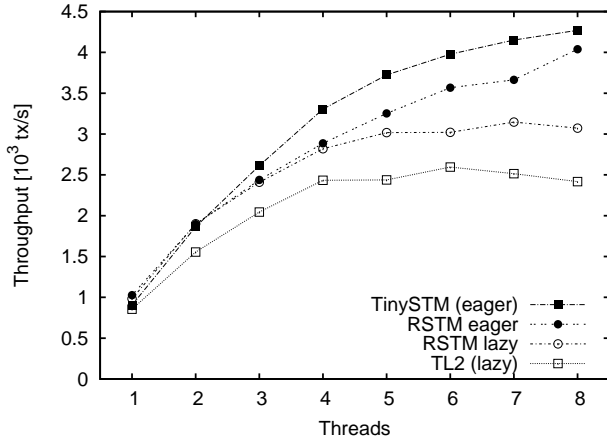


Figure 7. Throughput of eager and lazy conflict detection schemes in read-dominated STMBench7

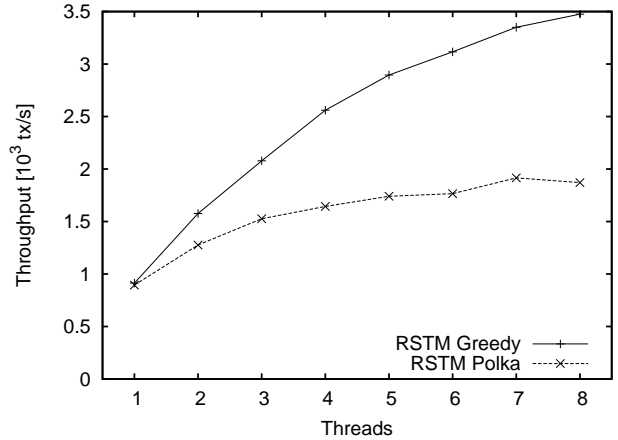


Figure 9. Throughput of Polka vs. Greedy (in RSTM) in the read-dominated workload of STMBench7

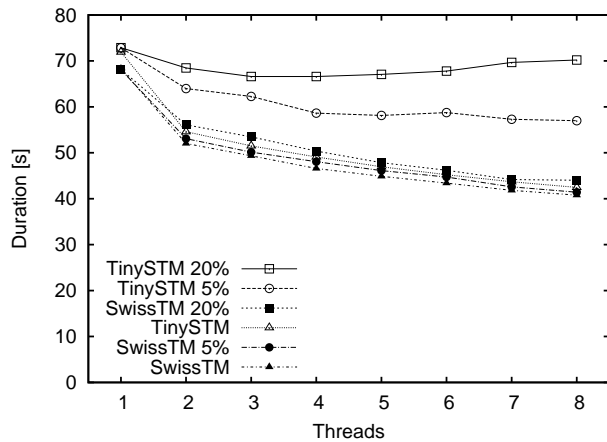


Figure 8. Execution time of SwissTM vs. TinySTM in “irregular” Lee-TM benchmark with *memory board* input

short transactions having a read/write conflict with longer ones to proceed, thus increasing parallelism.

Figure 2 suggests that the mixed eager-lazy scheme gives better performance than pure eager scheme which, in turn, outperforms lazy scheme (Figure 7). This does not say what kind of workloads benefit most from the mixed scheme, and what part of the performance boost of SwissTM can be attributed to its two-phase contention manager (and not to the mixed conflict detection).

To answer this question, we modified slightly the Lee-TM benchmark. The performance of the original Lee-TM does not seem to be influenced by the choice of a conflict detection and contention management schemes, because the transactions in Lee-TM are highly regular—they first read and then write. We introduce a small irregularity in Lee-TM by adding a single object O_c that every transaction reads at its start. A small ratio R of transactions (chosen randomly) also updates O_c , causing a read/write conflict with all the other transactions. The contention manager used by SwissTM does not provide a lot of benefit in this case, as the number of write/write conflicts introduced by this irregularity is not large (we keep R low in experiments).

Figure 8 compares the performance of TinySTM and SwissTM for R of 0%, 5%, and 20%. Due to its conflict detection scheme, the SwissTM performance degrades only slightly even when R is

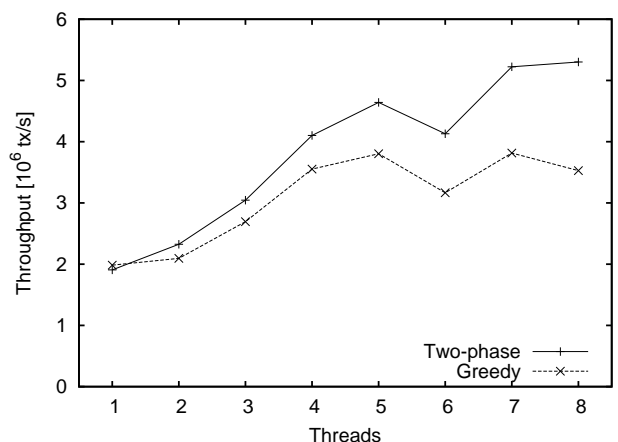


Figure 10. Throughput of the two-phase contention manager vs. Greedy (in SwissTM) with the red-black tree

relatively high (20%). Also, SwissTM continues to scale well as the number of threads increases. On the other hand, the performance of TinySTM performance degrades significantly even when R is only 5%, while with R of 20% it stops scaling already at three threads.

We conclude here that applications exhibiting regular access patterns benefit the most from lowering single-location access costs and are not significantly influenced by the conflict detection scheme itself. However, for applications where the access patterns introduce even small irregularities, especially those creating longer-lasting read/write conflicts, SwissTM’s optimistic approach yields significant benefits.

Contention management. The contention manager gets invoked in SwissTM only on write/write conflicts. The reasoning here is simple—read/write conflicts are detected only at commit time of the writer, and it makes little sense for the reader to abort the writer that is already committing. This is why the reader waits until the writer (quickly) commits, before attempting to revalidate its read set.

Figure 9 shows that Greedy performs better than Polka (which was shown previously to perform very well in a range of smaller scale benchmarks [32]) in larger-scale workloads. However, Greedy performs poorly with short transactions, because all

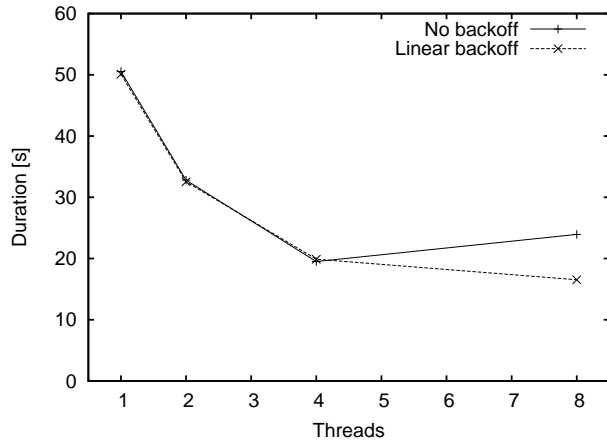


Figure 11. Execution time of backoff vs. no backoff (in SwissTM) in the STAMP *intruder* application

transactions increment a single shared counter at their start, which causes a lot of cache misses and significantly degrades performance and scalability (Figure 10). This problem is not noticeable with longer transactions as the overhead caused by cache misses is relatively small compared to the work of the transactions themselves. As shown in Figure 10, our two-phase contention manager overcomes this issue completely, improving both performance and scalability over Greedy. This is because it allows all short and read-only transactions to commit without incrementing the shared counter used by the Greedy algorithm, yet it provides the strong progress guarantees of Greedy for long ones.

It might seem beneficial to make transactions restart as soon as possible after conflicts that force them to rollback, as waiting just decreases the reaction time before the transaction re-executes. However, restarting immediately tends to increase contention on cache lines containing data that gets updated very frequently. Consequently, short back-offs after transaction rollbacks can improve performance. Figure 11 compares the performance of SwissTM in the STAMP *intruder* benchmark with and without the back-off scheme. (The *intruder* benchmark is a good example here, because it indeed contains a “hot spot”: a high number of transactions dequeue elements from a single queue.) The figure shows that immediately restarting transactions after rollback causes scalability problem with eight threads. A simple randomized linear (in the number of successive aborts) back-off scheme resolves the scalability issue.

Finally, we evaluate the influence of our two-phase contention manager on the overall performance of SwissTM. Figure 12 shows that the two-phase contention manager improves performance by as much as 16% in high-contention workloads. Its influence is lower in the read-dominated workload, which is not surprising given that this workload is characterized by a small number of write/write conflicts.

Locking granularity. An important implementation choice in an STM is its lock table configuration, in particular the size of the memory stripe that gets mapped to the same (lock table) entry. Increasing the size of memory stripes reduces locking and validation time, due to the data access locality, but increases abort rates by introducing false conflicts when the memory stripe becomes too large. The optimal value for this parameter is application specific and we searched for the best value across all benchmarks we used. Figure 13 depicts the average speedup (minus 1) of each logarithmic lock granularity compared to all the others at eight threads (32 bit word). The figure shows that the granularity of 2^4 bytes achieves

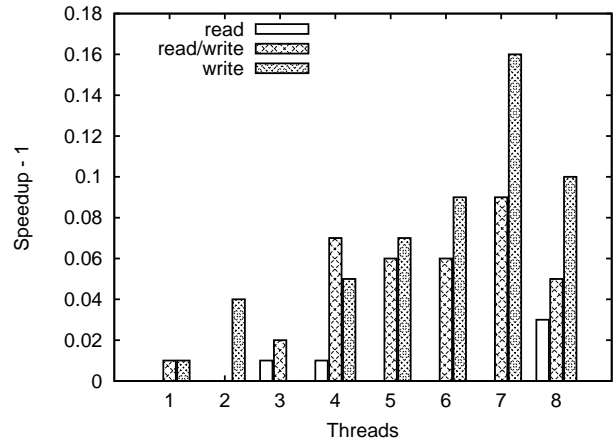


Figure 12. Speedup (with 1 subtracted) of the two-phase contention manager over the timid contention manager in SwissTM with STMBench7

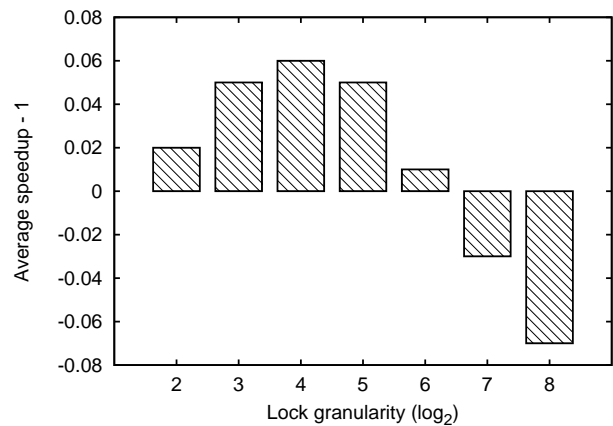


Figure 13. Average speedup (with one subtracted) of locking granularities from 2^2 to 2^8 compared to all other granularities

the best performance, with 2^3 and 2^5 being slightly slower. It is interesting to note that the commonly used sizes of one word (2^2) and one cache line (2^6) have performance of 4% and 5% lower on average than the one we select. We give a breakdown of these differences across benchmarks in Table 2.

6. Concluding Remarks

This paper presents SwissTM—an effective compilation of STM design choices for mixed workloads characterized by non-uniform, dynamic data structures and various transaction sizes. Those kinds of workloads are inherent to many applications that might be expected to significantly benefit from the STM paradigm and multi-core architectures. SwissTM significantly outperforms state-of-the-art STMs in precisely such workloads, while also delivering good performance in smaller-scale scenarios.

Not surprisingly, the design of SwissTM is a result of trial-and-error. We reported in the paper on various choices that might have seemed natural, but revealed inappropriate. Besides those, we also experimented with nested transactions (closed nesting) and multi-versioning, but we could not see a clear advantage of those

Benchmark	Difference in performance		
	2 ⁴ vs. 2 ²	2 ⁴ vs. 2 ⁶	2 ² vs. 2 ⁶
bayes	0.16	0.81	0.57
genome	0.13	-0.03	-0.14
intruder	0	-0.04	-0.04
kmeans-high	0.19	0.4	0.18
kmeans-low	0.14	0.05	-0.08
labyrinth	-0.12	-0.09	0.04
ssca2	0	0	0
vacation-high	0.14	-0.03	-0.15
vacation-low	0.12	-0.05	-0.15
yada	0	0.12	0.12
red-black tree	-0.01	0	0.01
Lee-TM memory	0.01	-0.03	-0.04
Lee-TM main	0.02	-0.01	-0.02
STMBench7 read	0	-0.02	-0.02
STMBench7 read-write	-0.01	-0.03	-0.02
STMBench7 write	-0.04	-0.06	-0.02
Average	0.05	0.06	0.01

Table 2. Comparing three different lock granularities

techniques in the considered workloads. Further experiments might be needed in this direction.

7. Availability

SwissTM is open-source and available from: [url_omitted_for_submission].⁶ The source code of benchmarks used in the paper can be downloaded from the same location.

References

- [1] M. Abadi, A. Birrell, T. Harris, and M. Isard. Semantics of transactional memory and automatic mutual exclusion. In *POPL*, 2008.
- [2] A.-R. Adl-Tabatabai, B. T. Lewis, V. Menon, B. R. Murphy, B. Saha, and T. Shpeisman. Compiler and runtime support for efficient software transactional memory. In *PLDI*, 2006.
- [3] M. Ansari, C. Kotselidis, K. Jarvis, M. Luján, C. Kirkham, and I. Watson. Lee-TM: A non-trivial benchmark for transactional memory. In *ICA3PP*, 2008.
- [4] The atomic_ops project. http://www.hpl.hp.com/research/linux/atomic_ops.
- [5] C. Blundell, H. Cain, M. M. Michael, P. Wu, S. Chiras, and S. Chatterjee. Software transactional memory: Why is it only a research toy? *ACM Queue*, Sept. 2008.
- [6] J. Cachopo and A. Rito-Silva. Versioned boxes as the basis for memory transactions. *Sci. Comput. Program.*, 63(2):172–185, 2006.
- [7] C. Cao Minh, J. Chung, C. Kozyrakas, and K. Olukotun. STAMP: Stanford transactional applications for multi-processing. In *IISWC*, 2008.
- [8] D. Dice, M. Herlihy, D. Lea, Y. Lev, V. Luchangco, W. Mesard, M. Moir, K. Moore, and D. Nussbaum. Applications of the adaptive transactional memory test platform. In *TRANSACT*, 2008.
- [9] D. Dice, O. Shalev, and N. Shavit. Transactional locking II. In *DISC*, 2006.
- [10] A. Dragojević, R. Guerraoui, and M. Kapalka. Dividing transactional memories by zero. In *TRANSACT*, 2008.
- [11] R. Ennals. Software transactional memory should not be obstruction-free. Technical Report IRC-TR-06-052, Intel Research Cambridge Tech Report, Jan. 2006.
- [12] K. P. Eswaran, J. N. Gray, R. A. Lorie, and I. L. Traiger. The notions of consistency and predicate locks in a database system. *Commun. ACM*, 19(11):624–633, 1976.
- [13] R. Guerraoui, M. Herlihy, and B. Pochon. Polymorphic Contention Management. In *DISC*, 2005.
- [14] R. Guerraoui, M. Herlihy, and B. Pochon. Toward a theory of transactional contention managers. In *PODC*, 2005.
- [15] R. Guerraoui and M. Kapalka. On the correctness of transactional memory. In *PPoPP*, 2008.
- [16] R. Guerraoui, M. Kapalka, and J. Vitek. STMBench7: A benchmark for software transactional memory. In *EuroSys*, 2007.
- [17] T. Harris, S. Marlow, S. Peyton-Jones, and M. Herlihy. Composable memory transactions. In *PPoPP*, 2005.
- [18] T. Harris, M. Plesko, A. Shinnar, and D. Tarditi. Optimizing memory transactions. In *PLDI*, 2006.
- [19] M. Herlihy, V. Luchangco, M. Moir, and W. Scherer. Software transactional memory for dynamic-sized data structures. In *PODC*, 2003.
- [20] M. Herlihy and J. E. B. Moss. Transactional memory: Architectural support for lock-free data structures. In *ISCA*, 1993.
- [21] J. R. Larus and R. Rajwar. *Transactional Memory*. Morgan&Claypool, 2007.
- [22] Y. Lev, M. Moir, and D. Nussbaum. PhTM: Phased transactional memory. In *TRANSACT*, 2007.
- [23] V. J. Marathe, M. F. Spear, C. Heriot, A. Acharya, D. Eisenstat, W. N. Scherer III, and M. L. Scott. Lowering the overhead of software transactional memory. In *TRANSACT*, 2006.
- [24] T. R. Pascal Felber and C. Fetzer. Dynamic performance tuning of word-based software transactional memory. In *PPoPP*, 2008.
- [25] T. Riegel, P. Felber, and C. Fetzer. A lazy snapshot algorithm with eager validation. In *DISC*, 2006.
- [26] RSTM home page. <http://www.cs.rochester.edu/research/synchronization/rstm>.
- [27] B. Saha, A.-R. Adl-Tabatabai, R. L. Hudson, C. Cao Minh, and B. Hertzberg. McRT-STM: a high performance software transactional memory system for a multi-core runtime. In *PPoPP*, 2006.
- [28] W. N. Scherer III and M. L. Scott. Contention management in dynamic software transactional memory. In *CSJP*, 2004.
- [29] N. Shavit and D. Touitou. Software transactional memory. In *PODC*, 1995.
- [30] M. F. Spear, V. J. Marathe, W. N. Scherer III, and M. L. Scott. Conflict detection and validation strategies for software transactional memory. In *DISC*, 2006.
- [31] T. Sweeney. The next mainstream programming language: a game developer’s perspective. Invited talk at POPL, 2006.
- [32] I. William N. Scherer and M. L. Scott. Advanced contention management for dynamic software transactional memory. In *PODC*, 2005.

⁶The url can be obtained by the committee chair on request.