

PMVERIFY: Robustness Verification for Checking Crash Consistency of Non-volatile Memory

Abstract

The emerging non-volatile memory (NVM) technologies provide competitive performance with DRAM and ensure data persistence in the event of system failure. However, it exhibits weak behaviour on hardware architectures in terms of the order in which stores are committed to NVMs, and therefore requires developers to manually flush pending writes. To ensure correctness of this error-prone task, it is crucial to develop a rigid method to check crash consistency of programs running on NVM devices. Most existing solutions are based on testing and rely on user input to dynamically detect such deficiencies. In this paper, we present a fully automated method to statically verify robustness, a newly established property for ensuring crash consistency of such programs. The method is based on the observation that, reachability of a post-crash non-volatile state under a given pre-crash execution can be reduced to validity of the pre-crash execution with additional ordering constraints under memory consistency model. Our robustness verification algorithm employs a search-based framework to explore all partial executions and states, and checks if any non-volatile state is reachable under certain pre-crash execution. Once a reachable non-volatile state is obtained, we further check its reachability under memory consistency model. The algorithm is implemented in a prototype tool PMVERIFY that leverages symbolic encoding of the program and utilizes an SMT solver to efficiently explore all executions and states. A dedicated theory solver is integrated into the DPLL(T) framework to optimize the robustness checking algorithm. Experiments on the PMDK example benchmark show that PMVERIFY is the first tool to establish robustness, and is competitive with the state-of-the-art dynamic tool, PSAN, in terms of robustness violation detection.

Keywords: persistent memory, non-volatile memory, robustness, program verification, crash consistency

1 Introduction

Non-volatile memory (a.k.a. NVM, or persistent memory) is a kind of non-conventional, byte-addressable storage device that preserves its content after a power failure [35, 36]. It enables direct access to persistent data using standard load and store instructions, and thus avoids the overhead of OS system calls. Due to its competitive performance with DRAM and guarantee of data persistence, it has been widely used in persistency-critical systems such as databases [4, 47, 49] and file systems [8, 13, 33, 42, 58–61, 65].

However, modern processors have write-back caches that induce non-determinism in the order stores are written to memory. Since cache systems are volatile, it may lead to data loss if some stores have not been committed to NVM when a crash happens. The exact order in which stores are written back to NVM, referred to as *persist order*, is constrained by the cache coherence protocol. Similar to memory consistency models which specify visibility order of memory operations, in recent works the Intel-x86 [10, 37, 51–53] and ARMv8 [10, 54] persistency models have been formalized which prescribes the persist order. Both architectures exhibit weak behaviours in terms of persist order.

As a simple example on Intel-x86, assume crash happens after executing the two instructions $a = 1$; $b = 1$. Upon recovery, it is possible to observe the non-volatile state $a = 0$; $b = 1$ (we assume 0 is the initial value of a). In general, persist order might differ from the order memory operations are made visible. Figure 1 shows a possible execution of these two instructions and relevant orders. Here the store $a = 1$ is issued and becomes visible first per program order, but remains in caches. On the contrary, the store $b = 1$ is issued later but leaves the cache before the system fails. The store $a = 1$ in the volatile cache is thus lost due to the crash.

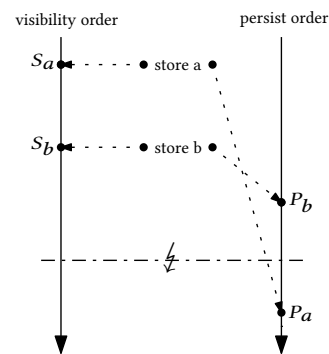


Figure 1. A possible visibility and persist order of two stores $a = 1$; $b = 1$; in a single thread. S_a, S_b are the points the stores are made visible to all threads, and P_a, P_b are when they are committed to NVM. ζ signifies system failure.

Overall, persistent programming is an error-prone task. It is the responsibility of the developers to avoid corruption of data residing on NVMs, since any inconsistency would persist across reboots. This necessitates a clear understanding of the persistency semantics. Although instructions (e.g. `clflush` and `clflushopt` on Intel-x86) have been introduced to constrain persist order of memory operations, the fact that stores are committed to NVMs in an out-of-order

manner can be counter-intuitive. The matter becomes even more intricate for a multi-threaded program.

To assist developers in correctly programming NVMs, researchers primarily pursue two approaches. On the one hand, high-level mechanisms such as *transactions* [6, 11, 18–20, 27, 54] and *locks* [5, 7, 26, 29, 44] have been developed to facilitate the development of NVM programs. However, these mechanisms often introduce significant overhead. On the other hand, another line of research focuses on enhancing the reliability of NVM programs [9, 12, 16, 23, 34, 43, 45, 46, 48, 55].

An important property, known as *crash consistency*, has been proposed to characterize the reliability of NVM programs. It ensures that the program state recovered from NVM after a system failure is always consistent, thereby enabling seamless resumption of program execution [55]. However, a considerable amount of existing tools require specifications provided by users for accurate bug detection. For instance, XFDetector [45] requires user annotation of commit variables to avoid false alarms, while PMTest [46] and PMDebugger [12] requires explicit annotation of ordering constraints in the program. [48] is able to prove correctness of NVM programs in terms of persistency invariant, a predicate that always holds on recovered state, but it is also restricted to the specifications provided by the user. The model checkers Jaaru [23] and Yat [43] do not require user input, but they only detect observable bugs, i.e. segmentation fault or assertion violation.

To circumvent the aforementioned difficulties in crash consistency checking, Gorjiara et al. [22] proposed a novel correctness criterion called *robustness*. Intuitively, a program is *robust* if the state recovered from NVM after system failure is guaranteed to be reachable under memory consistency model¹. For example, consider the program in Figure 1 and the observed post-crash state $a = 0; b = 1$. This state is not reachable if we ignore the existence of NVM devices and possible crashes, i.e. we only consider its normal executions. Therefore, the program is non-robust.

An advantage of robustness is that user annotation is no longer necessary for verification. In this setting, crash consistency checking can be separated into two steps: (1) prove the program is robust, and (2) verify program correctness under memory consistency model. The latter problem is well-studied and numerous methods for checking weak memory consistency exist in the literature [3, 15, 24, 28, 39, 57, 62], which could be reused on a robust program.

Since robustness acts as a bridge that reduces crash consistency to memory consistency, in this paper, we focus on developing a method for checking robustness of NVM programs. The tool PSAN developed in [22] employs a dynamic algorithm to sample execution traces from the program, and

checks these traces for robustness violation. While able to find robustness violation, it relies on test input generation and sampling and thus is inherently incapable of proving robustness. In contrast with PSAN, we propose a static method aimed at formally proving robustness.

Our method is based upon an observation that, the reachability of a post-crash non-volatile state under certain pre-crash execution can be reduced to validity of the execution with some additional ordering constraints. This enables us to efficiently identify reachable non-volatile states given a pre-crash execution, and check if they are also reachable under memory consistency model. The latter reachability checks are further optimized by shrinking its search space using the pre-crash execution.

Furthermore, we leverage a search method to explore all possible executions and (non-volatile) states and check their reachability. Apart from general-purpose search algorithms, some methods have been designed for efficient exploration of the vast search space in concurrent programs. These include stateless model checking algorithms with dynamic partial order reduction [1, 39–41] and SMT-based methods that encode the program and rely on constraint solving [3, 15, 17, 24, 57, 63, 64]. Our implementation opts for the latter method for exploration, which depends on a symbolic encoding of the input program, and a dedicated theory solver for robustness checking. The solver utilizes the emerging ordering consistency theory [24] for optimized validity checking used in the dual reachability checks and is incorporated into the DPLL(T) framework. Robustness violation is reported whenever we find a non-volatile state that is unreachable under memory consistency model, and we confirm robustness of the program if the exploration is exhaustive.

The proposed method has been implemented in a prototype tool called PMVERIFY. On 26 programs collected from the PMDK [27] `pmemobj` libraries, PMVERIFY is able to report 12 robustness violations and successfully proves robustness of one case. Compared to the dynamic model checking tool PSAN [22], our method finds 6 more violations while PSAN fails to prove robustness. Besides, on a set of 12 manually crafted robust programs, PMVERIFY is able to prove 6 of them, while PSAN is unable to prove any case.

In summary, our main contributions are:

1. We show that the reachability checking problem of a post-crash non-volatile state under a given pre-crash execution can be reduced to the well-studied validity checking problem of a concurrent execution (Section 3).
2. We propose a novel and efficient algorithm for checking robustness of all possible executions within a non-volatile memory program (Section 4). This algorithm is encapsulated as a dedicated theory solver and incorporated into the DPLL(T) framework (Section 5).

¹The definition of robustness in this paper is formulated differently from the original definition in [22], where it is defined using the notion of strong persistency model instead of reachability of recovered state.

3. The approach is implemented in a prototype tool, and we conduct experiments on PMDK benchmarks and a set of manually crafted robust programs. Evaluation results show our method is competitive with dynamic tool PSAN on robustness violation detection and outperforms on robustness verification (Section 6).

2 Preliminaries

2.1 x86 Persistency Model

In this paper, we focus on checking robustness of non-volatile memories on Intel-x86 platforms. The visibility order of memory operations is characterized by the standard x86-TSO model [56], while persist order is prescribed by Px86, a persistency model formalized in [53]. In this setting, a system typically employs a three-layer memory hierarchy per Px86 operational semantics: instructions are issued to thread-local store buffers first, then propagated to a global persistent buffer (write-back caches), from where stores are committed to NVM.

Cache line write-back instructions can be used to constrain persist order. The Intel-x86 architecture provides three such instructions: (1) cache line flush instruction `clflush`, (2) cache line optimized flush instruction `clflushopt`, and (3) cache line write back instruction `clwb`. All three of these instructions write back the content of a single cache line, but differ in how they could be reordered with other instructions. `clflush` instruction has stronger constraints and can only be reordered with loads, while `clflushopt` can be reordered with store, `clflush` and `clflushopt` instructions to other cache lines. `clwb` has the same semantics as `clflushopt` but does not invalidate the cache line, providing better performance. To further constrain the order, the memory barriers `mfence` and `sfence` can be used. `mfence` can not be reordered while `sfence` allows reordering with loads.

Table 1 summarizes the order between relevant instructions based on the standard x86-TSO model and Px86 semantics. Note only visibility order is characterized in the table, which roughly corresponds to the order in which instructions propagate from store buffers to the persistent buffer on Px86.

We can now define the persist order and reachability of non-volatile states:

Definition 1 (Persist Order). Given a fixed visibility order (defined later in Definition 5 as `hb`), the persist order, written `nvo`, is defined as a total order on all stores and flushes that satisfies the following two axioms [53]:

1. The visibility order and persist order coincide between stores to the same variable.
2. If a store is (visibility-)ordered before a flush to the same variable, then it must persist before any stores (visibility-)ordered after the flush.

At any point during program execution, only stores in a *prefix* of the events in `nvo` have persisted. In the case of system failure, these persisted stores in the prefix are safe and recoverable, which induce a *non-volatile state* s where for each location x , $s(x)$ equals the last store to x in the prefix. If a state s is induced by a prefix of some persist order `nvo` of the program that contains all flushes, s is said to be a *reachable* non-volatile state.

2.2 Program and Execution

2.2.1 Programs. We formulate a simple concurrent language for demonstration of our approach. It assumes a set of thread-local variables \mathcal{V}_l (written a, b, c etc.) and shared variables \mathcal{V}_p (written x, y, z etc.). All shared variables reside on non-volatile memory. For simplicity, flush operations work at the granularity of variables instead of cache lines. Let e represent an expression built from local variables, integers and arithmetic operators, an instruction i is then defined by the following grammar:

$$i ::= a = x \mid x = e \mid \text{fence} \mid \text{flush } x$$

Following [22] we ignore the differences between the flush instructions and assume a single `flush` operation. Likewise, we only consider a memory barrier `fence`. `flush` and `fence` exhibit semantics of `clflush` and `mfence` respectively as in Table 1. We note that our implementation supports all variants of flush operations and barriers (Section 6).

A thread consists of a sequence of instructions, and a (concurrent) program is the parallel composition of one or more threads. We use the symbol \parallel for parallel composition, and for each thread, we designate a thread identifier $\tau \in \text{Tid}$. Likewise, each instruction in a thread is associated with an event identifier $i \in \mathbb{N}$, i.e. the index in the sequence of instructions. For a given program, a state s is defined as a valuation of all shared variables, i.e. $s \in \mathcal{V}_p \rightarrow \mathbb{Z}$.

2.2.2 Event Order Graph. Similar to the standard declarative methods in the literature [3, 17, 24, 40, 64], we represent execution of a concurrent program by an *event order graph* (EOG). However, executions on non-volatile memory are slightly different from those on traditional DRAM, in that system failure might happen before a program execution finishes. Therefore, we distinguish between partial and total executions, and adjust the definition of EOGs accordingly. We first define a memory event:

Definition 2 (Event). An event e is a triple (τ, i, l) where $\tau \in \text{Tid}$ is a thread identifier, $i \in \mathbb{N}$ is an event identifier, and l is an event label that can be one of the following:

- $R(x, v)$, marking the event as a read event, where v is the value read from the shared variable x ,
- $W(x, v)$, marking the event as a write event, where v is the value stored to shared variable x ,

Table 1. The preserved program order of Intel-x86 instructions relevant to persistency. \times means two instructions can be reordered, while \checkmark means they are always ordered. CL means the pair of instructions is only ordered when on the same cache line.

		Later in Program Order					
		read	write	mfence	sfence	clflushopt	clflush
Earlier in Program Order	read	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark
	write	\times	\checkmark	\checkmark	\checkmark	CL	\checkmark
	mfence	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark
	sfence	\times	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark
	clflushopt	\times	\times	\checkmark	\checkmark	\times	CL
	clflush	\times	\checkmark	\checkmark	\checkmark	CL	\checkmark

- FL(x), marking the event as a flush event, where all pending stores to x are forced to persist in the order they are issued
- F, marking the event as a memory fence event, which prevents reordering of events before and after it.

Remark 1 (Notation). Given an event label l , the functions type, loc, and val returns the type (R, W, FL, F), location (x), value read or written (v) of l if applicable. Given an event e , the functions Tid, # and lab return the thread identifier, event identifier and event label respectively. The functions on event labels (type, loc etc.) are also lifted to events. For a program P , we write E_P for the set of memory events in P . The method to generate E_P is straightforward [37] by simulation of program execution, and we omit the details here. We abuse the symbols R, W, FL, F for the set of events with the corresponding label in E_P .

Given a relation r , we write r^+ for the transitive closure of r , and r^{-1} for its inverse. Given a relation A , $r|_A$ is r restricted to A . We write $r_1; r_2$ as the relation composition of the two relations r_1 and r_2 . For a set of events E , E_x is the subset of E restricted to events on variable x , i.e. $E_x = \{e \in E \mid \text{loc}(e) = x\}$. For any ordering relation r over E , we also write $e_1 \prec_r e_2$ for $(e_1, e_2) \in r$.

An EOG is then defined with respect to ordering relations over memory events:

Definition 3 (Event Order Graph). An event order graph $G = (E, E_0, \text{po}, \text{rf})$ consists of a set of events E and a subset of initialization events $E_0 \subseteq E$ containing a single write event to each shared variable. po, rf are relations over E where,

- $\text{po} \subseteq E \times E$ is the program order, a total order of events in each thread. Moreover, initialization events in E_0 are ordered before the other events in E . po can be derived syntactically from the program, i.e. $\text{po} = \{(e_1, e_2) \mid \text{Tid}(e_1) = \text{Tid}(e_2) \wedge \#e_1 < \#e_2\} \cup (E_0 \times E \setminus E_0)$.

- $\text{rf} \subseteq (E \cap W) \times (E \cap R)$ is the read-from relation between write and read events on the same variable. Intuitively, $(e_w, e_r) \in \text{rf}$ if e_r reads the value written by e_w . It is obvious that each read event should read from at most one write event, i.e. for any events $e_w^1, e_w^2 \in (E \cap W)$ and $e_r \in (E \cap R)$, $(e_w^1, e_r) \in \text{rf} \wedge (e_w^2, e_r) \in \text{rf} \rightarrow e_w^1 = e_w^2$.

For convenience, we use $G.x$ to refer to the element of G , where x can be E, E_0, po or rf . When the context is clear, we write E, po etc. directly.

An event order graph G represents a (total) *execution* of a concurrent program P , if $G.E$ equals E_P , and $G.\text{rf}$ assigns a write event to each read event. In this case, the execution finishes without being interrupted by a potential crash. However, not all executions are *valid*, or *consistent*, per the underlying memory consistency model that prescribes allowed visibility order of memory operations.

Each memory consistency model M essentially defines a predicate over executions, denoted $\text{cons}_M(\cdot)$, for the set of valid executions under M . The memory consistency model adopted in this paper is an extension of standard x86-TSO [56] for Intel-x86 platform which identifies a global *happens-before* order hb over all events. Unlike Sequential Consistency (SC), only *preserved program order* (ppo), a subset po that can not be reordered according to architecture specification, is included in hb .

Table 1 summarizes the ordering constraints of relevant instructions. It shows that writes and flushes might be reordered with later reads. The preserved program order ppo is then formally defined as:

$$\text{ppo} \triangleq \{(e_1, e_2) \in \text{po} \mid e_1 \in W \cup \text{FL} \rightarrow e_2 \notin R\}$$

To define happens-before order hb , we introduce coherence order (co):

Definition 4 (Coherence Order). Given an execution G , a coherence order $\text{co} \subseteq (G.E \cap W) \times (G.E \cap W)$ is the disjoint union of relations co_x for each shared variable $x \in \mathcal{V}_p$, where co_x is a strict total order on write events to x .

Orders within the same thread (called *internal* orders) are distinguished from those across different threads (called *external* orders), and we denote them with suffix i and e respectively. For instance, rf_i is the relation $\text{rf} \cap (\text{po} \cup \text{po}^{-1})$. Additionally, given a coherence order co , the *from-read* relation fr between a write event and a read event is derived as $\text{fr} \triangleq \text{rf}^{-1}; \text{co}$. Intuitively, if we have $(e_w, e_r) \in \text{rf}$ and $(e_w, e'_w) \in \text{co}$, e_r must happen before e'_w since e_r would read from e'_w otherwise.

Definition 5 (x86-TSO). An execution G is valid under x86-TSO, written $\text{cons}_{\text{TSO}}(G)$, if there is a coherence order co such that

1. $\text{hb} = (\text{ppo} \cup \text{rf}_e \cup \text{co} \cup \text{fr})^+$ is irreflexive,
2. $\text{fr}; \text{po}$ is irreflexive (per-location coherence)

Each valid execution G induces a (volatile) state s , where for each $x \in \mathcal{V}_p$, $s(x)$ equals the value written by the last write event to x in hb . If s is induced by some execution of the program P , it is said to be a *reachable* state of P under x86-TSO.

2.3 Crash Consistency and Robustness

Crash consistency is an essential property of programs running on non-volatile memories. Given that the system may crash at any time, it specifies that program execution can be correctly resumed from the recovered non-volatile state, as defined in Section 2.1. This essentially requires that the post-crash execution starting from the state does not terminate unexpectedly (e.g. segmentation faults or assertion violation) or cause data corruption. However, most tools in the literature rely on user annotation for crash consistency bug detection, and the few automatic tools only detect observable bugs. To tackle the problem, *robustness* is proposed [22] as a sufficient condition for crash consistency of lock-free programs:

Definition 6 (Robustness). A program P is *robust* iff all reachable non-volatile states of P , as defined in Section 2.1, are reachable under x86-TSO, as defined in Section 2.2.

In other words, the set of reachable non-volatile states is subsumed by the set of reachable states under x86-TSO. Crash consistency requires safe execution from any post-crash state. In this case, to prove crash consistency of a robust program, we only need to apply existing methods for ensuring correctness of a concurrent program under some weak memory consistency model, which is x86-TSO in our case. For a robust program, the problem in question is essentially reduced to the classical safety verification problem of concurrent programs. Furthermore, since consistency checking and proving robustness are decoupled from each other,

this method is fully automated, and user annotation is not needed.

3 Checking Reachability of a Non-volatile State

In this paper, we focus on proving robustness. Since robustness is a universal property over non-volatile states, it is necessary to explore all non-volatile states and check if all states are reachable per definitions in Section 2.1 and Section 2.2. In this section, we focus on how to observe a potential non-volatile state from the program and check reachability of the post-crash state given a fixed pre-crash execution.

3.1 Recovery Observer

To enumerate non-volatile states efficiently, we leverage recovery observer to instrument the program. Recovery observer is originally proposed in [50] as a hypothetical notion that atomically observes the entire content of the NVM. It is then adopted for verification of software performing file I/O [38]. Unsurprisingly, the semantics of I/O operations to storage devices are analogous to memory operations on NVM. In fact, it has been utilized later for persistent invariant checking [48].

Intuitively, recovery observer is a virtual thread that reads each shared variable. As the recovery observer acts as an additional thread, the reads in it interleave with other memory operations. By going through all possible interleaving of the threads, each read also iterates through all possible writes. It facilitates enumeration of states since we could rely on the rf relation of these reads for a proper post-crash state. Figure 2 shows an example program with recovery observer. The third thread is the recovery observer with a read to each shared variable in the program.

```

x = 1;      y = 2;
flush x;   flush y;
a = y;     || b = x;     || r1 = x;
x = a;     || y = b;     || r2 = y;
flush x;   flush y;
    
```

Figure 2. An example of recovery observer. $r1 = x$; and $r2 = y$ are not ordered.

To adopt recovery observer to our setting, we instrument the program with the virtual thread and introduce a dedicated reads-from order for reads in this thread. Formally, given a program P , its instrumented version is $P' = P \parallel P_r$, where P_r represents the recovery observer that contains an instruction $a_x = x$ for each $x \in \mathcal{V}_p$. Let REC be the set of events in the recovery observer P_r , i.e. $\text{REC} = \{e \in E_{P'} \mid \text{Tid}(e) = P_r\}$. To ensure they could observe all states of the program, these read events are not ordered with any other events, and in particular, they are not ordered with each

other. We thus have the following definition of instrumented execution:

Definition 7 (Instrumented Execution). An instrumented execution G^i of P is an execution of the instrumented program $P \parallel P_r$. In particular, it satisfies $G^i.\text{po} \cap (\text{REC} \times G^i.E) = \emptyset$ and $G^i.\text{po} \cap (G^i.E \times \text{REC}) = \emptyset$.

We define the *recovery read-from* relation rrf of G^i as the projection of the read-from relation to REC , i.e. $\text{rrf} \triangleq \text{rf}|_{(\text{W} \times \text{REC})}$. The relation induces an observed non-volatile state s_o such that for each shared variable $x \in \mathcal{V}_p$, $s_o(x)$ equals the store read by REC_x .

Once a non-volatile state s_o is observed by the recovery observer, the next step is to check reachability of s_o under G^i , which is elaborated in Section 3.2. If s_o is indeed a reachable non-volatile state, we then need to check if s_o is reachable under x86-TSO. While this is a well-studied problem and not the topic of this paper, we note that recovery observer can be tweaked for this purpose as well. Briefly speaking, we retain the instrumented execution G^i , but group the read events in REC together as an *atomic block*, i.e. we only allow it to read the whole memory simultaneously. In this way, the recovery observer now signifies an equivalent volatile state instead. The details are given in Section 4.

3.2 Reduction to Validity Checking

In this section, we check reachability of the observed non-volatile state s_o under a given execution. Reachability under Px86, which is a more general problem than ours, has been previously proved to be decidable [2], but no algorithmic method is given. To solve the problem, the key is to reduce it to an equivalent validity checking problem of a pre-crash execution that is augmented with additional ordering constraints.

To accomplish this, we leverage the *derived TSO propagation order* (dtpo) from [37] as a bridge between memory consistency and persistency. Given an instrumented execution G^i of the program, we have:

$$\text{dtpo} \triangleq \bigcup_{x \in \mathcal{V}_p} \text{FL}_x \times \{w \in \mathbb{W}_x \mid \exists w' \in \text{dom}(\text{rrf}), (w', w) \in \text{co}\}$$

dtpo orders any flush on the shared variable x before any store w to x that are co -ordered after the store w' to x read by rrf . Note that although it is derived from the persistency-related relation rrf , it characterizes visibility order between flushes and certain stores. The correctness of this derived order can be seen by the following argument: if the flush event is instead ordered after w , then all pending writes to x , including w , should be committed to the NVM. Since w happens after w' , it would overwrite w' , which contradicts the fact that w' is the last write to x that has persisted. Fig. 3 demonstrates a possible instrumented execution of the program in Fig. 2 where the state $x = 0; y = 1$ is observed. Since the observer reads the initial value of x , which is ordered

before the two stores $x = 1$ and $x = a$, two dtpo ordering constraints are induced.

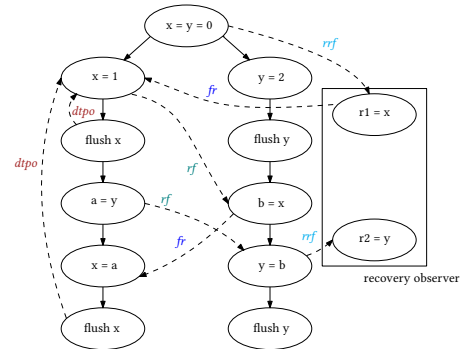


Figure 3. Visualization of an event order graph representing a possible instrumented execution of the program in Fig. 2, where the state $x = 0; y = 1$ is observed by the recovery observer. The execution is not valid under DPTSO per Definition 8 due to cycles induced by the dtpo order.

Our goal is to check if the observed non-volatile state s_o , as induced by rrf , is reachable under G^i . Note the execution G^i only characterizes visibility order, and we do not know the exact persist order nvo . A naive approach to check reachability is to enumerate all possible persist order, and check whether the last persisted store to each shared variable conforms with s_o . Nevertheless, it is redundant to consider all possible orders. Even if we prune the search by leveraging the two axioms for nvo (Definition 1), the exhaustive method is still inefficient.

What matters in reachability checking is the last persisted stores and the additional ordering constraints they generate (dtpo). Therefore, the crux is to check if G^i is a valid execution under a memory model augmented with dtpo , called DPTSO in [37]:

Definition 8 (DPTSO). An instrumented execution G^i is valid under DPTSO, written $\text{cons}_{\text{DPTSO}}(G^i)$, if there is a coherence order co such that

1. $\text{hb} = (\text{ppo} \cup \text{rf}_e \cup \text{co} \cup \text{fr} \cup \text{dtpo})^+$ is irreflexive,
2. $\text{fr}; \text{po}$ is irreflexive (per-location coherence)

The predicate $\text{cons}_{\text{DPTSO}}(\cdot)$ can be checked in an analogous way to checking $\text{cons}_{\text{TSO}}(\cdot)$. It is basically a cycle detection algorithm on a directed graph where the orders rf , co , dtpo etc. are regarded as edges. Since DPTSO only adds dtpo as a component in the happens-before relation, existing methods for checking validity of a concurrent execution could be easily adopted. For example, the execution depicted in Fig. 3 is not valid under DPTSO due to cycles introduced by the extra dtpo edges.

Theorem 1. Given an instrumented execution G^i valid under x86-TSO, the non-volatile state s_o induced by rrf is reachable under G^i iff $\text{cons}_{\text{DPTSO}}(G^i)$ holds.

661 *Proof.* The proof is elaborated in Appendix A. \square

663 4 Robustness Checking Algorithm

664 In this section, we discuss the general search-based frame-
 665 work for robustness checking. More specifically, we extend
 666 the algorithm to partial executions and introduce the overall
 667 exploration algorithm based on it.

668 Previously, reachability of states is defined with an as-
 669 sumption that all memory events have been propagated from
 670 the store buffer and made visible to all threads before the
 671 program terminates. This is in line with the Px86 and DPTSO
 672 models. However, to verify programs running on NVM, it is
 673 necessary to reason about crashes, and in particular when a
 674 system failure would occur. Therefore, the verification algo-
 675 rithm must take partial executions into account.

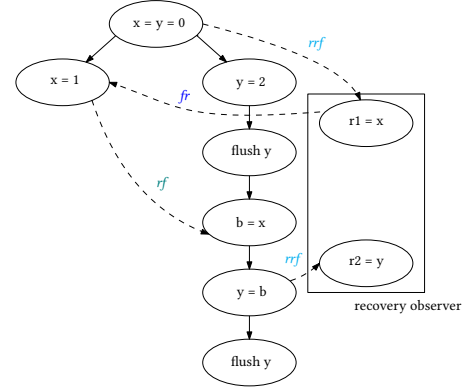
677 **Partial Executions.** In Section 2.2, we defined an exe-
 678 cution of a program P as an EOG that contains all memory
 679 events of a program, E_p , and assigns a value to each read
 680 event in rf . Likewise, a *partial execution* is an EOG that
 681 assigns a value to each read in rf , but it contains only a
 682 subset of all memory events E_p . However, the events must
 683 be *prefix-closed*. In other words, any event in the $porf$ -
 684 prefix of an event in the partial execution should also be con-
 685 tained in the event set. This requirement corresponds to the fact
 686 that when a system fails, only a prefix of the total execution
 687 has been propagated. A partial instrumented execution is
 688 defined similarly.

689 Take the total execution of the program in Fig. 2 as an
 690 example, as depicted in Fig. 3. It is shown previously that
 691 this execution is not valid under DPTSO, thus the state s_o is
 692 not reachable. Now consider the realistic scenario where the
 693 instructions $\text{flush } x$; $a = y$; $x = a$; $\text{flush } x$; are not
 694 propagated before the crash. In this case, we obtain a partial
 695 execution as visualized below:

696 It can be easily checked that this partial execution is valid
 697 under DPTSO, and thus the state $x = 0$; $y = 1$; is a reachable
 698 non-volatile state.

700 **Exploration Algorithm.** The reachability checking algo-
 701 rithm in Section 3 can be naturally lifted to partial executions
 702 since $dtpo$ is defined analogously on partial executions. An
 703 exploration algorithm can then search through all partial exe-
 704 cutions and states, and utilize the aforementioned reduction
 705 to check reachability of a state with a given partial execution.
 706 Whenever a reachable non-volatile state is found, robustness
 707 is checked locally first by checking reachability of this state
 708 under x86-TSO. Robustness of the whole program is proved
 709 if no violation is found when the exploration ends.

710 The presence of recovery observer has embedded a non-
 711 volatile state in an instrumented execution, thus the explo-
 712 ration method only needs to search through all instrumented
 713 partial executions. Besides, it allows that some events in REC
 714 are not contained in the partial execution, i.e. it induces a



716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770

Figure 4. A partial instrumented execution of the program in Fig. 2, where the state $x = 0$; $y = 1$; is observed by the recovery observer. No $dtpo$ orders are derived since flush events to x are not in this partial execution. This execution is valid under DPTSO.

736 partial state. This boosts performance and allows our robust-
 737 ness checking algorithm to have the flexibility of leveraging
 738 different search methods, from brute force searching to more
 739 advanced stateless model checking with dynamic partial
 740 order reduction, or simply relying on program encoding
 741 and constraint solving. Our robustness checking algorithm
 742 could be incorporated into any exploration method capable
 743 of search tree pruning. Therefore, we abstract away the de-
 744 tails and assume the exploration method provides the next
 745 and hasNext interface for exploration, and block interface
 746 to block a subset of partial instrumented executions. Sec-
 747 tion 5 will elaborate on this topic.

749 **Remark 2 (Notation).** We say the partial instrumented exe-
 750 cution G' is an *expansion* of partial instrumented execution
 751 G , written $G \prec G'$, if $G.E \subset G'.E$ and $G.rf \subset G'.rf$.

752 G' is an *alternation* of G , written $G' \simeq G$, if $G.E = G'.E$
 753 and $G.r rf = G'.r rf$ (other orders in $G.rf$ and $G'.rf$ might
 754 differ).

756 The overall algorithm framework is shown in Algorithm 1.
 757 The input program is first instrumented with recovery ob-
 758 server, then the exploration method takes over the search.
 759 Each time a partial execution G is yielded, a partial state is
 760 also generated. We first check if it is reachable under the
 761 current execution, i.e. if $\text{cons}_{\text{DPTSO}}(G)$ holds. If not, we make
 762 sure not to extend G and further explore its expansion, since
 763 an invalid execution with additional ordering constraints is
 764 still invalid under DPTSO (Line 10). This optimization can be
 765 implemented in most search methods. In depth-first search-
 766 ing, for instance, the search immediately backtracks to avoid
 767 further redundant exploration.

768 If G exhibits a reachable non-volatile state s_o , the next
 769 step is to check if s_o is reachable under the x86-TSO model,

Algorithm 1: Robustness Checking Algorithm
 Framework for Non-volatile Memories

```

771 input : A program  $P$  running on non-volatile
772         memory.
773 output: If  $P$  is robust.
774
775 1  $P' \leftarrow P \parallel P_r$ 
776 2 while hasNext( $P'$ ) do
777   3  $G \leftarrow \text{next}(P')$ 
778   4 if consDPTSO( $G$ ) holds then
779     5 foreach  $G' \preceq G$  do
780       6 if consTSO(atomic( $G'$ )) holds then
781         7  $\text{goto } 2$ ;
782       8 return false;
783     9 else
784       10  $\text{block}(\{G' \mid G \prec G'\})$ ;
785   11 return true
  
```

which typically involves another search over all total executions (Line 5). There are abundant algorithms for this task in the literature. In our algorithm, we take advantage of the recovery observer by regarding the read events in them as an atomic block and keeping the *rrf* orders. The formal definition of *atomic* is given in Appendix B due to space limit, but intuitively this allows the recovery observer to signify an equivalent volatile state. As an example, Fig. 5 shows how to check validity of an execution under x86-TSO with the help of recovery observer.

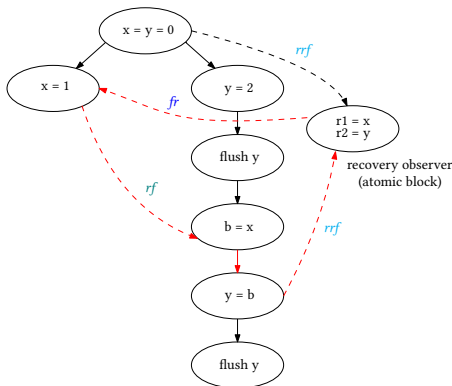


Figure 5. Visualization of validity checking of the execution in Fig. 4, where the recovery observer is regarded as an atomic block. It is invalid under x86-TSO due to the cycle in red.

While searching through total executions, we only alter other *rf* orders, keeping *G.E* the same. Essentially, enumerating alternations of *G* is enough for checking reachability under x86-TSO. The soundness of this method is shown in the following lemma:

Lemma 1. If for every alternation G' of G , $\text{cons}_{\text{TSO}}(G')$ does not hold, then the observed state s_o must be unreachable under x86-TSO.

Proof. The proof is given in Appendix A. \square

If the state s_o is proved unreachable under x86-TSO in this way, we report violation of robustness (Line 8). Note that s_o might be a partial state, but it is obvious that any total state that conforms with s_o is still invalid under x86-TSO. Therefore, we could always add *rrf* orders that read from the last store per *co*. This way validity under DPTSO is not affected. On the other hand, if s_o is reachable, once its reachability is proven we proceed with the exploration of executions (Line 7). If all executions have been explored and no violation is found, P has been proved robust (Line 11).

5 Integration with DPLL(T)

Algorithm 1 is parameterized by an exploration method. In this section, we instantiate our algorithm to leverage program encoding and constraint solving for this task.

5.1 Encoding

A complete encoding of a concurrent program should cover both functional program behaviours and possible interleaving of the threads, i.e. the orders between memory events. Following the standard encoding [3], the input program is firstly transformed into SSA form and the program event set E_p . The instructions in each thread, including the recovery observer in our case, are then naturally translated to atoms in first-order logic². As a simple example, the program in Fig. 2 is encoded as:

$$\begin{aligned}
 \rho_{ssa} &= x_0 = 0 \wedge y_0 = 0 && \text{(initial value)} \\
 &\wedge x_1 = 1 \wedge a = y_1 \wedge x_2 = a && \text{(first thread)} \\
 &\wedge y_2 = 2 \wedge b = x_3 \wedge y_3 = b && \text{(second thread)} \\
 &\wedge r_1 = x_4 \wedge r_2 = y_4 && \text{(recovery observer)}
 \end{aligned}$$

Note that it does not encode ordering relations. The read y_1 , for instance, could potentially read from y_0, y_2 and y_3 . To add ordering constraints to the encoding, we first model a partial execution by a predicate enabled (implemented as a Boolean variable) defined for every event, where $\text{enabled}(e)$ signifies e is included in the partial execution. Furthermore, each order relations used in x86-TSO or DPTSO are represented by Boolean variables explicitly. For instance, since the coherence order *co* is total, it is encoded by adding a Boolean variable $ws_{i,j}^x$ for each pair of stores x_i and x_j to the shared variable x . We have $(x_i, x_j) \in \text{co}$ iff $ws_{i,j}^x$ is assigned true. Additional *axioms* are included in the encoding that constrain its assignment:

²flush and fence operations are not included in the functional encoding of a program since they are irrelevant. However, they are still numbered and contribute to ordering constraints.

$$\begin{aligned}
 \rho_{i,j}^{\text{co}} &= \text{ws}_{i,j}^x \rightarrow \text{enabled}(x_i) \wedge \text{enabled}(x_j) && (\text{ws-cond}) \\
 &\wedge \text{ws}_{i,j}^x \rightarrow x_i \prec_{\text{co}} x_j && (\text{ws-order}) \\
 &\wedge (\text{enabled}(x_i) \wedge \text{enabled}(x_j)) \rightarrow \text{ws}_{i,j}^x \vee \text{ws}_{j,i}^x && (\text{ws-some})
 \end{aligned}$$

Similarly, we introduce for each variable x a Boolean variables $\text{rf}_{j,i}^x$ for any read x_i and store x_j and axioms for rf as follows:

$$\begin{aligned}
 \rho_{j,i}^{\text{rf}} &= \text{rf}_{j,i}^x \rightarrow \text{enabled}(x_i) \wedge \text{enabled}(x_j) \wedge x_j = x_i && (\text{rf-val}) \\
 &\wedge \text{rf}_{j,i}^x \rightarrow x_j \prec_{\text{rf}} x_i && (\text{rf-ord}) \\
 &\wedge \text{enabled}(x_i) \rightarrow \bigvee_{x_j \in W_x} \text{rf}_{j,i}^x && (\text{rf-some})
 \end{aligned}$$

Since fr can be derived from rf and co as discussed in Section 2, for each variable x , we introduce the following axiom for any two stores x_j, x_k and read x_i :

$$\rho_{j,i,k}^{\text{fr}} = \text{rf}_{j,i}^x \wedge \text{ws}_{j,k}^x \rightarrow r_{x_i} \prec_{\text{fr}} w_{x_k}$$

For dtpo , we need the flush event to be enabled. Thus, for each variable x , we introduce the following axiom for any two stores x_j, x_k , read event x_i in the recovery observer, and flush FL_q^x :

$$\rho_{q,j,i,k}^{\text{dtpo}} = \text{enabled}(\text{FL}_q^x) \wedge \text{rf}_{j,i}^x \wedge \text{ws}_{j,k}^x \rightarrow \text{FL}_q^x \prec_{\text{dtpo}} w_{x_k}$$

To ensure the prefix-closed property of a partial execution, encoding for each ordering relation requires the pair of events to be both enabled (e.g. the rule ws-cond and rf-val above), and an extra axiom is added to the encoding: for any two events e_1 and e_2 that are ordered by ppo , we have $\text{enabled}(e_2) \rightarrow \text{enabled}(e_1)$. The encoding of the program Ψ is then the conjunction of ρ_{ssa} and all axioms ($\rho_{i,j}^{\text{co}}, \rho_{j,i}^{\text{rf}}$ etc.) related to ordering constraints.

5.2 DPLL(T) and Exploration

The encoded formula of the program Ψ is solved by an SMT solver. While it searches for a model of the formula, variables in it are assigned values. In particular, the assignment of Boolean variables representing various order relations corresponds to a partial execution and state. Modern SMT solvers typically utilize the DPLL(T) framework. In the framework, formulas are in a combination of certain first-order background theories. Each background theory \mathcal{T} has a *theory solver* which decides \mathcal{T} -satisfiability of a conjunction of literals in \mathcal{T} . An overview of this framework is shown in Fig. 6.

In this framework, each atom in the given formula Ψ is first replaced by a Boolean variable, and the satisfiability of the resulting propositional formula $B(\Psi)$ is checked by

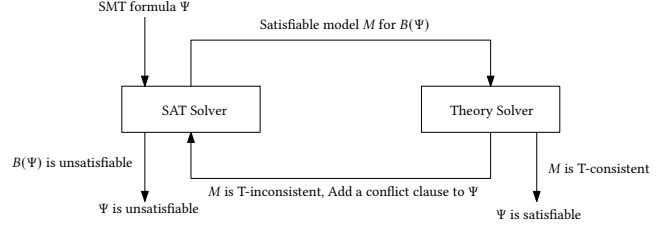


Figure 6. Overview of the DPLL(T) framework.

an SAT solver. If $B(\Psi)$ is unsatisfiable, so is Ψ . Otherwise, since $B(\Psi)$ is an over-approximation of Ψ , theory solvers are called to check if the model M returned by the SAT solver is compatible with the underlying background theories. The theory solver also returns a *conflict clause* to prevent the SAT solver from exploring the same assignment.

Following [24], while DPLL(T) controls the exploration, in our implementation each ordering constraints in the formula, such as $w_{x_i} \prec_{\text{ws}} w_{x_j}$ from ws-order above, are passed to a dedicated theory solver for robustness checking. Based on the solver for ordering consistency theory, the backend employs an incremental cycle detection algorithm for efficient checking validity under DPTSO. If the current partial state is reachable, we use the solver in [15] to check its reachability under x86-TSO. Otherwise, a conflict clause is generated and returned to the DPLL(T) framework which blocks further exploration of this partial instrumented execution.

6 Implementation and Evaluation

We have implemented our method in a prototype tool called PMVERIFY, expanding on Deagle [25], a concurrent program verification tool that supports weak memory consistency on top of the bounded model checker CBMC. We extend Deagle's frontend to recognize NVM programs using a selected set of APIs in `pmemobj` library from PMDK [27]. A dedicated backend for robustness checking is implemented to complement the default solver of Deagle for weak memory consistency checking.

To evaluate the efficacy of PMVERIFY, we collect the example programs that accompany the `pmemobj` library in PMDK as the benchmark. It contains 26 small to medium-sized programs (548 LOC on average) that implement simple algorithms and basic data structures on non-volatile memory, such as binary search and persistent lists.

As a comparison, we also run PSAN [22], the only robustness violation detection tool in the literature, on the same benchmark. PSAN is implemented based on the dynamic model checking tool JAARU [23] which observes the outcome of memory operations at runtime and checks for persistency bug on the observed trace. PSAN offers a model checking mode that exhaustively enumerates program traces, as well as a random mode that relies on sampled traces.

Table 2. Evaluation of PMVERIFY and PSAN on the PMDK pmemobj benchmarks (26 programs in total). The rows YES/No/UNKNOWN contain the total number of cases each tool proves to be robust, non-robust, or fails to give an answer. Unique No. is the number of cases only solved by this tool.

	PMVERIFY	PSAN
YES	1	0
No	12	6
UNKNOWN	13	20
Unique No.	7	0
Average Time	2768.42s	16.7s
Standard Deviation	1045.26s	9.98s

The experiments are conducted on an Intel® Core™ i5-10400 @ 2.90GHz CPU with 16GB memory. Timeout is set to two hours. For PSAN running in random mode, the maximum amount of sampled traces is set to 100,000.

6.1 Experimental Results on PMDK benchmark

Table 2 shows the overall results of the experiment. Out of 26 programs, PMVERIFY is able to solve 13. It successfully proves robustness of one of the programs (`manpage.c`), and detects robustness violation of 12 test programs. The results are manually checked to ensure soundness. PMVERIFY fails on the remaining programs mainly due to usage of some PMDK primitives that are not modeled by the frontend, such as `pmemobj_tx_add_range_direct`, and timeouts on one program it supports.

Compared to PSAN, our tool is able to find robustness violation of six more programs. Besides, since PSAN employs an incomplete dynamic method, it is only able to refute robustness instead of verifying the property. On the contrary, PMVERIFY is able to prove robustness of a program after all executions have been explored.

In terms of performance, Fig. 7 shows the accumulated solving time of PMVERIFY. It takes around 45 minutes to complete verification on average. PSAN, on the other hand, takes no more than a minute for the six cases it solves, which is, on average, 33 times shorter than PMVERIFY. We note that PSAN is similar to a testing tool and could very efficiently find potential bugs in the program due to its dynamic nature. However, it cannot verify robustness. PMVERIFY has adopted some optimization methods to improve the performance of exploration, but still faces the common state explosion problem. The exhaustive exploration is necessary to ensure completeness, at the cost of performance. Therefore, our method and PSAN can complement each other in robustness verification and bug detection.

Table 3. Evaluation results of PMVERIFY on 12 instrumented programs. PSAN is unable to prove robustness of any program.

Program	LOC	PMVerify Time (s)
<code>btree</code>	493	5468.95
<code>buffons_needle_problem</code>	432	timeout
<code>lists</code>	551	timeout
<code>pi</code>	570	3565.19
<code>examine_arttree</code>	6379	7021.70
<code>arttree</code>	1793	timeout
<code>fifo</code>	207	4677.78
<code>data_store</code>	5512	timeout
<code>mapcli</code>	742	timeout
<code>main</code>	195	1765.79
<code>reader</code>	95	timeout
<code>writer</code>	67	5905.49
Average	1420	4734.15

6.2 Evaluation on robust programs

The robust case PMVERIFY solves, `manpage.c`, is a simple program that opens a persistency memory pool and does nothing. In this section, to further demonstrate the ability of PMVERIFY to prove robustness, we manage to instrument each of the 12 non-robust programs to manually produce a set of robust programs.

More specifically, we insert a cache line flush instruction after each memory operation. In this way, the instrumented program is guaranteed to be robust. We then run both PMVERIFY and PSAN on this new benchmark.

Table 3 shows the results of PMVERIFY and PSAN running on the set of instrumented programs. PMVERIFY is able to prove robustness of six programs with an average running time of 4734.15 seconds, including the medium-sized programs `examine_arttree` and `data_store`. This shows the ability to scale to larger robust programs. On the other hand, PSAN is unable to prove robustness of any program.

Due to the increase in program size, the performance of PMVERIFY on this benchmark degrades by around half. We note that although adding a flush operation after every memory operation introduces considerable redundancy and increases the overall exploration space, the set of reachable non-volatile states is smaller because of stronger constraints. Since PMVERIFY checks for reachability under DPTSO model first, we can avoid later steps of checking x86-TSO consistency for some states. Therefore, the running time of PMVERIFY does not grow exponentially. In fact, all cases could be finished within three hours if we do not consider time limits,

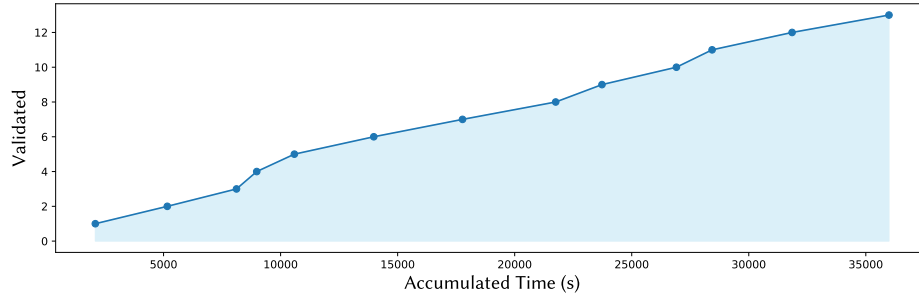


Figure 7. Accumulated solving time of PMVERIFY on the 13 solved case. 35989.46 seconds are spent in total.

with an average running time of 6712 seconds. This shows the efficacy of our tool PMVERIFY for verification of robustness.

7 Related Works

Persistent Models. The early studies on NVMs rely on certain persistency models, an extension of memory consistency models, to prescribe constraints on the persistence order. In [50], Pelly et al. classified these models into three categories: *strict persistency*, *epoch persistency* and *strand persistency*. The original definition of robustness in [22] is based on strict persistency, the strongest model where any recovered state is guaranteed to be an observable volatile state [14, 31].

Epoch persistency under sequential consistency is described in [50], while [32] proposes a persist barrier implementation that works on x86-TSO [56]. The first formal definition of epoch persistency is given by [30] under release consistency, and [52] formally describes operational and declarative semantics of epoch persistency under x86-TSO. StrandWeaver [21] implements strand persistency in hardware to minimally constrain persists to NVMs.

Recently, a line of work focused on formally defining the persistency model of hardware architecture. [54] develops PARMv8 model for ARMv8, followed by Px86 [53] for Intel-x86. Later, the PEx86 model [51] is proposed with formalized semantics of non-temporal stores. Alternative models such as DPTSO [37] and view-based models for Intel-x86 and Armv8 [10] are proposed to further develop these formalisms.

Memory consistency checking. The essential idea of multi-threaded program verification is to explore the possible executions caused by thread interleaving. [3] gives a framework for using partial order relations to model possible executions and encode program behaviors into a formula. Several works expand on this idea and rely on bounded model checking, such as lazy sequentialization [28] and a line of work that employs the scheduling constraints-based abstraction refinement method (SCAR) [62–64]. [17] proposes to

solve the difference logic-based ordering constraints more efficiently with DPLL(T) framework. [24] proposes an ordering consistency theory and integrates a dedicated theory solver into the DPLL(T) routine, which is extended to weak memory consistency in [15]. On the other hand, stateless model checking (SMC) methods enumerate all interleavings with respect to an equivalence class, i.e. a Mazurkiewicz trace. Several algorithms have been proposed to further weaken the ordering requirement and efficiently explore the search space [1, 39, 40].

Persistency Bug Detection. Several tools have been developed to assist persistent programming, including testing applications such as XFDetector [45], PMTest [46] and PMDebugger [12]. Yat [43] is a model checker that exhaustively explores all persistence orders and crash points. The model checker Jaaru [23] reduces search space by focusing on the last writes to each location. To our knowledge, the only automated verification tool for persistent memories is introduced in [48], which utilizes an SMT-based method to formally verify persistent invariants.

8 Conclusion

In this paper, we propose a novel approach to check robustness, a sufficient condition for crash consistency of lock-free programs running on non-volatile memories. Our algorithm employs a search method to explore all partial executions and non-volatile states, and check reachability of the state under the pre-crash execution. This is achieved by reducing the reachability checking problem to checking validity of an instrumented execution under an alternative model DPTSO. Our implementation is based on encoding the program into a SMT formula and constraint solving. It succeeds in robustness verification of a set of example programs in PMDK while the dynamic robustness violation tool PSAN fails. In terms of robustness bug detection, our tool also outperforms PSAN.

References

- [1] Parosh Aziz Abdulla, Stavros Aronis, Mohamed Faouzi Atig, Bengt Jonsson, Carl Leonardsson, and Konstantinos Sagonas. Stateless model

- checking for TSO and PSO. *Acta Informatica*, 54(8):789–818, December 2017.
- [2] Parosh Aziz Abdulla, Mohamed Faouzi Atig, Ahmed Bouajjani, K. Narayan Kumar, and Prakash Saivasan. Deciding reachability under persistent x86-TSO. *Proceedings of the ACM on Programming Languages*, 5(POPL):1–32, January 2021.
- [3] Jade Alglave, Daniel Kroening, and Michael Tautschnig. Partial Orders for Efficient Bounded Model Checking of Concurrent Software. In Natasha Sharygina and Helmut Veith, editors, *Computer Aided Verification*, Lecture Notes in Computer Science, pages 141–157, Berlin, Heidelberg, 2013. Springer.
- [4] Joy Arulraj and Andrew Pavlo. How to Build a Non-Volatile Memory Database Management System. In *Proceedings of the 2017 ACM International Conference on Management of Data*, SIGMOD ’17, pages 1753–1758, New York, NY, USA, May 2017. Association for Computing Machinery.
- [5] Hans-J. Boehm and Dhruva R. Chakrabarti. Persistence programming models for non-volatile memory. *ACM SIGPLAN Notices*, 51(11):55–67, June 2016.
- [6] Daniel Castro, Paolo Romano, and João Barreto. Hardware Transactional Memory meets memory persistency. *Journal of Parallel and Distributed Computing*, 130:63–79, August 2019.
- [7] Dhruva R. Chakrabarti, Hans-J. Boehm, and Kumud Bhandari. Atlas: Leveraging locks for non-volatile memory consistency. *ACM SIGPLAN Notices*, 49(10):433–452, October 2014.
- [8] Youmin Chen, Jiwu Shu, Jiaxin Ou, and Youyou Lu. HiNFS: A Persistent Memory File System with Both Buffering and Direct-Access. *ACM Transactions on Storage*, 14(1):4:1–4:30, April 2018.
- [9] Zhangyu Chen, Yu Hua, Yongle Zhang, and Luochangqi Ding. Efficiently detecting concurrency bugs in persistent memory programs. In *Proceedings of the 27th ACM International Conference on Architectural Support for Programming Languages and Operating Systems*, ASPLOS ’22, pages 873–887, New York, NY, USA, February 2022. Association for Computing Machinery.
- [10] Kyeongmin Cho, Sung-Hwan Lee, Azalea Raad, and Jeehoon Kang. Revamping hardware persistency models: View-based and axiomatic persistency models for Intel-x86 and Armv8. In *Proceedings of the 42nd ACM SIGPLAN International Conference on Programming Language Design and Implementation*, PLDI 2021, pages 16–31, New York, NY, USA, June 2021. Association for Computing Machinery.
- [11] Joel Coburn, Adrian M. Caulfield, Ameen Akel, Laura M. Grupp, Rajesh K. Gupta, Ranjit Jhala, and Steven Swanson. NV-Heaps: Making persistent objects fast and safe with next-generation, non-volatile memories. *ACM SIGARCH Computer Architecture News*, 39(1):105–118, March 2011.
- [12] Bang Di, Jiawen Liu, Hao Chen, and Dong Li. Fast, flexible, and comprehensive bug detection for persistent memory programs. In *Proceedings of the 26th ACM International Conference on Architectural Support for Programming Languages and Operating Systems*, ASPLOS ’21, pages 503–516, New York, NY, USA, April 2021. Association for Computing Machinery.
- [13] Subramanya R. Dulloor, Sanjay Kumar, Anil Keshavamurthy, Philip Lantz, Dheeraj Reddy, Rajesh Sankaran, and Jeff Jackson. System software for persistent memory. In *Proceedings of the Ninth European Conference on Computer Systems*, EuroSys ’14, pages 1–15, New York, NY, USA, April 2014. Association for Computing Machinery.
- [14] Per Ekemark, Yuan Yao, Alberto Ros, Konstantinos Sagonas, and Stefanos Kaxiras. TSOPEP: Efficient Coherence-Based Strict Persistency. In *2021 IEEE International Symposium on High-Performance Computer Architecture (HPCA)*, pages 125–138, February 2021.
- [15] Hongyu Fan, Zhihang Sun, and Fei He. Satisfiability Modulo Ordering Consistency Theory for SC, TSO, and PSO Memory Models. *ACM Transactions on Programming Languages and Systems*, 45(1):6:1–6:37, March 2023.
- [16] Xinwei Fu, Wook-Hee Kim, Ajay Paddayuru Shreepathi, Mohannad Ismail, Sunny Wadkar, Dongyoon Lee, and Changwoo Min. Witcher: Systematic Crash Consistency Testing for Non-Volatile Memory Key-Value Stores. In *Proceedings of the ACM SIGOPS 28th Symposium on Operating Systems Principles*, SOSP ’21, pages 100–115, New York, NY, USA, October 2021. Association for Computing Machinery.
- [17] Cunjing Ge, Feifei Ma, Jeff Huang, and Jian Zhang. SMT Solving for the Theory of Ordering Constraints. In Xipeng Shen, Frank Mueller, and James Tuck, editors, *Languages and Compilers for Parallel Computing*, pages 287–302, Cham, 2016. Springer International Publishing.
- [18] Kaan Genç, Michael D. Bond, and Guoqing Harry Xu. Crafty: Efficient, HTM-compatible persistent transactions. In *Proceedings of the 41st ACM SIGPLAN Conference on Programming Language Design and Implementation*, PLDI 2020, pages 59–74, New York, NY, USA, June 2020. Association for Computing Machinery.
- [19] Jerrin Shaji George, Mohit Verma, Rajesh Venkatasubramanian, and Pratap Subrahmanyam. Go-pmem: Native support for programming persistent memory in go. In *2020 USENIX Annual Technical Conference (USENIX ATC 20)*, pages 859–872. USENIX Association, July 2020.
- [20] Ellis Giles, Kshitij Doshi, and Peter Varman. Continuous checkpointing of HTM transactions in NVM. *ACM SIGPLAN Notices*, 52(9):70–81, June 2017.
- [21] Vaibhav Gogte, William Wang, Stephan Diestelhorst, Peter M. Chen, Satish Narayanasamy, and Thomas F. Wenisch. Relaxed Persist Ordering Using Strand Persistency. In *2020 ACM/IEEE 47th Annual International Symposium on Computer Architecture (ISCA)*, pages 652–665, May 2020.
- [22] Hamed Gorjiara, Weiyu Luo, Alex Lee, Guoqing Harry Xu, and Brian Demsky. Checking robustness to weak persistency models. In *Proceedings of the 43rd ACM SIGPLAN International Conference on Programming Language Design and Implementation*, PLDI 2022, pages 490–505, New York, NY, USA, June 2022.
- [23] Hamed Gorjiara, Guoqing Harry Xu, and Brian Demsky. Jaaru: Efficiently model checking persistent memory programs. In *Proceedings of the 26th ACM International Conference on Architectural Support for Programming Languages and Operating Systems*, ASPLOS ’21, pages 415–428, New York, NY, USA, April 2021. Association for Computing Machinery.
- [24] Fei He, Zhihang Sun, and Hongyu Fan. Satisfiability modulo ordering consistency theory for multi-threaded program verification. In *Proceedings of the 42nd ACM SIGPLAN International Conference on Programming Language Design and Implementation*, PLDI 2021, pages 1264–1279, New York, NY, USA, June 2021. Association for Computing Machinery.
- [25] Fei He, Zhihang Sun, and Hongyu Fan. Deagle: An SMT-based Verifier for Multi-threaded Programs (Competition Contribution). In Dana Fisman and Grigore Rosu, editors, *Tools and Algorithms for the Construction and Analysis of Systems*, pages 424–428, Cham, 2022. Springer International Publishing.
- [26] Terry Ching-Hsiang Hsu, Helge Brügger, Indrajit Roy, Kimberly Keeton, and Patrick Eugster. NVthreads: Practical Persistence for Multi-threaded Applications. In *Proceedings of the Twelfth European Conference on Computer Systems*, EuroSys ’17, pages 468–482, New York, NY, USA, April 2017. Association for Computing Machinery.
- [27] Intel Corporation. Persistent Memory Development Kit. <https://pmem.io/pmdl/>, 2023.
- [28] Omar Inverso, Ermenegildo Tomasco, Bernd Fischer, Salvatore La Torre, and Gennaro Parlato. Bounded Model Checking of Multi-threaded C Programs via Lazy Sequentialization. In Armin Biere and Roderick Bloem, editors, *Computer Aided Verification*, Lecture Notes in Computer Science, pages 585–602, Cham, 2014. Springer International Publishing.
- [29] Joseph Izraelevitz, Terence Kelly, and Aasheesh Kolli. Failure-Atomic Persistent Memory Updates via JUSTDO Logging. *ACM SIGARCH*

- 1321 *Computer Architecture News*, 44(2):427–442, March 2016.
- 1322 [30] Joseph Izraelevitz, Hammurabi Mendes, and Michael L. Scott. Lin-
- 1323 earizability of Persistent Memory Objects Under a Full-System-Crash
- 1324 Failure Model. In Cyril Gavoille and David Ilcinkas, editors, *Distributed*
- 1325 *Computing*, pages 313–327, Berlin, Heidelberg, 2016. Springer.
- 1326 [31] Jungi Jeong and Changhee Jung. PMEM-spec: Persistent memory
- 1327 speculation (strict persistency can trump relaxed persistency). In
- 1328 *Proceedings of the 26th ACM International Conference on Architectural*
- 1329 *Support for Programming Languages and Operating Systems*, ASPLOS
- 1330 '21, pages 517–529, New York, NY, USA, April 2021. Association for
- 1331 Computing Machinery.
- 1332 [32] Arpit Joshi, Vijay Nagarajan, Marcelo Cintra, and Stratis Viglas. Effi-
- 1333 cient persist barriers for multicores. In *Proceedings of the 48th Inter-*
- 1334 *national Symposium on Microarchitecture*, MICRO-48, pages 660–671,
- 1335 New York, NY, USA, December 2015. Association for Computing Ma-
- 1336 chinery.
- 1337 [33] Rohan Kadekodi, Se Kwon Lee, Sanidhya Kashyap, Taesoo Kim,
- 1338 Aasheesh Kolli, and Vijay Chidambaram. SplitFS: Reducing software
- 1339 overhead in file systems for persistent memory. In *Proceedings of*
- 1340 *the 27th ACM Symposium on Operating Systems Principles*, SOSP '19,
- 1341 pages 494–508, New York, NY, USA, October 2019. Association for
- 1342 Computing Machinery.
- 1343 [34] Tomasz Kapela. An introduction to pmemcheck (part 1) - bas-
- 1344 ics. [https://pmem.io/blog/2015/07/an-introduction-to-pmemcheck-](https://pmem.io/blog/2015/07/an-introduction-to-pmemcheck-part-1-basics/)
- 1345 [part-1-basics/](https://pmem.io/blog/2015/07/an-introduction-to-pmemcheck-part-1-basics/), July 2015.
- 1346 [35] T. Kawahara, K. Ito, R. Takemura, and H. Ohno. Spin-transfer torque
- 1347 RAM technology: Review and prospect. *Microelectronics Reliability*,
- 1348 52(4):613–627, April 2012.
- 1349 [36] Terence Kelly. Is persistent memory persistent? *Communications of*
- 1350 *the ACM*, 63(9):48–54, August 2020.
- 1351 [37] Artem Khyzha and Ori Lahav. Taming x86-TSO persistency. *Proceed-*
- 1352 *ings of the ACM on Programming Languages*, 5(POPL):1–29, January
- 1353 2021.
- 1354 [38] Michalis Kokologiannakis, Ilya Kaysin, Azalea Raad, and Viktor
- 1355 Vafeiadis. PerSeVerE: Persistency semantics for verification under ext4.
- 1356 *Proceedings of the ACM on Programming Languages*, 5(POPL):43:1–
- 1357 43:29, January 2021.
- 1358 [39] Michalis Kokologiannakis, Iason Marmanis, Vladimir Gladstein, and
- 1359 Viktor Vafeiadis. Truly stateless, optimal dynamic partial order reduc-
- 1360 tion. *Proceedings of the ACM on Programming Languages*, 6(POPL):49:1–
- 1361 49:28, January 2022.
- 1362 [40] Michalis Kokologiannakis and Viktor Vafeiadis. HMC: Model Check-
- 1363 ing for Hardware Memory Models. In *Proceedings of the Twenty-Fifth*
- 1364 *International Conference on Architectural Support for Programming*
- 1365 *Languages and Operating Systems*, ASPLOS '20, pages 1157–1171, New
- 1366 York, NY, USA, March 2020. Association for Computing Machinery.
- 1367 [41] Michalis Kokologiannakis and Viktor Vafeiadis. GenMC: A Model
- 1368 Checker for Weak Memory Models. In Alexandra Silva and K. Rus-
- 1369 tan M. Leino, editors, *Computer Aided Verification*, pages 427–440,
- 1370 Cham, 2021. Springer International Publishing.
- 1371 [42] Youngjin Kwon, Henrique Fingler, Tyler Hunt, Simon Peter, Emmett
- 1372 Witchel, and Thomas Anderson. Strata: A Cross Media File System.
- 1373 In *Proceedings of the 26th Symposium on Operating Systems Princi-*
- 1374 *ples*, SOSP '17, pages 460–477, New York, NY, USA, October 2017.
- 1375 Association for Computing Machinery.
- 1376 [43] Philip Lantz, Subramanya Dulloori, Sanjay Kumar, Rajesh Sankaran,
- 1377 and Jeff Jackson. Yat: A validation framework for persistent memory
- 1378 software. In *2014 USENIX Annual Technical Conference (USENIX ATC*
- 1379 *14)*, pages 433–438, Philadelphia, PA, June 2014. USENIX Association.
- 1380 [44] Qingrui Liu, Joseph Izraelevitz, Se Kwon Lee, Michael L. Scott, Sam H.
- 1381 Noh, and Changhee Jung. iDO: Compiler-Directed Failure Atomicity
- 1382 for Nonvolatile Memory. In *2018 51st Annual IEEE/ACM International*
- 1383 *Symposium on Microarchitecture (MICRO)*, pages 258–270, October
- 1384 2018.
- 1385 [45] Sihang Liu, Korakit Seemakhupt, Yizhou Wei, Thomas Wenisch,
- 1386 Aasheesh Kolli, and Samira Khan. Cross-Failure Bug Detection in Per-
- 1387 sistent Memory Programs. In *Proceedings of the Twenty-Fifth Interna-*
- 1388 *tional Conference on Architectural Support for Programming Languages*
- 1389 *and Operating Systems*, ASPLOS '20, pages 1187–1202, New York, NY,
- 1390 USA, March 2020. Association for Computing Machinery.
- 1391 [46] Sihang Liu, Yizhou Wei, Jishen Zhao, Aasheesh Kolli, and Samira
- 1392 Khan. PMTest: A Fast and Flexible Testing Framework for Persistent
- 1393 Memory Programs. In *Proceedings of the Twenty-Fourth International*
- 1394 *Conference on Architectural Support for Programming Languages and*
- 1395 *Operating Systems*, ASPLOS '19, pages 411–425, New York, NY, USA,
- 1396 April 2019. Association for Computing Machinery.
- 1397 [47] Virendra J. Marathe, Margo Seltzer, Steve Byan, and Tim Harris. Persi-
- 1398 stent memcached: Bringing legacy code to Byte-Addressable persistent
- 1399 memory. In *9th USENIX Workshop on Hot Topics in Storage and File*
- 1400 *Systems (HotStorage 17)*, Santa Clara, CA, July 2017. USENIX Associa-
- 1401 tion.
- 1402 [48] Iason Marmanis and Viktor Vafeiadis. SMT-Based Verification of Per-
- 1403 sistency Invariants of Px86 Programs. In Akash Lal and Stefano
- 1404 Tonetta, editors, *Verified Software. Theories, Tools and Experiments.*, Lec-
- 1405 ture Notes in Computer Science, pages 92–110, Cham, 2023. Springer
- 1406 International Publishing.
- 1407 [49] Diego Ongaro, Stephen M. Rumble, Ryan Stutsman, John Ousterhout,
- 1408 and Mendel Rosenblum. Fast crash recovery in RAMCloud. In *Proceed-*
- 1409 *ings of the Twenty-Third ACM Symposium on Operating Systems*
- 1410 *Principles*, SOSP '11, pages 29–41, New York, NY, USA, October 2011.
- 1411 Association for Computing Machinery.
- 1412 [50] Steven Pelley, Peter M. Chen, and Thomas F. Wenisch. Memory persi-
- 1413 stency. *ACM SIGARCH Computer Architecture News*, 42(3):265–276,
- 1414 June 2014.
- 1415 [51] Azalea Raad, Luc Maranget, and Viktor Vafeiadis. Extending Intel-x86
- 1416 consistency and persistency: Formalising the semantics of Intel-x86
- 1417 memory types and non-temporal stores. *Proceedings of the ACM on*
- 1418 *Programming Languages*, 6(POPL):1–31, January 2022.
- 1419 [52] Azalea Raad and Viktor Vafeiadis. Persistence semantics for weak
- 1420 memory: Integrating epoch persistency with the TSO memory model.
- 1421 *Proceedings of the ACM on Programming Languages*, 2(OOPSLA):137:1–
- 1422 137:27, October 2018.
- 1423 [53] Azalea Raad, John Wickerson, Gil Neiger, and Viktor Vafeiadis. Persi-
- 1424 stency semantics of the Intel-x86 architecture. *Proceedings of the ACM*
- 1425 *on Programming Languages*, 4(POPL):11:1–11:31, January 2020.
- 1426 [54] Azalea Raad, John Wickerson, and Viktor Vafeiadis. Weak persistency
- 1427 semantics from the ground up: Formalising the persistency seman-
- 1428 tics of ARMv8 and transactional models. *Proceedings of the ACM on*
- 1429 *Programming Languages*, 3(OOPSLA):135:1–135:27, October 2019.
- 1430 [55] Jinglei Ren, Jishen Zhao, Samira Khan, Jongmoo Choi, Yongwei Wu,
- 1431 and Onur Mutlu. ThyNVM: Enabling software-transparent crash
- 1432 consistency in persistent memory systems. In *Proceedings of the 48th*
- 1433 *International Symposium on Microarchitecture*, MICRO-48, pages 672–
- 1434 685, New York, NY, USA, December 2015. Association for Computing
- 1435 Machinery.
- 1436 [56] Peter Sewell, Susmit Sarkar, Scott Owens, Francesco Zappa Nardelli,
- 1437 and Magnus O. Myreen. X86-TSO: A rigorous and usable program-
- 1438 mer's model for x86 multiprocessors. *Communications of the ACM*,
- 1439 53(7):89–97, January 2010.
- 1440 [57] Zhihang Sun, Hongyu Fan, and Fei He. Consistency-preserving propa-
- 1441 gation for SMT solving of concurrent program verification. *Proceedings*
- 1442 *of the ACM on Programming Languages*, 6(OOPSLA2):158:929–158:956,
- 1443 October 2022.
- 1444 [58] Haris Volos, Sanketh Nalli, Sankarlingam Panneerselvam,
- 1445 Venkatanathan Varadarajan, Prashant Saxena, and Michael M.
- 1446 Swift. Aerie: Flexible file-system interfaces to storage-class memory.
- 1447 In *Proceedings of the Ninth European Conference on Computer Systems*,
- 1448 EuroSys '14, pages 1–14, New York, NY, USA, April 2014. Association
- 1449 for Computing Machinery.
- 1450

1431	for Computing Machinery.	1486
1432	[59] Xiaojian Wu and A. L. Narasimha Reddy. SCMFS: A file system for	1487
1433	storage class memory. In <i>Proceedings of 2011 International Conference</i>	1488
1434	<i>for High Performance Computing, Networking, Storage and Analysis</i> ,	1489
1435	SC '11, pages 1–11, New York, NY, USA, November 2011. Association	1490
1436	for Computing Machinery.	1491
1437	[60] Jian Xu and Steven Swanson. NOVA: A log-structured file system	1492
1438	for hybrid Volatile/Non-volatile main memories. In <i>14th USENIX</i>	1493
1439	<i>Conference on File and Storage Technologies (FAST 16)</i> , pages 323–338,	1494
1440	Santa Clara, CA, February 2016. USENIX Association.	1495
1441	[61] Jian Xu, Lu Zhang, Amirsaman Memaripour, Akshatha Gangadharaiah,	1496
1442	Amit Borase, Tamires Brito Da Silva, Steven Swanson, and Andy	1497
1443	Rudoff. NOVA-Fortis: A Fault-Tolerant Non-Volatile Main Memory	1498
1444	File System. In <i>Proceedings of the 26th Symposium on Operating Systems</i>	1499
1445	<i>Principles</i> , SOSP '17, pages 478–496, New York, NY, USA, October 2017.	1500
1446	Association for Computing Machinery.	1501
1447	[62] Liangze Yin, Wei Dong, Wanwei Liu, Yunchou Li, and Ji Wang. YOGAR-	1502
1448	CBMC: CBMC with Scheduling Constraint Based Abstraction Refine-	1503
1449	ment. In Dirk Beyer and Marieke Huisman, editors, <i>Tools and Algo-</i>	1504
1450	<i>rithms for the Construction and Analysis of Systems</i> , pages 422–426,	1505
1451	Cham, 2018. Springer International Publishing.	1506
1452	[63] Liangze Yin, Wei Dong, Wanwei Liu, and Ji Wang. Scheduling con-	1507
1453	straint based abstraction refinement for weak memory models. In	1508
1454	<i>Proceedings of the 33rd ACM/IEEE International Conference on Auto-</i>	1509
1455	<i>mated Software Engineering</i> , ASE '18, pages 645–655, New York, NY,	1510
1456	USA, September 2018. Association for Computing Machinery.	1511
1457	[64] Liangze Yin, Wei Dong, Wanwei Liu, and Ji Wang. On Scheduling	1512
1458	Constraint Abstraction for Multi-Threaded Program Verification. <i>IEEE</i>	1513
1459	<i>Transactions on Software Engineering</i> , 46(5):549–565, May 2020.	1514
1460	[65] Bohong Zhu, Youmin Chen, Qing Wang, Youyou Lu, and Jiwu Shu.	1515
1461	Octopus+: An RDMA-Enabled Distributed Persistent Memory File	1516
1462	System. <i>ACM Transactions on Storage</i> , 17(3):19:1–19:25, August 2021.	1517
1463		1518
1464		1519
1465		1520
1466		1521
1467		1522
1468		1523
1469		1524
1470		1525
1471		1526
1472		1527
1473		1528
1474		1529
1475		1530
1476		1531
1477		1532
1478		1533
1479		1534
1480		1535
1481		1536
1482		1537
1483		1538
1484		1539
1485		1540

A Proof

Theorem 2. Given an instrumented execution G^i valid under x86-TSO, the non-volatile state s_o induced by `rrf` is reachable under G^i iff $\text{cons}_{\text{DPTSO}}(G^i)$ holds.

Proof. To clarify, we use suffixes TSO and DPTSO to distinguish the respective `co` and `hb` order in Definition 5 and Definition 8. Assume that there are m shared variables $\mathcal{V}_p = \{x_1, \dots, x_m\}$ and $|G^i.E \cap W| + |G^i.E \cap FL| = n$, i.e. there are n store and flush events in total.

(\leftarrow): Suppose $\text{cons}_{\text{DPTSO}}(G^i)$ holds.

To prove s_o is reachable, let $\text{hb}_{\text{TSO}} = \text{hb}_{\text{DPTSO}}$ and $\text{co}_{\text{TSO}} = \text{co}_{\text{DPTSO}}$. We first arbitrarily construct a persist order `nvo` that satisfies the two axioms per Definition 1. The first axiom requires `nvo` to conform with the per-location store order of hb_{TSO} , which is exactly co_{TSO} in this case. Note `nvo` is a total order over all stores and flushes. Let the sequence of events induced by `nvo` be e_1, \dots, e_n , and for each variable x_i , let $s_o(x_i)$ equals the store e_{k_i} where $1 \leq k_i \leq n$.

Let e_f be the last flush event in the sequence. Consider $u = \max(k_1, k_2, \dots, k_m, f)$. Now we have the prefix $\bar{e} = e_1, \dots, e_u$ and try to adjust `nvo` such that \bar{e} induces s_o . In other words, it requires $\forall 1 \leq i \leq m, \forall k_i < j \leq u. \text{loc}(e_j) \neq x_i$ holds.

This is ensured by repeatedly reordering events in `nvo` while adhering to the two axioms. At each step, we pick a shared variable x_i for which the above condition does not hold and find the event e_p such that $\forall k_i < j \leq u. \text{loc}(e_j) = x_i \rightarrow j \leq p$, i.e. e_p is the last store to x_i in the range $[e_{k_i}, e_u]$. We then rearrange `nvo` so that e_p succeeds e_u . Apparently, the new persist order does not infringe the first axiom, since e_p is the last store in the range and no co_{TSO} order is violated by the reorder.

We now prove that the second axiom is not violated. From the above assumption, we know $(e_{k_i}, e_p) \in \text{co}_{\text{TSO}}$, so for any flush event e_q on x_i we also have $(e_q, e_p) \in \text{dtpo}$. Since $\text{cons}_{\text{DPTSO}}(G^i)$ holds, the first requirement of Definition 8 gives $(e_p, e_q) \notin (\text{ppo} \cup \text{rf} \cup \text{fr} \cup \text{co}_{\text{DPTSO}})^+$ which simplifies to $(e_p, e_q) \notin \text{hb}_{\text{TSO}}$. In this case, the premise of the second axiom does not hold, and thus e_p is not `nvo`-ordered with any flush events. e_p is therefore safe to be reordered.

Since the reorder of e_p in `nvo` does not violate the two axioms, after finite steps, the prefix \bar{e} must satisfy the aforementioned condition and induces s_o , thus we have proved s_o is reachable under G^i .

(\rightarrow): Suppose s_o is reachable, then for some $\text{co}_{\text{TSO}}, \text{hb}_{\text{TSO}}$, there is a persist order `nvo` and its prefix $\bar{e} = e_1, \dots, e_u$ that induces s_o . Let $\text{co}_{\text{DPTSO}} = \text{co}_{\text{TSO}}$. For each variable x_i , let $s_o(x_i)$ equals the store e_{k_i} where $1 \leq k_i \leq u$.

Assume some store e_p to x_i happens after e_{k_i} , then $(e_{k_i}, e_p) \in \text{co}_{\text{TSO}}$. By definition of a reachable state, we have $p > u$. Note that the prefix contains all flush events. Therefore, all flushes e_q to x_i must be `nvo`-ordered before e_p . This entails that $(e_q, e_p) \in \text{hb}_{\text{TSO}}$, otherwise the first axiom of `nvo` is violated.

Now consider the validity of G^i under DPTSO. Requirement (2) of Definition 8 follows directly from Definition 5. Suppose that requirement (1) is violated. Since G^i is valid under x86-TSO, there must be a flush event e_q , store event e_{k_i} and e_p on x_i such that $(e_q, e_p) \in \text{dtpo}$ and $(e_p, e_q) \in \text{hb}_{\text{DPTSO}}$. From the reasoning above, we have $(e_q, e_p) \in \text{hb}_{\text{TSO}}$.

We now consider the path from e_p to e_q on the directed graph. If no `dtpo` edge is on the path, then we have $(e_p, e_q) \in \text{hb}_{\text{TSO}}$. If there are one or more `dtpo` edge on the directed path, we note that any $(e'_p, e'_q) \in \text{dtpo}$ entails $(e'_p, e'_q) \in \text{hb}_{\text{TSO}}$ by the argument above as well. Therefore we also have $(e_p, e_q) \in \text{hb}_{\text{TSO}}$. In either case, it contradicts the assumption that G^i is valid under x86-TSO. We hereby prove requirement (1) holds, i.e. $\text{cons}_{\text{DPTSO}}(G^i)$ holds. \square

Lemma 2. If for every alternation G' of G , $\text{cons}_{\text{TSO}}(G')$ does not hold, then the observed state s_o must be unreachable under x86-TSO.

Proof. Suppose the state is reachable under x86-TSO and let G'' be an execution that induces this state, thus $\text{cons}_{\text{TSO}}(G'')$ holds. Since x86-TSO (Definition 5) only requires acyclicity of certain orders, by restricting $G''.E$ and the ordering constraints of G'' on G , we can always construct an execution G' such that $\text{cons}_{\text{TSO}}(G')$ also holds. G' is an alternation of G , thus contradicts with the assumption. \square

B Atomic Block

Intuitively, regarding the recovery observer as atomic forces all read events in it to happen at the same time.

Definition 9. Given a partial instrumented execution G' , regarding the recovery observer as atomic block yield the partial execution $\text{atomic}(G') = (E', E'_0, \text{po}', \text{rf}')$ such that

- $E' = (G'.E \setminus \text{REC}) \cup \{e_r\}$
- $E'_0 = E_0$
- $\text{po}' = \text{po}$
- $\text{rf}' = \text{rf} \setminus \text{rrf} \cup \{(e_1, e_r) \mid \exists e_2 \in \text{REC}. (e_1, e_2) \in \text{rrf}\}$

where e_r is a fresh event with $\text{Tid}(e_r) = P_r$.