

The Semantics of Progress in Lock-Based Transactional Memory

Rachid Guerraoui Michał Kapalka

EPFL, Switzerland

{rachid.guerraoui, michal.kapalka}@epfl.ch

Abstract

Transactional memory (TM) is a promising paradigm for concurrent programming. Whereas the number of TM implementations is growing, however, little research has been conducted to precisely define TM semantics, especially their progress guarantees. This paper is the first to formally define the progress semantics of lock-based TMs, which are considered the most effective in practice.

We use our semantics to reduce the problems of reasoning about the correctness and computability power of lock-based TMs to those of simple try-lock objects. More specifically, we prove that checking the progress of any set of transactions accessing an arbitrarily large set of shared variables can be reduced to verifying a simple property of each individual (logical) try-lock used by those transactions. We use this theorem to determine the correctness of state-of-the-art lock-based TMs and highlight various configuration ambiguities. We also prove that lock-based TMs have consensus number 2. This means that, on the one hand, a lock-based TM cannot be implemented using only read-write memory, but, on the other hand, it does not need very powerful instructions such as the commonly used compare-and-swap.

We finally use our semantics to formally capture an inherent trade-off in the performance of lock-based TM implementations. Namely, we show that the space complexity of every lock-based software TM implementation that uses invisible reads is at least exponential in the number of objects accessible to transactions.

Categories and Subject Descriptors D.1.3 [Programming Techniques]: Concurrent Programming; F.2.2 [Analysis of Algorithms and Problem Complexity]: Nonnumerical Algorithms and Problems

General Terms Theory, Algorithms, Verification

Keywords Transactional memory, lock, try-lock, consensus number, impossibility, lower bound, reduction, semantics

1. Introduction

Multi-core processors are predicted to be common in home computers, laptops, and maybe even smoke detectors. To exploit the power of modern hardware, applications will need to become increasingly parallel. However, writing scalable concurrent programs is hard and error-prone with traditional locking techniques. On the one hand, coarse-grained locking throttles parallelism and causes

lock contention. On the other hand, fine-grained locking is usually an engineering challenge, and as such is not suitable for use by the masses of programmers.

Transactional memory (TM) (Herlihy and Moss 1993) is a promising technique to facilitate concurrent programming while delivering comparable performance to fine-grained locking implementations. In short, a TM allows concurrent threads of an application to communicate by executing lightweight, in-memory *transactions*. A transaction accesses shared data and then either commits or aborts. If it commits, its operations are applied to the shared state *atomically*. If it aborts, however, its changes to the shared data are lost and never visible to other transactions.

While a large number of TM implementations have been proposed so far, there is still no precise and complete description of the *semantics* of a TM. Indeed, a correctness criterion for TM, called *opacity*, has been proposed (Guerraoui and Kapalka 2008a), and the progress properties of *obstruction-free* TM implementations have been defined (Guerraoui and Kapalka 2008c). However, opacity is only concerned with safety—it does not specify when transactions need to commit. (For example, a TM that aborts every transaction could trivially ensure opacity.) Moreover, TM implementations that are considered effective (Ennals 2006), e.g., TL2 (Dice et al. 2006), TinySTM (Felber et al. 2008), a version of RSTM (Marathe et al. 2006), BartokSTM (Harris et al. 2006), or McRT-STM (Adl-Tabatabai et al. 2006) are not obstruction-free. They internally use locking, in order to reduce the overheads of TM mechanisms, and do not ensure obstruction-freedom, which inherently precludes the use of locks.

Lock-based TMs do ensure some progress for transactions, for otherwise nobody would use them. However, this has never been precisely defined. The lack of such a definition hampers the portability of applications that use lock-based TMs, and makes it difficult to reason formally about their correctness or to establish whether any performance limitation is inherent or simply an artifact of a specific implementation.

This paper defines the progress semantics of lock-based TMs. We do so by introducing a new property, which we call *strong progressiveness*,¹ and which stipulates the two following requirements.

1. A transaction that encounters no *conflict* must be able to commit. (Basically, a conflict occurs when two or more concurrent transactions access the same transactional variable and at least one of those accesses is not read-only.)
2. If a number of transactions have only a “simple” conflict, i.e., on a single transactional variable, then at least one of them must be able to commit.

The former property captures the common intuition about the progress of any TM (see (Scott 2006)). The second property ensures that conflicts that are easy to resolve do not cause all con-

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

POPL’09, January 18–24, 2009, Savannah, Georgia, USA.
Copyright © 2009 ACM 978-1-60558-379-2/09/01...\$5.00

¹ We call it “strong” by opposition to a weaker form of progressiveness that we also introduce in this paper.

flicting transactions to be aborted. This is especially important when non-transactional accesses to shared variables are encapsulated inside unit transactions to ensure strong atomicity (Blundell et al. 2006). Strong progressiveness, together with opacity and operation-level wait-freedom,² is ensured by state of the art lock-based implementations, such as TL2, TinySTM, RSTM, BartokSTM, and McRT-STM.³

We use our strong progressive semantics to reduce the problems of reasoning about the correctness and computability power of lock-based TMs to those of simple *try-lock* objects (Scott and Scherer III 2001; Jayanti 2003). We first show that proving strong progressiveness of a set of transactions accessing any number of shared variables can be reduced to proving a simple property of every individual logical try-lock that protects those variables. Basically, we prove that if it is possible to say which parts of a TM algorithm can be viewed as logical try-locks (in a precise sense we define in the paper), and if those logical try-locks are *strong*, then the TM is strongly progressive. Intuitively, a try-lock is strong if it guarantees that among processes that compete for the unlocked try-lock, one always acquires the try-lock (most try-locks in the literature that are implemented from compare-and-swap or test-and-set are strong). We illustrate our reduction approach on state-of-the-art lock-based TMs. We formally establish and prove their correctness while highlighting some of their configurations that, maybe unexpectedly, violate the progress semantics.

Then, still using the try-lock reduction, we show that a lock-based TM has *consensus number 2* in the parlance of (Herlihy 1991). The consensus number is a commonly used metric for the computational power of a shared-memory abstraction, and is expressed as the maximum number of processes that can solve a non-trivial agreement problem (namely consensus (Herlihy 1991)) in a wait-free manner using this abstraction. The fact that a lock-based TM has consensus number 2 means that such a TM cannot be implemented using only read-write memory instructions, but, on the other hand, powerful instructions such as compare-and-swap are not necessary to implement a lock-based TM.

In fact, we give an implementation of a lock-based TM using read-write and test-and-set instructions. This implementation might be interesting in its own right when compare-and-swap instructions are not available or simply too expensive. Interestingly, we highlight an alternative semantics we call *weak progressiveness* which enables TMs with consensus number 1 and can thus be implemented using only read-write memory. Intuitively, weak progressiveness requires only that a transaction that encounters no conflicts commits. This might be considered a viable alternative to strong progressiveness for “lightweight” lock-based implementations.

We finally use our progress semantics to determine an inherent trade-off between the required memory and the latency of reads in lock-based TMs. This trade-off impacts the performance and/or progress guarantees of a TM but it was never formally established, precisely because of the lack of any precise semantics. We show that the space complexity of every lock-based TM that uses the *invisible reads* strategy⁴ is at least exponential in the number of variables available to transactions. This might seem surprising, since it

² Wait-freedom (Herlihy 1991) requires threads executing operations on transactional data within transactions to make progress independently, i.e., without waiting for each other. Maybe surprisingly, this property can easily be ensured by lock-based TMs.

³ The source code of the implementations of BartokSTM and McRT-STM is not publicly available. We could thus verify strong progressiveness of those TMs only from their algorithm descriptions in (Harris et al. 2006) and (Adl-Tabatabai et al. 2006), respectively.

⁴ With invisible reads, the reading of transactional variables is performed optimistically, without any (shared or exclusive) locking or updates to

is not obvious that modern lock-based TMs have non-linear space complexity. The exponential (or, in fact, unbounded) complexity comes from the use of timestamps that determine version numbers of shared variables. TM implementations usually reserve a constant-size word for each version number (which gives linear space complexity). However, an overflow can happen and has to be handled in order to guarantee correctness (opacity). As we explain in Section 6.3, this requires (a) limiting the progress of transactions when overflow occurs and (b) preventing read-only transactions from being completely invisible. Concretely speaking, our result means that efficient TM implementations (the ones that use invisible reads) must either intermittently (albeit very rarely) violate progress guarantees, or use unbounded timestamps.

Summary of contributions. To summarize, this paper contributes to the understanding of TM design and implementations by presenting the first precise semantics of a large class of popular TMs—lock-based ones. We precisely define the progress semantics of such TMs and propose reduction approaches to simplify their verification and computational study. We also use our semantics to study their inherent performance bottlenecks.

Roadmap. The rest of the paper proceeds as follows. First, in Section 2, we describe the basic model and terminology used to state our semantics and prove our results. Then, in Section 3, we define the progress semantics of lock-based TMs. In Section 4, we show how to simplify the verification of strong progressiveness. Next, in Sections 5 and 6, we establish the fundamental power and limitations of lock-based TMs. We also discuss in those sections the impact of weakening progress properties. Finally, in Section 7, we discuss possible extensions of the results presented in this paper.

Related work. It is worth noting that there has been an attempt to describe the overall semantics of TMs (Scott 2006) (including lock-based ones). However, the approach taken there is very low-level—the properties are defined with respect to specific TM protocols and strategies. Our approach is more general: we define semantics that is implementation-agnostic and that is visible through the public interface of a TM to a user. We also show how this semantics can be verified.

There have also been other attempts to describe the semantics of a TM, e.g., in (Vitek et al. 2004; Jagannathan et al. 2005; Abadi et al. 2008; Moore and Grossman 2008; Menon et al. 2008). Those papers, however, focus on safety, i.e., serializability. In (Moore and Grossman 2008) there is a notion of progress, but it refers to deadlock-freedom of the whole system (i.e., making sure that at least one thread can execute a step at any given time) rather than progress of individual transactions.

2. Preliminaries

2.1 Shared Objects and their Implementations

We consider an asynchronous shared memory system of n processes (threads) p_1, \dots, p_n that communicate by executing operations on (shared) objects. (At the hardware level, a shared object is simply a word in shared memory with the instructions supported by a given processor, e.g., read, write, or compare-and-swap.) An example of a very simple shared object is a *register*,⁵ which exports only *read* and *write* operations. Operation *read* returns the current state (value) of the register, and *write*(v) sets the state of the register to value v . Hence, a register provides the basic read-write memory semantics.

shared state. Invisible reads are used by most TM implementations and considered crucial for good performance in read-dominated workloads.

⁵ Note that we use here the term “register” in its distributed computing sense: a read-write abstraction.

Consider a single run of any algorithm. A *history* is a sequence of invocations and responses of operations that were executed by processes on (shared) objects in this run. A history of an object x is a history that contains only operations executed on x . (Note here that we assume that events executed in a given run can be totally ordered by their execution time; events that are issued at the same time, e.g., on multi-processor systems, can be ordered arbitrarily.)

An object x may be implemented either directly in hardware, or from other, possibly more primitive, objects, which we call *base objects*. If I_x is an implementation of an object x , then an *implementation history* of I_x is a sequence of (1) invocations and responses of operations on x , and (2) corresponding operations on base objects (called *steps*) that were executed by I_x (i.e., by processes executing I_x) in some run. Hence, intuitively, a history of an object x represents what happened in some run at the (public) interface of x . An implementation history, in addition, shows what steps the implementation of x executed in response to the operations invoked on x .

In algorithms, for simplicity, we assume that base objects such as registers and test-and-set objects are atomic, i.e., linearizable (Herlihy and Wing 1990). That is, operations on those objects appear (to the application) as if they happened instantaneously at some unique point in time between their invocation and response events. (For example, in Java, a “volatile” variable is an atomic register, while an object of class `AtomicInteger` is an atomic object that supports operations such as *get*, *set*, *incrementAndGet*, etc.)

However, assuming a weaker memory model does not impact our results: the progress properties we define do not rely on atomicity, strong try-lock objects are not linearizable, and atomic registers of any size can be implemented out of 1-bit safe (the most primitive) registers (Lamport 1986).

If E is an (implementation) history, then $E|p_i$ denotes the restriction of E to events (including steps) executed by process p_i , and $E|x$ denotes the restriction of E to events on object x and steps of the implementation of x . We assume that processes execute operations on objects sequentially. That is, in every restriction $E|p_i$ of an (implementation) history E , no two operations and no two steps overlap.

We focus on object implementations that are *wait-free* (Herlihy 1991). Intuitively, an implementation I_x of an object x is wait-free if a process that invokes an operation on x is never blocked indefinitely long inside the operation, e.g., waiting for other processes. Hence, processes can make progress independently of each other. More precisely:

DEFINITION 1. *An implementation I_x of an object x is wait-free, if whenever any process p_i invokes an operation on x , p_i returns from the operation within a finite number of its own steps.*

2.2 Transactional Memory (TM)

A TM enables processes to communicate by executing transactions. For simplicity, we will say that a transaction T performs some action, meaning that the process executing T performs this action within the transactional context of T . A transaction T may perform operations on *transactional variables*, which we call *t-variables* for short. For simplicity, we assume that every t-variable x supports only two operations: *read* that returns the current state (value) of x , and *write(v)* that sets the state of x to value v . We discuss in Section 7 what changes when t-variables are arbitrary objects, i.e., objects that have operations beyond *read* and *write* (e.g., *incrementAndGet*). Note, however, that most existing TMs either provide only read-write t-variables (e.g., word-based TMs), or effectively treat all operations on t-variables as reads and writes (e.g., without exploiting the commutativity relations between non-read-only operations).

Each transaction has its own unique identifier, e.g., T_1, T_2 , etc. A transaction T_k may access (read or write) any number of t-variables. Then, T_k may either *commit* or *abort*. We assume that once T_k commits or aborts T_k does not perform any further actions. In this sense, restarting a transaction T_k (i.e., the computation T_k was supposed to perform) is considered in our model as a different transaction (with a different identifier).

We can treat a TM as an object with the following operations:

- $tread_k(x)$ and $twrite_k(x, v)$ that perform, respectively, a *read* or a *write(v)* operation on a t-variable x within a transaction T_k ,
- $tryC_k$ that is a request to commit transaction T_k ,
- $tryA_k$ that is a request to abort transaction T_k .

Each of the above operations can return a special value A_k that indicates that the operation has failed and the respective transaction T_k has been aborted. Operation $tryC_k$ returns value C_k if committing T_k has been successful. Operation $tryA_k$ always returns A_k (i.e., it always succeeds in aborting transaction T_k).

The above operations of a TM, in some form, are either explicitly used by a programmer (e.g., in TL2, TinySTM, RSTM), or inserted by a TM-aware compiler (e.g., in McRT-TM, Bartok-STM). Even if the compiler is responsible for inserting those operations, the programmer must specify which blocks of code are parts of transactions, and retains full control of what operations on which t-variables those transactions perform. Hence, in either case, this TM interface is visible to a programmer, and so are properties defined with respect to this interface.

If H is an (implementation) history of a TM object, then $H|T_k$ denotes the restriction of H to only events of transaction T_k . We say that a transaction T_k is in a history H , and write $T_k \in H$, if $H|T_k$ is a non-empty sequence.

Let H be any history and T_k be any transaction in H . We say that T_k is *committed* in H , if H contains response C_k of operation $tryC_k$. We say that T_k is *aborted* in H , if H contains response A_k of any TM operation.

We say that a transaction T_k *follows* a transaction T_i in a history H , if T_i is committed or aborted in H and the first event of T_k in H follows the last event of T_i in H . If neither T_k follows T_i in H , nor T_i follows T_k in H , then we say that T_i and T_k are *concurrent* in H .

We assume that every transaction itself is sequential. That is, for every history H of a TM and every transaction $T_k \in H$, $H|T_k$ is a sequence of non-overlapping TM operations. Clearly, operations of different transactions can overlap. We also assume that each transaction is executed by a single process, and that each process executes only one transaction at a time (i.e., transactions at the same process are never concurrent).

2.3 Try-Locks

All lock-based TMs we know of use (often implicitly) a special kind of locks, usually called *try-locks* (Scott and Scherer III 2001). Intuitively, a try-lock is an object that provides mutual exclusion (like a lock), but does not block processes indefinitely. That is, if a process p_i requests a try-lock L but L is already acquired by a different process, p_i is returned the information that its request failed instead of being blocked waiting until L is released.

Try-locks keep the TM implementation simple and avoid deadlocks. Moreover, if any form of fairness is needed, it is provided at a higher level than at the level of individual locks—then more information about a transaction can be used to resolve conflicts and provide progress. Ensuring safety and progress can be effectively separate tasks.

More precisely, a try-lock is an object with the following operations:

1. *trylock*, that returns *true* or *false*; and
2. *unlock*, that always returns *ok*.

Let L be any try-lock. If a process p_i invokes *trylock* on L and is returned *true*, then we say that p_i has *acquired* L . Once p_i acquires L , we say that (1) p_i holds L until p_i invokes operation *unlock* on L , and (2) L is *locked* until p_i returns from operation *unlock* on L . (Hence, L might be locked even if no process holds L —when some process that was holding L is still executing operation *unlock* on L .)

Every try-lock L guarantees the following property, called *mutual exclusion*: no two processes hold L at the same time.

For simplicity, we assume that try-locks are not reentrant. That is, a process p_i may invoke *trylock* on a try-lock L only when p_i does not hold L . Conversely, p_i may invoke *unlock* on L only when p_i holds L .

Intuitively, we say that a try-lock L is *strong* if whenever several processes compete for L , then one should be able to acquire L . This property corresponds to deadlock-freedom, livelock-freedom, or progress (Raynal 1986) properties of (blocking) locks.

DEFINITION 2. We say that a try-lock L is strong, if L ensures the following property, in every run: if L is not locked at some time t and some process invokes operation *trylock* on L at t , then some process acquires L after t .

While there exists a large number of lock implementations, only a few are try-locks or can be converted to try-locks in a straightforward way. The technical problems of transforming a queue (blocking) lock into a try-lock are highlighted in (Scott and Scherer III 2001). It is trivial to transform a typical TAS or TATAS lock (Raynal 1986) into a strong try-lock (e.g., Algorithm 4 in Section 5.2).

3. Progress of a Lock-Based TM

Lock-based TMs are TM implementations that use (internally) mutual exclusion to handle some phases of a transaction. Most of them use some variant of the two-phase locking protocol, well-known in the database world (Eswaran et al. 1976).

From the user’s perspective, however, the choice of the mechanism used internally by a TM implementation is not very important. What is important is the semantics the TM manifests on its public interface, and the time/space complexities of the implementation. If those properties are known, then the designer of a lock-based TM is free to choose the techniques that are best for a given hardware platform, without the fear of breaking existing applications that use a TM.

As we already mentioned, the correctness criterion for TMs, including lock-based ones, is usually *opacity* (Guerraoui and Kapalka 2008a). This property says, intuitively, that (1) committed transactions should appear as if they were executed sequentially, in an order that agrees with their real-time ordering, (2) no transaction should ever observe the modifications to shared state done by aborted or live transactions, and (3) all transactions, including aborted and live ones, should always observe a consistent state of the system. The first two properties correspond, roughly, to the classical database properties: strict serializability (Papadimitriou 1979) and the strongest variant of recoverability (Hadzilacos 1988), respectively. The last property is unique to TMs, and needs to be ensured to prevent unexpected crashes or incorrect behavior of applications that use a TM.

However, opacity is not enough. A TM that always aborts every transaction, or that blocks transactions infinitely long, could ensure opacity and still be useless from the user’s perspective. In this section, we define the *progress* properties of a lock-based TM. These involve individual operations of transactions, where it is

typical to require *wait-freedom*, and entire transactions, for which we will require our notion of *strong progressiveness*.

3.1 Liveness of TM Operations

If a process p_i invokes an operation (*tread*, *twrite*, *tryC*, or *tryA*) on a TM, we expect that p_i eventually gets a response from the operation. The response might be the special value A_k that informs p_i that its current transaction T_k has been aborted.

We assume that each implementation of a TM is a wait-free object. That is, a TM ensures wait-freedom on the level of its operations. This property is indeed ensured by many current lock-based TMs: if a transaction T_k encounters a conflict, T_k is immediately aborted and the control is returned to the process executing T_k .

Note that a TM may use a *contention manager* to decide what to do in case of a conflict. A contention manager is a logically external module that can reduce contention by delaying or aborting some of the transactions that conflict. In principle, a contention manager could make transactions wait for each other, in which case wait-freedom would be violated. However, such contention managers change the progress properties of a TM significantly and as such should be considered separately.

Operation wait-freedom may also be violated periodically by some TM mechanisms that handle overflows. While those can be unavoidable, as we discuss in Section 6.3, they are executed very rarely. Moreover, one can easily predict when they could start. In this sense, wait-freedom can be guaranteed except for some short periods that can be signalled in advance to processes by, e.g., setting a global flag.

3.2 Progress of Transactions

Intuitively, a transaction makes progress when it commits. One would like most transactions to commit, except those that were explicitly requested by the application to be aborted (using a *tryA* operation of a TM). However, a TM may be often forced to abort some transactions when the conflicts between them cannot be easily resolved. We will call such transactions *forcefully aborted*. The *strong progressiveness* property we introduce here defines when precisely a transaction can be forcefully aborted.

Intuitively, strong progressiveness says that (1) if a transaction has no *conflict* then it cannot be forcefully aborted, and (2) if a group of transactions conflict on a single t-variable, then not all of those transactions can be forcefully aborted. Roughly speaking, two or more transactions conflict if they access the same t-variable in a conflicting way, i.e., if at least one of those accesses is a write operation. (It is worth noting that the notion of a conflict can be easily generalized to t-variables with arbitrary operations, and to arbitrary mappings between t-variables and locks that may allow *false* conflicts. We discuss this in Section 7.)

Strong progressiveness is not the strongest possible progress property. The strongest one, which requires that no transaction is ever forcefully aborted, cannot be implemented without throttling significantly the parallelism between transactions, and is thus impractical in multi-processor systems.

Strong progressiveness, however, still gives a programmer the following important advantages. First, it guarantees that if two independent subsystems of an application do not share any memory locations (or t-variables), then their transactions are completely isolated from each other (i.e., a transaction executed by a subsystem A does not cause a transaction in a subsystem B to be forcefully aborted). Second, it avoids “spurious” aborts: the cases when a transaction can abort are strictly defined. Third, it ensures global progress for single-operation transactions, which is important when non-transactional accesses to t-variables are encapsulated into transactions in order to ensure strong atomicity (Blundell et al. 2006). Finally, it ensures that processes are able to eventu-

ally communicate via transactions (albeit in a simplified manner—through a single t-variable at a time). Nevertheless, one can imagine many other reasonable progress properties, for which strong progressiveness can be a good reference point.

More precisely, let H be any history of a TM and T_k be any transaction in H . We say that T_k is *forcefully aborted* in H , if T_k is aborted in H and there is no invocation of operation $tryA_k$ in H . We denote by $WSet_H(T_k)$ and $RSet_H(T_k)$ the sets of t-variables on which T_k executed, respectively, a *write* or a *read* operation in H . We denote by $RWSet_H(T_k)$ the union of sets $RSet_H(T_k)$ and $WSet_H(T_k)$, i.e., the set of t-variables accessed (read or written) by T_k in history H . We say that two transactions T_i and T_k in H *conflict on a t-variable x* , if (1) T_i and T_k are concurrent in H , and (2) either x is in $WSet_H(T_k)$ and in $RWSet_H(T_i)$, or x is in $WSet_H(T_i)$ and in $RWSet_H(T_k)$. We say that T_k *conflicts with a transaction T_i in H* if T_i and T_k conflict in H on some t-variable.

Let H be any history, and T_i be any transaction in H . We denote by $CVar_H(T_i)$ the set of t-variables on which T_i conflicts with any other transaction in history H . That is, a t-variable x is in $CVar_H(T_i)$ if there exists a transaction $T_k \in H$, $k \neq i$, such that T_i conflicts with T_k on t-variable x .

Let Q be any subset of the set of transactions in a history H . We denote by $CVar_H(Q)$ the union of sets $CVar_H(T_i)$ for all $T_i \in Q$.

Let $CTrans(H)$ be the set of subsets of transactions in a history H , such that a set Q is in $CTrans(H)$ if no transaction in Q conflicts with a transaction *not* in Q . In particular, if T_i is a transaction in a history H and T_i does not conflict with any other transaction in H , then $\{T_i\} \in CTrans(H)$.

DEFINITION 3. A TM implementation M is *strongly progressive*, if in every history H of M the following property is satisfied: for every set $Q \in CTrans(H)$, if $|CVar_H(Q)| \leq 1$, then some transaction in Q is not forcefully aborted in H .

4. Verifying Strong Progressiveness

Verifying that a given TM implementation M ensures a given property P might often be difficult as one has to reason about a large number of histories involving an arbitrary number of transactions accessing an arbitrary number of t-variables. This complexity is greatly reduced if one can reduce the verification task to some small subset of histories of M , e.g., involving a limited number of t-variables or transactions. This approach has been used, e.g., in (Guerraoui et al. 2008) to automatically check opacity, obstruction-freedom, and lock-freedom of TMs that feature certain symmetry properties.

In this section, we show how to reduce the problem of proving strong progressiveness of histories with arbitrary numbers of transactions and t-variables to proving a simple property of each individual (logical) try-lock used in those histories. Basically, we show that if a TM implementation M uses try-locks, or if one can assign “logical” try-locks to some parts of the algorithm of M , and if each of those try-locks is strong, then M ensures strong progressiveness. Unlike in (Guerraoui et al. 2008), we do not assume any symmetry properties of a TM. Our result is thus complementary to that of (Guerraoui et al. 2008), not only because it concerns a different property, but also because it uses a different approach.

Our reduction theorem is general as it encompasses lock-based TMs that use invisible reads, i.e., in which readers of a t-variable are not visible to other transactions, as well as those that use visible ones. We show also how the theory presented here can be used to prove strong progressiveness of TL2, TinySTM, RSTM, and McRT-STM. Finally, we point out one of the ambiguities of ensuring strong progressiveness with visible reads.

4.1 Reduction Theorem

Let M be any TM implementation, and E be any implementation history of M . Let E' be any implementation history that is obtained from E by inserting into E any number of invocations and responses of operations of a try-lock L_x for every t-variable x . We say that E' is a *strong try-lock extension* of E , if the following conditions are satisfied in E' :

STLE1. For every t-variable x , $E'|L_x$ is a valid history of a strong try-lock object;

STLE2. For every process p_i and every t-variable x , if, at some time t , p_i invokes *trylock* on L_x or p_i holds L_x , then p_i executes at t in E' a transaction T_k such that $x \in WSet_{E'}(T_k)$;

STLE3. For every process p_i and every transaction $T_k \in E'|p_i$, if T_k is forcefully aborted in E' , then either (1) p_i while executing T_k is returned *false* from every operation *trylock* on some try-lock L_x , or (2) there is a t-variable $x \in RSet_{E'}(T_k)$, such that some process other than p_i holds L_x at some point while p_i executes T_k but before T_k acquires L_x (if at all).

THEOREM 4. For any TM implementation M , if there exists a strong try-lock extension of every implementation history of M , then M is strongly progressive.

Proof. Assume, by contradiction, that there exists a TM implementation M , such that some implementation history E of M has a strong try-lock extension E' , but E violates strong progressiveness. This means that there is a set Q in $CTrans(E)$, such that $|CVar_E(Q)| \leq 1$ and every transaction in Q is forcefully aborted in E . Recall that Q is a subset of transactions, such that no transaction in Q has a conflict with a transaction outside of Q .

Assume first that $CVar_E(Q) = \emptyset$. But then no transaction in set Q has a conflict, and so, by STLE1–2, no transaction in Q can fail to acquire a try-lock, or read a t-variable x such that try-lock L_x is held by a concurrent transaction. Hence, by STLE3, no transaction in Q can be forcefully aborted—a contradiction.

Let x be the t-variable that is the only element of set $CVar_E(Q)$. Note first that if a transaction T_k in Q invokes operation *trylock* on some try-lock L_y (where y is a different t-variable than x) then, by STLE2, no other transaction concurrent to T_k invokes *trylock* on L_y or reads t-variable y . This is because no transaction in Q conflicts on a t-variable different than x .

Assume first that no transaction in set Q acquires try-lock L_x . But then, by STLE1–3, no transaction in Q can be forcefully aborted—a contradiction.

Let T_k be the first transaction from set Q to acquire try-lock L_x . By STLE3, and because T_k is forcefully aborted, there is a transaction T_i that holds L_x after T_k starts and before T_k acquires L_x . Clearly, by STLE2, x must be in $WSet_E(T_i)$, and so T_i must be in set Q . But then T_i acquires L_x before T_k —a contradiction with the assumption that T_k is the first transaction from set Q to acquire L_x . \square

4.2 Examples

We show here how our reduction theorems can be used to prove the strong progressiveness of TL2, TinySTM, RSTM (one of its variants), and McRT-STM. None of those TM implementations explicitly use try-locks, and so we need to show which parts of their algorithms correspond to operations on “logical” try-locks for respective t-variables. We assume the use of a simple contention manager that makes each transaction that encounters a conflict abort itself. Such a contention manager (possibly with a back-off protocol) is usually the default one in word-based TMs. We also assume that the mapping between t-variables and locks is a one-to-

one function (which is the default in RSTM). This assumption is revisited in Section 7.

TL2. This TM uses commit-time locking and deferred updates. That is, locking and updating t-variables is delayed until the commit time of transactions. The TL2 algorithm is roughly the following (for a process p_i executing a transaction T_k):

1. When T_k starts, p_i reads the *read timestamp* of T_k from a global counter C .
2. If T_k reads a t-variable x , p_i checks whether x is not locked and whether the version number of x is lower or equal to the read timestamp of T_k . If any of those conditions is violated then T_k is aborted.
3. Once T_k invokes $tryC_k$, p_i first tries to lock all t-variables that were written to by T_k . Locking of a t-variable x is done by executing a compare-and-swap (CAS) operation on a memory word $w(x)$ that contains, among other information, a *locked* flag. If p_i successfully changes the *locked* flag from *false* to *true*, then p_i becomes the exclusive owner of x and can update x . If CAS fails, however, T_k is aborted.
4. Once all t-variables written to by T_k are locked, p_i atomically increments and reads the value of the global counter C . The read value is the *write timestamp* of T_k .
5. Next, p_i validates transaction T_k by checking, for every t-variable x read by T_k , whether x is not locked by a transaction other than T_k and whether the version number of x is lower or equal to the read timestamp of T_k . Again, if any of those conditions is violated then transaction T_k is aborted (and its locks released).
6. Then, p_i updates all the states of the locked t-variables with the values written by T_k and the write timestamp of T_k .
7. Finally, T_k releases all the locked t-variables.

It is easy to assign logical try-locks to the above algorithm of TL2, i.e., to build a try-lock extension of every implementation history E of TL2. Basically, we put an invocation and response of operation $trylock$ on a try-lock L_x around any CAS operation that operates on the *locked* flag of any t-variable x . The response is *true* if CAS succeeds, and *false* otherwise. We also put an invocation and response of operation $unlock$ on L_x around the write operation that sets the *locked* flag of x to *false*. It is straightforward to see that this way we indeed obtain a valid try-lock extension of any implementation history E of TL2:

1. Property STLE1 is ensured because a CAS on a word $w(x)$ can fail only when some other CAS on $w(x)$ already succeeded, and once a CAS on $w(x)$ succeeds, no other CAS on $w(x)$ can succeed until the *locked* flag is reset. Hence, the single CAS operation indeed implements a strong try-lock.
2. Property STLE2 is ensured because a transaction T_i invokes CAS on a word $w(x)$ only when (1) T_i wrote to t-variable x , and (2) T_i is in its commit phase.
3. To prove that TL2 ensures property STLE3, consider any forcefully aborted transaction T_k executed by some process p_i (in some implementation history E of TL2). Assume first that a CAS operation executed by T_k (i.e., by p_i while executing T_k) on some word $w(x)$ fails. But then (1) T_k could not have locked try-lock L_x before, and (2) T_k is immediately aborted afterwards. Hence, property STLE3 is trivially ensured. This means that T_k reads some t-variable x and either (1) $w(x)$ has the *locked* flag set to *true* when T_k reads x (and $w(x)$ is not locked by T_k), or (2) the version number of x is larger than the read timestamp of T_k . In case (1) property STLE3 is trivially en-

sured. Assume then case (2). This means that some transaction T_m that has a write timestamp greater than the read timestamp of T_k wrote to x either (a) before T_k read x , or (b) after T_k read x and before T_k locked $w(x)$. But then T_m must have acquired its write timestamp, while holding try-lock L_x , after T_k acquired its read timestamp and before T_k locked L_x (if at all). Hence, STLE3 is ensured.

We thus obtain the following theorem:

THEOREM 5. *TL2 (with a one-to-one t-variable to try-lock mapping) is strongly progressive.*

TinySTM. There are two major differences with TL2. First, TinySTM locks a t-variable x already inside any *write* operation on x , i.e., locking is not delayed until the commit time of transactions. Second, if a transaction T_k reads a t-variable x that has a version number higher than the read timestamp of T_k , then T_k tries to validate itself to avoid being aborted, instead of aborting itself immediately. TinySTM uses CAS for locking, in the same way as TL2. Hence, we can insert the invocations and responses of operations on logical try-locks into any implementation history of TinySTM in the same way as for TL2.

It is worth noting, however, that the overflow handling mechanism, which can be turned on at compile time, breaks strong progressiveness in very long histories. As we discuss in Section 6.3, this mechanism is necessary to overcome the complexity lower bound and still guarantee correctness. However, strong progressiveness is still ensured in histories with the number of transactions lower than the maximum value of the t-variable version number, or between version number overflows.

THEOREM 6. *TinySTM (with the overflow handling mechanism turned off, and with a one-to-one t-variable to try-lock mapping) is strongly progressive.*

RSTM. This TM is highly configurable: currently there are four TM backends to choose from, and each has a number of configuration options. The two backends that are of interest here are *LLT* and *RedoLock*. LLT is virtually identical to TL2. RedoLock has object-level lock granularity. That is, transactions conflict if they access (in a conflicting way) the same object, not necessarily the same memory location (i.e., t-variables in RSTM are objects, not single memory words as in TL2 and TinySTM). However, the algorithm of RedoLock is, depending on the configuration option, similar to either TL2 or TinySTM. The main difference is in the validation heuristic that decides when a transaction needs to validate its read set, but this does not impact strong progressiveness (the heuristic does not by itself abort any transaction—it just determines when to validate the read set of a transaction). Like in TL2 and TinySTM, RedoLock uses CAS for locking, and so the same technique as for TL2 and TinySTM can be used to prove that RSTM with RedoLock backend is strongly progressive.

THEOREM 7. *RSTM with the RedoLock backend is strongly progressive.*

McRT-STM. The algorithm of McRT-STM (as described in (Adl-Tabatabai et al. 2006)) is essentially the same as the one of TinySTM, except that McRT-STM does not validate reads until the commit time of a transaction (and so the timestamp-based read validation technique is not necessary). McRT-STM also does not handle timestamp overflows. Hence, as McRT-STM uses CAS for locking, it is immediate that McRT-STM is strongly progressive.

THEOREM 8. *McRT-STM is strongly progressive.*

Visible reads. It may seem that the simplest way of implementing a strongly progressive TM that uses visible reads is to use read-write try-locks. Then, if a transaction T_i wants to read a t-variable x , T_i must first acquire a shared (read) try-lock on x , and if T_i wants to write to x , T_i must acquire an exclusive (write) try-lock on x . However, this simple algorithm does not ensure strong progressiveness, even if the read-write try-locks are (in some sense) strong. Consider transactions T_i and T_k that read a t-variable x . Clearly, both transactions acquire a shared lock on x . But then, if both T_i and T_k want to write to x , it may happen that both get aborted. This is because a transaction T_k cannot acquire an exclusive try-lock on x if any other transaction holds a shared try-lock on x .

A simple way to implement a strongly progressive TM with invisible reads is to use (standard) try-locks. Then, only the writing to a t-variable x requires acquiring a try-lock on x . A transaction that wants to read x simply adds itself to the list of readers of x (if the try-lock of x is unlocked). This list, however, is not used to implement a read-write try-lock semantics, but to allow a transaction that writes to x to invalidate and abort all the current readers of x . Such a TM can be verified to be strongly progressive using our reduction theorem. A separate reduction theorem, based on read-write try-locks, is thus not necessary, and would probably be incorrect (trying to provide such a theorem allowed us to discover this ambiguity).

5. The Power of a Lock-Based TM

In this section, we use our semantics to determine the computational power of lock-based TMs. We use the notion of *consensus number* (Herlihy 1991) as the metric of power of an object. The consensus number of an object x is defined as the maximum number of processes for which one can implement a wait-free *consensus* object using any number of instances of x (i.e., objects of the same type as x) and registers. A consensus object, intuitively, allows processes to agree on a single value chosen from the values those processes have proposed. More formally, a consensus object implements a single operation: *propose*(v). When a process p_i invokes *propose*(v), we say that p_i *proposes* value v . When p_i is returned value v' from *propose*(v), we say that p_i *decides* value v' . Every consensus object ensures the following properties in every execution: (1) no two processes decide different values (agreement), and (2) every value decided is a value proposed by some process (validity).

According to (Herlihy 1991), if an object x has consensus number k , then one cannot implement x using objects with consensus number lower than k . For example, a queue and test-and-set have consensus number 2, and so they cannot be implemented from only registers (which have consensus number 1).

We prove here that the consensus number of a strongly progressive TM is 2. We do so in the following way. First, we prove that a strongly progressive TM is computationally equivalent to a strong try-lock. That is, one can implement a strongly progressive TM from (a number of) strong try-locks and registers, and vice versa. Second, we determine that the consensus number of a strong try-lock is 2.

The equivalence to a strong try-lock is interesting in its own right. It might also help proving further impossibility results as a strong try-lock is a much simpler object to reason about than a lock-based TM.

5.1 Equivalence between Lock-Based TMs and Try-Locks

To prove that a strongly progressive TM is (computationally) equivalent to a strong try-lock, we exhibit two algorithms: Algorithm 1 that implements a strong try-lock from a strongly progressive TM object and a shared memory register, and Algorithm 2 that imple-

Algorithm 1: An implementation of a strong try-lock from a strongly progressive TM object (k is a unique transaction identifier generated for every operation call)

uses: M —TM object, x_1, x_2, \dots —binary t-variables,
 V —register

initially: $x_1, x_2, \dots = false, V = 1$

```

1 operation trylock
2    $v \leftarrow V.read;$ 
3    $locked \leftarrow M.tread_k(x_v);$ 
4   if  $locked = A_k$  or  $locked = true$  then return false;
5    $s \leftarrow M.twrite_k(x_v, true);$ 
6   if  $s = A_k$  then return false;
7    $s \leftarrow M.tryC_k;$ 
8   if  $s = A_k$  then return false;
9   else return true;

10 operation unlock
11    $v \leftarrow V.read;$ 
12    $V.write(v + 1);$ 
13   return ok;

```

ments a strongly progressive TM from a number of strong try-locks and registers. Both algorithms are not meant to be efficient or practical: their sole purpose is to prove the equivalence result.

The intuition behind Algorithm 1 is the following. We use an unbounded number of binary t-variables x_1, x_2, \dots (each initialized to *false*) and a single register V holding an integer (initialized to 1). If the value of V is v , then the next operation (*trylock* or *unlock*) will use t-variable x_v . If x_v equals *true*, then the lock is locked. A process p_i acquires the lock when p_i manages to execute a transaction T_k that changes the value of x_v from *false* to *true*. Then, p_i releases the lock by incrementing the value of register V , so that $x_{v'} = false$ where v' is the new value of V . (Note that incrementing V in two steps is safe here, as only one process—the one that holds the lock—may execute lines 11–12 at a time.) The implemented try-lock is strong because whenever several processes invoke *trylock*, at least one of those processes will commit its transaction (as the TM is strongly progressive) and acquire the try-lock.

Due to space constraints, the following lemma is proved in the extended version of this paper (Guerraoui and Kapalka 2008b).

LEMMA 9. *Algorithm 1 implements a strong try-lock.*

The intuition behind Algorithm 2 is the following. We use a typical two-phase locking scheme with eager updates, optimistic (invisible) reads, and incremental validation (this can be viewed as a simplified version of TinySTM that explicitly uses strong try-locks). Basically, whenever a transaction T_i invokes operation *write* on a t-variable x for the first time, T_i acquires the corresponding try-lock L_x (line 13) and marks x as locked (line 21). Then, T_i may update the state of x in $TVar[x]$ any number of times. The original state of x is saved by T_i in $oldval[x]$, so that if T_i aborts then all the updates of t-variables done by T_i can be rolled back (line 39). If, at any time, T_i fails to acquire a try-lock, T_i aborts. This ensures freedom from deadlocks.

If T_i invokes operation *read* on a t-variable y that T_i has not written to before, T_i reads the current value of y (line 2) and *validates* itself (function *validate*). Validation ensures that none of the t-variables that T_i read so far has changed or has been locked, thus preventing T_i from having an inconsistent view of the system. If validation fails, T_i is aborted. Because values written to any t-variable are not guaranteed to be unique, and because, in our

Algorithm 2: An implementation of a strongly progressive TM from strong try-locks and registers

uses: L_x —strong try-lock (for each t-variable x),
 $TVar$ —array of registers (other variables are local)
initially: $TVar[x] = (0, 0, false)$ for each t-variable x ,
 $rset = wset = \emptyset$

```

1 operation  $tread_k(x)$ 
2    $(v, ts, locked) \leftarrow TVar[x].read;$ 
3   if  $x \in wset$  then return  $v$ ;
4   if  $x \notin rset$  then
5      $readts[x] \leftarrow ts;$ 
6      $rset \leftarrow rset \cup \{x\};$ 
7   if (not  $validate$ ) or  $locked$  then
8      $abort;$ 
9     return  $A_k$ ;
10  return  $v$ ;

11 operation  $twrite_k(x, v)$ 
12  if  $x \notin wset$  then
13     $locked \leftarrow L_x.trylock;$ 
14    if not  $locked$  then
15       $abort;$ 
16      return  $A_k$ ;
17   $(v', ts, locked) \leftarrow TVar[x].read;$ 
18  if  $x \notin wset$  then
19     $oldval[x] \leftarrow v';$ 
20     $wset \leftarrow wset \cup \{x\};$ 
21   $TVar[x].write(v, ts, true);$ 
22  return  $ok$ ;

23 operation  $tryC_k$ 
24  if not  $validate$  then
25     $abort;$ 
26    return  $A_k$ ;
27  for  $x \in wset$  do
28     $(v, ts, locked) \leftarrow TVar[x].read;$ 
29     $TVar[x].write(v, ts + 1, false);$ 
30     $L_x.unlock;$ 
31   $wset \leftarrow rset \leftarrow \emptyset;$ 
32  return  $C_k$ ;

33 operation  $tryA_k$ 
34   $abort;$ 
35  return  $A_k$ ;

36 function  $abort$ 
37  for  $x \in wset$  do
38     $(v, ts, locked) \leftarrow TVar[x].read;$ 
39     $TVar[x].write(oldval[x], ts, false);$ 
40     $L_x.unlock;$ 
41   $wset \leftarrow rset \leftarrow \emptyset;$ 

42 function  $validate$ 
43  for  $x \in rset$  do
44     $(v, ts, locked) \leftarrow TVar[x];$ 
45    if ( $locked$  and  $x \notin wset$ ) or  $ts \neq readts[x]$  then
46      return  $false$ ;
47  return  $true$ ;

```

Algorithm 3: An implementation of wait-free consensus from a strong try-lock in a system of 2 processes (code for process $p_i, i = 1, 2$)

uses: L —strong try-lock, V_1, V_2 —registers

```

1 operation  $propose(v)$ 
2    $V_i.write(v);$ 
3    $locked \leftarrow L.trylock;$ 
4   if  $locked$  then return  $v$ ;
5   else return  $V_{(3-i)}.read;$ 

```

simplified model, a try-lock does not have an operation that would read its state, we store with the state of each t-variable x a (unique) timestamp (version number) of x and a $locked$ flag that is set to $true$ if x is being written to by some transaction. The timestamps and $locked$ flags are used for validation.

To commit a transaction T_i , the algorithm first validates T_i (line 24). Then, for each t-variable x written to by T_i , the timestamp of x is incremented by 1, the $locked$ flag of x is set to $false$ (line 29), and finally the try-lock L_x of x is unlocked (line 30). Aborting T_i requires rolling back all the updates done by T_i (line 39) and unlocking all the try-locks acquired by T_i (line 40).

LEMMA 10. *Algorithm 2 implements a strongly progressive TM.*

Due to space constraints, the proof of Lemma 10 is omitted. It can be found in the extended version of this paper (Guerraoui and Kapalka 2008b). Note that strong progressiveness of Algorithm 2 is trivial to verify using our reduction theorem, because every implementation history of this algorithm is its own try-lock extension (i.e., it ensures properties STLE1–3).

From Lemma 9 and Lemma 10, we immediately obtain the following result (recall that an object x is (computationally) equivalent to an object y , if y can be implemented from any number of instances of x and registers, and x can be implemented from any number of instances of y and registers):

THEOREM 11. *Every strongly progressive TM is equivalent to a strong try-lock.*

5.2 Consensus Number of Strong Try-Locks

To prove that the consensus number of a strong try-lock is 2, we show that (1) a strong try-lock can implement consensus in a system of 2 processes, and (2) there is no algorithm that implements consensus using (any number of) strong try-locks and registers in a system of 3 (or more) processes.

Algorithm 3 shows an implementation of consensus for two processes (p_1 and p_2) using a single strong try-lock (L) and two registers (V_1 and V_2). The process p_i that acquires L is the winner: the value proposed by p_i , and written by p_i to register V_i , is decided by both p_1 and p_2 . Because L is a strong try-lock, if both processes concurrently execute operation $propose$, at least one of them acquires L . Because no process ever unlocks L , at most one process acquires L . Hence, exactly one process is the winner.

Due to space constraints, the following lemma is proved in the extended version of this paper (Guerraoui and Kapalka 2008b).

LEMMA 12. *Algorithm 3 implements wait-free consensus in a system of 2 processes.*

To prove that there is no algorithm that implements consensus using strong try-locks and registers in a system of 3 (or more) processes, we show in Algorithm 4 that a strong try-lock can be

Algorithm 4: An implementation of a strong try-lock from a test-and-set object

uses: S —test-and-set object
initially: $S = false$

```

1 operation trylock
2    $locked \leftarrow S.test\text{-}and\text{-}set;$ 
3   return  $\neg locked;$ 

4 operation unlock
5    $S.reset;$ 

```

implemented from a single test-and-set object.⁶ Because a test-and-set object has consensus number 2, the algorithm proves that a strong try-lock cannot have consensus number higher than 2. Note that the presented algorithm is a non-blocking version of a simple and well-known TAS lock (Raynal 1986). The following lemma is thus trivial to verify:

LEMMA 13. *Algorithm 4 implements a strong try-lock.*

From Lemma 12 and Lemma 13, we immediately obtain the following result:

THEOREM 14. *A strong try-lock has consensus number 2.*

Hence, by Theorem 11 and Theorem 14, the following theorem holds:

THEOREM 15. *Every strongly progressive TM has consensus number 2.*

COROLLARY 16. *There is no algorithm that implements a strongly progressive TM using only registers.*

5.3 Weakening Strong Progressiveness

Interestingly, nailing down precisely the progress property of a lock-based TM also helps consider alternative semantics and their impacts. We discuss here how one has to weaken the progress semantics of a lock-based TM so that it could be implemented with registers only. We define a property called *weak progressiveness* that enables (lightweight) TM implementations with consensus number 1.

Intuitively, a TM is weakly progressive if it can forcefully abort a transaction T_i only if T_i has a conflict with another transaction. More precisely:

DEFINITION 17. *A TM implementation M is weakly progressive, if in every history H of M the following property is satisfied: if a transaction $T_i \in H$ is forcefully aborted, then T_i conflicts with some transaction in H .*

We correlate this notion with the concept of a *weak try-lock*: a try-lock which operation *trylock* executed by a process p_i may always return *false* if another process is concurrently executing *trylock* on the same try-lock object. That is, p_i is guaranteed to acquire a weak try-lock L only if L is not locked and no other process tries to acquire L at the same time. More precisely:

DEFINITION 18. *We say that a try-lock L is weak if L has the following property: if a process p_i invokes *trylock* on L at some time t , L is not locked at t , and no process other than p_i executes operation *trylock* on L at time t or later, then p_i is returned *true*.*

⁶ A test-and-set object has two operations: *test-and-set*, which atomically reads the state of the object, sets the state to *true*, and returns the state read, and *reset*, which sets the state to *false*.

Algorithm 5: An implementation of a weak try-lock using registers (code for process p_i)

uses: $R[1, \dots, n]$ —array of registers
initially: $R[k] = 0$ for $k = 1, \dots, n$

```

1 operation trylock
2    $s \leftarrow getsum;$ 
3   if  $s > 0$  then return false;
4    $R[i].write(1);$ 
5    $s \leftarrow getsum;$ 
6   if  $s = 1$  then return true;
7    $R[i].write(0);$ 
8   return false;

9 operation unlock
10   $R[i].write(0);$ 
11  return ok;

12 function getsum
13   $s \leftarrow 0;$ 
14  for  $k = 1$  to  $n$  do  $s \leftarrow s + R[k].read;$ 
15  return  $s;$ 

```

While we do not know of any existing implementation of a weak try-lock, such an implementation can be easily obtained from several well-known (blocking) mutual exclusion algorithms, e.g., those proposed in (Lamport 1985) that ensure at least the *shutdown safety* property introduced in the same paper.

An example implementation of a weak try-lock using only registers, similar in concept to some of the lock implementations in (Lamport 1985), is given in Algorithm 5. The intuition behind the algorithm is the following. If a process p_i invokes operation *trylock* on a try-lock L implemented by the algorithm, p_i first checks whether any other process holds L (lines 2–3). If not, p_i announces that it wants to acquire L by setting register $R[i]$ to 1 (line 4). Then, p_i checks whether it is the only process that wants to acquire L (lines 5–6). If yes, then p_i acquires L (returns *true*). Otherwise, p_i resets $R[i]$ back to 0 (so that future invocations of *trylock* may succeed) and returns *false*. Clearly, if two processes execute *trylock* in parallel, then both can reach line 6. However, then at least one of them must observe that more than one register in array L is set to 1, and return *false*.

LEMMA 19. *Algorithm 5 implements a weak try-lock.*

Proof. Denote Algorithm 5 by A , and by L —a try-lock object implemented by A . First, it is straightforward to see that A is wait-free: it does not have any loops or waiting statements and all base objects used by A are wait-free.

Assume, by contradiction, that A does not ensure mutual exclusion. Hence, there is an implementation history E of A in which some two processes, say p_i and p_k , hold L at some time t . Consider only the latest *trylock* operations of p_i and p_k on L before t . Both of those operations must have returned *true*. Process p_i observes that $R[k] = 0$ in line 5, and so p_i reads $R[k]$ before p_k writes 1 to $R[k]$ in line 4. Hence, p_k reads $R[i]$ (in line 5) after p_i writes 1 to $R[i]$. Thus, p_k reads that $R[i]$ and $R[k]$ equal 1 and returns *false* in line 6—a contradiction.

It is easy to see that, for any process p_i , if $R[i] = 1$ then either p_i holds L or p_i is executing operation *trylock* on L . Hence, if a process p_i returns *false* from *trylock*, then either L is held by another process or another process is executing *trylock* concurrently with p_i . This means that L is a weak try-lock. \square

From Lemma 19, we obtain the following result:

THEOREM 20. *A weak try-lock has consensus number 1.*

It is straightforward to see that using weak try-locks instead of strong ones in the TM implementation shown in Algorithm 2 gives a TM that ensures weak progressiveness. Hence, by Theorem 20, we immediately prove the following result:

THEOREM 21. *Every weakly progressive TM has consensus number 1.*

6. Performance Trade-Off

We prove that the space complexity of every weakly (and, a fortiori, strongly) progressive TM that uses invisible reads is at least exponential with the number of t-variables available to transactions. The invisible reads strategy is used by a majority of TM implementations (Dice et al. 2006; Marathe et al. 2006; Harris et al. 2006; Adl-Tabatabai et al. 2006; Felber et al. 2008) as it allows efficient optimistic reading of t-variables. Intuitively, if invisible reads are used, a transaction that reads a t-variable does not write any information to base objects. Hence, many processors can concurrently execute transactions that read the same t-variables, without invalidating each other’s caches and causing high traffic on the inter-processor bus. However, transactions that update t-variables do not know whether there are any concurrent transactions that read those variables.

6.1 Semantics of Invisible Reads

We state our lower bound result assuming a simplified definition of the notion of invisible reads. This is sufficient for our lower bound proof, and is in agreement with what is ensured by various TM implementations (Dice et al. 2006; Marathe et al. 2006; Felber et al. 2008). Intuitively, we say that a TM implementation M uses invisible reads, if it does not modify the state of any base object when processing a *read* operation on any t-variable.

We capture this more precisely using the notion of a configuration. A *configuration* is the state of all base objects at a given point in time. Assuming that the initial state of base objects is fixed, and that base objects are deterministic, the configuration after any implementation history E can be precisely determined.

Let E be any implementation history of a TM. We define an *operation execution* of a process p_i in E to be any pair of (a) an invocation of operation *tread* or *twrite* and (b) the subsequent response of this operation in the sub-history $E|p_i$. If e is an operation execution of some process p_i in E , then every step in $E|p_i$ between the invocation and the response of e is said to be *corresponding to e*.

DEFINITION 22. *A TM implementation M uses invisible reads if, for every implementation history E of M , no step corresponding to an execution of operation *tread* in E changes the configuration.*

6.2 The Lower Bound

The *size* of a t-variable or a base object x can be defined as the number of distinct, reachable states of x . In particular, if x is a t-variable or a register object, then the size of x is the number of values that can be written to x . For example, the size of a 32-bit register is 2^{32} .

THEOREM 23. *Every weakly progressive TM implementation that uses invisible reads and provides to transactions N_s t-variables of size K_s uses $\Omega(K_s^{N_s}/K_b)$ base objects of size K_b .*

Due to space constraints, we give here only a sketch of the proof. The full proof can be found in the companion technical report (Guerraoui and Kapalka 2008b).

Proof. (sketch) Let M be any weakly progressive TM implementation that uses invisible reads and provides to transactions N_s t-variables x_1, \dots, x_{N_s} of size K_s . Basically, the proof consists of constructing a set of implementation histories E_1, \dots, E_L of M , such that the following conditions are satisfied.

1. The number L of those histories is $\Omega(K_s^{N_s})$.
2. Each history E_k is characterized by a pair of a base object b and a state w of b , such that (roughly speaking) base object b is never in state w within histories E_{k+1}, \dots, E_L .

This implies that M must use at least L/K_b base objects of size K_b (i.e., L base object-value pairs, each characterizing one implementation history E_k). We explain the construction of this set of histories in the following paragraphs.

Consider an implementation history E of M in which process p_1 executes a number of transactions, each of which writes to every t-variable. In parallel, process p_2 tries to execute in E a single transaction T_2 that reads every t-variable (i.e., takes a snapshot of all N_s t-variable values). Because (a) T_2 executes only *read* operations and (b) M uses invisible reads, process p_2 , while executing T_2 , cannot change the state of any base object until the commit time of T_2 . Hence, p_1 has no way to determine that p_2 is concurrently executing some transaction—in this sense T_2 is effectively invisible to p_1 .

Let Q be the set of configurations after every transaction of p_1 . If p_2 executes T_2 entirely between some two transactions of p_1 , then p_2 observes some configuration c from set Q . In this case, p_2 cannot abort T_2 , because there is no transaction concurrent to T_2 (and we assume M to be weakly progressive). Transaction T_2 must then return the values (v_1, \dots, v_{N_s}) written to t-variables by the latest transaction of p_1 executed before T_2 . Conversely, if p_2 determines that the current configuration is c , then p_2 cannot abort T_2 and has to return to T_2 values (v_1, \dots, v_{N_s}) . In this sense, the tuple (v_1, \dots, v_{N_s}) of t-variable values *corresponds* to configuration c , and vice versa. (Note that there may be many configurations in set Q that correspond to a given tuple of values, but at most one tuple of values can correspond to a given configuration.)

Process p_2 , in each step, can only read the state of a single base object. However, between any two steps of p_2 , process p_1 may execute any number of transactions and almost entirely change the configuration. (This is because the system is asynchronous and process p_1 is not aware of T_2 until the commit time of T_2 .) Hence, if every step of p_2 falls in between two transactions executed by p_1 , then p_2 , while executing T_2 , may observe (almost) any combination of configurations from set Q . But opacity requires that T_2 is returned values from one of the tuples written in E by the transactions of process p_1 . (Note that opacity requires that T_2 is returned a consistent snapshot of t-variable values regardless of whether T_2 eventually commits or aborts). That is, T_2 cannot be returned values from any tuple s that is not written by transactions of p_1 in E . Therefore, no combination of configurations from set Q can correspond to such an “illegal” tuple s .

Let $S = \{s_1, s_2, \dots, s_L\}$ be the set of all possible tuples of t-variable values ($L = K_s^{N_s}$). That is, each value s_i in S is a tuple (v_1, \dots, v_{N_s}) , where each v_k is a value of t-variable x_k . Let S_i denote the subset $\{s_i, \dots, s_L\}$ of S . The crux of the proof is to show that there exist a sequence of sets $Q_L \subset Q_{L-1} \subset \dots \subset Q_1$ that have the following properties:

1. Every element of Q_1 is a configuration just after some transaction of p_1 in some implementation history E of M in which p_1 executes a sequence of transactions, each writing to all t-variables.
2. Every set Q_k contains at least one configuration that corresponds to tuple s_k of t-variable values.

- For every configuration c_i in any set Q_k , we can find a non-empty sequence G_i of transactions, each writing to t-variables the values from a tuple in set S_k , such that if the system is in configuration c_i and process p_1 executes the sequence G_i , then the system is back to configuration c_i . That is, we can get from any configuration in set Q_k , which corresponds to a tuple from set S_k , back to the same configuration by executing only transactions that write tuples from set S_k .

The main reason why we can prove that such a sequence of sets exists is the fact that the number of configurations is finite (if it is infinite, Theorem 23 trivially holds). This means that if process p_1 executes infinitely many transactions, each writing values from a tuple in a set S_k , then some configurations (at least one corresponding to each tuple in S_k) must appear infinitely many times.

The properties of sets Q_1, \dots, Q_L imply that, for each set Q_k , we can construct an implementation history E_k of M , in which process p_1 executes transactions that write tuples from set S_k and process p_2 , while executing transaction T_2 that reads all t-variables, observes an *arbitrary* combination c of configurations from set Q_k (at least until the commit time of T_2). Basically, whenever p_2 is about to access a base object b , we let p_1 execute a sequence of transactions (each writing a tuple from S_k) that results in a configuration in which the state of b is the same as the state of b in (virtual) configuration c .

Clearly, a combination c of configurations from set Q_k cannot correspond to any tuple from set $S - S_k$. Otherwise, T_2 would read a snapshot of t-variable values that was never written by any transaction of p_1 in E_k , which would violate opacity. Therefore, we can show that, for each set S_k , there must be a pair (b_{l_k}, w_{l_k}) , where b_{l_k} is a base object and w_{l_k} its state, such that no configuration in set Q_{k+1} has b_{l_k} in state w_{l_k} . Hence, M needs $L = K_s^{N_s}$ such pairs, i.e., $K_s^{N_s} / K_b$ base objects. \square

6.3 Overcoming the Lower Bound

Our lower bound relies on three properties of a TM: weak progressiveness, operation wait-freedom, and invisible reads. It could seem that weakening (reasonably) any of those properties would allow overcoming the lower bound. We explain (informally) in the following paragraphs why this is not the case, and what has to be done in order for the lower bound not to hold.

Consider the following progress property, which is strictly weaker than weak progressiveness: if a transaction T_i is forcefully aborted, then there must be a transaction concurrent to T_i . We say that a TM that ensures this property is *non-trivial*—indeed, this seems like a basic requirement for a TM. However, non-trivial TMs do not overcome the complexity bound if they ensure operation wait-freedom and use invisible reads. Basically, in the proof of Theorem 23, transactions executed by processes p_1 and p_2 are not aware of any concurrent transactions, and so they will not be forcefully aborted in a non-trivial TM.

Consider the following liveness property, which we call *termination*: if a process p_i invokes an operation on a TM object, then p_i eventually returns from the operation. Clearly, termination is strictly weaker than wait-freedom. Consider a TM that ensures weak progressiveness and termination, and that uses invisible reads. Again, the complexity lower bound holds for such a TM: as in the proof of Theorem 23 process p_1 and p_2 is not aware of the operations executed by the other process, no process can block waiting for the other one to execute steps. Hence, in the particular execution used in the proof, termination would be sufficient.

Assume now that we allow a TM that uses invisible reads to update the state of a constant number of base objects in the first operation of every transaction, even if this operation is a *read*. We say then that such a TM uses *weak invisible reads*. Hence, each

transaction is allowed to announce its start. This means that, in the proof of Theorem 23, process p_1 can be aware of transaction T_2 executed by process p_2 . However, if the transactions executed by p_1 and p_2 access not all but almost all t-variables (all except for a constant number), then p_2 would not be allowed (in general) to forcefully abort its transactions, as there would be no guarantee that there is a conflict between those transactions and transaction T_2 .

This means that to overcome the lower bound we need to weaken more than one property of the TM. For example, TinySTM can be compiled with an option to enable a mechanism that handles timestamp overflows. (Without such a mechanism TinySTM can violate opacity in very long executions, as can TL2 or the LLT backend of RSTM.) Then, TinySTM uses weak invisible reads and may periodically violate both strong progressiveness and operation wait-freedom. Roughly speaking, once a transaction overflows a version number of a t-variable x , all transactions that access x are aborted, and all transactions that start afterwards are blocked on a barrier. Once there is no running transaction, the version number of x can be reset and transactions can proceed. This means that TinySTM ensures strong progressiveness and operation wait-freedom between timestamp overflows, but when an overflow happens the TM becomes only non-trivial and its operation-level liveness is reduced to termination.

7. Concluding Remarks

The two major assumptions we made in this paper were that t-variables support only *read* and *write* operations, and that the mapping between t-variables and corresponding try-locks is a one-to-one relation. We discuss here how those assumptions can be relaxed (at the price of increasing the complexity of the model and definitions). We also discuss the problem of model checking TMs for strong progressiveness.

Arbitrary t-variables. Object-based TMs support t-variables of arbitrary type. However, most of them classify all the operations of t-variables as either read-only or update ones. In those cases, there is no need to extend our simplified model, because read-only operations are effectively *reads*, and update operations are effectively pairs of *reads* and *writes*.

We can, however, imagine a TM that exploits the commutativity relations between some operations of t-variables of any type. In this case, one can extend the model of a TM to allow for arbitrary operations on t-variables, and redefine the notion of a conflict. Indeed, operations that commute should not conflict. Consider for example a counter object and its operation *inc* that increments the counter and does not return any meaningful value. It is easy to see that there is no real conflict between transactions that concurrently invoke operation *inc* on the same counter: the order of those operations does not matter and is not known to transactions (it would be, however, if *inc* returned the current value of the counter).

Once the notion of a conflict is defined, our definitions of progress properties remain correct even for t-variables with arbitrary operations. If we assume that a TM must support t-variables with operations *read* and *write* (in addition to other t-variables), then also the consensus number and complexity lower bound results hold for those TMs. However, the question of how to verify strong progressiveness of TM implementations with arbitrary t-variables is an open problem.

Arbitrary t-variable to try-lock mappings. Many lock-based TMs employ a hash function to map a t-variable to the corresponding try-lock. It may thus happen that a false conflict occurs between transactions that access disjoint sets of t-variables, and so, a priori, strong progressiveness might be violated. However, it is easy to take the hash function h of a TM implementation M into account in

the definition of strong progressiveness. Basically, when a transaction T_i reads or writes a t-variable x in a history H of M , we add to, respectively, the read set ($RSet_H(T_i)$) or the write set ($WSet_H(T_i)$) of T_i not only x , but also every t-variable y such that $h(x) = h(y)$. The definition of a conflict hence also takes into account false conflicts between transactions, and the strong progressiveness property can be ensured by M . (Such a property could be called *h-based strong progressiveness*.) It is important to note, however, that the hash function must be known to a user of a TM, or even provided by the user. Otherwise, strong progressiveness (and, for that matter, any other property that relies on the notion of a conflict) would no longer be visible, and very meaningful, to a user.

Model checking. While our reduction theorem simplifies proving strong progressiveness of a TM implementation, it might still be difficult to verify this property in an automatic manner. Indeed, even when verifying histories from the perspective of individual try-locks, we have to deal with an unbounded number of states. A solution would be to propose a reduction theorem along the lines of (Guerraoui et al. 2008), assuming that a TM implementation has certain symmetry properties. Two problems arise then. First, one has to express those properties in the fine-grained model we use ((Guerraoui et al. 2008) assumes operations like *validate* or *commit* to be atomic). Second, one has to prove that a given TM implementation ensures those properties, which is not always trivial (e.g., properties P6 and P7 in (Guerraoui et al. 2008)). Both problems remain open.

Acknowledgments

We would like to thank Leaf Petersen, Benjamin Pierce, and the anonymous reviewers for their invaluable help in improving the contents and the presentation of this paper. We would also like to thank Aleksandar Dragojević, Vincent Gramoli, Seth Gilbert, and Jan Vitek for their helpful comments and discussions.

References

- Martin Abadi, Andrew Birrell, Tim Harris, and Michael Isard. Semantics of transactional memory and automatic mutual exclusion. In *Proceedings of the 35th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL)*, 2008.
- Ali-Reza Adl-Tabatabai, Brian T. Lewis, Vijay Menon, Brian R. Murphy, Bratin Saha, and Tatiana Shpeisman. Compiler and runtime support for efficient software transactional memory. In *Proceedings of the ACM SIGPLAN 2006 Conference on Programming Language Design and Implementation (PLDI)*, 2006.
- Colin Blundell, E Christopher Lewis, and Milo M. K. Martin. Subtleties of transactional memory atomicity semantics. *IEEE Computer Architecture Letters*, 5(2), 2006.
- Dave Dice, Ori Shalev, and Nir Shavit. Transactional locking II. In *Proceedings of the 20th International Symposium on Distributed Computing (DISC)*, 2006.
- Robert Ennals. Software transactional memory should not be obstruction-free. Technical Report IRC-TR-06-052, Intel Research Cambridge Tech Report, Jan 2006.
- Kapali P. Eswaran, Jim N. Gray, Raymond A. Lorie, and Irving L. Traiger. The notions of consistency and predicate locks in a database system. *Commun. ACM*, 19(11):624–633, 1976.
- Pascal Felber, Torvald Riegel, and Christof Fetzer. Dynamic performance tuning of word-based software transactional memory. In *Proceedings of the 13th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming (PPoPP)*, Feb 2008.
- Rachid Guerraoui and Michał Kapalka. On the correctness of transactional memory. In *Proceedings of the 13th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming (PPoPP)*, 2008a.
- Rachid Guerraoui and Michał Kapalka. The semantics of progress in lock-based transactional memory. Technical Report LPD-REPORT-2008-015, EPFL, October 2008b.
- Rachid Guerraoui and Michał Kapalka. On obstruction-free transactions. In *Proceedings of the 20th ACM Symposium on Parallelism in Algorithms and Architectures (SPAA)*. ACM, June 2008c.
- Rachid Guerraoui, Thomas Henzinger, Barbara Jobstmann, and Vasu Singh. Model checking transactional memories. In *Proceedings of the ACM SIGPLAN 2008 Conference on Programming Language Design and Implementation (PLDI)*, 2008.
- Vassos Hadzilacos. A theory of reliability in database systems. *Journal of the ACM*, 35(1):121–145, 1988.
- Tim Harris, Mark Plesko, Avraham Shinnar, and David Tarditi. Optimizing memory transactions. In *Proceedings of the ACM SIGPLAN 2006 Conference on Programming Language Design and Implementation (PLDI)*, 2006.
- Maurice Herlihy. Wait-free synchronization. *ACM Transactions on Programming Languages and Systems*, 13(1):124–149, January 1991.
- Maurice Herlihy and J. Eliot B. Moss. Transactional memory: Architectural support for lock-free data structures. In *Proceedings of the 20th Annual International Symposium on Computer Architecture*, pages 289–300, May 1993.
- Maurice Herlihy and Jeannette M. Wing. Linearizability: a correctness condition for concurrent objects. *ACM Transactions on Programming Languages and Systems*, 12(3):463–492, June 1990.
- Suresh Jagannathan, Jan Vitek, Adam Welc, and Antony Hosking. A transactional object calculus. *Science of Computer Programming*, 57(2):164–186, 2005.
- Prasad Jayanti. Adaptive and efficient abortable mutual exclusion. In *Proceedings of the 22nd Annual ACM Symposium on Principles of Distributed Computing (PODC)*, 2003.
- Leslie Lamport. The mutual exclusion problem—part II: Statement and solutions. *Journal of the ACM*, 33(2), 1985.
- Leslie Lamport. On interprocess communication—part II: Algorithms. *Distributed Computing*, 1(2), 1986.
- Virendra J. Marathe, Michael F. Spear, Christopher Heriot, Athul Acharya, David Eisenstat, William N. Scherer III, and Michael L. Scott. Lowering the overhead of software transactional memory. In *1st ACM SIGPLAN Workshop on Transactional Computing (TRANSACT)*, Jun 2006.
- Vijay Menon, Steven Balensiefer, Tatiana Shpeisman, Ali-Reza Adl-Tabatabai, Richard L. Hudson, Bratin Saha, and Adam Welc. Practical weak-atomicity semantics for java stm. In *Proceedings of the 20th Annual Symposium on Parallelism in Algorithms and Architectures (SPAA)*, 2008.
- Katherine F. Moore and Dan Grossman. High-level small-step operational semantics for transactions. In *Proceedings of the 35th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL)*, 2008.
- Christos H. Papadimitriou. The serializability of concurrent database updates. *Journal of the ACM*, 26(4):631–653, 1979.
- Michel Raynal. *Algorithms for Mutual Exclusion*. The MIT Press, 1986.
- Michael L. Scott. Sequential specification of transactional memory semantics. In *1st ACM SIGPLAN Workshop on Transactional Computing (TRANSACT)*, 2006.
- Michael L. Scott and William N. Scherer III. Scalable queue-based spin locks with timeout. In *Proceedings of the 8th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming (PPoPP)*, 2001.
- Jan Vitek, Suresh Jagannathan, Adam Welc, and Antony Hosking. A semantic framework for designer transactions. In *Proceedings of the European Symposium on Programming (ESOP)*, March 2004.