

On the Semantics of Snapshot Isolation

Azalea Raad¹, Ori Lahav², and Viktor Vafeiadis¹

¹ MPI-SWS, Germany

² Tel Aviv University, Israel

{azalea,viktor}@mpi-sws.org

orilahav@tau.ac.il

Abstract. Snapshot isolation (SI) is a standard transactional consistency model used in databases, distributed systems and software transactional memory (STM). Its semantics is formally defined both declaratively as an acyclicity axiom, and operationally as a concurrent algorithm with memory bearing timestamps.

We develop two simpler equivalent operational definitions of SI as lock-based reference implementations that do not use timestamps. Our first locking implementation is prescient in that requires *a priori* knowledge of the data accessed by a transaction and carries out transactional writes eagerly (in-place). Our second implementation is non-prescient and performs transactional writes lazily by recording them in a local log and propagating them to memory at commit time. Whilst our first implementation is simpler and may be better suited for developing a program logic for SI transactions, our second implementation is more practical due to its non-prescience. We show that both implementations are sound and complete against the declarative SI specification and thus yield equivalent operational definitions for SI.

We further consider, for the first time *formally*, the use of SI in a context with racy non-transactional accesses, as can arise in STM implementations of SI. We introduce *robust snapshot isolation* (RSI), an adaptation of SI with similar semantics and guarantees in this mixed setting. We present a declarative specification of RSI as an acyclicity axiom and analogously develop two operational models as lock-based reference implementations (one eager, one lazy). We show that these operational models are both sound and complete against the declarative RSI model.

1 Introduction

Transactions are the *de facto* synchronisation mechanism in databases and geo-replicated distributed systems, and are thus gaining adoption in the shared-memory setting via *software transactional memory* (STM) [33,20]. In contrast to other synchronisation mechanisms, transactions readily provide atomicity, isolation, and consistency guarantees for sequences of operations, allowing programmers to focus on the high-level design of their systems.

However, providing these guarantees comes at a significant cost. As such, various transactional consistency models in the literature trade off consistency

guarantees for better performance. At nearly the one end of the spectrum, we have *serialisability* [28], which requires transactions to appear to have been executed in some total order. Serialisability provides strong guarantees, but is widely considered too expensive to implement. The main problem is that two conflicting transactions (e.g. one reading from and one updating the same datum) cannot both execute and commit in parallel.

Consequently, most major databases, both centralised (e.g. Oracle and MS SQL Server) and distributed [15,32,29], have opted for a slightly weaker model called *snapshot isolation* (SI) [7] as their default consistency model. SI has much better performance than serialisability by allowing conflicting transactions to execute concurrently and commit successfully as long as they do not have a write-write conflict. This in effect allows reads of SI transactions to read from an earlier memory snapshot than the one affected by their writes, and permits the *write skew anomaly* [11] depicted in Fig. 1. Besides this anomaly, however, SI is essentially the same as serialisability: Cerone et al. [11] provide a widely applicable condition under which SI and serialisability coincide for a given set of transactions. For these reasons, SI has also started gaining adoption in the generic programming language setting via STM implementations [1,16,8,26,25] that provide SI semantics for their transactions.

The formal study of SI, however, has so far not accounted for the more general STM setting in which both transactions and uninstrumented non-transactional code can access the same memory locations. While there exist two equivalent definitions of SI—one declarative in terms of an acyclicity constraint [10,11] and one operational in terms of an optimistic multi-version concurrency control algorithm [7]—neither definition supports *mixed-mode* (i.e. both transactional and non-transactional) accesses to the same locations. Extending the definitions to do so is difficult for two reasons: (1) the operational definition attaches a timestamp to every memory location, which heavily relies on the absence of non-transactional accesses; and (2) there are subtle interactions between the transactional implementation and the weak memory model underlying the non-transactional accesses.

In this article, we address these limitations of SI. We develop two simple lock-based reference implementations for SI that do not use timestamps. Our first implementation is *prescient* [19] in that it requires *a priori* knowledge of the data accessed by a transaction, and performs transactional writes *eagerly* (in-place). Our second implementation is non-prescient and carries out transactional writes *lazily* by first recording them in a local log and subsequently propagating them to memory at commit time. Our first implementation is simpler and may be better suited for understanding and developing a program logic for SI transactions, whilst our second implementation is more practical due to its non-prescience. We show that both implementations are sound and complete against the declarative SI specification and thus yield equivalent operational definitions for SI.

We then extend both our eager and lazy implementations to make them robust under uninstrumented non-transactional accesses, and characterise declaratively

the semantics we obtain. We call this extended model *robust snapshot isolation* (RSI) and show that it gives reasonable semantics with mixed-mode accesses.

To provide SI semantics, instead of timestamps, our implementations use *multiple-readers-single-writer* (MRSW) locks. They acquire locks in reader mode to take a snapshot of the memory locations accessed by a transaction and then promote the relevant locks to writer mode to enforce an ordering on transactions with write-write conflicts. As we discuss in §4, the equivalence of the RSI implementation and its declarative characterisation depends heavily upon the axiomatisation of MRSW locks: here, we opted for the weakest possible axiomatisation that does not order any concurrent reader lock operations and present an MRSW lock implementation that achieves this.

Outline In §2 we present an overview of our contributions by describing our reference implementations for both SI and RSI. In §3 we define the declarative framework for specifying STM programs. In §4 we present the declarative SI specification against which we demonstrate the soundness and completeness of our SI implementations. In §5 we formulate a declarative specification for RSI and demonstrate the soundness and completeness of our RSI implementations. We discuss related and future work in §6.³

2 Background and Main Ideas

As noted earlier, the key challenge in specifying STM transactions lies in accounting for the interactions between mixed-mode accesses to the same data. One simple approach is to treat each non-transactional access as a singleton mini-transaction and to provide *strong isolation* [27,9], i.e. full isolation between transactional and non-transactional code. This, however, requires *instrumenting* non-transactional accesses to adhere to same access policies as transactional ones (e.g. acquiring the necessary locks), which incurs a substantial performance penalty for non-transactional code. A more practical approach is to enforce isolation only amongst transactional accesses, an approach known as *weak isolation* [27,9], adopted by the relaxed transactions of C++ [2].

As our focus is on STMs with SI guarantees, instrumenting non-transactional accesses is not feasible. In particular, as we expect many more non-transactional accesses than transactional ones, we do not want to incur any performance degradation on non-transactional code when executed in parallel with transactional code. As such, we opt for an STM with SI guarantees under *weak isolation*. Under weak isolation, however, transactions with explicit abort instructions are problematic as their intermediate state may be observed by non-transactional code. As such, weakly isolated STMs (e.g. C++ relaxed transactions [2]) often forbid explicit aborts altogether. Throughout our development we thus make two simplifying assumptions: (1) transactions are not nested; and (2) there are no explicit abort instructions, following the example of weakly isolated relaxed

³ A full version of this article is available at [31].

| | | | | | |
|--|--|---|--|--|--|
| T1: $[a := x; //0 \parallel x := a + 1;$ (LU) Lost Update SI: ✗ | T2: $[b := x; //0 \parallel x := b + 1;$ | T1: $[a := x; //0 \parallel y := 1;$ (WS) Write Skew SI: ✓ | T2: $[b := y; //0 \parallel x := 1;$ | T1: $[y := 1;$ T3: $[a := x; //0 \parallel$ (WS2) Write Skew Variant SI: ✓ | T2: $[b := y; //0 \parallel x := 1;$ |
| T1: $[x := 1; \parallel y := 1;$ T3: $[a := y; //2 \parallel$ (LU2) Lost Update Variant SI: ✗ | T2: $[b := x; //0 \parallel y := 2;$ | $x := 1;$ T1: $[a := z; \parallel b := y; //0 \parallel$ (SBT) Store Buffering RSI: ✓ | $y := 1;$ T2: $[c := z; \parallel d := x; //0$ | $x := 1;$ $y := 1;$ T2: $[a := y; //1 \parallel b := x; //0$ (MPT) Message Passing RSI: ✗ | |

Fig. 1: Litmus tests illustrating transaction anomalies and their admissibility under SI and RSI. In all tests, initially, $x = y = z = 0$. The $//v$ annotation next to a read records the value read.

transactions of C++. As we describe later in §2.3, it is straightforward to lift the latter restriction (2) for our lazy implementations.

For non-transactional accesses, we naturally have to pick some consistency model. For simplicity and uniformity, we pick the release/acquire (RA) subset of the C++ memory model [6,23], a well-behaved platform-independent memory model, whose compilation to x86 requires no memory fences.

Snapshot Isolation (SI) The initial model of SI in [7] is described informally in terms of a multi-version concurrent algorithm as follows. A transaction T proceeds by taking a *snapshot* S of the shared objects. The execution of T is then carried out locally: read operations query S and write operations update S . Once T completes its execution, it attempts to *commit* its changes and succeeds *only if* it is not *write-conflicted*. Transaction T is write-conflicted if another *committed* transaction T' has written to a location also written to by T , since T recorded its snapshot. If T fails the conflict check it aborts and may restart; otherwise, it commits its changes, and its changes become visible to all other transactions that take a snapshot thereafter.

To realise this, the shared state is represented as a series of *multi-versioned* objects: each object is associated with a history of several versions at different *timestamps*. In order to obtain a snapshot, a transaction T chooses a *start-timestamp* t_0 , and reads data from the committed state as of t_0 , ignoring updates after t_0 . That is, updates committed after t_0 are invisible to T . In order to commit, T chooses a *commit-timestamp* t_c larger than any existing start- or commit-timestamp. Transaction T is deemed write-conflicted if another transaction T' has written to a location also written to by T *and* the commit-timestamp of T' is in the execution interval of T ($[t_0, t_c]$).

2.1 Towards an SI Reference Implementation without Timestamps

While the SI description above is suitable for understanding SI, it is not useful for integrating the SI model in a language such as C/C++ or Java. From a programmer’s perspective, in such languages the various threads directly access the *uninstrumented* (single-versioned) shared memory; they do not access their own instrumented snapshot at a particular timestamp, which is loosely related to the snapshots of other threads. Ideally, what we would therefore like is an equivalent description of SI in terms of accesses to uninstrumented shared memory and a synchronisation mechanism such as locks.

In what follows, we present our first lock-based reference implementation for SI that does not rely on timestamps. To do this, we assume that the locations accessed by a transaction can be statically determined. Specifically, we assume that each transaction T is supplied with its *read set*, RS , and *write set*, WS , containing those locations read and written by T , respectively (a static over-approximation of these sets suffices for soundness.). As such, our first reference implementation is *prescient* [19] in that it requires *a priori* knowledge of the locations accessed by the transaction. Later in §2.3 we lift this assumption and develop an SI reference implementation that is *non-prescient* and similarly does not rely on timestamps.

Conceptually, a candidate implementation of transaction T would (1) obtain a snapshot of the locations read by T ; (2) lock those locations written by T ; (3) execute T *locally*; and (4) unlock the locations written. The snapshot is obtained via `snapshot(RS)` in Fig. 3 where the values of locations in RS are recorded in a local array \mathbf{s} . The local execution of T is carried out by executing $\langle T \rangle$ in Fig. 3, which is obtained from T by (i) modifying read operations to read locally from the snapshot in \mathbf{s} , and (ii) updating the snapshot after each write operation. Note that the snapshot must be obtained *atomically* to reflect the memory state at a particular instance (*cf.* start-timestamp). An obvious way to ensure the snapshot atomicity is to lock the locations in the read set, obtain a snapshot, and unlock the read set. However, as we must allow for two transactions *reading* from the same location to execute in parallel, we opt for *multiple-readers-single-writer* (MRSW) locks.

Let us now try to make this general pattern more precise. As a first attempt, consider the implementation in Fig. 2a written in a simple while language, which releases all the reader locks at the end of the snapshot phase before acquiring any writer locks. This implementation is unsound as it admits the lost update (LU) anomaly in Fig. 1 disallowed under SI [11]. To understand this, consider a scheduling where T_2 runs between lines 3 and 4 of T_1 in Fig. 2a, which would result in T_1 having read a stale value. The problem is that the writer locks on WS are acquired too late, allowing two conflicting transactions to run concurrently. To address this, writer locks must be acquired early enough to pre-empt the concurrent execution of write-write-conflicting transactions. Note that locks have to be acquired early even for locations only written by a transaction to avoid exhibiting a variant of the lost update anomaly (LU2).

| | | |
|--|---|---|
| <pre> 1. for (x ∈ RS) lock_r x 2. snapshot(RS); 3. for (x ∈ RS) unlock_r x 4. for (x ∈ WS) lock_w x 5. (T); 6. for (x ∈ WS) unlock_w x </pre> <p style="text-align: center;">(a)</p> | <pre> 1. for (x ∈ WS) lock_w x; 2. for (x ∈ RS \ WS) lock_r x 3. snapshot(RS); 4. for (x ∈ RS \ WS) unlock_r x 5. (T); 6. for (x ∈ WS) unlock_w x </pre> <p style="text-align: center;">(b)</p> | <pre> 1. for (x ∈ RS ∪ WS) lock_r x 2. snapshot(RS); 3. for (x ∈ RS ∪ WS) { 4. if (x ∈ WS) promote x 5. else unlock_r x; } 6. (T); 7. for (x ∈ WS) unlock_w x </pre> <p style="text-align: center;">(c)</p> |
| <p>Sound: ✗ allows (LU), (LU2)</p> | <p>Sound: ✓ Complete: ✗ disallows (WS)</p> | <p>Sound: ✓ Complete: ✗ disallows (WS2)</p> |

Fig. 2: Candidate SI implementations of transaction T given read/write sets RS,WS

As such, our second candidate implementation in Fig. 2b brings forward the acquisition of writer locks. Whilst this implementation is sound (and disallows lost update), it nevertheless disallows behaviours deemed valid under SI such as the write skew anomaly (WS) in Fig. 1, and is thus incomplete. The problem is that such early acquisition of writer locks not only pre-empts concurrent execution of *write-write-conflicting* transactions, but also those of *read-write-conflicting* transactions (e.g. WS) due to the exclusivity of writer locks.

To remedy this, in our third candidate implementation in Fig. 2c we first acquire weaker reader locks on all locations in RS or WS, and later *promote* the reader locks on WS to exclusive writer ones, while releasing the reader locks on RS. The promotion of a reader lock signals its intent for exclusive ownership and awaits the release of the lock by other readers before claiming it exclusively as a writer. To avoid deadlocks, we further assume that $RS \cup WS$ is ordered so that locks are promoted in the same order by all threads.

Although this implementation is “more complete” than the previous one, it is still incomplete as it disallows certain behaviour admitted by SI. In particular, consider a variant of the write skew anomaly (WS2) depicted in Fig. 1, which is admitted under SI, but not admitted by this implementation.

To understand why this is admitted by SI, recall the operational SI model using timestamps. Let the domain of timestamps be that of natural numbers \mathbb{N} . The behaviour of (WS2) can be achieved by assigning the following execution intervals for T1: $[t_0^{T1}=2, t_c^{T1}=2]$; T2: $[t_0^{T2}=1, t_c^{T2}=4]$; and T3: $[t_0^{T3}=3, t_c^{T3}=3]$. To see why the implementation in Fig. 2c does not admit the behaviour in (WS2), let us assume without loss of generality that x is ordered before y . Upon executing lines 3-5, a) T1 promotes y ; b) T2 promotes x and then c) releases the reader lock on y ; and d) T3 releases the reader lock on x . To admit the behaviour in (WS2), the release of y in (c) must occur before the promotion of y in (a) since otherwise T2 cannot read 0 for y . Similarly, the release of x in (d) must occur before its promotion in (b). On the other hand, since T3 is executed by the same thread after T1, we know that (a) occurs before (d). This however leads to circular execution: (b)→(c)→(a)→(d)→(b), which cannot be realised.

| | |
|--|---|
| <pre> 0. LS := ∅; 1. for (x ∈ RS ∪ WS) lock_r x 2. snapshot(RS); 3. for (x ∈ RS \ WS) unlock_r x 4. for (x ∈ WS) { 5. if (can-promote x) LS.add(x) 6. else { 7. for (x ∈ LS) unlock_w x 8. for (x ∈ WS \ LS) unlock_r x 9. goto line 0 } 10. } 11. (T); 12. for (x ∈ WS) unlock_w x </pre> | <pre> snapshot(RS) \triangleq for (x ∈ RS) s_x := x snapshot_{RSI}(RS) \triangleq start: for (x ∈ RS) s_x := x for (x ∈ RS) { if (s_x != x) goto start } (a:=x) \triangleq a:=s_x (x:=a) \triangleq x:=a; s_x:=a (S₁; S₂) \triangleq (S₁); (S₂) (while(e) S) \triangleq while(e) (S) ... and so on ... </pre> |
|--|---|

Fig. 3: SI implementation of transaction T given RS, WS; the code in blue ensures deadlock avoidance. The RSI implementation (§5) is obtained by replacing snapshot on line 2 with snapshot_{RSI}.

To overcome this, in our final candidate execution in Fig. 3 (ignoring the code in blue), after obtaining a snapshot, we *first* release the reader locks on RS, and *then* promote the reader locks on WS, rather than simultaneously in one pass. As we demonstrate in §4, the implementation in Fig. 3 is both *sound and complete* against its declarative SI specification.

Avoiding Deadlocks As two distinct reader locks on x may simultaneously attempt to promote their locks, promotion is done on a ‘first-come-first-served’ basis to avoid *deadlocks*. A call to `can-promote x` by reader *r* thus returns a boolean denoting either (i) successful promotion (true); or (ii) failed promotion as another reader *r'* is currently promoting a lock on x (false). In the latter case, *r* must release its reader lock on x to ensure the successful promotion of x by *r'* and thus avoid deadlocks. To this end, our implementation in Fig. 3 includes a deadlock avoidance mechanism (code in blue) as follows. We record a list LS of those locks on the write set that have been successfully promoted so far. When promoting a lock on x succeeds (line 5), the LS is extended with x. On the other hand, when promoting x fails (line 6), all those locks promoted so far (i.e. in LS) as well as those yet to be promoted (i.e. in WS \ LS) are released and the transaction is restarted.

Remark 1. Note that the deadlock avoidance code in blue does not influence the correctness of the implementation in Fig. 3, and is merely included to make the reference implementation more realistic. In particular, the implementation without the deadlock avoidance code is both sound and complete against the SI specification, provided that the conditional `can-promote` call on line 5 is replaced by the blocking `promote` call.

Avoiding Over-Synchronisation due to MRSW Locks Consider the store buffering program (SBT) shown in Fig. 1. If, for a moment, we ignore transactional accesses, our underlying memory model (RA)—as well as all other weak memory models—allows the annotated weak behaviour. Intuitively, placing the two transactions that only *read* z in (SBT) should still allow the weak behaviour since the two transactions do not need to synchronise in any way. Nevertheless, most MRSW lock implementations forbid this outcome because they use a single global counter to track the number of readers that have acquired the lock, which inadvertently also synchronises the readers with one another. As a result, the two read-only transactions act as memory fences forbidding the weak outcome of (SBT). To avoid such synchronisation, in the technical appendix [31] we provide a different MRSW implementation using a separate location for each thread so that reader lock acquisitions do not synchronise.

To keep the presentation simple, we henceforth assume an abstract specification of a MRSW lock library providing operations for acquiring/releasing reader/writer locks, as well as promoting reader locks to writer ones. We require that (1) calls to writer locks (to acquire, release or promote) *synchronise* with all other calls to the lock library; and (2) writer locks provide *mutual exclusion* while held. We formalise these notions in §4. These requirements do not restrict synchronisation between *two read* lock calls: two read lock calls may or may not synchronise. Synchronisation between read lock calls is relevant only for the completeness of our RSI implementation (handling mixed-mode code); for that result, we further require that (3) read lock calls not synchronise.

2.2 Handling Racy Mixed-Mode Accesses

Let us consider what happens when data accessed by a transaction is modified concurrently by an uninstrumented atomic non-transactional write. Since such writes do not acquire any locks, the snapshots taken may include values written by non-transactional accesses. The result of the snapshot then depends on the order in which the variables are read. Consider the (MPT) example in Fig. 1. In our implementation, if in the snapshot phase y is read before x , then the annotated weak behaviour is not possible because the underlying model (RA) disallows this weak “message passing” behaviour. If, however, x is read before y , then the weak behaviour is possible. In essence, this means that the SI implementation described so far is of little use when there are races between transactional and non-transactional code. Technically, our SI implementation violates *monotonicity* with respect to wrapping code inside a transaction. The weak behaviour of the (MPT) example is disallowed by RA if we remove the transaction block T2, and yet it is exhibited by our SI implementation with the transaction block.

To get monotonicity under RA, it suffices for the snapshots to read the variables in the same order they are accessed by the transactions. Since a static calculation of this order is not always possible, following [30], we achieve this by reading each variable twice. In more detail, our `snapshotRSI` implementation in Fig. 3 takes *two* snapshots of the locations read by the transaction, and checks that they both return the same values for each location. This ensures

that every location is read both before and after every other location in the transaction, and hence all the high-level happens-before orderings in executions of the transactional program are also respected by its implementation. As we demonstrate in §5, our RSI implementation is both *sound and complete* against our proposed declarative semantics for RSI. There is however one caveat: since equality of values is used to determine whether the two snapshots agree, we will miss cases where different non-transactional writes to a location write the same value. In our formal development (see §5), we thus assume that if multiple non-transactional writes write the same value to the same location, they cannot race with the same transaction. Note that this assumption cannot be lifted without instrumenting non-transactional writes, and thus impeding performance substantially. That is, to lift this restriction we must instead replace every non-transactional *write* $x := v$ with `lock_w x; x := v; unlock_w x`.

2.3 Non-Prescient Reference Implementations without Timestamps

Recall that the SI and RSI implementations in §2.1 are prescient in that they require knowledge of the read and write sets of transactions beforehand. In what follows we present alternative SI and RSI implementations that are *non-prescient*.

Non-Prescient SI Reference Implementation In Fig. 4 we present a *lazy* lock-based reference implementation for SI. This implementation is *non-prescient* and does not require *a priori* knowledge of the read set \mathbf{RS} and the write set \mathbf{WS} . Rather, the \mathbf{RS} and \mathbf{WS} are computed on the fly as the execution of the transaction unfolds. As with the SI implementation in Fig. 3, this implementation does not rely on timestamps and uses MRSW locks to synchronise concurrent accesses to shared data. As before, the implementation consults a local *snapshot* at \mathbf{s} for read operations. However, unlike the eager implementation in Fig. 3 where transactional writes are performed *in-place*, the implementation in Fig. 4 is *lazy* in that it logs the writes in the local array \mathbf{s} and propagates them to memory at commit time, as we describe shortly.

Ignoring the code in blue, the implementation in Fig. 4 proceeds with initialising \mathbf{RS} and \mathbf{WS} with \emptyset (line 1); it then populates the local snapshot array at \mathbf{s} with initial value \perp for each location \mathbf{x} (line 2). It then executes $\langle \mathbf{T} \rangle$ which is obtained from \mathbf{T} as follows. For each *read* operation $\mathbf{a} := \mathbf{x}$ in \mathbf{T} , first the value of $\mathbf{s}[\mathbf{x}]$ is inspected to ensure it contains a snapshot of \mathbf{x} . If this is not the case (i.e. $\mathbf{x} \notin \mathbf{RS} \cup \mathbf{WS}$), a reader lock on \mathbf{x} is acquired, a snapshot of \mathbf{x} is recorded in $\mathbf{s}[\mathbf{x}]$, and the read set \mathbf{RS} is extended with \mathbf{x} . The snapshot value in $\mathbf{s}[\mathbf{x}]$ is subsequently returned in \mathbf{a} . Analogously, for each write operation $\mathbf{x} := \mathbf{a}$, the \mathbf{WS} is extended with \mathbf{x} , and the written value is lazily logged in $\mathbf{s}[\mathbf{x}]$. Recall from our candidate executions in Fig. 2 that to ensure implementation correctness, for each written location \mathbf{x} , the implementation must first acquire a reader lock on \mathbf{x} , and subsequently promote it to a writer lock. As such, for each write operation in \mathbf{T} , the implementation first checks if a reader lock for \mathbf{x} has been acquired (i.e. $\mathbf{x} \in \mathbf{RS} \cup \mathbf{WS}$) and obtains one if this is not the case.

| | |
|--|--|
| <pre> 0. LS:=∅; 1. RS:=∅; WS:=∅; 2. for (x∈Locs) s[x]:=⊥ 3. ⟨T⟩; 4. for (x∈RS\WS) unlock_r x 5. for (x∈WS) { 6. if (can-promote x) LS.add(x) 7. else { 8. for (x∈LS) unlock_w x 9. for (x∈WS\LS) unlock_r x 10. goto line 0 } } 11. for (x∈WS) x:=s[x] 12. for (x∈WS) unlock_w x </pre> | <pre> ⟨a:=x⟩ ≜ if (x∉RS∪WS) { lock_r x; RS.add(x); s[x]:=x; } a:=s[x]; ⟨x:=a⟩ ≜ if (x∉RS∪WS) lock_r x; WS.add(x); s[x]:=a; ⟨S₁;S₂⟩ ≜ ⟨S₁⟩;⟨S₂⟩ ⟨while(e) S⟩ ≜ while(e) ⟨S⟩ ... and so on... </pre> |
|--|--|

Fig. 4: Non-prescient SI implementation of transaction T with RS and WS computed on the fly; the code in blue ensures deadlock avoidance.

Once the execution of $\langle T \rangle$ is completed, the implementation proceeds to *commit* the transaction. To this end, the reader locks on RS are released (line 4), reader locks on WS are promoted to writer ones (line 6), the writes logged in \mathbf{s} are propagated to memory (line 11), and finally the writer locks on WS are released (line 12). As we demonstrate later in §4, the implementation in Fig. 4 is both sound and complete against the declarative SI specification.

Note that the implementation in Fig. 4 is optimistic in that it logs the writes performed by the transaction in the local array \mathbf{s} and propagates them to memory at commit time, rather than performing the writes *in-place* as with its pessimistic counterpart in Fig. 3. As before, the code in blue ensures deadlock avoidance and is identical to its counterpart in Fig. 3. As before, this deadlock avoidance code does not influence the correctness of the implementation and is merely included to make the reference implementation more practical.

Non-Prescient RSI Reference Implementation In Fig. 5 we present a *lazy* lock-based reference implementation for RSI. As with its SI counterpart, this implementation is non-prescient and computes the RS and WS on the fly. As before, the implementation does not rely on timestamps and uses MRSW locks to synchronise concurrent accesses to shared data. Similarly, the implementation consults the local *snapshot* at \mathbf{s} for read operations, whilst logging write operations lazily in a *write sequence* at \mathbf{wseq} , as we describe shortly.

Recall from the RSI implementation in §2.1 that to ensure snapshot validity, each location is read twice to preclude intermediate non-transactional writes. As such, when writing to a location x , the initial value read (recorded in \mathbf{s}) must not be *overwritten* by the transaction to allow for subsequent validation of the snapshot. To this end, for each location x , the snapshot array \mathbf{s} contains a *pair* of values, (r, c) , where r denotes the snapshot value (initial value read), and c denotes the current value which may have overwritten the snapshot value.

Recall that under weak isolation, the intermediate values written by a transaction may be observed by non-transactional reads. For instance, given the

| | |
|---|---|
| <pre> 0. LS:=∅; 1. RS:=∅; WS:=∅; wseq:= []; 2. for (x∈Locs) s[x]:= (⊥, ⊥) 3. $\langle\mathbb{T}\rangle$; 4. for (x∈RS) {(r,-):=s[x]; 5. if (x!=r) { //read x again 6. for (x∈RS∪WS) unlock_r x 7. goto line 0 } } 8. for (x∈RS\WS) unlock_r x 9. for (x∈WS) { 10. if (can-promote x) LS.add(x) 11. else { 12. for (x∈LS) unlock_w x 13. for (x∈WS\LS) unlock_r x 14. goto line 0 } } 15. for ((x,v)∈wseq) x:=v 16. for (x∈WS) unlock_w x </pre> | <pre> (a:=x) \triangleq if (x∉RS∪WS) { lock_r x; RS.add(x); r:=x; s[x]:= (r,r); } (-,c):=s[x]; a:=c; (x:=a) \triangleq if (x∉RS∪WS) lock_r x WS.add(x); (r,-):=s[x]; s[x]:= (r,a); wseq:=wseq++[(x,a)]; (S₁;S₂) \triangleq (S₁);(S₂) (while(e) S) \triangleq while(e) (S) ... and so on... </pre> |
|---|---|

Fig. 5: Non-prescient RSI implementation of transaction \mathbb{T} with \mathbf{RS} and \mathbf{WS} computed on the fly; the code in blue ensures deadlock avoidance.

\mathbb{T} : $[x := 1; x := 2 \parallel a := x]$ program, the non-transactional read $a := x$, may read either 1 or 2 for x . As such, at commit time, it is not sufficient solely to propagate the last written value (in program order) to each location (e.g. to propagate only the $x := 2$ write in the example above). Rather, to ensure implementation completeness, one must propagate all written values to memory, in the order they appear in the transaction body. To this end, we track the values written by the transaction as a (FIFO) *write sequence* at location \mathbf{wseq} , containing items of the form (x, v) , denoting the location written (x) and the associated value (v).

Ignoring the code in blue, the implementation in Fig. 5 initialises \mathbf{RS} and \mathbf{WS} with \emptyset , initialises \mathbf{wseq} as an empty sequence $[]$ (line 1), and populates the local snapshot array \mathbf{s} with initial value (\perp, \perp) for each location x (line 2). It then executes $\langle\mathbb{T}\rangle$, obtained from \mathbb{T} in an analogous manner to that in Fig. 4. For every read $\mathbf{a}:=\mathbf{x}$ in $\langle\mathbb{T}\rangle$, the current value recorded for x in \mathbf{s} (namely c when $\mathbf{s}[x]$ holds $(-, c)$) is returned in \mathbf{a} . Dually, for every write $\mathbf{x}:=\mathbf{a}$ in $\langle\mathbb{T}\rangle$, the current value recorded for x in \mathbf{s} is updated to \mathbf{a} , and the write is logged in the write sequence \mathbf{wseq} by appending (x, \mathbf{a}) to it.

Upon completion of $\langle\mathbb{T}\rangle$, the snapshot in \mathbf{s} is *validated* (lines 4-7). Each location x in \mathbf{RS} is thus read again and its value is compared against the snapshot value in $\mathbf{s}[x]$. If validation fails (line 5), the locks acquired are released (line 6) and the transaction is restarted (line 7).

If validation succeeds, the transaction is committed: the reader locks on \mathbf{RS} are released (line 8), the reader locks on \mathbf{WS} are promoted (line 10), the writes in \mathbf{wseq} are propagated to memory in FIFO order (line 15), and finally the writer locks on \mathbf{WS} are released (line 16).

As we show in §5, the implementation in Fig. 5 is both sound and complete against our proposed declarative specification for RSI. As before, the code in blue ensures deadlock avoidance; it does not influence the implementation correctness and is merely included to make the implementation more practical.

Supporting Explicit Abort Instructions It is straightforward to extend the lazy implementations in Fig. 4 and Fig. 5 to handle transactions containing explicit abort instructions. More concretely, as the effects (writes) of a transaction are logged locally and are not propagated to memory until commit time, upon reaching an `abort` in $\langle T \rangle$ no roll-back is necessary, and one can simply release the locks acquired so far and return. That is, one can extend $\langle \cdot \rangle$ in Fig. 4 and Fig. 5, and define $\langle \text{abort} \rangle \triangleq \text{for } (x \in \text{RS} \cup \text{WS}) \text{ unlock_r } x; \text{ return.}$

3 A Declarative Framework for STM

We present the notational conventions used in the remainder of this article, and describe a general framework for declarative concurrency models. Later in this article, we present SI, its extension with non-transactional accesses, and their lock-based implementations as instances of this general definition.

Notation Given a relation r on a set A , we write $r^?$, r^+ and r^* for the reflexive, transitive and reflexive-transitive closure of r , respectively. We write r^{-1} for the inverse of r ; $r|_A$ for $r \cap (A \times A)$; $[A]$ for the identity relation on A , i.e. $\{(a, a) \mid a \in A\}$; $\text{irreflexive}(r)$ for $\nexists a. (a, a) \in r$; and $\text{acyclic}(r)$ for $\text{irreflexive}(r^+)$. Given two relations r_1 and r_2 , we write $r_1; r_2$ for their (left) relational composition, i.e. $\{(a, b) \mid \exists c. (a, c) \in r_1 \wedge (c, b) \in r_2\}$. Lastly, when r is a strict partial order, we write $r|_{\text{imm}}$ for the *immediate* edges in r : $\{(a, b) \in r \mid \nexists c. (a, c) \in r \wedge (c, b) \in r\}$.

Assume finite sets of *locations* LOC ; *values* VAL ; *thread identifiers* TID , and *transaction identifiers* TXID . We use x, y, z to range over locations, v over values, τ over thread identifiers, and ξ over transaction identifiers.

Definition 1 (Events). *An event is a tuple $\langle n, \tau, \xi, l \rangle$, where $n \in \mathbb{N}$ is an event identifier, $\tau \in \text{TID} \uplus \{0\}$ is a thread identifier (0 is used for initialisation events), $\xi \in \text{TXID} \uplus \{0\}$ is a transaction identifier (0 is used for non-transactional events), and l is an event label that takes one of the following forms:*

- A memory access label: $\text{R}(x, v)$ for reads; $\text{W}(x, v)$ for writes; and $\text{U}(x, v_r, v_w)$ for updates.
- A lock label: $\text{RL}(x)$ for reader lock acquisition; $\text{RU}(x)$ for reader lock release; $\text{WL}(x)$ for writer lock acquisition; $\text{WU}(x)$ for writer lock release; and $\text{PL}(x)$ for reader to writer lock promotion.

We typically use a, b , and e to range over events. The functions tid , tx , lab , typ , loc , val_r and val_w respectively project the thread identifier, transaction identifier, label, type (in $\{\text{R}, \text{W}, \text{U}, \text{RL}, \text{RU}, \text{WL}, \text{WU}, \text{PL}\}$), location, and read/written values of an event, where applicable. We assume only reads and writes are used in transactions ($\text{tx}(a) \neq 0 \implies \text{typ}(a) \in \{\text{R}, \text{W}\}$).

Given a relation r on events, we write r_{loc} for $\{(a, b) \in r \mid \text{loc}(a) = \text{loc}(b)\}$. Analogously, given a set A of events, we write A_x for $\{a \in A \mid \text{loc}(a) = x\}$.

Definition 2 (Execution graphs). *An execution graph, G , is a tuple of the form $(E, \text{po}, \text{rf}, \text{mo}, \text{lo})$, where:*

- E is a set of events, assumed to contain a set E_0 of initialisation events, consisting of a write event with label $\mathbb{W}(x, 0)$ for every $x \in \text{Loc}$. The sets of read events in E is denoted by $\mathcal{R} \triangleq \{e \in E \mid \text{typ}(e) \in \{\text{R}, \text{U}\}\}$; write events by $\mathcal{W} \triangleq \{e \in E \mid \text{typ}(e) \in \{\text{W}, \text{U}\}\}$; update events by $\mathcal{U} \triangleq \mathcal{R} \cap \mathcal{W}$; and lock events by $\mathcal{L} \triangleq \{e \in E \mid \text{typ}(e) \in \{\text{RL}, \text{RU}, \text{WL}, \text{WU}, \text{PL}\}\}$. The sets of reader lock acquisition and release events, \mathcal{RL} and \mathcal{RU} , writer lock acquisition and release events, \mathcal{WL} and \mathcal{WU} , and lock promotion events \mathcal{PL} are defined analogously. The set of transactional events in E is denoted by \mathcal{T} ($\mathcal{T} \triangleq \{e \in E \mid \text{tx}(e) \neq 0\}$); and the set of non-transactional events is denoted by \mathcal{NT} ($\mathcal{NT} \triangleq E \setminus \mathcal{T}$).
- $\text{po} \subseteq E \times E$ denotes the ‘program-order’ relation, defined as a disjoint union of strict total orders, each ordering the events of one thread, together with $E_0 \times (E \setminus E_0)$ that places the initialisation events before any other event. We assume that events belonging to the same transaction are ordered by po , and that any other event po -between them also belongs to the same transaction.
- $\text{rf} \subseteq \mathcal{W} \times \mathcal{R}$ denotes the ‘reads-from’ relation, defined between write and read events of the same location with matching read and written values; it is total and functional on reads, i.e. every read is related to exactly one write.
- $\text{mo} \subseteq \mathcal{W} \times \mathcal{W}$ denotes the ‘modification-order’ relation, defined as a disjoint union of strict total orders, each ordering the write events on one location.
- $\text{lo} \subseteq \mathcal{L} \times \mathcal{L}$ denotes the ‘lock-order’ relation, defined as a disjoint union of strict orders, each of which (partially) ordering the lock events to one location.

In the context of an execution graph $G=(E, \text{po}, \text{rf}, \text{mo}, \text{lo})$ —we often use “ G .” as a prefix to make this explicit—the ‘same-transaction’ relation, $\text{st} \in \mathcal{T} \times \mathcal{T}$, is the equivalence relation given by $\text{st} \triangleq \{(a, b) \in \mathcal{T} \times \mathcal{T} \mid \text{tx}(a) = \text{tx}(b)\}$. Given a relation $r \subseteq E \times E$, we write $r_{\mathcal{T}}$ for lifting r to transaction classes: $r_{\mathcal{T}} \triangleq \text{st}; (r \setminus \text{st}); \text{st}$. For instance, when $(w, r) \in \text{rf}$, w is a transaction ξ_1 event and r is a transaction ξ_2 event, then all all events in ξ_1 are $\text{rf}_{\mathcal{T}}$ -related to all events in ξ_2 . Analogously, we write $r_{\mathcal{I}}$ to restrict r to its *intra-transactional* edges (within a transaction): $r_{\mathcal{I}} \triangleq r \cap \text{st}$; and write $r_{\mathcal{E}}$ to restrict r to its *extra-transactional* edges (outside a transaction): $r_{\mathcal{E}} \triangleq r \setminus \text{st}$. Lastly, the ‘reads-before’ relation is defined by $\text{rb} \triangleq (\text{rf}^{-1}; \text{mo}) \setminus [E]$. Intuitively, rb relates a read r to all writes w that are mo -after the write r reads from; i.e. when $(w', r) \in \text{rf}$ and $(w', w) \in \text{mo}$, then $(r, w) \in \text{rb}$. In the transactional literature, this is known as the *anti-dependency* relation [4,3].

Execution graphs of a given program represent traces of shared memory accesses generated by the program. The set of execution graphs associated with programs written in our while language can be straightforwardly defined by induction over the structure of programs as in e.g. [35]. Each execution of a program P has a particular program *outcome*, prescribing the final values of local variables in each thread. In this initial stage, the execution outcomes are almost unrestricted as there are very few constraints on the rf , mo and lo relations. Such restrictions and thus the permitted outcomes of a program are determined by defining the set of *consistent* executions, which is defined separately for each model we consider. Given a program P and a model M , the set $\text{outcomes}_M(P)$ collects the outcomes of every M -consistent execution of P .

4 Snapshot Isolation (SI)

We present a declarative specification of SI and demonstrate that the SI implementations presented in Fig. 3 and Fig. 4 are both sound and complete with respect to the SI specification.

In [11] Cerone and Gotsman developed a declarative specification for SI using dependency graphs [4,3]. Below we adapt their specification to the notation of §3. As with [11], throughout this section, we take *SI execution graphs* to be those in which $E = \mathcal{T} \subseteq (\mathcal{R} \cup \mathcal{W}) \setminus \mathcal{U}$. That is, the SI model handles transactional code only, consisting solely of read and write events (excluding updates).

Definition 3 (SI consistency [11]). *An SI execution $G = (E, \text{po}, \text{rf}, \text{mo}, \text{lo})$ is SI-consistent if the following conditions hold:*

- $\text{rf}_I \cup \text{mo}_I \cup \text{rb}_I \subseteq \text{po}$ (INT)
- $\text{acyclic}((\text{po}_T \cup \text{rf}_T \cup \text{mo}_T); \text{rb}_T?)$ (EXT)

Informally, (INT) ensures the consistency of each transaction internally, while (EXT) provides the synchronisation guarantees among transactions. In particular, we note that the two conditions together ensure that if two read events in the same transaction read from the same location x , and no write to x is po-between them, then they must read from the same write (known as ‘internal read consistency’).

Next, we provide an alternative equivalent formulation of SI-consistency which will serve as the basis of our extension with non-transactional accesses in §5.

Proposition 1. *An SI execution $G = (E, \text{po}, \text{rf}, \text{mo}, \text{lo})$ is SI-consistent if and only if INT holds and the ‘SI-happens-before’ relation $\text{si-hb} \triangleq (\text{po}_T \cup \text{rf}_T \cup \text{mo}_T \cup \text{si-rb})^+$ is irreflexive, where $\text{si-rb} \triangleq [\mathcal{R}_E]; \text{rb}_T; [\mathcal{W}]$ and $\mathcal{R}_E \triangleq \{r \mid \exists w. (w, r) \in \text{rf}_E\}$.*

Proof. The full proof is given in the technical appendix [31].

Intuitively, SI-happens-before orders events of different transactions; this order is due to either the program order (po_T), or synchronisation enforced by the implementation ($\text{rf}_T \cup \text{mo}_T \cup \text{si-rb}$). By contrast, events of the same transaction are unordered, as the implementation may well execute them in a different order (in particular, by taking a snapshot, it executes external reads before the writes).

In more detail, the rf_T corresponds to transactional synchronisation due to *causality*, i.e. when one transaction T_2 observes an effect of an earlier transaction T_1 . The inclusion of rf_T ensures that T_2 cannot read from T_1 without observing its *entire* effect. This in turn ensures that transactions exhibit ‘all-or-nothing’ behaviour: they cannot mix-and-match the values they read. For instance, if T_1 writes to both x and y , transaction T_2 may not read x from T_1 but read y from an earlier (in ‘happens-before’ order) transaction T_0 .

The mo_T corresponds to transactional synchronisation due to *write-write conflicts*. Its inclusion enforces write-conflict-freedom of SI transactions: if T_1 and T_2 both write to x via events w_1 and w_2 such that $(w_1, w_2) \in \text{mo}$, then T_1 must commit before T_2 , and thus its entire effect must be visible to T_2 .

To understand **si-rb**, first note that \mathcal{R}_E denotes the *external* transactional reads (i.e. those reading a value written by another transaction). That is, the \mathcal{R}_E are the read events that get their values from the transactional snapshot phases. By contrast, internal reads (those reading a value written by the same transaction) happen only after the snapshot is taken. Now let there be an **rb_T** edge between two transactions, T_1 and T_2 . This means there exist a read event r of T_1 and a write event w of T_2 such that $(r, w) \in \mathbf{rb}$; i.e. there exists w' such that $(w', r) \in \mathbf{rf}$ and $(w', w) \in \mathbf{mo}$. If r reads internally (i.e. w' is an event in T_1), then T_1 and T_2 are conflicting transactions and as accounted by **mo_T** described above, all events of T_1 happen before those of T_2 . Now, let us consider the case when r reads externally (w' is not in T_1). From the timestamped model of SI, there exists a start-timestamp $t_0^{T_1}$ as of which the T_1 snapshot (all its external reads including r) is recorded. Similarly, there exists a commit-timestamp $t_c^{T_2}$ as of which the updates of T_2 (including w) are committed. Moreover, since $(r, w) \in \mathbf{rb}$ we know $t_0^{T_1} < t_c^{T_2}$ (otherwise r must read the value written by w and not w'). That is, we know all events in the snapshot of T_1 (i.e. all external reads in T_1) happen before all writes of T_2 .⁴

We use the declarative framework in §3 to formalise the semantics of our implementation. Here, our programs include only non-transactional code, and thus *implementation execution graphs* are taken as those in which $\mathcal{T} = \emptyset$. Furthermore, we assume that locks in implementation programs are used in a *well-formed* manner: the sequence of lock events for *each location*, in each thread (following po), should match (a prefix of) the regular expression $(\mathbf{RL} \cdot \mathbf{RU} \mid \mathbf{WL} \cdot \mathbf{WU} \mid \mathbf{RL} \cdot \mathbf{PL} \cdot \mathbf{WU})^*$. For instance, a thread never releases a lock, without having acquired it earlier in the program. As a consistency predicate on execution graphs, we use the C11 release/acquire consistency augmented with certain constraints on lock events.

Definition 4. *An implementation execution graph $G = (E, \text{po}, \mathbf{rf}, \mathbf{mo}, \mathbf{lo})$ is RA-consistent if the following hold, where $\mathbf{hb} \triangleq (\text{po} \cup \mathbf{rf} \cup \mathbf{lo})^+$ denotes the ‘RA-happens-before’ relation:*

- $\forall x. \forall a \in \mathcal{WL}_x \cup \mathcal{WU}_x \cup \mathcal{PL}_x, b \in \mathcal{L}_x. a = b \vee (a, b) \in \mathbf{lo} \vee (b, a) \in \mathbf{lo}$ (WSYNC)
- $[\mathcal{WL} \cup \mathcal{PL}]; (\mathbf{lo} \setminus \text{po}); [\mathcal{L}] \subseteq \text{po}; [\mathcal{WU}]; \mathbf{lo}$ (WEX)
- $[\mathcal{RL}]; (\mathbf{lo} \setminus \text{po}); [\mathcal{WL} \cup \mathcal{PL}] \subseteq \text{po}; [\mathcal{RU} \cup \mathcal{PL}]; \mathbf{lo}$ (RSHARE)
- acyclic($\mathbf{hb}_{loc} \cup \mathbf{mo} \cup \mathbf{rb}$) (ACYC)

The (WSYNC) states that write lock calls (to acquire, release or promote) *synchronise* with all other calls to the same lock.

The next two constraints ensure the ‘single-writer-multiple-readers’ paradigm. In particular, (WEX) states that write locks provide *mutual exclusion* while held: any lock event l of thread τ **lo**-after a write lock acquisition or promotion event l' of another thread τ' , is **lo**-after a subsequent write lock release event u of τ' (i.e. $(l', u) \in \text{po}$ and $(u, l) \in \mathbf{lo}$). As such, the lock cannot be acquired (in read or write mode) by another thread until it has been released by its current owner.

The (RSHARE) analogously states that once a thread acquires a lock in read mode, the lock cannot be acquired in write mode by other threads until it has

⁴ By taking **rb_T** instead of **si-rb** in Prop. 1 one obtains a characterisation of *serialisability*.

either been released, or promoted to a writer lock (and subsequently released) by its owner. Note that this does not preclude other threads from simultaneously acquiring the lock in read mode. In the technical appendix [31] we present two MRSW lock implementations that satisfy the conditions outlined above.

The last constraint (**ACYC**) is that of C11 RA consistency [23], with the **hb** relation extended with **lo**.

Remark 2. Our choice of implementing the SI STMs on top of the RA fragment is purely for presentational convenience. Indeed, it is easy to observe that execution graphs of $\langle P \rangle$ are data race free, and thus, **ACYC** could be replaced by any condition that implies $\forall x. ([\mathcal{W}_x]; (\text{po} \cup \text{lo})^+; [\mathcal{W}_x]; (\text{po} \cup \text{lo})^+; [\mathcal{R}_x]) \cap \text{rf} = \emptyset$ and that is implied by $\text{acyclic}(\text{po} \cup \text{rf} \cup \text{lo} \cup \text{mo} \cup \text{rb})$. In particular, the C11 non-atomic accesses or sequentially consistent accesses may be used.

We next show that our SI implementations in Fig. 3 and Fig. 4 are *sound and complete* with respect to the declarative specification given above. The proofs are non-trivial and the full proofs are given in the technical appendix [31].

Theorem 1 (Soundness and completeness). *Let P be a transactional program; let $\langle P \rangle_E$ denote its eager implementation as given in Fig. 3 and $\langle P \rangle_L$ denote its lazy implementation as given in Fig. 4. Then:*

$$\text{outcomes}_{\text{SI}}(P) = \text{outcomes}_{\text{RA}}(\langle P \rangle_E) = \text{outcomes}_{\text{RA}}(\langle P \rangle_L)$$

Proof. The full proofs for both implementations is given in the technical appendix [31].

Stronger MRSW Locks As noted in §2, for both (prescient and non-prescient) SI implementations our soundness and completeness proofs show that the same result holds for a stronger lock specification, in which reader locks synchronise as well. Formally, this specification is obtained by adding the following to Def. 4:

$$- \forall x. \forall a, b \in \mathcal{RL}_x \cup \mathcal{RU}_x. a = b \vee (a, b) \in \text{lo} \vee (b, a) \in \text{lo} \quad (\text{RSYNC})$$

Soundness of this stronger specification ($\text{outcomes}_{\text{RA}}(\langle P \rangle_x) \subseteq \text{outcomes}_{\text{SI}}(P)$ for $x \in \{E, L\}$) follows immediately from Thm. 1. Completeness ($\text{outcomes}_{\text{SI}}(P) \subseteq \text{outcomes}_{\text{RA}}(\langle P \rangle_x)$ for $x \in \{E, L\}$), however, is more subtle, as we need to additionally satisfy (**RSYNC**) when constructing **lo**. While we can do so for SI, it is essential for the completeness of our RSI implementations that reader locks not synchronise, as shown by (**SBT**) in §2.

In the technical appendix [31] we present two MRSW lock implementations *sound* against the **lo** conditions in Def. 4. Additionally, the first implementation is *complete* against the conditions of Def. 4 augmented with (**RSYNC**), whilst the second is complete against the conditions of Def. 4 alone.

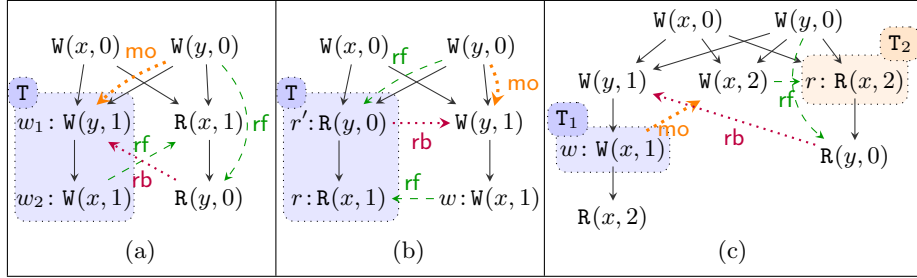


Fig. 6: RSI-inconsistent executions due to (a) rsi-po; (b) $[\mathcal{N}\mathcal{T}]$; rf; st; (c) $(\text{mo}; \text{rf})_{\mathcal{T}}$

5 Robust Snapshot Isolation (RSI)

We explore the semantics of SI STMs in the presence of non-transactional code with *weak isolation* guarantees (see §2). We refer to this model as *robust snapshot isolation* (RSI), due to its ability to provide SI guarantees between transactions even in the presence of non-transactional code. We propose the first declarative specification of RSI programs and develop two lock-based reference implementations that are both *sound and complete* against our proposed specification.

A Declarative Specification of RSI STMs We formulate a declarative specification of RSI semantics by adapting the SI semantics in Prop. 1 to account for non-transactional accesses. To specify the abstract behaviour of RSI programs, *RSI execution graphs* are taken to be those in which $\mathcal{L} = \emptyset$. Moreover, as with SI graphs, RSI execution graphs are those in which $\mathcal{T} \subseteq (\mathcal{R} \cup \mathcal{W}) \setminus \mathcal{U}$. That is, RSI transactions comprise solely read and write events, excluding updates.

Definition 5 (RSI consistency). *An execution $G = (E, \text{po}, \text{rf}, \text{mo}, \text{lo})$ is RSI-consistent iff INT holds and acyclic($\text{rsi-hb}_{\text{loc}} \cup \text{mo} \cup \text{rb}$), where $\text{rsi-hb} \triangleq (\text{rsi-po} \cup \text{rsi-rf} \cup \text{mo}_{\mathcal{T}} \cup \text{si-rb})^+$ is the ‘RSI-happens-before’ relation, with $\text{rsi-po} \triangleq (\text{po} \setminus \text{po}_{\mathcal{I}}) \cup [\mathcal{W}]$; $\text{po}_{\mathcal{I}}; [\mathcal{W}]$ and $\text{rsi-rf} \triangleq (\text{rf}; [\mathcal{N}\mathcal{T}]) \cup ([\mathcal{N}\mathcal{T}]; \text{rf}; \text{st}) \cup \text{rf}_{\mathcal{T}} \cup (\text{mo}; \text{rf})_{\mathcal{T}}$.*

As with SI and RA, we characterise the set of executions admitted by RSI as graphs that lack cycles of certain shapes. To account for non-transactional accesses, similar to RA, we require $\text{rsi-hb}_{\text{loc}} \cup \text{mo} \cup \text{rb}$ to be acyclic (recall that $\text{rsi-hb}_{\text{loc}} \triangleq \{(a, b) \in \text{rsi-hb} \mid \text{loc}(a) = \text{loc}(b)\}$). The RSI-happens-before relation rsi-hb includes both the synchronisation edges enforced by the transactional implementation (as in si-hb), and those due to non-transactional accesses (as in hb of the RA consistency). The rsi-hb relation itself is rather similar to si-hb . In particular, the $\text{mo}_{\mathcal{T}}$ and si-rb subparts can be justified as in si-hb ; the difference between the two lies in rsi-po and rsi-rf .

To justify rsi-po , recall from §4 that si-hb includes $\text{po}_{\mathcal{T}}$. The rsi-po is indeed a strengthening of $\text{po}_{\mathcal{T}}$ to account for non-transactional events: it additionally includes (i) po to and from non-transactional events; and (ii) po between two *write* events in a transaction. We believe (i) comes as no surprise to the reader;

for (ii), consider the execution graph in Fig. 6a, where transaction T is denoted by the dashed box labelled T, comprising the write events w_1 and w_2 . Removing the T block (with w_1 and w_2 as non-transactional writes), this execution is deemed inconsistent, as this weak “message passing” behaviour is disallowed in the RA model. We argue that the analogous transactional behaviour in Fig. 6a must be similarly disallowed to maintain monotonicity with respect to wrapping non-transactional code in a transaction (see Thm. 3). As in SI, we cannot include the entire po in rsi-hb because the write-read order in transactions is not preserved by the implementation.

Similarly, rsi-rf is a strengthening of rf_T to account for non-transactional events: in the absence of non-transactional events rsi-rf reduces to $\text{rf}_T \cup (\text{mo}; \text{rf})_T$ which is contained in si-hb. The $\text{rf}; [\mathcal{NT}]$ part is required to preserve the ‘happens-before’ relation for non-transactional code. That is, as rf is included in the hb relation of underlying memory model (RA), it is also included in rsi-hb.

The $[\mathcal{NT}]; \text{rf}; \text{st}$ part asserts that in an execution where a read event r of transaction T reads from a non-transactional write w , the snapshot of T reads from w and so all events of T happen after w . Thus, in Fig. 6b, r' cannot read from the overwritten initialisation write to y .

For the $(\text{mo}; \text{rf})_T$ part, consider the execution graph in Fig. 6c where there is a write event w of transaction T_1 and a read event r of transaction T_2 such that $(w, r) \in \text{mo}; \text{rf}$. Then, transaction T_2 must acquire the read lock of $\text{loc}(w)$ after T_1 releases the writer lock, which in turn means that every event of T_1 happens before every event of T_2 .

Remark 3. Recall that our choice of modelling SI and RSI STMs in the RA fragment is purely for presentational convenience (see Remark 2). Had we chosen a different model, the RSI consistency definition (Def. 5) would largely remain unchanged, with the exception of $\text{rsi-rf} \triangleq (\text{sw}; [\mathcal{NT}]) \cup ([\mathcal{NT}]; \text{sw}; \text{st}) \cup \text{rf}_T \cup (\text{mo}; \text{rf})_T$, where in the highlighted changes the rf relation is replaced with sw, denoting the ‘synchronises-with’ relation. As in the RA model $\text{sw} \triangleq \text{rf}$, we have inlined this in Def. 5.

SI and RSI Consistency We next demonstrate that in the absence of non-transactional code, the definitions of SI-consistency (Prop. 1) and RSI-consistency (Def. 5) coincide. That is, for all executions G , if $G.\mathcal{NT} = \emptyset$, then G is SI-consistent if and only if G is RSI-consistent.

Theorem 2. *For all executions G , if $G.\mathcal{NT} = \emptyset$, then:*

$$G \text{ is SI-consistent} \iff G \text{ is RSI-consistent}$$

Proof. The full proof is given in the technical appendix [31].

Note that the above theorem implies that for all transactional programs P , if P contains no non-transactional accesses, then $\text{outcomes}_{\text{SI}}(P) = \text{outcomes}_{\text{RSI}}(P)$.

RSI Monotonicity We next prove the *monotonicity* of RSI when wrapping non-transactional events into a transaction. That is, wrapping a block of non-transactional code inside a new transaction does not introduce additional behaviours. More concretely, given a program P , when a block of non-transactional code in P is wrapped inside a new transaction to obtain a new program P_T , then $\text{outcomes}_{\text{RSI}}(P_T) \subseteq \text{outcomes}_{\text{RSI}}(P)$. This is captured in the theorem below, with its full proof given in the technical appendix [31].

Theorem 3 (Monotonicity). *Let P_T and P be RSI programs such that P_T is obtained from P by wrapping a block of non-transactional code inside a new transaction. Then:*

$$\text{outcomes}_{\text{RSI}}(P_T) \subseteq \text{outcomes}_{\text{RSI}}(P)$$

Proof. The full proof is given in the technical appendix [31].

Lastly, we show that our RSI implementations in §2 (Fig. 3 and Fig. 5) are sound and complete with respect to Def. 5. This is captured in the theorem below. The soundness and completeness proofs are non-trivial; the full proofs are given in the technical appendix [31].

Theorem 4 (Soundness and completeness). *Let P be a program that possibly mixes transactional and non-transactional code. Let $(P)_E$ denote its eager RSI implementation as given in Fig. 3 and $(P)_L$ denote its lazy RSI implementation as given in Fig. 5. If for every location x and value v , every RSI-consistent execution of P contains either (i) at most one non-transactional write of v to x ; or (ii) all non-transactional writes of v to x are happens-before-ordered with respect to all transactions accessing x , then:*

$$\text{outcomes}_{\text{RSI}}(P) = \text{outcomes}_{\text{RA}}((P)_E) = \text{outcomes}_{\text{RA}}((P)_L)$$

Proof. The full proofs for both implementations are given in the technical appendix [31].

6 Related and Future Work

Much work has been done in formalising the semantics of weakly consistent *database transactions* [7,34,3,4,10,12,11,13,14,18], both operationally and declaratively. On the operational side, Berenson et al. [7] gave an operational model of SI as a multi-version concurrent algorithm. Later, Sovran et al. [34] described and operationally defined the *parallel snapshot isolation* model (PSI), as a close relative of SI with weaker guarantees.

On the declarative side, Adya et al. [3,4] introduced *dependency graphs* (similar to execution graphs of our framework in §3) for specifying transactional semantics and formalised several ANSI isolation levels. Cerone et al. [10,12] introduced *abstract executions* and formalised several isolation levels including SI and PSI. Later in [11], they used dependency graphs of Adya to develop equivalent SI

and PSI semantics; recently in [13], they provided a set of algebraic laws for connecting these two declarative styles.

To facilitate client-side reasoning about the behaviour of database transactions, Gotsman et al. [18] developed a proof rule for proving invariants of client applications under a number of consistency models.

Recently, Kaki et al. [21] developed a program logic to reason about transactions under ANSI SQL isolation levels (including SI). To do this, they formulated an operational model of such programs (parametric in the isolation level). They then proved the soundness of their *logic* with respect to their proposed operational model. However, the authors did not establish the *soundness* or *completeness* of their *operational model* against existing formal semantics, e.g. [11]. The lack of the completeness result means that their proposed operational model may exclude behaviours deemed valid by the corresponding declarative models. This is a particular limitation as possibly many valid behaviours cannot be shown correct using the logic and is thus detrimental to its usability.

By contrast, transactional semantics in the STM setting with mixed (both transactional and non-transactional) accesses is under-explored on both operational and declarative sides. Recently, Dongol et al. [17] applied execution graphs [5] to specify *serialisable* STM programs under weak memory models. Raad et al. [30] formalised the semantics of PSI STMs declaratively (using execution graphs) and operationally (as lock-based reference implementations). Neither work, however, handles the semantics of SI STMs with weak isolation guarantees.

Finally, Khyzha et al. [22] formalise the sufficient conditions on STMs and their programs that together ensure strong isolation. That is, non-transactional accesses can be viewed as singleton transactions (transactions containing single instructions). However, their conditions require *serialisability* for fully transactional programs, and as such, RSI transactions do not meet their conditions. Nevertheless, we conjecture that a DRF guarantee for strong atomicity, similar to [22], may be established for RSI. That is, if all executions of a fully transactional program have no races between singleton and non-singleton transactions, then it is safe to replace all singleton transactions by non-transactional accesses.

In the future, we plan to build on the work presented here by developing reasoning techniques that would allow us to verify properties of STM programs. This can be achieved by either extending existing program logics for weak memory, or developing new ones for currently unsupported models. In particular, we can reason about the SI models presented here by developing custom proof rules in the existing program logics for RA such as [24,35].

Acknowledgements We thank the VMCAI reviewers for their constructive feedback. The first author was supported in part by a European Research Council (ERC) Consolidator Grant for the project “RustBelt”, under the European Union’s Horizon 2020 Framework Programme (grant agreement number 683289). The second author was supported by the Israel Science Foundation (grant number 5166651), and by Len Blavatnik and the Blavatnik Family foundation.

References

1. The Clojure Language: Refs and Transactions, <http://clojure.org/refs>
2. Technical specification for C++ extensions for transactional memory (2015), <http://www.open-std.org/jtc1/sc22/wg21/docs/papers/2015/n4514.pdf>
3. Adya, A.: Weak consistency: A generalized theory and optimistic implementations for distributed transactions. Ph.D. thesis, MIT (1999)
4. Adya, A., Liskov, B., O’Neil, P.: Generalized isolation level definitions. In: Proceedings of the 16th International Conference on Data Engineering. pp. 67–78 (2000)
5. Alglave, J., Maranget, L., Tautschnig, M.: Herding cats: Modelling, simulation, testing, and data mining for weak memory. *ACM Trans. Program. Lang. Syst.* **36**(2), 7:1–7:74 (2014)
6. Batty, M., Owens, S., Sarkar, S., Sewell, P., Weber, T.: Mathematizing C++ concurrency. In: Proceedings of the 38th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages. pp. 55–66 (2011)
7. Berenson, H., Bernstein, P., Gray, J., Melton, J., O’Neil, E., O’Neil, P.: A critique of ANSI SQL isolation levels. In: Proceedings of the 1995 ACM SIGMOD International Conference on Management of Data. pp. 1–10 (1995)
8. Bieniusa, A., Fuhrmann, T.: Consistency in hindsight: A fully decentralized STM algorithm. In: Proceedings of the 2010 IEEE International Symposium on Parallel and Distributed Processing, IPDPS 2010. pp. 1 – 12 (2010)
9. Blundell, C., C. Lewis, E., M. K. Martin, M.: Deconstructing transactions: The subtleties of atomicity. In: 4th Annual Workshop on Duplicating, Deconstructing, and Debunking (2005)
10. Cerone, A., Bernardi, G., Gotsman, A.: A framework for transactional consistency models with atomic visibility. In: Proceedings of the 26th International Conference on Concurrency Theory. pp. 58–71 (2015)
11. Cerone, A., Gotsman, A.: Analysing snapshot isolation. In: Proceedings of the 2016 ACM Symposium on Principles of Distributed Computing. pp. 55–64 (2016)
12. Cerone, A., Gotsman, A., Yang, H.: Transaction chopping for parallel snapshot isolation. In: Proceedings of the 29th International Symposium on Distributed Computing - Volume 9363. pp. 388–404 (2015)
13. Cerone, A., Gotsman, A., Yang, H.: Algebraic laws for weak consistency. In: CONCUR (2017)
14. Crooks, N., Pu, Y., Alvisi, L., Clement, A.: Seeing is believing: A client-centric specification of database isolation. In: Proceedings of the ACM Symposium on Principles of Distributed Computing. pp. 73–82. PODC ’17, ACM, New York, NY, USA (2017). <https://doi.org/10.1145/3087801.3087802>, <http://doi.acm.org/10.1145/3087801.3087802>
15. Daudjee, K., Salem, K.: Lazy database replication with snapshot isolation. In: Proceedings of the 32Nd International Conference on Very Large Data Bases. pp. 715–726 (2006)
16. Dias, R.J., Distefano, D., Seco, J.a.C., Lourenço, J.a.M.: Verification of snapshot isolation in transactional memory java programs. In: Proceedings of the 26th European Conference on Object-Oriented Programming. pp. 640–664. ECOOP’12, Springer-Verlag, Berlin, Heidelberg (2012). <https://doi.org/10.1007/978-3-642-31057-7-28>, <http://dx.doi.org/10.1007/978-3-642-31057-7-28>
17. Dongol, B., Jagadeesan, R., Riely, J.: Transactions in relaxed memory architectures. *Proc. ACM Program. Lang.* **2**(POPL), 18:1–18:29 (Dec 2017). <https://doi.org/10.1145/3158106>, <http://doi.acm.org/10.1145/3158106>

18. Gotsman, A., Yang, H., Ferreira, C., Najafzadeh, M., Shapiro, M.: 'cause i'm strong enough: Reasoning about consistency choices in distributed systems. In: Proceedings of the 43rd Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages. pp. 371–384. POPL '16, ACM, New York, NY, USA (2016). <https://doi.org/10.1145/2837614.2837625>, <http://doi.acm.org/10.1145/2837614.2837625>
19. Harris, T., Larus, J., Rajwar, R.: Transactional Memory, 2Nd Edition. Morgan and Claypool Publishers, 2nd edn. (2010)
20. Herlihy, M., Moss, J.E.B.: Transactional memory: Architectural support for lock-free data structures. In: Proceedings of the 20th Annual International Symposium on Computer Architecture. pp. 289–300 (1993)
21. Kaki, G., Nagar, K., Najafzadeh, M., Jagannathan, S.: Alone together: Compositional reasoning and inference for weak isolation. Proc. ACM Program. Lang. **2**(POPL), 27:1–27:34 (Dec 2017). <https://doi.org/10.1145/3158115>, <http://doi.acm.org/10.1145/3158115>
22. Khyzha, A., Attiya, H., Gotsman, A., Rinetzky, N.: Safe privatization in transactional memory. In: Proceedings of the 23rd ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming. pp. 233–245 (2018)
23. Lahav, O., Giannarakis, N., Vafeiadis, V.: Taming release-acquire consistency. In: Proceedings of the 43rd Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages. pp. 649–662 (2016)
24. Lahav, O., Vafeiadis, V.: Owicki-Gries reasoning for weak memory models. In: Proceedings, Part II, of the 42Nd International Colloquium on Automata, Languages, and Programming - Volume 9135. pp. 311–323 (2015)
25. Litz, H., Cheriton, D., Firoozshahian, A., Azizi, O., Stevenson, J.P.: SI-TM: Reducing transactional memory abort rates through snapshot isolation. SIGPLAN Not. pp. 383–398 (2014)
26. Litz, H., Dias, R.J., Cheriton, D.R.: Efficient correction of anomalies in snapshot isolation transactions. ACM Trans. Archit. Code Optim. **11**(4), 65:1–65:24 (Jan 2015). <https://doi.org/10.1145/2693260>
27. Martin, M., Blundell, C., Lewis, E.: Subtleties of transactional memory atomicity semantics. IEEE Comput. Archit. Lett. **5**(2), 17–17 (2006)
28. Papadimitriou, C.H.: The serializability of concurrent database updates. J. ACM **26**(4), 631–653 (Oct 1979). <https://doi.org/10.1145/322154.322158>, <http://doi.acm.org/10.1145/322154.322158>
29. Peng, D., Dabek, F.: Large-scale incremental processing using distributed transactions and notifications. In: Proceedings of the 9th USENIX Conference on Operating Systems Design and Implementation. pp. 251–264 (2010)
30. Raad, A., Lahav, O., Vafeiadis, V.: On parallel snapshot isolation and release/acquire consistency. In: Proceedings of the 27th European Symposium on Programming (2018), to appear
31. Raad, A., Lahav, O., Vafeiadis, V.: The technical appendix for this paper is available at <https://arxiv.org/abs/1805.06196> (2018)
32. Serrano, D., Patino-Martinez, M., Jimenez-Peris, R., Kemme, B.: Boosting database replication scalability through partial replication and 1-copy-snapshot-isolation. In: Proceedings of the 13th Pacific Rim International Symposium on Dependable Computing. pp. 290–297 (2007)
33. Shavit, N., Touitou, D.: Software transactional memory. In: Proceedings of the Fourteenth Annual ACM Symposium on Principles of Distributed Computing. pp. 204–213 (1995)

34. Sovran, Y., Power, R., Aguilera, M.K., Li, J.: Transactional storage for geo-replicated systems. In: Proceedings of the Twenty-Third ACM Symposium on Operating Systems Principles. pp. 385–400 (2011)
35. Vafeiadis, V., Narayan, C.: Relaxed separation logic: A program logic for C11 concurrency. In: Proceedings of the 2013 ACM SIGPLAN International Conference on Object Oriented Programming Systems Languages & Applications. pp. 867–884 (2013)