

Software Verification for Weak Memory via Program Transformation^{*}

Jade Alglave^{1,2}, Daniel Kroening², Vincent Nimal², and Michael Tautschnig^{2,3}

¹ University College London

² University of Oxford

³ Queen Mary, University of London

dedicated to the memory of Kohei Honda

Abstract Multiprocessors implement weak memory models, but program verifiers often assume *Sequential Consistency* (SC), and thus may miss bugs due to weak memory. We propose a sound transformation of the program to verify, enabling SC tools to perform verification w.r.t. weak memory. We present experiments for a broad variety of models (from x86-TSO to Power) and a vast range of verification tools, quantify the additional cost of the transformation and highlight the cases when we can drastically reduce it. Our benchmarks include work-queue management code from PostgreSQL.

1 Introduction

Current multi-core architectures such as Intel’s x86, IBM’s Power or ARM implement *weak memory models* for performance reasons, allowing optimisations such as *instruction reordering*, *store buffering* or *write atomicity relaxation* [3]. These models make concurrent programming and debugging extremely challenging, because the execution of a concurrent program might not be an interleaving of its instructions, as would be the case on a Sequentially Consistent (SC) architecture [21]. As an instance, the lock-free signalling code in the open-source database PostgreSQL failed regression tests on a PowerPC cluster, due to the memory model. We study this bug in detail in Sec. 5.

This observation highlights the crucial need for weak memory aware verification. Yet, most existing work assume SC, hence might miss bugs specific to weak memory. Recent work addresses the design or the adaptation of existing methods and tools to weak memory [25,29,17,13,23,11,2], but often focuses on one specific model or cannot handle the write atomicity relaxation of Power/ARM: generality remains a challenge.

Since we want to avoid writing one tool per architecture of interest, we propose a unified method. Given a program analyser handling SC concurrency for C programs, we *transform its input* to simulate the possible non-SC behaviours of the program whilst executing the program on SC. Essentially, we augment our programs with arrays to simulate (on SC) the buffering and caching scenarios due to weak memory.

^{*} Supported by ERC project 280053, EPSRC project EP/G026254/1 and the Semiconductor Research Corporation (SRC) under task 2269.002.

The verification problem for weak memory models is known to be hard (e.g. non-primitive recursive for TSO), if not undecidable (e.g. for RMO-like models) [9]. This means that we cannot design a *complete* verification method. Yet, we can achieve *soundness*, by implementing our tools in tandem with the design of a proof, and by stressing our tools with test cases reflecting subtle points of the proof.

We also aim for an effective and unified verification setup, where one can easily plug a tool of choice. This paper meets these objectives by making three new contributions:

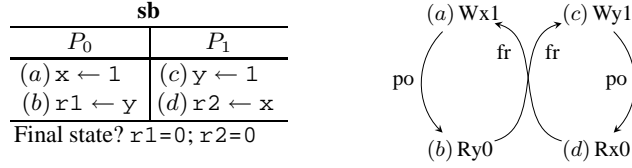
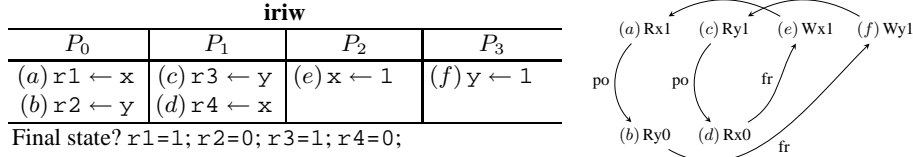
1. To design our transformation, we define in Sec. 3 an abstract state machine that we prove (in the Coq proof assistant) equivalent to the framework of [8] (recalled in Sec. 2). We also explain how this equivalence proof allows us to design a drastically improved transformation with a speed-up of more than two orders of magnitude.
2. Sec. 4 describes our implementation, highlighting the generality of our approach: we support a broad variety of models (x86/TSO, PSO, RMO and Power) and all concurrency-aware program analysers for C programs (cf. experiments below).
3. Sec. 5 details our experiments. i) We systematically validate our implementation w.r.t. our theoretical study with 555 *litmus tests* exercising weak memory artefacts. We study the overhead and validate the viability of our transformation using Blender [20], CheckFence [13], ESBMC [14], MMChecker [17], Poirot [1], SatAbs [15], and Threader [16]. ii) We verify an excerpt of the relational database software PostgreSQL, which has a bug specific to Power. iii) Our transformation easily scales to systems code from the Linux kernel or the Apache HTTP server, and also industrial code.

We provide the source and documentation of our tools, our benchmarks, experimental reports, Coq proofs and their typeset sketches online: www.cprover.org/wmm/

Related Work We focus here on the *verification* problem, i.e., detecting the behaviours that are buggy, not all the non-SC ones. This problem is non-primitive recursive for TSO [9]. It is undecidable if read/write or read/read pairs can be reordered, as in RMO-like models [9]. Forbidding *causal loops* restores decidability; relaxing write atomicity makes the problem undecidable again [10].

Existing solutions use various bounds over the objects of the model [11,19], over-approximate the possible program behaviours [20,18], or relinquish termination [22]. For TSO, [2] presents a sound and complete solution. We present a provably sound method that allows to lift any SC method or tool to a large spectrum of weak memory models, ranging from x86 to Power. We build an operational model; [24] presented such a model, but theirs is restricted to TSO. Given the undecidability of the problem, we cannot provide completeness, as we focus on soundness. We do not use any bound in our theoretical model (Sec. 3), but our implementation uses finite buffers (Sec. 4).

Our approach also reduces the amount of instrumentation in a provably sound manner. Unlike [11], we only instrument selected shared memory accesses. For TSO this would follow immediately from [12], but we generalise to models such as Power.


Figure 1. Store Buffering (**sb**)

Figure 2. Independent Reads of Independent Writes (**iriw**)

2 Context: Axiomatic Memory Model

In an operational view, weak memory effects occur as follows: A processor can commit a write first to a store buffer, then to a cache, and finally to memory. When a write hits the memory, all the processors agree on its value. But while the write is in transit through store buffers and caches, a read can occur before the value is actually available to all processors from the memory.

To describe such scenarios, we use the framework of [8], which provably embraces several (*weak*) architectures: SC [21], Sun TSO (i.e. the x86 model [24]), PSO and RMO, Alpha, and a fragment of Power. At the core of this framework we use *relations* over *read and write memory events*. We introduce this framework on *litmus tests*, as shown in Fig. 1. The left-hand side of the figure shows a multi-threaded program. The shared variables x and y are initialised to zero. A store instruction (e.g. $x \leftarrow 1$ on P_0) gives rise to a write event ($(a)Wx1$), and a load (e.g. $r1 \leftarrow y$ on P_0) to a read event ($(b)Ry0$). The property of interest is whether there exists an execution of the program such that the final state is $r1=0$ and $r2=0$. To determine this, we study the *event graph*, given on the right-hand side of the figure. An architecture allows an execution when it represents a *global happens-before* order over all processors. A cycle in an event graph is a violation of global happens before, unless the architecture relaxes any of the relations contributing to this cycle. Thus, if the graph has a cycle, we check if the architecture *may relax* some relations. Such a relaxation makes the graph acyclic, which implies that the architecture allows the final state.

In SC, nothing is relaxed, thus the cycle in Fig. 1 forbids the execution. On the other hand, x86 relaxes the program order (po in Fig. 1) between writes and reads, thus the forbidding cycle no longer exists, and the given final state can be observed.

Formalisation An *event* is a read or a write memory access, composed of a unique identifier, a direction R for read or W for write, a memory address, and a value. We represent each instruction by the events it issues. In Fig. 2, we associate the store $x \leftarrow 1$ on processor P_2 with the event $(e)Wx1$. We define two utility functions on events:

$\text{proc}(e)$ returns the processor executing the event e , and $\text{addr}(e)$ yields the address of a read or write event e .

A set of events \mathbb{E} and their program order po form an *event structure* $E \triangleq (\mathbb{E}, \text{po})$. po is a per-processor total order over the events of \mathbb{E} . We write $\text{dp} \subseteq \text{po}$ for the relation that models the *dependencies* between instructions, e.g. an *address dependency* occurs when computing the address of a load or store from the value of a preceding load.

We represent the *communication* between processors leading to the final state via an *execution witness* $X \triangleq (\text{ws}, \text{rf})$, which consists of two relations over the events. First, the *write serialisation* ws is a per-address total order on writes which models the *memory coherence* widely assumed by modern architectures. It links a write w to any write w' to the same address that hits the memory after w . Second, the *read-from* relation rf links a write w to a read r such that r reads the value written by w .

Given a pair of writes $(w', w) \in \text{ws}$ and a read-from pair $(w', r) \in \text{rf}$, we are to complete global happens before: w' happens before w by ws and r reads from w' by rf . Thus r is to happen before w , as otherwise it would have to read from w . To that aim, we derive the *from-read* relation fr from ws and rf . A read r is in fr with a write w when the write w' from which r reads hit the memory before w did. Formally, we have: $(r, w) \in \text{fr} \triangleq \exists w', (w', r) \in \text{rf} \wedge (w', w) \in \text{ws}$.

In Fig. 2, the specified outcome corresponds to the execution on the right if each memory location initially holds 0. If $r1=1$ in the end, the read (a) obtained its value from the write (e) on P_2 , hence $(e, a) \in \text{rf}$. If $r2=0$ in the end, the read (b) obtained its value from the initial state, thus before the write (f) on P_3 , hence $(b, f) \in \text{fr}$. Similarly, we have $(f, c) \in \text{rf}$ from $r3=1$, and $(d, e) \in \text{fr}$ from $r4=0$.

Relaxed or safe We model the scenario of reads to occur in advance, as described at the beginning of this section, by some subrelation of the read-from rf being *relaxed*, i.e. not included in global happens before. When a processor can read from its own store buffer [3] (the typical TSO/x86 scenario), we relax the internal read-from rfi . When two processors P_0 and P_1 can communicate privately via a cache (a case of *write atomicity* relaxation [3]), we relax the external read-from rfe , and call the corresponding write *non-atomic*. This is the main particularity of Power or ARM, and cannot happen on TSO/x86. Some program-order pairs may be relaxed (e.g. write-read pairs on x86, and all but dp ones on Power), i.e. only a subset of po is guaranteed to occur in this order. This subset constitutes *preserved program order*, ppo .

When a relation may not be relaxed, we call it *safe*. Architectures provide special *fence* (or *barrier*) instructions to prevent weak behaviours. Following [8], the relation $\text{fence} \subseteq \text{po}$ induced by a fence is *non-cumulative* when it only orders certain pairs of events surrounding the fence, i.e. fence is *safe*. The relation fence is *cumulative* when it additionally makes writes atomic, e.g. by flushing caches. In our axiomatic model, this amounts to making sequences of external read-from and fences ($\text{rfe}; \text{fence}$ or $\text{fence}; \text{rfe}$) safe, even though rfe alone would not be safe for the architecture. We denote the union of fence and the additional cumulativity by ab .

Architectures An *architecture* A determines the set safe_A of the relations safe on A , i.e. the relations embedded in global happens before. Following [8], we always consider the write serialisation ws and the from-read relation fr safe. SC relaxes nothing, i.e. rf and

po are safe. TSO authorises the reordering of write-read pairs and store buffering but nothing else. Fences are safe by design, thus $\text{ab} \subseteq \text{safe}_A$.

Finally, an execution (E, X) is *valid* on A when the three following conditions hold.

1. SC holds per address, i.e. the communication and the program order for accesses with same address po-loc are compatible: $\text{uniproc}(E, X) \triangleq \text{acyclic}(\text{ws} \cup \text{rf} \cup \text{fr} \cup \text{po-loc})$.
2. Values do not come out of thin air, i.e. there is no causal loop: $\text{thin}(E, X) \triangleq \text{acyclic}(\text{rf} \cup \text{dp})$.
3. There exists a linearisation of events in global happens before, i.e. the safe relations do not form a cycle: $\text{ghb}(E, X) \triangleq \text{acyclic}((\text{ws} \cup \text{rf} \cup \text{fr} \cup \text{po}) \cap \text{safe}_A)$.

Formally:

$$\text{valid}_A(E, X) \triangleq \text{uniproc}(E, X) \wedge \text{thin}(E, X) \wedge \text{ghb}(E, X)$$

3 Simulating Weak Behaviours on SC

We develop a provably correct instrumentation strategy for programs. To this end, we first give an operational description of memory models in terms of an *abstract state machine* (Sec. 3.1). We then show in Sec. 3.3 the equivalence of the axiomatic model of Sec. 2 and the abstract machine. We explain in Sec. 3.4 how this equivalence proof guides our instrumentation strategy.

3.1 Abstract machine

We define a non-deterministic state machine that reads a sequence of *labels*. The machine has a designated bad state \perp , and all other states of the machine represent system configurations, i.e. the memory, write buffers, and the set of pending reads. We write addr , evt , and rln for the types of memory addresses, events and relations, respectively.

Definition 1 (State). A state of the machine is either \perp or a triple (m, b, rs) , where

- the memory $(m : \text{addr} \rightarrow \text{evt})$ maps a memory address ℓ to a write to ℓ ;
- the write buffer $(b : \text{rln evt})$ is a total order over writes to the same address; the buffer has a special symbol \perp_b , placed before all events in the buffer;
- the read set $(rs : \text{set evt})$ is a set of read events.

We have a single set of reads, but one totally ordered buffer per address. Existing formalisations [24,11] use per-thread buffers, whereas our buffers are solely per-address objects. This allows us to model not only store buffering (which per-thread objects would allow), but also caching scenarios (fully non-atomic stores) as exhibited by **irw+dps**, i.e. the **irw** test of Fig. 2 with dependencies between the reads on P_0 and P_1 to prevent their reordering.

The machine performs transitions depending on *delay* and *flush* labels. Intuitively, a delay label pushes an object in the write buffer or read set. A flush label makes it exit the write buffer or read set. The details of transitions are described below.

Definition 2 (Label). For a write event w , $d(w(w))$ denotes its delay label, and $f(w(w))$ its flush label. For a read event r , its delay label (with direction r , read) is denoted by $d(r(w, r))$, and its flush is denoted by $f(r(w, r))$.

$$\begin{aligned}
\text{updm}(\mathbf{m}, w) &\triangleq x \mapsto \text{if } \text{addr}(x) = \text{addr}(w) \text{ then } w \text{ else } x \\
\text{updb}(\mathbf{b}, w) &\triangleq \mathbf{b} \cup \{(w_1, w_2) \mid w_1 = \perp_{\mathbf{b}} \vee ((\perp_{\mathbf{b}}, w_1) \in \mathbf{b} \wedge \text{addr}(w_1) = \text{addr}(w)) \wedge \\
&\quad w_2 = w\} \\
\text{updrs}(\mathbf{rs}, r) &\triangleq \mathbf{rs} \cup \{r\} \\
\text{delb}(\mathbf{b}, w) &\triangleq \{(w_1, w_2) \mid (w_1, w_2) \in \mathbf{b} \wedge w_1 \neq w \wedge w_2 \neq w\} \\
\text{delrs}(\mathbf{rs}, r) &\triangleq \{e \mid e \in \mathbf{rs} \wedge e \neq r\} \\
\text{last}(\mathbf{b}, w) &\triangleq (\neg(\exists w', (\perp_{\mathbf{b}}, w') \in \mathbf{b}) \wedge w = \perp_{\mathbf{b}}) \vee \\
&\quad ((\exists w', (\perp_{\mathbf{b}}, w') \in \mathbf{b}) \wedge (\perp_{\mathbf{b}}, w) \in \mathbf{b} \wedge \neg(\exists w', (w', w) \in \mathbf{b})) \\
\text{rfm}(\mathbf{m}, \mathbf{b}, w) &\triangleq w = \mathbf{m}(\text{addr}(r)) \wedge \text{rr}(\mathbf{b}, \{w \mid (w, r) \in \text{po-loc}\}) = \emptyset
\end{aligned}$$

$$\begin{array}{c}
\text{WRITE TO BUFFER} \\
\top \\
\hline
s \xrightarrow{d(w)} (\mathbf{m}, \text{updb}(\mathbf{b}, w), \mathbf{rs})
\end{array}
\qquad
\begin{array}{c}
\text{DELAY READ} \\
\top \\
\hline
s \xrightarrow{d(r(w, r))} (\mathbf{m}, \mathbf{b}, \text{updrs}(\mathbf{rs}, r))
\end{array}$$

$$\begin{array}{c}
\text{READ FROM SET} \\
r \in \mathbf{rs} \wedge \quad (\text{R1}) \\
\mathbf{rs} \cap \{r \mid (r, w) \in \text{dp}\} = \emptyset \wedge \quad (\text{R2}) \\
\text{rr}(\mathbf{b}, \{e \mid (e, r) \in \text{ppo} \cup \text{ab}\}) = \emptyset \wedge \quad (\text{R3}) \\
\mathbf{rs} \cap \{e \mid (e, w) \in \text{ppo} \cup \text{ab}\} = \emptyset \wedge \quad (\text{R4}) \\
\mathbf{rs} \cap \{r \mid (r, w) \in \text{po-loc}\} = \emptyset \wedge \quad (\text{R5}) \\
\text{rfm}(\mathbf{m}, \mathbf{b}, w) \vee \quad (\text{R5}) \\
\text{last}(\text{rr}(\mathbf{b}, \{e \mid \text{addr}(e) = \ell\}), w) \quad (\text{W4}) \quad (w \neq \mathbf{m}(\text{addr}(r)) \wedge w \in \mathbf{b} \wedge \text{visible}(w, r)) \quad (\text{R6}) \\
\hline
s \xrightarrow{f(w)} (\text{updm}(\mathbf{m}, w), \text{delb}(\mathbf{b}, w), \mathbf{rs}) \qquad s \xrightarrow{f(r(w, r))} (\mathbf{m}, \mathbf{b}, \text{delrs}(\mathbf{rs}, r))
\end{array}$$

Figure 3. The abstract machine

A set L of labels is well-formed w.r.t. an event structure E when: in $d(w(w))$ or $f(w(w))$, w is a write of E ; in $d(r(w, r))$ or $f(r(w, r))$, w is a write of E and r a read of E , both with the same address; any event of E has a unique corresponding flush label in L ; when a flush label belongs to L , so does its delay counterpart.

Transitions We write $s \xrightarrow{l} s'$ to denote that the machine can make a transition from state s to state s' reading label l . Let the machine be in a state $(\mathbf{m}, \mathbf{b}, \mathbf{rs})$. Given a label, the machine performs transitions from one state to another if the conditions described below are fulfilled. Otherwise, the machine transitions to \perp (it gets stuck).

In Fig. 3, we give the formal definition of the transitions of our machine. We need to define a few auxiliary functions, also formally defined in Fig. 3. We update the memory with a write w via $\text{updm}(\mathbf{m}, w)$, a buffer with a write w via $\text{updb}(\mathbf{b}, w)$, and a set with a read r via $\text{updrs}(\mathbf{rs}, r)$. We delete a write w from a buffer via $\text{delb}(\mathbf{b}, w)$ and we delete a read r from a set via $\text{delrs}(\mathbf{rs}, r)$. We write $\text{rr}(R, S)$ for the restriction of a relation R to a set S , i.e. $\{(x, y) \mid (x, y) \in R \wedge x \in S \wedge y \in S\}$. We pick the last write to an address ℓ of a buffer via $\text{last}(\mathbf{b}, w)$. In prose, the transitions are as follows. To

avoid ambiguity in wording, we write “r-before” or “r-after” to express before or after w.r.t. the relation r.

- *Write to buffer*: a write $d(w(w))$ to address ℓ can always enter the buffer \mathbf{b} , taking its place \mathbf{b} -after all the writes to ℓ that are already in \mathbf{b} .
- *Delay read*: a read $d(r(w, r))$ can always enter the read set \mathbf{rs} .
- *Write from buffer to memory*: a write $f(w(w))$ to address ℓ exits the buffer \mathbf{b} and updates the memory at ℓ if:
 - there is no event e in the buffer nor in the read set which is $\mathbf{ppo} \cup \mathbf{ab}$ -before w (Conditions (W1) and (W2));
 - and there is no read from ℓ in the buffer which is \mathbf{po} -before w (Cond. (W3));
 - and there is no write to ℓ in the buffer which is \mathbf{b} -before w (Condition (W4)).
- *Read from set*: a read $f(r(w, r))$ from ℓ (Condition (R1)) exits the read set if:
 - there is no read in the read set that is \mathbf{dp} -before w (Condition (R2));
 - and there is no event in the buffer or in the read set that is $\mathbf{ppo} \cup \mathbf{ab}$ -before r (Conditions (R3) and (R4));
 - and either w is in memory, and there is no write to ℓ in the buffer that is \mathbf{po} -before r (Condition (R5));
 - or if w is not in memory, w is in the buffer and is *visible to* r (a notion defined below) (Condition (R6)).

To define a write w as *visible to a read* r , we need a few auxiliary functions. We define the part of the buffer visible to a read r as follows: $\mathbf{b}_r \triangleq \{w \mid (\perp_{\mathbf{b}}, w) \in \mathbf{b} \wedge ((\mathbf{rfi} \subseteq \mathbf{safe}_A) \Rightarrow \text{proc}(w) = \text{proc}(r)) \wedge ((\mathbf{rfe} \subseteq \mathbf{safe}_A) \Rightarrow \text{proc}(w) \neq \text{proc}(r))\}$. Now, w is visible to r when:

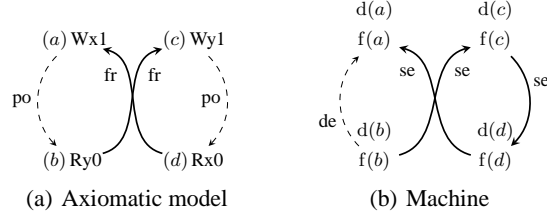
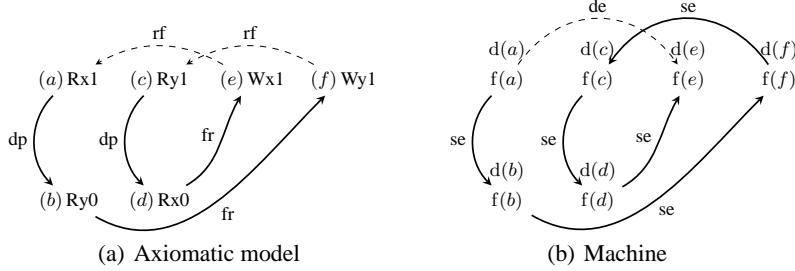
- w and r share the same address ℓ ;
- w is in the part of the buffer visible to r , namely if \mathbf{rfi} (resp. \mathbf{rfe}) is safe then w cannot be on the same (resp. a different) thread as r ($w \in \mathbf{b}_r$);
- w is \mathbf{b} -before the first write w_a to ℓ that is \mathbf{po} -after r ;
- w is equal to, or \mathbf{b} -after, the last write w_b to ℓ that is \mathbf{po} -before r .

All states except \perp are accepting states. Thus, the abstract machine accepts a sequence p of labels l_0, l_1, \dots if there is a sequence of states s_0, s_1, \dots such that $s_i \xrightarrow{l_i} s_{i+1}$ and $s_i \neq \perp$ for all i .

Definition 3 (Accepting sequence). A sequence p is a total order over L compatible with the program order, i.e. for two events $(x, y) \in \mathbf{po}$, their delay labels appear in the same order in p . It is accepting iff the sequence p is accepted by the abstract machine.

3.2 Illustration using examples

We illustrate the machine by revisiting the **sb** test of Fig. 1 for TSO and the **iriw** test of Fig. 2 for Power. Fig. 4 and 5 reproduce on the left the event graphs from Fig. 1 and 2. On the right, they show the counterparts in the abstract machine. We explain the labels on the arrows in the next section (§“From the axiomatic model to the machine”). We use the following graphical conventions. In the axiomatic world (i.e. on the left of our

**Figure 4.** Revisiting **sb** on TSO with our machine**Figure 5.** Revisiting **iriw+dps** on Power with our machine

figures), we reflect a pair that an architecture relaxes by a dashed arrow. For example, in the **sb** test of Fig. 4 on TSO, the write-read pairs (a, b) and (d, c) can be relaxed. Likewise, in the **iriw+dps** test of Fig. 5 on Power, the read-from pairs (e, a) and (f, c) can be relaxed (as opposed to the read-read pairs (a, b) on P_0 and (c, d) on P_1 , which are safe because of dependencies).

In any given execution, the abstract machine may choose to relax any pair that is not safe. Such pairs are depicted with a dashed arrow. Pairs that the machine does not relax are depicted with a thick arrow.

In Fig. 1, the pairs (a, b) on P_0 and (c, d) on P_1 are relaxed on TSO. Our machine may simulate the behaviour permitted on TSO by following the scenario in Fig. 4(b), which corresponds to the path $d(a) \rightarrow d(b) \rightarrow d(c) \rightarrow d(d) \rightarrow f(b) \rightarrow f(c) \rightarrow f(d) \rightarrow f(a)$. In the figure, the label “se” corresponds to a safe exit, and “de” to a delay exit, which are formalised below. The machine delays all events w.r.t. program order. In this scenario, the machine chooses to relax the pairs (a, b) by flushing the read b before the write a , ensuring that the registers $r1$ and $r2$ hold 0 in the end.

In Fig. 2, assume dependencies between the reads on P_0 and P_1 , so that (a, b) on P_0 and (c, d) on P_1 are safe on Power. Yet (e, a) and (f, c) may be relaxed on Power, because Power has non-atomic writes. Our machine may simulate the weak behaviour exhibited on Power by following Fig. 5(b), which corresponds to the path $d(e) \rightarrow d(a) \rightarrow f(a) \rightarrow d(b) \rightarrow f(b) \rightarrow d(f) \rightarrow f(f) \rightarrow d(c) \rightarrow f(c) \rightarrow d(d) \rightarrow f(d) \rightarrow f(e)$. Since (a, b) and (c, d) are safe on Power, our machine flushes a before b (resp. c before d). Since $(b, f) \in \text{fr}$ (resp. $(d, e) \in \text{fr}$), which is always safe, the machine flushes b before f (resp. d before e), ensuring that b and d read from memory, thus $r2$

and $r4$ hold 0 in the end. Finally, in this scenario, the machine chooses to relax the pairs (e, a) by flushing a before e , ensuring that $r1$ and $r3$ hold the value 1 in the end.

3.3 Equivalence of the axiomatic model and the abstract machine

We now prove the equivalence of the axiomatic model of Sec. 2 and the machine defined in Sec. 3.1. We first show that we can build an execution valid in the axiomatic model from any path of labels accepted by the machine (Thm. 1). We then show that we can build a path of labels accepted by the machine from any execution that is valid in axiomatic model (Thm. 2).

Thm. 1 (From the machine to the axiomatic model). *Let E be an event structure and L be a set of labels well-formed w.r.t. E . Then there exists an execution witness valid for E , if there is an accepting sequence p over L .*

Let $\text{ptoX}(p, L)$ denote the execution witness of Thm. 1. Recall from Sec. 2 that an execution witness is a pair of write serialisation and read-from map. Intuitively, we build these as follows. The write serialisation gathers the pairs of writes to the same address according to the order of their flushed parts in the accepting sequence p : $\{(w_1, w_2) \mid \text{addr}(w_1) = \text{addr}(w_2) \wedge (f(w(w_1)), f(w(w_2))) \in p\}$. For the read-from map, we simply gather the pairs given by the labels of L : $\{(w, r) \mid \text{addr}(w) = \text{addr}(r) \wedge f(r(w, r)) \in L\}$.

Proof (Thm. 1). We need to show that $(E, \text{ptoX}(p, L))$ passes the uniproc, thin and ghb checks. The three proofs follow the same lines, thus we focus on the first for brevity.

The execution passes the uniproc check iff for all $(x, y) \in \text{po-loc}$, we do not have $(y, x) \in \text{rf} \cup \text{fr} \cup \text{ws} \cup (\text{ws}; \text{rf}) \cup (\text{fr}; \text{rf})$ [4, App. A]. By contradiction take $(x, y) \in \text{po-loc}$ and $(y, x) \in \text{rf} \cup \text{fr} \cup \text{rf}$. We proceed by case disjunction over $(y, x) \in \text{rf} \cup \text{fr} \cup \text{ws} \cup (\text{ws}; \text{rf}) \cup (\text{fr}; \text{rf})$. We write ℓ for the address shared by x and y .

If $(y, x) \in \text{rf}$, $f(r(y, x))$ is in L . Since p is accepting, the Read from set transition on $f(r(y, x))$ does not block. Hence y is in memory, or y is in the buffer and visible to x . If y is in memory, y has been flushed, i.e. the Write from buffer to memory transition on $f(w(y))$ did not block. Hence there is no read from ℓ po-before y in the set. Yet $(x, y) \in \text{po-loc}$, and x is still in the set when y is in memory, a contradiction. If y is in the buffer and visible to x , y is in the buffer before the first write to ℓ po-after x . Yet, $(x, y) \in \text{po-loc}$, a contradiction.

For brevity, we present only the rf case; all the other cases are similar, using the premises of the rules of the machine. For example the $(y, x) \in \text{ws}$ case uses the Write from buffer to memory rule, in particular the fact that y exits the buffer if there is no write to ℓ before it in the buffer; yet x is still in there. The $(y, x) \in \text{fr}$ case uses the Read from set rule, in particular the fact that if the write w from which x reads is in memory, then there is no write to ℓ po-before y in the buffer; yet x is in there. If w is in the buffer, we use the fact that w is equal to, or in the buffer after, the last write to ℓ po-before x , which will block the flush of w , a contradiction. \square

For the other direction, we first build labels from the events of E . We augment our events with directions: a write w becomes $w(w)$ and r becomes $r(w, r)$, where $(w, r) \in \text{rf}$. Then we *split* an augmented event e into its delayed part $d(e)$, and its flushed part $f(e)$. We write $\text{labels}(E, X)$ for the labels built from the events of E .

Then we form the *delay pairs* of (E, X) , as follows. We build the relation ndelay over the events of E , such that: $((\text{ws} \cup \text{rf} \cup \text{fr}) \cap \text{safe}_A) \subseteq \text{ndelay}$; ndelay is transitive;

ndelay is irreflexive; if $(x, y) \notin \text{ndelay}$ then $(y, x) \in \text{ndelay}$. The delay pairs are the pairs (x, y) of events of E that are not in ndelay .

Given (E, X) and a choice of delay pairs, we build an accepting path p as follows, with e , e_1 , and e_2 denoting augmented events:

Delay before flush we always delay an event e before we flush it, i.e. $(d(e), f(e)) \in p$;
Enter $(e_1, e_2) \in \text{po}$ enter the buffer or set in this order, i.e. $(d(e_1), d(e_2)) \in p$;
Rf a write enters before we flush a read from it, i.e. $(d(e_1), f(e_2)) \in p$ if $(e_1, e_2) \in \text{rf}$;
Safe Exit $(e_1, e_2) \in \text{ndelay}$ are flushed in the same order, i.e. $(f(e_1), f(e_2)) \in p$.
Delay Exit $(e_1, e_2) \notin \text{ndelay}$ are flushed in the opposite order, i.e. $(f(e_2), f(e_1)) \in p$.

Reconsider Fig. 4(b) and 5(b). We omit the arrows corresponding to the first three cases to ease the reading of the figures. In Fig. 4(b), we chose (a, b) to be a delay pair, hence we flush them b before a , following the delay exit rule. On the contrary, (b, c) , (c, d) and (d, a) are not delay pairs, hence we flush b before c , c before d and d before a , following the safe exit rule. The same explanation applies in Fig. 5 to the pair (e, a) being delayed, and (a, b) , (f, c) , (c, d) and (d, e) being safe.

We build $X_{\text{top}}(E, X, \text{ndelay})$ as above. As ndelay is transitive and irreflexive, $X_{\text{top}}(E, X, \text{ndelay})$ is acyclic. Hence the transitive closure $(X_{\text{top}}(E, X, \text{ndelay}))^+$ is a partial order of the labels. Any linearisation $\text{lin}((X_{\text{top}}(E, X, \text{ndelay}))^+)$ of this transitive closure forms an actual path, which we show accepting when 1. X is valid 2. this linearisation has finite prefixes, in which case we say that (E, X) has finite prefixes:

Thm. 2 (From the axiomatic model to the machine). *For any valid execution (E, X) with finite prefixes, there is an accepting path p over labels L well-formed w.r.t. E .*

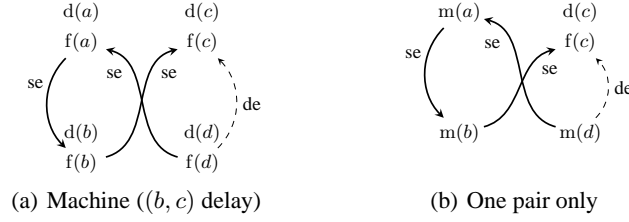
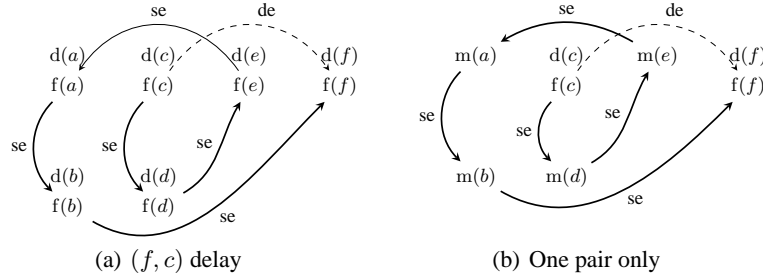
Proof. We need to show that no transition can block the machine. The Write to buffer and Delay read transitions are trivial since they can never block.

For the Write from buffer to memory case, suppose as a contradiction that the transition blocks on a write w to an address ℓ . If there is e $\text{ppo} \cup \text{ab}$ -before w in the buffer or the set, (e, w) cannot be a delay pair (because ppo and ab are safe), i.e. should be flushed in order, contradicting the presence of e in the buffer or the set. Otherwise, there is in the set a read r from ℓ po -before w . Therefore (r, w) is in fr , thus safe, hence cannot be a delay pair, and the same argument applies. Finally, if there is a write w' to ℓ before w in the buffer; one can show that (w', w) is in ws , hence w' should be flushed before w , a contradiction.

For the Read from set case, suppose as a contradiction that the transition blocks on a read (w, r) with address ℓ . If there is a read r' dp -before w in the set, one can show that r' should be flushed before r , and r should be flushed before r' (i.e. a thin-air cycle in X), a contradiction. If there is an event $\text{ppo} \cup \text{ab}$ -before r in the buffer or the set, the reasoning is the same as above in the write case. If w is in memory and there is a write to ℓ po -before r in the buffer, we create a uniproc cycle, a contradiction. If w is in the buffer and not visible to r , there are two cases. Either w is not on a thread whose buffer r can read w.r.t. A , in which case (w, r) do not form a delay pair and should be flushed in this order, contradicting the presence of w in the buffer. Or w is in the buffer after the first write to ℓ po -after r (or before the last write to ℓ po -before r), in which case we create a uniproc cycle. \square

3.4 Instrumentation

Thm. 2 leaves freedom in the instrumentation strategy. We can exploit the choice of delay pairs and the choice of the linearisation of $X_{\text{top}}(E, X)$ in order to reduce the overhead of running or verifying an instrumented program.

**Figure 6.** Choices for instrumenting **sb** for TSO**Figure 7.** Choices for instrumenting **iriw+dps** for Power

Choice of delay pairs The conditions on the **ndelay** relation restrict the choice of delay pairs. We have to put at least all the safe pairs into **ndelay**, by the first condition.

Since **ndelay** is transitive and irreflexive, it is acyclic. An execution (E, X) presents a cycle iff it is not SC (if it is SC, all pairs are safe and there is no cycle). [7, Thm.1] shows that an execution is valid on A but not on SC iff it contains *critical cycles*⁴. Thus we can put all pairs in **ndelay**, except one unsafe pair per critical cycle, which corresponds to the last condition over **ndelay**.

In Fig. 4(b), we build an accepting path corresponding to the axiomatic execution of Fig. 4(a) by choosing the unsafe pair (a, b) on the cycle to be a delay. In Fig. 6(a), we choose the unsafe pair (c, d) . Similarly for Fig. 5(a), we can build an accepting path corresponding to the axiomatic execution of Fig. 5(a) by choosing e.g. (e, a) as delay (cf. Fig. 5(b)). In Fig. 7(a), we choose (f, c) as delay.

Our examples are symmetric, thus the choice of which pair to delay should not make a difference. In Fig. 1, (a, b) and (c, d) are write-read pairs. Similarly in Fig. 2, (e, a) and (f, c) are of the same nature, namely **rfe** pairs. For asymmetric examples, the

⁴ We recall here the definition of [7]. Two events (x, y) are *competing*, written $(x, y) \in \text{cmp}$, if they are from distinct processors, to the same address, and at least one of them is a write (e.g. in Fig. 2, the read (a) from x on P_0 and the write (e) to x on P_2). A cycle $\sigma \subseteq \text{cmp} \cup \text{po}$ is critical when it is not a cycle in $(\text{cmp} \cup (\text{ppo} \cap \text{safe}_A))^+$ and it satisfies the two following properties: **(i)** Per processor, there are at most two memory accesses (x, y) on this processor and $\text{addr}(x) \neq \text{addr}(y)$. **(ii)** For a given memory address x , there are at most three accesses relative to x , and these accesses are from distinct processors $((w, w') \in \text{cmp}, (w, r) \in \text{cmp}, (r, w) \in \text{cmp}$ or $\{(r, w), (w, r')\} \subseteq \text{cmp}$). Fig. 2, shows a critical cycle of **iriw** on Power.

chosen delayed pair can make a crucial difference (cf. Sec. 5), if the instrumentation of one pair causes more execution or verification time overhead than the other.

Choice of the linearisation Thm. 2 accepts any linearisation of $(X_{\text{top}}(E, X, \text{ndelay}))^+$. Yet, some require less instrumentation than others. Consider Fig. 6(a) and (b): in both we choose to delay the pair (c, d) . On the left, we can pick any interleaving (compatible with X_{top}) of the delayed and flushed events to instantiate Thm. 2, e.g. $d(a) \rightarrow d(b) \rightarrow d(c) \rightarrow d(d) \rightarrow f(b) \rightarrow f(d) \rightarrow f(c) \rightarrow f(a)$.

On the right, we write $m(e)$ when the delayed and flushed part of an event happen without intervening events in between. Observe that in this case, the event e occurs w.r.t. memory: if it is a read, it reads from the memory; if it is a write, it writes to memory. In Fig. 6(b), we pick a particular interleaving, namely the one where all events are w.r.t. memory, except for the event c . This interleaving requires to instrument only one instruction, as opposed to all of them on the left.

Similarly in Fig. 7(a) and (b), we choose in both cases to delay the pair (f, c) . On the left, we instrument all instructions. On the right, we instrument only the pair (f, c) .

4 Implementation

4.1 Overview

We implemented the transformation technique of Sec. 3. Our tool reads a concurrent C program, possibly with inline assembly `m fence`, `sync`, or `lwsync` instructions (cf. Sec. 2). It generates a new concurrent C program augmented with C equivalents of write buffers and read sets of Sec. 3.1. The transformation proceeds in three main steps:

1. We devise an *abstract event structure*, as defined below, the concretisation of which amounts to all event structures (cf. Sec. 2) of the program.
2. Given an architecture, we identify potential critical cycles in this structure.
3. We instrument unsafe pairs in the cycle, as described in Sec. 3.4.

The resulting program is then passed to any SC program analyser.

The first two steps guide the program transformation of the third step, in order to reduce the overhead for subsequent verification. As our experiments confirm (Sec. 5), we drastically improve verification performance over instrumenting all instructions.

4.2 Abstract event structures

As described in Sec. 3, we can choose to delay only one pair per critical cycle. To do so, all critical cycles need to be identified first. Sec. 2 defines cycles over events and event structures, which use concrete addresses and values, and thus correspond to concrete execution traces. As the enumeration of all traces is infeasible, we compute a conservative, over-approximate set of possible cycles using static analysis. In this program analysis we introduce *abstract events*, which summarise all concrete events that have the same process identifier, program counter, direction and memory address. We extend the definition of event structure to *abstract event structures*, which are identical except that they use abstract events.

Statements to abstract events The derivation of an abstract event structure from a non-branching multi-threaded program is straight-forward. For each thread, decompose each statement into abstract events, extracting all writes or reads of shared memory. For an assignment to a location designated by a pointer variable, consider the example $*(&x+z) = y;$, where $\&x$ denotes the address of x and $*p$ the value held at address p . We first read y , then read z and finally we write to the object pointed to by $\&x+z$, which is determined using an alias analysis⁵. If the precision of the alias analysis is insufficient to determine the object, we assume that this write can target any of the objects in the program.

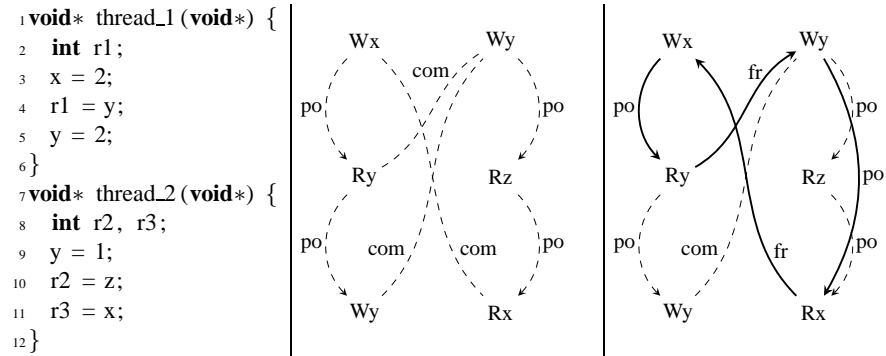


Figure 8. The program on the left contains an **sb** cycle (cf. Fig. 1). We build the abstract event graph in the middle, and indeed detect the cycle in the graph, on the right.

Abstract event graph In order to devise **SC** cycles that become critical cycles on a weaker architecture, we look for cycles in $ws \cup fr \cup rf \cup po$ (definition of **SC**, [5, Thm. 3]). Abstract events in each thread are ordered by program order, po , which we derive as described below. As we do not use concrete values, we compute over-approximations of the relations ws , rf and fr . We further abstract from directed edges and use undirected edges in these over-approximations. We call the abstract event structure equipped with over-approximations of ws , rf and fr an *abstract event graph*. We compute the over-approximations as follows:

- the internal rf , fr and ws pairs (relating two events on the same thread) are already covered by po edges;
- the external rf , fr and ws pairs (relating two events from different threads) are abstracted by undirected external communications, denoted by com , and relate any pair of write-read, read-write or write-write between two distinct threads.

Fig. 8 depicts this first step in the middle, which is the resulting abstract event graph of the program shown on the left-hand side. A concretisation of the abstract event graph may yield critical cycles. Fig. 8 shows an example of a critical cycle on the right-hand

⁵ The alias analysis we use is known to be sound for the weak architectures we consider [6].

side. Whether this cycle can be fully concretised to an execution witness, filling in concrete values in all abstract events, is left as task to a verification back end.

Control flow To build an abstract event graph for branching programs, we consider the if-then-else branches, loops and function calls. Functions are analysed as if they were inlined, thus recursion is not handled. For if-then-else, po in the abstract event graph follows both of the branches separately, and then joins at the end of the condition. For loops or backward jumps and given a pair $(x, y) \in \text{po}$, the back-edge may render x reachable from y as well. We thus include copies of x and y in the abstract event graph, such that (y, x) in po if such a back-edge exists. By [7] it suffices to use a single copy, as a critical cycle does not require more than two events in program order per thread.

The analysis proceeds in a forward manner along the control-flow graph of a given program. For each statement recorded in a node of the control-flow graph, the abstract events are computed. When preserved program order is defined via dp (cf. Sec. 2), possible dependencies between abstract events are recorded as well.

4.3 Detecting critical cycles

Given the abstract event graph of a program, we need to compute an over-approximate set of critical cycles. To increase scalability of this procedure, we first identify all strongly connected components (SCCs) in the graph using Tarjan’s 1972 algorithm [27], which is linear in the size of the abstract event graph. The detection of critical cycles can then be performed in parallel and independently for each SCC, as no cycle can span multiple SCCs. The SCCs also offer first insights about the program under test: two distinct SCCs will refer to two parts of the code that are independently accessing and updating shared memory.

Detecting all the critical cycles in an SCC Our cycle computation is based on Tarjan’s 1973 algorithm [28]. The abstract event graph, however, does not encode the transitive closure of po . Thus, we first extract *candidate cycles* by picking at most two abstract events per thread, which are guaranteed to be (transitively) linked by program order. For each candidate cycle we then perform additional filtering, as such a cycle need not be critical: a candidate is guaranteed to be *not* critical if it does not contain any *unsafe pair* for the given architecture, or is a cycle in *uniproc* or *thin-air*. All of these checks need to be performed a-posteriori for a complete cycle.

Tarjan’s original algorithm is worst-case exponential in the number of vertices (abstract events), and our subsequent filtering adds additional complexity. To deal with this complexity, we soundly limit the exploration using properties of critical cycles, such as all program-order pairs per address in a critical cycle being one of write-write, read-write, write-read or read-write-read [4].

4.4 Selecting and instrumenting delay pairs

The above cycle detection yields candidates for unsafe pairs of abstract events to be delayed in each cycle. Following Sec. 3.4, we instrument one pair to delay per cycle.

We may select these pairs arbitrarily, but we describe below a weighted instrumentation that decidedly reduces verification time, as we show in Sec. 5.

We first normalise the program such that all shared memory accesses appear in assignments only; any reads in branching conditions or function call parameters are moved to temporary variables as follows: $\mathbf{if}(\phi(x)) \dots; \mapsto \text{tmp} = \phi(x); \mathbf{if}(\text{tmp}) \dots;$ for an expression ϕ over a shared memory address x . In the following, we thus restrict ourselves to assignment statements.

For each memory address x of events in unsafe pairs we introduce an array $\mathbf{b}(x)$. In addition to the properties described in Sec. 3.1, we also keep track of the originating thread of the write to x . We introduce an additional pointer for each local variable reading from a shared memory address, i.e. an r such that $r = x$; In a pair to delay, in one of the critical cycles or after, we equip r with a pointer $\mathbf{rs}(r)$, which implements the read set of Sec. 3.1. We now describe the instrumentation of writes, then reads. To soundly over-approximate all possible behaviours, all instrumented operations are guarded by $\mathbf{if}(\ast)$, expressing non-deterministic choice.

Instrumenting writes We implement here the two operations associated to the weak-memory effects of a write w , as defined in Sec. 3.1: (1) delaying a write, $\mathbf{d}(w(w))$, by appending to the buffer, and (2) flushing a write, $\mathbf{f}(w(w))$, removing it from the buffer. A delayed write amounts to appending an element to the array:

$x = \text{smthg}; \mapsto \mathbf{if}(\ast) \mathbf{b}(x).\text{push}(\text{smthg}, \text{thread.number}); \mathbf{else} x = \text{smthg};$

According to Sec. 3.1, each delay is accompanied by a flush. Yet the point in time when the flush happens is not determined. We would thus need to add non-deterministic flush instructions at each statement in the program. This transformation would make the program highly non-deterministic, and very hard for a model checker to analyse. Therefore, we insert flushes only where they might have an effect, i.e. before each potential read from the address that was written to, and make them flush a non-deterministic number of writes in FIFO-manner. The function *take* implements the semantics of “write from buffer to memory” of Fig. 3 on C arrays for a non-deterministic number of elements, and returns the resulting in-memory value at address x .

$\text{smthg} = x; \mapsto \mathbf{if}(\ast) x = \mathbf{b}(x).\text{take}(\text{thread.number}); \text{smthg} = x;$

Instrumenting reads Here we are to implement the two operations for reads: delaying a read $\mathbf{d}(r(w, r))$ and reading from the set, $\mathbf{f}(r(w, r))$. We delay a read by recording the memory address to be read from. Note that, given our program normalisation, our reads manifest as assignments to local variables. For a local variable $r1$, we delay the read of x as follows:

$r1 = x; \mapsto \mathbf{if}(\ast) \mathbf{rs}(r1) = \&x; \mathbf{else} r1 = x;$

For flushing the read, considerations analogous to the write case are made: we flush non-deterministically upon an actual read (then of $r1$) only, instead of every program point. The flush dereferences the address previously recorded:

$r2 = r1; \mapsto \mathbf{if}(\mathbf{rs}(r1) \neq 0 \ \&\& \ \ast) \{ r1 = \ast\mathbf{rs}(r1); \mathbf{rs}(r1) = 0; \} r2 = r1;$

Input: the edges to instrument E , the cycles C_j
Problem: minimise $\sum_{e_i \in E} \mathbf{d}(e_i) * x_i$
s.t. $\forall j, \sum_{e_i \in C_j \cap E} x_i \geq 1$ (ensures soundness)
where
 e_i is a pair to potentially instrument,
 x_i is a Boolean variable stating whether we instrument e_i ,
and $\mathbf{d}()$ is the cost of an instrumentation.
Output: the x_i , stating which pairs to instrument

Figure 9. Mixed integer programming problem to choose the pairs to instrument

4.5 Weighted selection of unsafe pairs

Above, we selected an arbitrary unsafe pair per cycle, as this suffices to reveal all weak-memory effects (cf. Sec. 3). We do observe, however, that the choice of pairs has a strong effect on verification time. We thus assign an empirically devised cost \mathbf{d} to candidate pairs. With our implementation, we chose $\mathbf{d}(\text{poW}^*)=1$ (pairs in program order where the first event is a write), $\mathbf{d}(\text{poRW})=2$ (read-write pairs in program order), $\mathbf{d}(\text{rfe})=2$ (write-read pairs on different threads), $\mathbf{d}(\text{poRR})=3$ (read-read pairs in program order). Given a set E of pairs to delay in the graph with critical cycles C_j , we solve the mixed integer programming problem of Fig. 9. Our experiments show that this encoding yields a speedup of 26% over all architectures with an SC bounded model-checker.

5 Experimental Results

We exercised our method and measured its cost using 8 tools. We considered 5 ANSI-C model checkers: a bounded model checker based on CBMC; SatAbs, a verifier based on predicate abstraction, using Boom as the model checker for the Boolean program; ESBMC, a bounded model checker; Threader, a thread-modular verifier; and Poirot, which implements a context-bounded translation to sequential programs. These tools cover a broad spectrum of symbolic algorithms for verifying SC programs. We also experimented with Blender, CheckFence, and MMChecker. We ran our experiments on Linux 2.6.32 64-bit machines with 3.07 GHz (only Poirot was run on a Windows system). Further details on the results are available on our web page.

Validation First, we systematically validate our setup using 555 litmus tests exposing weak memory artefacts (e.g. instruction reordering, store buffering, write atomicity relaxation) in isolation. The diy tool automatically generates x86, Power and ARM assembly programs implementing an idiom that cannot be reached on SC, but can be reached on a given model. For example, **sb**

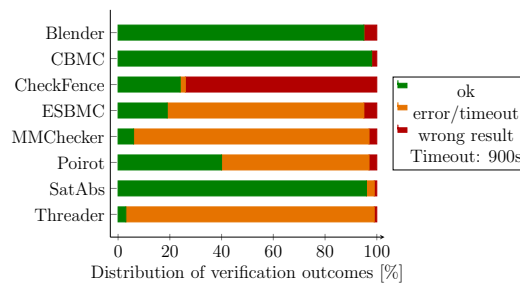
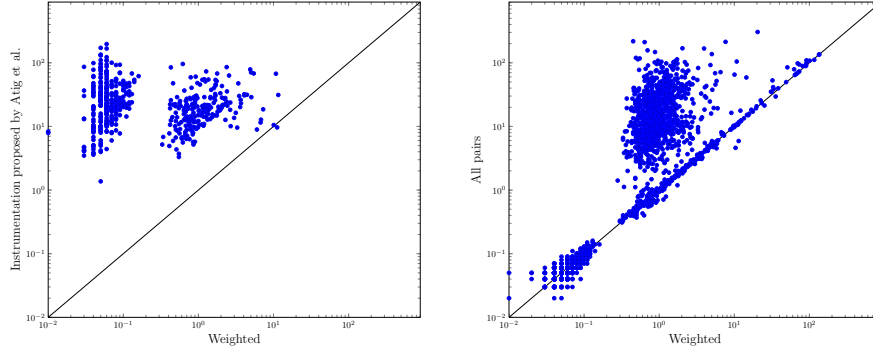


Figure 10. All tools on all litmus tests and models



(a) All accesses [11] vs. weighted selection (b) All pairs vs. weighted selection

Figure 11. Comparison of verification times of CBMC (seconds) for different instrumentations

(Fig. 1) exhibits store buffering, thus the final state can be reached on any weak model, from TSO to Power.

Each litmus test comes with an assertion that models the SC violation exercised by the test, e.g. the outcomes of Fig. 1 and 2. Thus, verifying a litmus test amounts to checking whether the model under scrutiny can reach the specified outcome. We then convert these tests automatically into C code, leading to programs of 48 lines on average, involving 2 to 4 threads.

These examples provide assurance that we soundly implement the theory of Sec. 3: we verify each test w.r.t. SC, i.e. without transformation, then w.r.t. TSO, PSO, RMO, and Power. Despite the tests being small, they provide challenging concurrent idioms to verify. Fig. 10 compares the tools on all tests and models. Most tools, with the exception of Blender, CBMC and SatAbs, time out or give wrong results on a vast majority of tests. Blender only expectedly fails on tests involving `lwsync` fences; CBMC and SatAbs return spurious results in 1.5% of the tests, caused by the over-approximation in the implementation of our instrumentation.

Fig. 11 compares the verification time using CBMC over all litmus families (e.g. `rfe` tests exercise store atomicity, `podwr` tests exercise the write-read reordering) for different instrumentation options. First, with the restriction to TSO, Fig. 11(a) compares the instrumentation of all shared memory accesses proposed in [11] to the weighted transformation (Sec. 4.5). On average, we observe a more than 300-fold speedup in verification time. In addition, the reduced instrumentation also yields 246 fewer spurious results. We also quantify the specific benefit of the weighted selection of pairs in Fig. 11(b). We compare the cost of the instrumentation of all pairs on critical cycles with that of the weighted transformation (Sec. 4.5) for all models, tools and tests. The average speedup over all models and tests is still more than one order of magnitude. We give the detailed results for all experiments online.

We also verified several TSO examples that have been used in the literature (details are online). Note that these examples in fact only exhibit idioms already covered by our litmus tests (e.g. Dekker corresponds to the `sb` test of Fig. 1). Furthermore, we applied the instrumentation to code taken from the Read-Copy-Update algorithm in the Linux

kernel and scheduling code in the Apache HTTP server, as well as industrial code from IBM. We observe that the instrumentation tool completes even on such code of up to 28,000 lines in less than 1 second, and in 32 seconds on IBM’s code. We now study one real-life example in detail, an excerpt of the relational database software PostgreSQL.

Worker Synchronization in PostgreSQL Mid 2011, PostgreSQL developers observed that a regression test occasionally failed on a multi-core PowerPC system.⁶ The test implements a protocol passing a token in a ring of processes. Further analysis drew the attention to an interprocess signalling mechanism. It turned out that the code had already been subject to an inconclusive discussion in late 2010.⁷

```

1 #define WORKERS 2
2 volatile _Bool latch[WORKERS];
3 volatile _Bool flag[WORKERS];
4 void worker(int i)
5 { while(! latch[i]);
6   for(;;)
7   { assert(! latch[i] || flag[i]);
8     latch[i] = 0;
9     if(flag[i])
10    { flag[i] = 0;
11      flag[(i+1)%WORKERS] = 1;
12      latch[(i+1)%WORKERS] = 1; }
13   while(! latch[i]); } }

```

Listing 1. Token passing in pgsq1.c

The code in Listing 1 is an inlined version of the problematic code, with an additional assertion in line 7. Each element of the array “latch” is a Boolean variable stored in shared memory to facilitate interprocess communication. Each working process waits to have its latch set and then expects to have work to do (from line 9 onwards). Here, the work consists of passing around a token via the array “flag”. Once the process is done with its work, it passes the token on (line 11), and sets the latch of the process the token was passed to (line 12).

Starvation seemingly cannot occur: when a process is woken up, it has work to do (has the token). Yet, the PostgreSQL developers observed that the wait in line 13 (which in the original code is bounded in time) would time out, thus signalling starvation of the ring of processes. The developers identified the memory model of the platform as possible culprit: it was assumed that the processor would at times delay the write in line 11 until after the latch had been set.

We transform the code of Listing 1 for two workers under Power. The event graphs show two idioms: **lb** (load buffering) and **mp** (message passing), in Fig. 12 and 13. The code fragments on the left-hand side give the corresponding line numbers in Listing 1.

The **lb** idiom contains the two *if* statements controlling the access to both critical sections. Since the **lb** idiom is yet unimplemented by Power machines (despite being allowed by the architecture [26]), we believe that this is not the bug observed by the PostgreSQL developers. Yet, it might lead to actual bugs on future machines.

In contrast, the **mp** case is commonly observed on Power machines (e.g. 1.7G/167G on Power 7 [26]). The **mp** case arises in the PostgreSQL code by the combination of some writes in the critical section of the first worker, and the access to the critical section of the second worker; the relevant code lines are in Fig. 13.

We first check the fully transformed code with SatAbs. After 21.34 seconds, SatAbs provides a counterexample (given online), where we first execute the first worker up to line 13. All accesses are w.r.t. memory, except at lines 11 and 12, where the values 0

⁶ <http://archives.postgresql.org/pgsql-hackers/2011-08/msg00330.php>

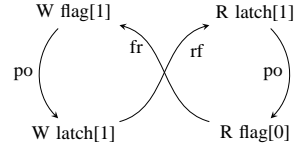
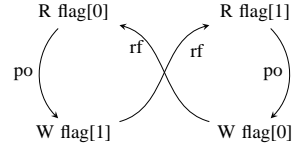
⁷ <http://archives.postgresql.org/pgsql-hackers/2010-11/msg01575.php>

pgsql (lb)	
Worker 0	Worker 1
(9) if (flag[0])	(9) if (flag[1])
(11) flag[1]=1;	(11) flag[0]=1;
Observed: flag[0]=1; flag[1]=1	

Figure 12. An lb idiom detected in `pgsql.c`

pgsql (mp)	
Worker 0	Worker 1
(11) flag[1]=1;	(5) while (!latch[1]);
(12) latch[1]=1;	(9) if (flag[1])
Observed: latch[1]=1; flag[1]=0	

Figure 13. An mp idiom detected in `pgsql.c`



and 1 are stored into the buffers of flag[0] and flag[1]. Then the second worker starts, reading the updated value 1 of latch[1]. It exits the blocking while (line 5) and reaches the assertion. Here, latch[1] still holds 1, and flag[1] still holds 0, as Worker 0 has not yet flushed the write waiting in its buffer. Thus, the condition of the *if* is not true, the critical section is skipped, and the program arrives at line 13, without having authorised the next worker to enter the critical section, and loops forever.

As **mp** can arise on Power e.g. because of non-atomic writes, we know by Sec. 3.4 that we only need to transform one **rfe** pair of the cycle, and relaunch the verification. SatAbs spends 1.29 seconds to check it (and finds a counterexample, as previously).

PostgreSQL developers discussed fixes, but only committed comments to the code base, as it remained unclear whether the intended fixes were appropriate. We proposed a provably correct patch solving both **lb** and **mp**. After discussion with the developers⁸, we improved it to meet the developers' desire to maintain the current API. The final patch introduces two `lwsync` barriers: after line 8 and before line 12.

6 Conclusion

We have presented a provably sound method to verify concurrent software w.r.t. weak memory. Our contribution allows to lift SC methods and tools to a wide range of weak memory models (from x86 to Power), by means of program transformation.

Our approach crucially relies on the definition of a generic operational model equivalent to the axiomatic one of [8]. We do not favour any style of model in particular, but we highlight the importance of the availability of several equivalent mathematical styles to model semantics as intricate as weak memory. In addition, operational models are often the style of choice in the verification community; we contribute here to the vocabulary to tackle the verification problem w.r.t. weak memory.

Our extensive experiments and in particular the PostgreSQL bug demonstrate the practicability of our approach from several different perspectives. First, we confirmed a known bug (**mp**), and validated the fix proposed by the developers, including an evaluation of different synchronisation options. Second, we found an additional idiom (**lb**), which will cause a bug on future Power machines; our fix repairs it already.

⁸ <http://archives.postgresql.org/pgsql-hackers/2012-03/msg01506.php>

References

1. <http://research.microsoft.com/en-us/projects/poirot/>
2. Abdulla, P., Atig, M.F., Chen, Y., Leonardsson, C., Rezine, A.: Counter-example guided fence insertion under TSO. In: TACAS (2012)
3. Adve, S.V., Gharachorloo, K.: Shared Memory Consistency Models: A Tutorial. *IEEE Computer* 29, 66–76 (1995)
4. Alglave, J.: A Shared Memory Poetics. Ph.D. thesis, Université Paris 7 and INRIA (2010)
5. Alglave, J.: A Formal Hierarchy of Weak Memory Models. In: FMSD (2012)
6. Alglave, J., Kroening, D., Lugton, J., Nimal, V., Tautschnig, M.: Soundness of data flow analyses for weak memory models. In: APLAS (2011)
7. Alglave, J., Maranget, L.: Stability in weak memory models. In: CAV (2011)
8. Alglave, J., Maranget, L., Sarkar, S., Sewell, P.: Fences in Weak Memory Models. In: CAV (2010)
9. Atig, M.F., Bouajjani, A., Burckhardt, S., Musuvathi, M.: On the verification problem for weak memory models. In: POPL (2010)
10. Atig, M.F., Bouajjani, A., Burckhardt, S., Musuvathi, M.: What’s decidable about weak memory models? In: ESOP (2012)
11. Atig, M.F., Bouajjani, A., Parlato, G.: Getting rid of store-buffers in the analysis of weak memory models. In: CAV (2011)
12. Bouajjani, A., Meyer, R., Moehlmann, E.: Deciding robustness against total store ordering. In: ICALP (2011)
13. Burckhardt, S., Alur, R., Martin, M.K.: Checkfence: Checking consistency of concurrent data types on relaxed memory models. In: PLDI (2007)
14. Cordeiro, L., Fischer, B.: Verifying multi-threaded software using SMT-based context-bounded model checking. In: ICSE. pp. 331–340. ACM (2011)
15. Donaldson, A., Kaiser, A., Kroening, D., Wahl, T.: Symmetry-aware predicate abstraction for shared-variable concurrent programs. In: CAV (2011)
16. Gupta, A., Popeea, C., Rybalchenko, A.: Threader: A constraint-based verifier for multi-threaded programs. In: CAV (2011)
17. Huynh, T., Roychoudhury, A.: A memory model sensitive checker for C#. In: FM (2006)
18. Jin, H., Yavuz-Kahveci, T., Sanders, B.A.: Java memory model-aware model checking. In: TACAS (2012)
19. Kuperstein, M., Vechev, M., Yahav, E.: Automatic inference of memory fences. In: FMCAD (2010)
20. Kuperstein, M., Vechev, M., Yahav, E.: Partial-Coherence Abstractions for Relaxed Memory Models. In: PLDI (2011)
21. Lamport, L.: How to Make a Correct Multiprocess Program Execute Correctly on a Multiprocessor. *IEEE Trans. Comput.* 46(7), 779–782 (1979)
22. Linden, A., P.Wolper: A verification-based approach to memory fence insertion in relaxed memory systems. In: SPIN (2011)
23. Owens, S.: Reasoning about the Implementation of Concurrency Abstractions on x86-TSO. In: ECOOP (2010)
24. Owens, S., Sarkar, S., Sewell, P.: A better x86 memory model: x86-TSO. In: TPHOL (2009)
25. Park, S., Dill, D.: An executable specification, analyzer and verifier for RMO. In: SPAA (1995)
26. Sarkar, S., Sewell, P., Alglave, J., Maranget, L., Williams, D.: Understanding Power multi-processors. In: PLDI (2011)
27. Tarjan, R.: Depth-first search and linear graph algorithms. *SIAM J. Comput.* (1972)
28. Tarjan, R.: Enumeration of the elementary circuits of a directed graph. *SIAM J. Comput.* (1973)
29. Yang, Y., Gopalakrishnan, G., Lindstrom, G.: Memory model sensitive data race analysis. In: ICFEM (2004)