

# Persistent Owicki-Gries Reasoning

A Program Logic for Reasoning about Persistent Programs on Intel-x86

AZALEA RAAD, MPI-SWS and Imperial College London, Germany and United Kingdom

ORI LAHAV, Tel Aviv University, Israel

VIKTOR VAFEIADIS, MPI-SWS, Germany

The advent of non-volatile memory (NVM) technologies is expected to transform how software systems are structured fundamentally, making the task of *correct* programming significantly harder. This is because ensuring that memory stores persist in the correct order is challenging, and requires low-level programming to flush the cache at appropriate points. This has in turn resulted in a noticeable *verification gap*.

To address this, we study the verification of NVM programs, and present *Persistent Owicki-Gries* (POG), the first program logic for reasoning about such programs. We prove the soundness of POG over the recent Intel-x86 model, which formalises the out-of-order persistence of memory stores and the semantics of the Intel cache line flush instructions. We then use POG to verify several programs that interact with NVM.

CCS Concepts: • **Theory of computation** → **Concurrency; Semantics and reasoning**; • **Software and its engineering** → **General programming languages**.

Additional Key Words and Phrases: non-volatile memory, program logic, x86-TSO, consistency, persistency

## ACM Reference Format:

Azalea Raad, Ori Lahav, and Viktor Vafeiadis. 2020. Persistent Owicki-Gries Reasoning: A Program Logic for Reasoning about Persistent Programs on Intel-x86. *Proc. ACM Program. Lang.* 4, OOPSLA, Article 151 (November 2020), 28 pages. <https://doi.org/10.1145/3428219>

## 1 INTRODUCTION

The emergence of *non-volatile memory* (NVM) technologies [Kawahara et al. 2012; Lee et al. 2009; Strukov et al. 2008] is expected to revamp the structure of modern software: NVM provides storage persistency across power failures with performance close to that of traditional (volatile) memory. As such, programs that require persistency of their data (e.g. databases) can achieve orders of magnitude lower latency by storing their data on NVM rather than on hard disks. It is therefore believed that NVM (a.k.a. persistent memory) will supplant RAM in the near future, thanks to its durable yet competitive performance. This belief is backed by widespread industrial support. Specifically, the two major architectures, ARMv8 and Intel-x86 which together account for almost 100% of the desktop and mobile market, have extended their official specifications to support persistent programming [Arm 2018; Intel 2019]. Intel has further (1) manufactured its own line of NVM, Optane technology [Intel 2019], with an extended academic study evaluating its performance [Izraelevitz et al. 2019]; and (2) released open-source NVM libraries such as PMDK [Intel 2015].

To describe the behaviour of programs under NVM, Intel has introduced a *persistency model* for their x86 architecture [Intel 2019], describing the order in which memory stores may persist to

---

Authors' addresses: Azalea Raad, MPI-SWS and Imperial College London, Germany and United Kingdom, [azalea@imperial.ac.uk](mailto:azalea@imperial.ac.uk); Ori Lahav, Tel Aviv University, School of Computer Science, 69978, Tel Aviv, Israel, [orilahav@tau.ac.il](mailto:orilahav@tau.ac.il); Viktor Vafeiadis, MPI-SWS, Saarland Informatics Campus, Germany, [viktor@mpi-sws.org](mailto:viktor@mpi-sws.org).

---

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

© 2020 Copyright held by the owner/author(s).

2475-1421/2020/11-ART151

<https://doi.org/10.1145/3428219>

NVM. This model was formalised by Raad et al. [2020], wherein they extended the x86-TSO model [Sewell et al. 2010] (with thread-local buffers to model the delayed propagation of writes to other threads) with a global persistent buffer to model the out-of-order propagation of stores to NVM.

Although NVM research has grown rapidly in recent years in both persistency semantics [Condit et al. 2009; Gogte et al. 2018; Izraelevitz et al. 2016; Joshi et al. 2015; Kolli et al. 2017, 2016; Raad and Vafeiadis 2018; Raad et al. 2020, 2019] and algorithms/libraries that exploit NVM [Friedman et al. 2018; Nawab et al. 2017; Zuriel et al. 2019], there has been little work on *verifying* such artifacts. To our knowledge, the existing work [Friedman et al. 2018; Nawab et al. 2017; Raad and Vafeiadis 2018; Raad et al. 2020, 2019; Zuriel et al. 2019] offer low-level proofs about the correctness of small persistent algorithms and often make simplifying assumptions. In particular, they all work at the level of *traces* rather than at the level of *program syntax*, while [Friedman et al. 2018; Nawab et al. 2017; Zuriel et al. 2019] further assume sequential consistency as their concurrency model.

This is a significant omission because developing correct persistent data structures is rather error-prone. Since memory stores are typically persisted out of order, one has to use special low-level instructions for flushing the cache in order to ensure a correct persist ordering. Moreover, such persistent implementations are virtually impossible to test and debug, as one would have to use custom hardware to simulate crashes and check correct recovery from them.

To close this gap, we consider the *formal verification* of (multi-threaded) programs running over NVM. To this end, we adapt the well-known Owicki-Gries (OG) proof system [Owicki and Gries 1976], and develop POG (Persistent OG), the *first* program logic for reasoning about NVM programs. We show that the POG proof system is *sound* with respect to the Intel-x86 persistency semantics.

We develop POG over Intel-x86 for several reasons. First, Intel-x86 is a ubiquitous architecture with a formally-defined persistency model [Raad et al. 2020]. Indeed, excluding academic models proposed as proofs of concept, the only real-world persistency models currently available are those of low-level architectures, i.e. those of ARMv8 and Intel-x86 [Raad et al. 2020, 2019], and no existing mainstream programming language such as C/C++ currently has a formal persistency model. Second, existing research on language-level persistency models (e.g. [Gogte et al. 2020; Kolli et al. 2017]) suggests that similar persistency primitives to those of Intel-x86 are considered at the language level, and will likely be lifted to higher-levels. As such, the reasoning principles of POG will be useful in the future higher-level persistency models. Lastly, POG presents the first formal framework for reasoning about persistency primitives and their behaviour *abstractly*, i.e. at the program syntax level, rather than delving into all possible program executions and reasoning at the trace level. Given the complexity of the Intel-x86 persistency model, even verifying simple examples is non-trivial, especially when done at the trace level, and we believe that our syntax-level proof rules in POG help simplify such proofs significantly.

**Challenges.** Developing the reasoning principles of POG over the Intel-x86 persistency model involves two main challenges: (1) dealing with weak memory consistency, i.e. with the thread-local FIFO buffers of Intel-x86; and (2) dealing with weak persistency, i.e. the persistent buffer of Intel-x86, which allows for stores to persist *asynchronously* and *out of order*.

To address the first challenge, we base our program logic on a variant of OG, named OGRA [Lahav and Vafeiadis 2015], proved sound under release-acquire consistency (a memory model weaker than x86-TSO) and thus also sound under the Intel-x86 model as far as consistency is concerned.

To address the second challenge, we develop an *intermediate* operational model of Intel-x86 persistency,  $\text{Ix86}_{\text{sim}}$ , which forgoes the persistent buffer altogether and operates on the original x86-TSO model (i.e. with only the thread-local buffers). We show that our  $\text{Ix86}_{\text{sim}}$  model correctly captures the effect of the Intel-x86's persistency buffer. That is, the possible outcomes of a program under the Intel-x86 persistency model with two types of buffers are the same as those under  $\text{Ix86}_{\text{sim}}$ .

Our next challenge in designing the POG reasoning principles is modelling the behaviour of explicit persist instructions on Intel-x86, **flush** and **flush<sub>opt</sub>**, which when executed persist all pending writes on a given cache line to the memory. As we describe shortly in §2, **flush<sub>opt</sub>** instructions offer weaker ordering constraints and may be reordered with respect to other instructions more freely. As such, their effect may not take place at the intended program point, making it more difficult to reason about their persistency behaviour. To keep the POG reasoning principles simple, we devise POG to focus only on the stronger **flush** instructions. We then provide a mechanism to extend our POG reasoning to weaker **flush<sub>opt</sub>** instructions. More concretely, we present a transformation that allows us in most cases to rewrite a program using **flush<sub>opt</sub>** to an *equivalent* program using **flush**. We can then use POG to reason about such programs by first using our transformation to replace **flush<sub>opt</sub>** with **flush**, and then using POG to verify the transformed program.

Although our main contribution is POG, we remark that our  $\text{Ix86}_{\text{sim}}$  model and our transformation are also valuable contributions *per se*. Specifically,  $\text{Ix86}_{\text{sim}}$  may serve as the input to automated verification tools for concurrency, e.g. model checkers, especially those that already support x86-TSO [Abdulla et al. 2015; Clarke et al. 2004; Huang and Huang 2016]. Our persistency-preserving program transformation can be used to optimise code (e.g. at compile time) to replace **flush** with **flush<sub>opt</sub>**.

**Contributions and Outline.** Our contributions (detailed in §2) are as follows: (1) in §3 we present POG, the *first* program logic for verifying persistency guarantees; (2) in §4 we use POG to verify several examples; (3) in §5 we present our  $\text{Ix86}_{\text{sim}}$  model and show that it faithfully captures Intel-x86 persistency; (4) in §6 we show the POG is sound with respect to the  $\text{Ix86}_{\text{sim}}$  model; (5) in §7 we present our transformation for rewriting programs with **flush<sub>opt</sub>** to equivalent ones with **flush**, and show that this transformation is sound. We discuss related and future work in §8.

**Additional Material.** The proofs of all theorems stated in this article are given in full in the accompanying technical appendix available at <http://plv.mpi-sws.org/pog/>.

## 2 OVERVIEW

### 2.1 Px86<sub>sim</sub> at a Glance

Memory *consistency* models typically describe the permitted behaviours of programs by constraining the *volatile memory order*, i.e. the order in which memory writes are made visible to other threads. Analogously, memory *persistency* models [Pelley et al. 2014] describe the permitted behaviours of programs upon recovering from a crash (e.g. a power failure) by defining a *persistent memory order*, i.e. the order in which writes are committed to persistent memory. To distinguish between the two memory orders, memory *stores* are differentiated from memory *persists*. The former denotes the process of making a write visible to other threads, whilst the latter denotes the process of committing writes durably to persistent memory.

Raad et al. [2020] recently developed the Px86 (‘persistent x86’) models, formalising the persistency semantics of the Intel-x86 architecture. As they noted, the Intel manual [Intel 2019] is ambiguous at times and allows for weaker behaviours than originally intended. They thus formulated two persistency models: (1)  $\text{Px86}_{\text{man}}$ , which reflects the behaviour outlined in the manual; and (2)  $\text{Px86}_{\text{sim}}$ , which is a strengthening of  $\text{Px86}_{\text{man}}$  and captures the architectural intent. As  $\text{Px86}_{\text{sim}}$  reflects the architectural intent, in this article we focus on the  $\text{Px86}_{\text{sim}}$  model.

The  $\text{Px86}_{\text{sim}}$  model follows a *buffered, relaxed* persistency model. Under a buffered model, memory persists occur *asynchronously* [Condit et al. 2009; Izraelevitz et al. 2016; Joshi et al. 2015]: they are buffered in a queue to be committed to persistent memory at a future time. This way, persists occur after their corresponding stores and as prescribed by the persistent memory order, while allowing the execution to proceed ahead of persists. As such, after recovering from a crash, only a *prefix*

$x := 1;$ $y := 1;$ (a)	$x := 1;$ <b>flush</b> $x'$ ; $y := 1;$ (b)	$x := 1;$ <b>flush<sub>opt</sub></b> $x'$ ; $y := 1;$ (c)	$x := 1;$ <b>flush<sub>opt</sub></b> $x'$ ; <b>sfence</b> ; $y := 1;$ (d)	$x := 1;$ <b>flush</b> $x'$ ; $y := 1;$	$a := y;$ <b>if</b> ( $a=1$ ) $z := 1;$ (e)
$\checkmark: x, y \in \{0, 1\}$	$\checkmark: y=1 \Rightarrow x=1$	$\checkmark: x, y \in \{0, 1\}$	$\checkmark: y=1 \Rightarrow x=1$	$\checkmark: z=1 \Rightarrow x=1$	

Fig. 1. Example P<sub>x86<sub>sim</sub></sub> programs and possible  $x, y, z$  values upon recovery; in all examples  $x, y, z$  are locations in persistent memory,  $a$  is a (local) register,  $x, x' \in X$  ( $x, x'$  are in cache line  $X$ ),  $y, z \notin X$ , and initially  $x=y=z=0$ . Replacing the **sfence** instruction in (d) with **mfence** or an atomic RMW yields the same result. Similarly, replacing **flush**  $x'$  in (e) with **flush<sub>opt</sub>**  $x'$ ;  $c$  yields the same result when  $c$  is an **sfence/mfence** or an RMW.

of the persistent memory order may have successfully persisted. Under relaxed persistency, the volatile and persistent memory orders may disagree: the order in which the writes are made visible to other threads may differ from the order in which they are persisted.

The relaxed and buffered persistency of P<sub>x86<sub>sim</sub></sub> is demonstrated in Fig. 1a. If a crash occurs during (or after) the execution of this program, at crash time either write may or may not have persisted and thus  $x, y \in \{0, 1\}$  upon recovery. The relaxed nature of P<sub>x86<sub>sim</sub></sub> allows for surprising behaviours that are not possible during normal (non-crashing) executions. In particular, the two writes cannot be reordered under Intel-x86 and thus at no point during the normal execution of this program  $x=0, y=1$  is observable. Nevertheless, in case of a crash it is possible under P<sub>x86<sub>sim</sub></sub> to observe  $x=0, y=1$  after recovery. This is due to the relaxed persistency of P<sub>x86<sub>sim</sub></sub>: the store order ( $x$  before  $y$ ) is separate from the persist order ( $y$  before  $x$ ). Under the P<sub>x86<sub>sim</sub></sub> model the writes may be persisted (1) in any order, when they are on distinct locations; or (2) in the volatile memory order, when they are on the same location. That is, for each location, its store and persist orders coincide.

Intel-x86 provides explicit *persist* instructions, **flush**  $x$ , **flush<sub>opt</sub>**  $x$  and **wb**  $x$ , in order to afford more control over when pending writes are persisted. When executed, these instructions *asynchronously* persist all pending writes on all locations in the cache line of  $x$  [Intel 2019]. That is, when location  $x$  is in cache line  $X$ , written  $x \in X$ , an explicit persist on  $x$  persists all pending writes on all locations  $x' \in X$ . As noted by Raad et al. [2020], **flush** instructions are the strongest of the three in terms of their constraints on instruction reordering, whereas **flush<sub>opt</sub>** and **wb** are equally weak and have the same specification, with **wb** providing better performance than **flush<sub>opt</sub>**. That is, **flush<sub>opt</sub>** and **wb** are indistinguishable under P<sub>x86<sub>sim</sub></sub>. As such, in the remainder of our discussion we focus on **flush** and **flush<sub>opt</sub>** instructions and describe their behaviour via several examples.

The persistency behaviour of **flush** is illustrated in Fig. 1b: given  $x, x' \in X$ , executing **flush**  $x'$  persists the earlier write on  $X$  (i.e.  $x := 1$ ). As such, if a crash occurs during the execution of this program and  $y=1$  upon recovery, then  $x=1$ . That is, if  $y := 1$  has executed and persisted before the crash, then so must the earlier  $x := 1$ ; **flush**  $x'$ . This is guaranteed by the ordering constraints on **flush**: **flush** instructions are ordered with respect to both earlier (in program order) and later writes. Hence, **flush**  $x'$  in Fig. 1b cannot be reordered with respect to  $x := 1$  or  $y := 1$ . As such, upon recovery  $y=1 \Rightarrow x=1$ . Note that **flush**  $x$  persists  $X$  *asynchronously*: its execution does not block until  $X$  is persisted; rather, the execution proceeds and  $X$  is made persistent at a future point.

In contrast to **flush**, **flush<sub>opt</sub>** instructions provide weaker ordering guarantees in relation to writes, in that they are only ordered with respect to *earlier writes* on the *same* cache line. This is illustrated in the example of Fig. 1c obtained from that in Fig. 1b by replacing **flush**  $x'$  with **flush<sub>opt</sub>**  $x'$ . Unlike in Fig. 1b, the **flush<sub>opt</sub>**  $x'$  instruction is not ordered with respect to the later write

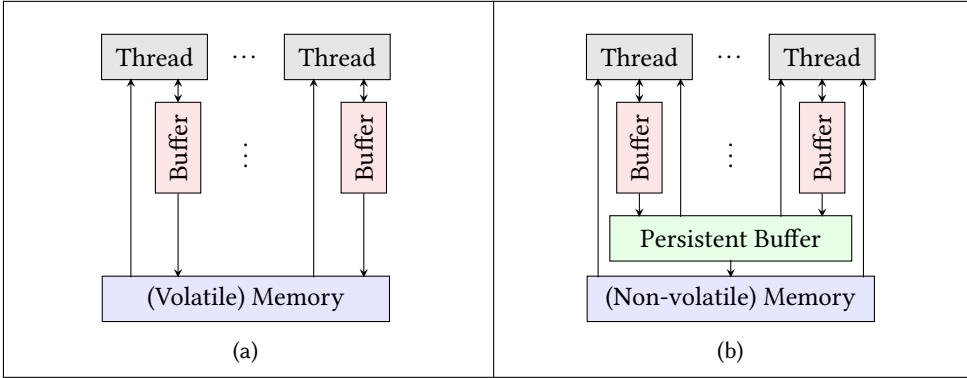


Fig. 2. Storage subsystems of x86-TSO (a) and Px86<sub>sim</sub> (b) as depicted in [Raad et al. 2020]

( $y := 1$ ) and may thus be reordered after it. As such,  $\mathbf{flush}_{\text{opt}} x'$  may not execute at the intended program point (after  $x := 1$  and before  $y := 1$ ) and thus may not guarantee the intended persist ordering. That is, unlike in Fig. 1b, the  $\mathbf{flush}_{\text{opt}} x'$  instructions does not guarantee that the  $x := 1$  write persists before  $y := 1$ , and is thus possible to observe  $x=0 \wedge y=1$  upon recovery.

In order to prevent such reorderings and to strengthen the ordering constraints between  $\mathbf{flush}_{\text{opt}}$  and later instructions, one can use either *fence* instructions, namely **sfence** (store fence) and **mfence** (memory fence), or atomic *read-modify-write* (RMW) instructions such as compare-and-set (CAS) and fetch-and-add (FAA). More concretely, **sfence**, **mfence** and RMW instructions are ordered with respect to all (both earlier and later)  $\mathbf{flush}_{\text{opt}}$ ,  $\mathbf{flush}$  and write instructions, and can be used to prevent reorderings such as that in Fig. 1c. This is illustrated in the example of Fig. 1d obtained from Fig. 1c by inserting an **sfence** after  $\mathbf{flush}_{\text{opt}}$ . Unlike in Fig. 1c, the intervening **sfence** ensures that  $\mathbf{flush}_{\text{opt}}$  in Fig. 1d is ordered with respect to  $y := 1$  and cannot be reordered after it, thus ensuring that  $x := 1$  persists before  $y := 1$  (i.e.  $y=1 \Rightarrow x=1$  upon recovery), as in Fig. 1b.

The example in Fig. 1e illustrates how persist orderings can be imposed on the writes of different threads using *message passing*. Note that the program in the left thread of Fig. 1e is that of Fig. 1b. A message is passed from thread  $\tau_1$  to  $\tau_2$  when  $\tau_2$  reads a value written by  $\tau_1$ . For instance, if the right thread in Fig. 1e reads 1 from  $y$  (written by the left thread), then the left thread passes a message to the right thread. Under Intel-x86 message passing ensures that the instruction writing the message and all those ordered before it (e.g.  $x := 1$ ;  $\mathbf{flush} x$ ;  $y := 1$ ) are executed (ordered) before the instruction reading it (e.g.  $a := y$ ). As such, since  $x := 1$ ;  $\mathbf{flush} x'$  is executed before  $a := y$ , and  $z := 1$  is executed after  $a := y$  when  $a=1$ , we know  $x := 1$ ;  $\mathbf{flush} x'$  is executed before  $z := 1$ . Consequently, if upon recovery  $z=1$  (i.e.  $z := 1$  has persisted before the crash), then  $x=1$  ( $x := 1$ ;  $\mathbf{flush} x'$  must have also persisted before the crash). As before, replacing  $\mathbf{flush} x'$  in Fig. 1e with the  $\mathbf{flush}_{\text{opt}} x'$ ;  $c$  sequence yields the same result upon recovery when  $c$  is an **sfence/mfence** or an RMW instruction.

Lastly, observe that  $\mathbf{flush}/\mathbf{flush}_{\text{opt}}$  instructions impose a particular *persist ordering*. Given  $x \in X$ , all writes on  $X$  ordered before  $\mathbf{flush} x/\mathbf{flush}_{\text{opt}} x$  persist before all instructions (regardless of their cache line) ordered after  $\mathbf{flush} x$ . For instance, since  $x := 1$  in Fig. 1b is ordered before  $\mathbf{flush} x'$ , and  $\mathbf{flush} x'$  is ordered before  $y := 1$ , the  $x := 1$  write is guaranteed to persist before  $y := 1$ . Similarly, as  $x := 1$  in Fig. 1d is ordered before  $\mathbf{flush}_{\text{opt}} x'$  which is in turn ordered before  $y := 1$  (thanks to the intervening **sfence**), the  $x := 1$  write is guaranteed to persist before  $y := 1$ .

**The Operational Px86<sub>sim</sub> Model.** Raad et al. [2020] developed their operational Px86<sub>sim</sub> model as an extension of the x86-TSO model by Sewell et al. [2010]. As illustrated in Fig. 2a, each thread in x86-TSO is connected to the (volatile) memory via a FIFO buffer. When a thread writes value  $v$  to location  $x$ , it records it in its buffer as the  $\langle x, v \rangle$  entry. When a thread reads from  $x$ , it first consults its own buffer. If it contains buffered writes for  $x$ , the thread reads the last such buffered write; otherwise, it consults the memory. Threads can debuffer their writes by propagating them (in FIFO order) to the memory at non-deterministic times. Additionally, the execution of a memory fence **mfence** drains the buffer of the executing thread.

To model the buffered persists of Px86<sub>sim</sub>, Raad et al. [2020] extended the x86-TSO storage subsystem with a *persistent buffer* as depicted in Fig. 2b, containing those writes that are pending to be persisted to the (non-volatile) memory. As with the memory, the persistent buffer is accessible by all threads. However, while the memory is non-volatile, the persistent buffer is volatile and its contents are lost upon a crash. When writes in the thread-local buffer are debuffered, they are propagated to the persistent buffer, denoting the store of the write (i.e. when the write becomes visible to other threads). Pending writes in the persistent buffer are in turn debuffered and propagated to the memory at non-deterministic times, denoting the persist of the write (i.e. when the write is committed durably to memory). The execution of reads accordingly traverses this hierarchy: when reading from  $x$ , the thread first inspects its own local buffer for the last write to  $x$  when such a write exists; otherwise, it consults the persistent buffer for the last store to  $x$  if such a store exists; otherwise, it reads  $x$  from the memory. Recall that the writes on distinct locations may persist in any order, whereas the writes on the same location persist in the store order. To capture this, the persistent buffer is modelled as a queue, where the pending writes on each location are propagated in the FIFO queue order, while those on different locations are propagated in an arbitrary order.

## 2.2 Eliminating **flush<sub>opt</sub>** Instructions via Program Transformation

As discussed above, the **flush<sub>opt</sub>** instructions provide weaker ordering constraints and can be reordered e.g. with respect to writes on different cache lines. While this flexibility may in certain cases lead to better performance by affording the compiler more optimisation opportunities through reordering, it significantly complicates the task of reasoning about persistency behaviours. In particular, the weak ordering constraints on **flush<sub>opt</sub>** can lead to unintended persistency behaviours such as that in Fig. 1c (where it is possible to observe  $x=0 \wedge y=1$  upon recovery), and to ensure correct persistency one must thus account for all possible such reorderings. It is therefore simpler to limit our reasoning to programs that solely use the stronger **flush** instructions.

As such, in order to support reasoning about the weaker **flush<sub>opt</sub>** instructions while simultaneously keeping our reasoning principles simple, we first (1) devise a mechanism that in most cases allows us to transform programs using **flush<sub>opt</sub>** to *equivalent* programs that *only* use **flush**; and then (2) design our POG reasoning principles for programs that only use **flush** instructions.

More concretely, for step (1) we note that the main use-case of **flush<sub>opt</sub>** (in which using **flush<sub>opt</sub>** rather than **flush** may prove advantageous for performance) prescribes a particular programming *pattern*. We then show that given a program  $C$  that uses **flush<sub>opt</sub>** in this pattern, one can transform  $C$  to a program  $C'$  that uses *only* **flush** instructions, such that  $C$  and  $C'$  have *equivalent persistency behaviours*, in that they yield the same values for all memory locations upon crash recovery.

Note that our intent through this transformation is not to forgo **flush<sub>opt</sub>** instructions altogether; rather, this transformation merely allows us to extend our reasoning to programs that use **flush<sub>opt</sub>** by considering equivalent programs using **flush**, while keeping our reasoning principles simple.

We present the formal details of this transformation in §7; in the remainder of this section and in §3-§6 we thus focus on programs using only **flush** (and not **flush<sub>opt</sub>**) instructions.

### 2.3 An Intermediate Operational Model for $Px86_{sim}$

Our goal in this paper is to devise a program logic for reasoning about the persistency behaviour of programs under  $Px86_{sim}$ . Although program logics are typically built over operational models that manipulate the underlying state, such operational models often operate on states that comprise the memory alone, without intermediate caches such as those of thread-local and persistent buffers in  $Px86_{sim}$ . This is because a large number of such logics operate under *sequential consistency* (SC) [Lamport 1979], while the presence of such caches introduces weak behaviours absent under SC. To remedy this, recent research [Lahav and Vafeiadis 2015; Sieczkowski et al. 2015; Svendsen et al. 2018; Turon et al. 2014; Vafeiadis and Narayan 2013] demonstrates how to reason about the weak behaviours introduced by e.g. thread-local buffers (see §2.4). However, no existing work currently supports the challenging task of reasoning about the persistency behaviour of programs. The difficulty of such reasoning is further compounded when considering the buffered behaviour of persists due to e.g. the persistent buffer of  $Px86_{sim}$ .

To streamline the task of devising a program logic for  $Px86_{sim}$ , we first (1) develop an *intermediate* operational semantics,  $Ix86_{sim}$ , that forgoes the persistent buffer, while emulating all valid  $Px86_{sim}$  behaviours; and then (2) devise a program logic for persistency reasoning over  $Ix86_{sim}$ . We proceed with an intuitive account of our  $Ix86_{sim}$  model; we briefly describe our program logic later in §2.4.

As discussed above, in our  $Ix86_{sim}$  model we forgo the persistent buffer altogether, thus operating on the x86-TSO storage system in Fig. 2a, with the volatile memory replaced with a non-volatile one. For simplicity, let us begin by assuming that **flush** instructions are executed *synchronously*. We later lift this assumption and describe how we handle the asynchronous behaviour of **flush**.

Recall that under  $Px86_{sim}$  the store and persist orders may disagree; i.e. the order in which writes in thread buffers are debuffed may differ from the order in which they are debuffed from the persistent buffer. As such, when forgoing the persistent buffer, additional care is required to preserve such weak behaviours. To see this, let us return to Fig. 1a, where  $x := 1$  is always store-ordered before  $y := 1$ , while  $y := 1$  may be persist-ordered before  $x := 1$ . We can then model the store order as in x86-TSO: upon executing each write the thread adds it to its local buffer, and non-deterministically debuffs its entries in the FIFO order, thus ensuring  $x := 1$  is store-ordered before  $y := 1$ . However, without the additional persistent buffer, we can no longer model the out-of-order persists.

To remedy this, for each location  $x$  we record two versions: (1) the ‘*volatile*’ version, written  $x_v$ , tracking the latest observable value of  $x$ ; and (2) the ‘*synchronously-persisted*’ (‘synchronous’) version, written  $x_s$ , tracking the latest persisted value of  $x$  provided that **flush** instructions are executed synchronously. Memory instructions (e.g. writes) are then carried out on volatile versions, leaving the synchronous versions untouched. Moreover, the volatile versions may non-deterministically propagate to the corresponding synchronous versions, modelling the notion that writes may be committed to memory at non-deterministic times. Similarly, since we assume **flush** instructions execute synchronously, given  $x \in X$ , executing **flush**  $x$  copies  $x'_v$  to  $x'_s$ , for all  $x' \in X$ .

Let us return to Fig. 1a and write  $x=v$  to denote that the latest value observable for  $x$  is  $v$ ; i.e. either the thread buffer contains no  $x$  entries and the value of  $x$  in memory is  $v$ , or the latest  $x$  entry in the thread buffer is  $\langle x, v \rangle$ . Similarly, let us write  $x \in \{v_1, v_2\}$  for  $x=v_1 \vee x=v_2$ . Therefore:

- (i) we assume that initially  $x_v=x_s=y_v=y_s=0$ ;
- (ii) after executing  $x := 1$  we have:  $x_v=1 \wedge x_s \in \{0, 1\} \wedge y_v=y_s=0$ ;
- (iii) upon subsequently executing  $y := 1$  we have:  $x_v=1 \wedge x_s \in \{0, 1\} \wedge y_v=1 \wedge y_s \in \{0, 1\}$ .

Note that since initially  $x_s=0$  and the value of  $x_v$  may be copied to  $x_s$  non-deterministically, we must account for this propagation in (ii) and thus we have  $x_s \in \{0, 1\}$ ; similarly for  $y_s$  in (iii). As

such, at all program points ((i)–(iii)) we have  $x_s, y_s \in \{0, 1\}$ . That is, if a crash occurs at any point, upon recovery we have  $x_s, y_s \in \{0, 1\}$ , thus emulating the desired behaviour in Fig. 1a.

We can analogously emulate the behaviour of Fig. 1b:

- (a) we assume that initially  $x_v=x_s=y_v=y_s=0$ ;
- (b) after executing  $x := 1$  we have:  $x_v=1 \wedge x_s \in \{0, 1\} \wedge y_v=y_s=0$ ;
- (c) after executing **flush**  $x$  we have:  $x_v=x_s=1 \wedge y_v=y_s=0$ ;
- (d) upon subsequently executing  $y := 1$  we have:  $x_v=x_s=1 \wedge y_v=1 \wedge y_s \in \{0, 1\}$ .

That is, executing **flush**  $x$  (c) copies  $x_v$  to  $x_s$ , and thus we have  $y_s=1 \Rightarrow x_s=1$  at all program points.

**Modelling the Asynchronous Behaviour of flush.** As we demonstrated above, tracking the synchronous version of locations (e.g.  $x_s$ ) allows us to capture the necessary persist orderings (e.g. that  $x_s$  persists before  $y_s$ ). However, as **flush** instructions execute asynchronously, synchronous versions do not accurately capture the memory state upon a crash. For instance, if a crash occurs after **flush**  $x'$  is executed (but not yet fully completed) in Fig. 1b, unlike what we wrote in (c) and (d) above,  $x$  may not necessarily contain 1. To address this, for each location  $x$  we record a third, *persisted* version, written  $x_p$ , denoting the latest persisted value of  $x$  (without assuming **flush** instructions are executed synchronously).

Intuitively, there is a chronological order on different versions of a location  $x$ :  $x_p \rightarrow x_s \rightarrow x_v$ , in that while  $x_p$  reflects the last persisted value of  $x$  in memory,  $x_s$  and  $x_v$  denote later updates on  $x$  that are yet to be persisted, with  $x_v$  describing the latest such update. As discussed,  $x_v$  may non-deterministically be copied to  $x_s$ , allowing  $x_s$  to catch up with  $x_v$ . Intuitively, this amounts to a pending write on  $x$  (in the persistent buffer) being committed to memory. Analogously, the effect of asynchronous **flush** instructions may take effect at non-deterministic times. We may thus be inclined to copy  $x_s$  to  $x_p$  non-deterministically, allowing  $x_p$  to catch up with  $x_s$ . However, this is too weak. To see this, let us extend the (a)–(d) steps of Fig. 1b with  $x_p$  and  $y_p$  constraints (highlighted):

- (a')  $x_v=x_s=x_p=y_v=y_s=y_p=0$
- (b')  $x_v=1 \wedge x_s, x_p \in \{0, 1\} \wedge (x_p=1 \Rightarrow x_s=1) \wedge y_v=y_s=y_p=0$
- (c')  $x_v=x_s=1 \wedge x_p \in \{0, 1\} \wedge y_v=y_s=y_p=0$
- (d')  $x_v=x_s=1 \wedge y_v=1 \wedge x_p, y_s, y_p \in \{0, 1\} \wedge (y_p=1 \Rightarrow y_s=1)$

First, note that  $(x_p=1 \Rightarrow x_s=1)$  in (b') captures the chronological order between  $x_s$  and  $x_p$ :  $x_s$  may be copied to  $x_p$  and thus if  $x_p=1$  then  $x_s=1$ ; similarly for  $(y_p=1 \Rightarrow y_s=1)$  in (d'). Second, note that  $x_p$  and  $y_p$  are our main interest as they denote the latest persisted values of  $x$  and  $y$ . We are hence interested in establishing  $y_p=1 \Rightarrow x_p=1$  at all program points, thus modelling the desired behaviour in Fig. 1b. This is, however, not the case as (d') allows  $y_p=1 \wedge x_p=0$ .

To remedy this, we require that the non-deterministic copying of synchronous values to persisted ones be carried out *for all locations* at once, and not a single location. In doing so, we ensure that  $y_p=1 \Rightarrow x_p=1$  holds at (d'), as desired. Intuitively, this captures the global persist orderings imposed by **flush**. In particular, recall that when  $x \in X$ , all  $X$  writes ordered before **flush**  $x$  persist before all instructions ordered after **flush**  $x$ . As such, upon propagating a write to the persistent memory, i.e. copying some  $y_s$  to  $y_p$ , we must ensure that the effects of each prior **flush**  $x$  instruction is completed in that its preceding writes on  $X$  have also reached the memory. This amounts to (simultaneous) copying of  $x_s$  to  $x_p$  for  $x \in X$  (since each prior **flush** on  $X$  has copied  $x_v$  to  $x_s$  for  $x \in X$ ).

In summary, in our  $\text{Ix86}_{\text{sim}}$  model: (1) memory operations on  $x$  (reads/writes) manipulate  $x_v$ ; (2) when  $x \in X$ , **flush**  $x$  copies  $x'_v$  to  $x'_s$  for every  $x' \in X$ ; (3)  $x_v$  may be copied to  $x_s$  non-deterministically; and (4)  $\overline{x_s}$  may be point-wise copied to  $\overline{x_p}$  non-deterministically, where  $\overline{x_s}$  denotes *all* synchronous locations and  $\overline{x_p}$  denotes the corresponding persistent ones.



## 2.4 POG: Persistent Owicki-Gries Reasoning

Having forgone the need for persistent buffers through our  $\text{Ix86}_{\text{sim}}$  operational model, our next goal is to devise a program logic for persistency reasoning over  $\text{Ix86}_{\text{sim}}$  which operates on the x86-TSO storage system in Fig. 2a. As mentioned in §2.3, our next challenge is accounting for the weak behaviours caused by the thread-local buffers in x86-TSO. Fortunately, existing work [Lahav and Vafeiadis 2015; Sieczkowski et al. 2015; Svendsen et al. 2018; Turon et al. 2014; Vafeiadis and Narayan 2013] demonstrate how existing program logics for SC can be adapted to reason about such weak behaviours. Here, we follow the simple approach of [Lahav and Vafeiadis 2015], which demonstrated how Owicki-Gries (OG) reasoning [Owicki and Gries 1976] can be adapted to reason about the release-acquire fragment of the C11 memory model [Lahav et al. 2017].

**OG Reasoning.** As in Hoare logic, the basic constructs in OG are *Hoare triples* of the form  $\{P\} C \{Q\}$ , where  $P$  and  $Q$  are assertions (sets of states) describing the pre- and post-condition of program  $C$ . OG reasoning extends the proof rules of Hoare logic with a rule to reason about concurrent programs of the form  $C_1 \parallel C_2$ , which allows one to compose the verified programs  $C_1$  and  $C_2$  into a verified concurrent program, provided that the two proofs are *non-interfering*:

$$\frac{\{P_1\} C_1 \{Q_1\} \quad \{P_2\} C_2 \{Q_2\} \quad \text{two proofs are non-interfering}}{\{P_1 \wedge P_2\} C_1 \parallel C_2 \{Q_1 \wedge Q_2\}}$$

As such, OG is often deemed non-compositional as it refers to non-interference of proof outlines that cannot be checked based solely on the two input triples. However, as demonstrated in [Lahav and Vafeiadis 2015], presenting OG in the *rely-guarantee* (RG) [Jones 1983] style allows compositional reasoning. In this presentation, Hoare triples are interpreted under an *RG context*,  $\langle \mathcal{R}; \mathcal{G} \rangle$ . The *rely* component,  $\mathcal{R}$ , comprises a set of assertions assumed to be *stable* under memory updates carried out by the environment (i.e. other threads). The *guarantee* component,  $\mathcal{G}$ , in turn comprises a set of *guarded updates* that the thread may perform. A guarded update is of the form  $\langle x, e, P \rangle$ , stating that when the program state satisfies the guard  $P$ , then the thread may update  $x$  to  $e$ .

An RG-style OG triple is of the form  $\langle \mathcal{R}; \mathcal{G} \rangle \vdash \{P\} C \{Q\}$ , stating that: (1) every terminating run of  $C$  from a state in  $P$  results in a state in  $Q$ ; (2)  $C$  updates the state in accordance with  $\mathcal{G}$  while satisfying the prescribed guards; and (3) the same holds when  $C$  is run in parallel with any program  $C'$  whose updates preserve the assertions in  $\mathcal{R}$ . We can then rewrite the parallel composition rule in the RG-style as shown in Fig. 3 (PAR), requiring that the RG contexts of the proofs be non-interfering. That is, every update of  $C_1$  in  $\mathcal{G}_1$  preserves every assertion in  $\mathcal{R}_2$ , and vice versa.

**Invariant-based Reasoning.** A main advantage of Hoare logic and its descendants such as OG is their support for *compositional* reasoning: once we verify the behaviour of a small program  $C$ , we can use its specification to verify larger programs that use  $C$ . For instance, having established  $\{P\} C \{Q\}$ , we can then use it to verify the larger program  $C' \triangleq C; C''$ . That is, it suffices to find  $Q'$  such that  $\{Q\} C'' \{Q'\}$ , as we can then compose the two triples to derive  $\{P\} C' \{Q'\}$ . In other words, this allows us to treat  $C$  as a black box, jumping over its body and assuming  $Q$  at the end.

However, this is no longer the case in the persistent setting: as the execution of  $C$  may crash, we cannot assume that its execution completed successfully and that  $Q$  describes the state on its completion. Rather, we must account for the possibility of  $C$  crashing at any point during its execution. As such, the state upon returning from  $C$  (either due to a crash or after successful completion) is that described by the *disjunction* (union) of states at each program point. To keep our presentation simple, we typically define an *invariant*,  $I$ , that holds at all program points, and could simply be defined as the disjunction of all states. Intuitively,  $I$  corresponds to the persistency behaviour we seek to establish, e.g.  $y=1 \Rightarrow x=1$  in Fig. 1b. At each program point we can still

strengthen  $I$  by adding additional conjuncts. However, when  $C$  is used in the larger context of  $C'$ , we can simply treat its specification as  $\{I\} C \{I\}$ . Note that this is a generalisation of the approach in [Chen et al. 2015], where Hoare triples are of the form  $\{P\} C \{Q\}\{I\}$ , with  $Q$  denoting the non-crashing postcondition that holds once  $C$  executes successfully, and  $I$  denoting the crashing postcondition that holds in case of a crash, and is itself typically defined as a disjunction of postconditions at each point. As such, the overall postcondition of  $C$  is described by  $Q \vee I$ .

**Stability.** Recall that under  $\text{Ix86}_{\text{sim}}$ ,  $x_v$  may be non-deterministically copied to  $x_s$  for each location, while  $x_s$  may in turn be non-deterministically copied to  $x_p$  for all locations at once. As such, we stipulate that the assertions used in our proofs be *stable* with respect to these propagations. In particular, we require that for all  $x: P \Rightarrow P[x_v/x_s]$ ; i.e. if  $P$  holds beforehand, it should still hold after  $x_v$  is propagated to  $x_s$ . Analogously, we require that for all  $x: P \Rightarrow P[x_s/x_p]$ .

**POG Reasoning.** We develop the POG ('persistent Owicki-Gries') logic for reasoning about persistency behaviours of programs under  $\text{Px86}_{\text{sim}}$ . We formulate POG as an extension of OG, and build it over our  $\text{Ix86}_{\text{sim}}$  model. We present the formal details of POG in §3. Here we introduce POG by verifying the example in Fig. 1b. Later in §4 we verify several other examples in POG.

Recall that our goal is to devise an invariant  $I$  that holds at all program points in Fig. 1b. In particular, we are interested in establishing  $I \triangleq y_p=1 \Rightarrow x_p=1$ , as shown in Fig. 1b. As discussed, memory operations (e.g. writes) on  $x$  in  $\text{Ix86}_{\text{sim}}$  manipulate the  $x_v$  location. This is reflected in the (WRITE) rule in Fig. 3, stating that when executing  $x := e$ , the postcondition  $Q$  is obtained from the precondition  $P$  by substituting the new value  $e$  for  $x_v$ . Analogously, the (FLUSH) proof rule in Fig. 3 states that when executing **flush**  $x$ , the postcondition  $Q$  is obtained from the precondition  $P$  by copying  $x'_v$  to  $x'_s$  for every  $x'$  in the cache line of  $x$ . As such, assuming that initially all locations hold 0 and that  $x$  and  $x'$  are in the same cache line, we can verify the program in Fig. 1b as follows:

$$\begin{aligned} & \{I \wedge x_v=x_s=x_p=0 \wedge y_v=y_s=y_p=0\} \\ & \quad x := 1; \\ & \{I \wedge x_v=1 \wedge x_s, x_p \in \{0,1\} \wedge y_v=y_s=y_p=0\} \\ & \quad \quad \mathbf{flush} \ x'; \\ & \{I \wedge x_v=x_s=1 \wedge x_p \in \{0,1\} \wedge y_v=y_s=y_p=0\} \\ & \quad \quad y := 1; \\ & \{I \wedge x_v=x_s=1 \wedge y_v=1 \wedge x_p, y_s, y_p \in \{0,1\}\} \end{aligned}$$

Note that as the program in Fig. 1b is sequential, we define the RG context as  $\langle \top; \top \rangle$ ; i.e. the environment must preserve all assertions ( $\mathcal{R}=\top$ ), and the program may perform any assignment ( $\mathcal{G}=\top$ ). For simplicity, we have elided the RG context above. Observe that although  $x := 1$  solely manipulates  $x_v$  and in its precondition we have  $x_s=x_p=0$ , after its execution we weakened the postcondition to  $x_s, x_p \in \{0,1\}$ . This is to ensure that the postcondition is *stable* since the value of  $x_v$  (i.e. 1) may non-deterministically propagate to  $x_s$  (and subsequently from  $x_s$  to  $x_p$ ), as discussed above. Moreover, when analogously stabilising the post-condition of  $y := 1$ , we must propagate both  $x_s$  to  $x_p$  and  $y_s$  to  $y_p$  *simultaneously*, thus ruling out  $x_p=0, y_p=1$  (as  $x_s=1$ ) and establishing  $I$ .

### 3 THE POG PROGRAM LOGIC

**POG Language.** The POG language given below is that of  $\text{Px86}_{\text{sim}}$  in [Raad et al. 2020] excluding **flush**<sub>opt</sub>. We assume a finite set  $\text{Loc}$  of memory locations; a finite set  $\text{REG}$  of (local) registers; a finite set  $\text{VAL}$  of values; a finite set  $\text{TID}$  of thread identifiers; any standard interpreted language for expressions containing registers and values; and a finite set  $\text{CL}$  of cache lines partitioning locations ( $\text{Loc} = \uplus \text{CL}$ ). We use  $x, y, z$  as metavariables for locations;  $a, b, c$  for registers;  $v$  for values;  $\tau$  for

$$\begin{array}{c}
\frac{}{\langle\{P\}; \emptyset\rangle \vdash \{P\} \text{ skip } \{P\}} \text{ (SKIP)} \quad \frac{P \Rightarrow Q[e/a]}{\langle\{P, Q\}; \emptyset\rangle \vdash \{P\} a := e \{Q\}} \text{ (ASSIGN)} \quad \frac{P \Rightarrow Q[x_v/a]}{\langle\{P, Q\}; \emptyset\rangle \vdash \{P\} a := x \{Q\}} \text{ (READ)} \\
\\
\frac{P \Rightarrow Q[e/x_v]}{\langle\{P, Q\}; \{\langle x_v, e, P \rangle\}\rangle \vdash \{P\} x := e \{Q\}} \text{ (WRITE)} \quad \frac{P \Rightarrow Q[x_v/a][x_v+e/x_v]}{\langle\{P, Q\}; \{\langle x_v, x_v+e, P \rangle\}\rangle \vdash \{P\} a := \text{FAA}(x, e) \{Q\}} \text{ (FAA)} \\
\\
\frac{P \wedge x_v \neq e_1 \Rightarrow Q[x_v/a] \quad P \wedge x_v = e_1 \Rightarrow Q[e_1/a][e_2/x_v]}{\langle\{P, Q\}; \{\langle x_v, e_2, P \wedge x_v = e_1 \rangle\}\rangle \vdash \{P\} a := \text{CAS}(x, e_1, e_2) \{Q\}} \text{ (CAS)} \quad \frac{x \in X \quad X = \{x^1 \dots x^n\} \quad P \Rightarrow Q[x_v^1/x_s^1 \dots x_v^n/x_s^n]}{\langle\{P, Q\}; \emptyset\rangle \vdash \{P\} \text{ flush } x \{Q\}} \text{ (FLUSH)} \\
\\
\frac{\langle\mathcal{R}_1; \mathcal{G}_1\rangle \vdash \{P\} c_1 \{R\} \quad \langle\mathcal{R}_2; \mathcal{G}_2\rangle \vdash \{R\} c_2 \{Q\}}{\langle\mathcal{R}_1 \cup \mathcal{R}_2; \mathcal{G}_1 \cup \mathcal{G}_2\rangle \vdash \{P\} c_1; c_2 \{Q\}} \text{ (SEQ)} \quad \frac{\langle\mathcal{R}; \mathcal{G}\rangle \vdash \{P \wedge e \neq 0\} c_1 \{Q\} \quad \langle\mathcal{R}; \mathcal{G}\rangle \vdash \{P \wedge e = 0\} c_2 \{Q\}}{\langle\mathcal{R} \cup \{P\}; \mathcal{G}\rangle \vdash \{P\} \text{ if } (e) \text{ then } c_1 \text{ else } c_2 \{Q\}} \text{ (ITE)} \\
\\
\frac{P \wedge e = 0 \Rightarrow Q \quad \langle\mathcal{R}; \mathcal{G}\rangle \vdash \{P \wedge e \neq 0\} c \{P\}}{\langle\mathcal{R} \cup \{Q\}; \mathcal{G}\rangle \vdash \{P\} \text{ while}(e) c \{Q\}} \text{ (WHILE)} \quad \frac{P \Rightarrow P' \quad \mathcal{R}' \subseteq \mathcal{R} \quad \mathcal{G}' \subseteq \mathcal{G} \quad Q' \Rightarrow Q \quad \langle\mathcal{R}'; \mathcal{G}'\rangle \vdash \{P'\} C \{Q'\}}{\langle\mathcal{R} \cup \{P, Q\}; \mathcal{G}\rangle \vdash \{P\} C \{Q\}} \text{ (CONSEQ)} \\
\\
\frac{\langle\mathcal{R}_1; \mathcal{G}_1\rangle \vdash \{P_1\} C_1 \{Q_1\} \quad \langle\mathcal{R}_2; \mathcal{G}_2\rangle \vdash \{P_2\} C_2 \{Q_2\} \quad \langle\mathcal{R}_1; \mathcal{G}_1\rangle \text{ and } \langle\mathcal{R}_2; \mathcal{G}_2\rangle \text{ are non-interfering}}{\langle\mathcal{R}_1 \cup \mathcal{R}_2 \cup \{Q\}; \mathcal{G}_1 \cup \mathcal{G}_2\rangle \vdash \{P_1 \wedge P_2\} C_1 || C_2 \{Q\}} \text{ (PAR)} \\
\text{fr}(\mathcal{R}_1, C_1) \cap \text{wr}(C_2) = \emptyset \quad \text{fr}(\mathcal{R}_2, C_2) \cap \text{wr}(C_1) = \emptyset \\
\\
\frac{P \Rightarrow P' \quad \langle\top; \top\rangle \vdash \{P'\} C \{Q'\} \quad \langle\top; \top\rangle \vdash \{R\} C_{\text{rec}} \{Q'\} \quad Q' \Rightarrow Q \quad Q' \Rightarrow \exists \bar{v}. R[\bar{v}/a][\bar{x}_p/\bar{x}_s][\bar{x}_p/\bar{x}_v]}{\langle\text{rec}\rangle \vdash \{P\} C \{Q\}} \text{ (REC)} \\
C_{\text{rec}} \vdash \{P\} C \{Q\}
\end{array}$$

Fig. 3. POG proof rules with the implicit assumption that the pre- and post-conditions are stable (Def. 2).

thread identifiers;  $e$  for expressions; and  $X$  for cache lines.

$$\begin{aligned}
\text{SCOM} \ni c ::= & \text{ skip } \mid \text{ while}(e) c \mid \text{ if } (e) \text{ then } c_1 \text{ else } c_2 \mid c_1; c_2 \mid a := e \mid a := x \mid x := e \\
& \mid \text{ flush } x \mid \text{ sfence } \mid \text{ mfence } \mid a := \text{CAS}(x, e_1, e_2) \mid a := \text{FAA}(x, e) \\
\text{COM} \ni C \triangleq & \text{TID} \xrightarrow{\text{fin}} \text{SCOM}
\end{aligned}$$

The sequential fragment of the language is given by the  $c$  grammar and includes the standard constructs of **skip**, loops, conditionals and sequential composition, as well as local variable assignment ( $a := e$ ), memory read from location  $x$  ( $a := x$ ), memory write to  $x$  ( $x := e$ ), memory persist (**flush**  $x$ ), store fence (**sfence**), memory fence (**mfence**) and atomic RMW (read-modify-write) instructions. The RMW instruction  $a := \text{CAS}(x, e_1, e_2)$  denotes the atomic ‘compare-and-swap’, where the value of  $x$  is compared against  $e_1$ : if the values match then  $x$  is set to  $e_2$  and 1 is returned in  $a$ ; otherwise  $x$  is left unchanged and 0 is returned in  $a$ . Analogously,  $a := \text{FAA}(x, e)$  denotes the atomic ‘fetch-and-add’, where  $x$  is incremented by  $e$  and its old value is returned in  $a$ . Lastly, we model a multi-threaded program  $C$  as a function mapping each thread to its (sequential) program. We write  $C = c_1 || \dots || c_n$  when  $\text{dom}(C) = \{\tau_1 \dots \tau_n\}$  and  $C(\tau_i) = c_i$ , and write  $C_1 || C_2$  for  $C_1 \uplus C_2$ . We lift a sequential program  $c$  to a program in  $\text{COM}$  and simply write  $c$  for  $C \triangleq [\tau \mapsto c]$ .

**Instrumented Locations.** Recall from §2 that in order to reason about the persistency behaviours of programs, we instrument the memory to track three separate versions for each memory location. To this end, we define the set of *instrumented locations* as  $\text{ILoc} \triangleq \{x_v, x_s, x_p \mid x \in \text{Loc}\}$ ,

and define an instrumented memory,  $\text{IM} : \text{IMEM}$ , as a finite map from instrumented locations to values:  $\text{IMEM} : \text{ILOC} \xrightarrow{\text{fin}} \text{VAL}$ . As discussed in §2, for each memory location  $x$ :

- (1)  $\text{IM}(x_v)$  denotes the *volatile* value of  $x$ , i.e. the value observed for  $x$  during the execution;
- (2)  $\text{IM}(x_s)$  denotes the *synchronously persisted* value of  $x$ , i.e. the value observed for  $x$  after a crash, had the **flush** instructions executed synchronously; and
- (3)  $\text{IM}(x_p)$  denotes the *persistent* value of  $x$ , i.e. the value observed for  $x$  after a crash.

Intuitively,  $\text{IM}(x_p)$  reflects the current value of  $x$  in memory, while  $\text{IM}(x_v)$  and  $\text{IM}(x_s)$  record additional information that enables persistent reasoning, as discussed in §2.

**Assertions.** The POG assertions represent sets of states in our  $\text{Ix86}_{\text{sim}}$  semantics. Our assertion language is that of first order logic with *equality* and i) three constant symbols  $x_v, x_s, x_p$  for each location  $x$ ; ii) a constant symbol  $a$  for each register; and iii) a constant symbol  $v$  for each value.

### 3.1 The POG Proof System

The *POG triples* are of the form:  $\langle \mathcal{R}; \mathcal{G} \rangle \vdash \{P\} C \{Q\}$ , stating that: (1) the persistent parts of  $Q$  (describing persistent versions of variables, e.g.  $x_p$ ) are *invariant* throughout the execution of  $C$ ; (2) every *terminating* run of  $C$  from a state in  $P$  results in a state in  $Q$ ; (3)  $C$  updates the state in accordance with  $\mathcal{G}$  while satisfying the prescribed guards; and (4) the above holds when  $C$  is run in parallel with any program  $C'$ , provided that the  $C'$  updates preserve the assertions in  $\mathcal{R}$ . That is, while the persistent parts of  $Q$  hold at *all points* during the execution of  $C$ , those parts describing register values and volatile/synchronous versions hold only when  $C$  terminates successfully.

**POG Proof Rules.** We present the *POG proof rules* in Fig. 3. Ignoring the (REC) rule at the bottom, most rules remain largely unchanged from their OG counterparts and are merely adapted to RG-style as in [Lahav and Vafeiadis 2015]. Intuitively, the pre- and post-conditions of triples are accumulated in the rely component  $\mathcal{R}$  to ensure they remain stable (invariant) under the updates performed by other threads. Conversely, each time a thread updates a location  $x$  via (WRITE), (CAS) and (FAA), this is recorded in its guarantee component  $\mathcal{G}$  with the corresponding precondition. Moreover, the (READ), (WRITE), (CAS) and (FAA) rules accessing memory location  $x$  have been accordingly adjusted to access the latest value of  $x$ , namely that in  $x_v$ .

As discussed in §2.3, when  $x \in X$ , **flush**  $x$  propagates the latest persist-pending value of each  $x' \in X$  to memory. This is reflected in the  $P \Rightarrow Q[x_v^1/x_s^1 \dots x_v^n/x_s^n]$  premise of the (FLUSH) rule.

The (PAR) rule describes the concurrent execution  $C_1 \parallel C_2$ , where the non-interference premise ensures that  $C_1$  and  $C_2$  do not interfere with one another. Intuitively,  $C_1$  and  $C_2$  are non-interfering iff each update performed by  $C_1$  preserves the state assumptions of  $C_2$ , and vice versa. Put formally, when  $C_1$  and  $C_2$  are respectively run under contexts  $\langle \mathcal{R}_1; \mathcal{G}_1 \rangle$  and  $\langle \mathcal{R}_2; \mathcal{G}_2 \rangle$ , then every update  $\langle x, e, P \rangle$  of  $C$  in  $\mathcal{G}_1$  must preserve every assertion  $R$  in  $\mathcal{R}_2$ , i.e.  $P \wedge R \Rightarrow R[e/x_v]$ ; and vice versa.

**Definition 1.** The tuples  $\langle \mathcal{R}_1; \mathcal{G}_1 \rangle$  and  $\langle \mathcal{R}_2; \mathcal{G}_2 \rangle$  are *non-interfering* iff:

- for all  $\langle x_v, e, P \rangle \in \mathcal{G}_1$  and all  $R \in \mathcal{R}_2$ :  $R \wedge P \Rightarrow R[e/x_v]$
- for all  $\langle x_v, e, P \rangle \in \mathcal{G}_2$  and all  $R \in \mathcal{R}_1$ :  $R \wedge P \Rightarrow R[e/x_v]$

**Stability.** Recall from §2.4 that we require assertions to be stable against volatile-to-synchronous and synchronous-to-persistent version propagations. This is formalised in Def. 2.

**Definition 2.** An assertion  $P$  is *stable*, written  $\text{stable}(P)$ , iff  $P \Rightarrow P[\overline{x_s/x_p}]$  and  $\forall x_s. P \Rightarrow P[x_v/x_s]$ .

**Reasoning about Crash Recovery.** The POG triples discussed thus far are of the form  $\langle \mathcal{R}; \mathcal{G} \rangle \vdash \{P\} C \{Q\}$ , where  $Q$  describes the state once  $C$  has *terminated* (i.e. without crashing). However, as discussed in §2, the execution of  $C$  may crash (e.g. due to power loss), at which point the volatile and

synchronous versions (e.g.  $x_v$  and  $x_s$ ) are lost while the persistent versions (e.g.  $x_p$ ) are preserved, and the execution is resumed by running a *recovery program*. As such, in anticipation of a possible crash at any point, we require the persistent parts of  $Q$  to be *invariant* throughout the execution.

In order to reason about programs in the presence of crashes, we present POG *recovery* triples of the form  $C_{\text{rec}} \vdash \{P\} C \{Q\}$ , stating that every run of  $C$  from a state in  $P$  either: (1) terminates successfully (without crashing) in a state in  $Q$ ; or (2) crashes and its execution is resumed by repeatedly running the recovery program  $C_{\text{rec}}$  until  $C_{\text{rec}}$  terminates successfully in a state in  $Q$ . That is, the execution of  $C_{\text{rec}}$  may itself crash and is rerun repeatedly until it terminates successfully. In other words, the volatile (and synchronous) values in the postcondition of POG triples describe the post-states *only if*  $C$  does not crash, whereas the volatile (and synchronous) values in the postcondition of recovery triples always describe the post-states after termination (either after  $C$  terminates successfully, or after  $C$  crashes and  $C_{\text{rec}}$  is run repeatedly until it terminates successfully).

The recovery rule (REC) is given at the bottom of Fig. 3. The  $\langle \top; \top \rangle \vdash \{P'\} C \{Q'\}$  premise ensures that executing  $C$  from a state in  $P'$  either terminates successfully in a state in  $Q'$ , or it crashes and the recovery  $C_{\text{rec}}$  is run thereafter from a state in  $R$ , obtained from  $Q'$  by resetting register values and replacing volatile/synchronous versions with persistent ones ( $Q' \Rightarrow \exists \bar{v}. R[\bar{v}/\bar{a}][x_p/x_s][x_p/x_v]$ ), as they are lost upon a crash. The existential quantification of  $\bar{v}$  assigns (havocs) arbitrary values to local registers  $\bar{a}$  upon recovery, ensuring that  $R$  makes no assumptions about the post-crash values of registers. The  $\langle \top; \top \rangle \vdash \{R\} C_{\text{rec}} \{Q'\}$  in turn ensures that executing  $C_{\text{rec}}$  from  $R$  either terminates successfully in  $Q'$ , or it crashes and is rerun from  $R$ . Put together, as  $P \Rightarrow P'$  and  $Q' \Rightarrow Q$ , this ensures that executing  $C$  (under  $C_{\text{rec}}$ ) from  $P$  eventually terminates successfully in  $Q$ . Lastly, the RG contexts  $\langle \top; \top \rangle$  ensure that  $C$  and  $C_{\text{rec}}$  are run as closed programs (i.e. not in parallel with another program). In §4 we present an example of using (REC) to reason about recovery.

**Fence Proof Rules.** Note that our proof system in Fig. 3 does not include rules for **sfence** and **mfence**. For **sfence**, we can simply extend our rules with:  $\langle \{P\}; \emptyset \rangle \vdash \{P\} \text{sfence} \{P\}$ . That is, **sfence** acts as a no-op in our  $\text{Ix86}_{\text{sim}}$  model and has the same specification as (SKIP). This is caused by two factors. First recall that as discussed in §2.3, under  $\text{Ix86}_{\text{sim}}$  a **flush**  $x$  instruction with  $x \in X$  copies  $x'_v$  to  $x'_s$  for each  $x' \in X$ . That is,  $\text{Ix86}_{\text{sim}}$  *eliminates* each **flush** instruction on  $x \in X = \{x^1 \cdots x^n\}$ , and simply treats it as a series of writes on  $x_s^1 \cdots x_s^n$ . Second, as noted by Raad et al. [2020], in the absence of **flush/flush<sub>opt</sub>** instructions, **sfence** instructions behave as no-ops (**skip**) and impose no additional ordering constraints. As such, since  $\text{Ix86}_{\text{sim}}$  eliminates **flush** instructions (and treats them as writes) and excludes **flush<sub>opt</sub>** instructions by design, **sfence** is a no-op under  $\text{Ix86}_{\text{sim}}$ . Nevertheless, as discussed in §2.1 and detailed later in in §7, **sfence** instructions can be used to enforce additional ordering constraints on **flush<sub>opt</sub>** instructions, allowing us in most cases to transform programs using **flush<sub>opt</sub>** to those using **flush** instead. We therefore opt to include **sfence** in the POG programming language to facilitate such transformations.

For **mfence**, we can derive reasoning principles by treating them as RMW instructions (as in Fig. 7 of [Lahav and Vafeiadis 2015]). More concretely, we can treat **mfence** instructions as RMWs (e.g. **FAA**) on a designated location  $f$ , which enforces a global order on all **mfence** instructions.

## 4 EXAMPLES

We use the POG proof rules to verify several representative examples.

**Example 1.** We begin with the concurrent example in Fig. 1e, with its proof sketch given in Fig. 4a. The RG contexts of the left and right threads are given respectively by  $\langle \mathcal{R}_1; \mathcal{G}_1 \rangle$  and  $\langle \mathcal{R}_2; \mathcal{G}_2 \rangle$ , defined in Fig. 4a. As discussed in §3, the rely of each thread contains the assertions used in its proof outline (e.g.  $P_1$  in  $\mathcal{R}_1$ ), while the guarantee contains a guarded assignment for each write performed by the thread (e.g.  $\langle P_2, y_v, 1 \rangle$  in  $\mathcal{G}_1$ ). The invariant we are interested in establishing is

$I \triangleq (x_s=1 \Rightarrow x_v=1) \wedge (z_s=1 \Rightarrow z_v=1)$ $\wedge (z_p=1 \Rightarrow z_s=1) \wedge (y_v=1 \Rightarrow x_s=1)$ $\wedge (z_v=1 \Rightarrow x_s=1) \wedge (z_p=1 \Rightarrow x_p=1)$ $\mathcal{R}_1 \triangleq \{I, P_1, P_2\} \quad \mathcal{R}_2 \triangleq \{I, Q_1, Q_2\}$ $\mathcal{G}_1 \triangleq \{\langle I, x_v, 1 \rangle, \langle P_2, y_v, 1 \rangle\} \quad \mathcal{G}_2 \triangleq \{\langle Q_2, z_v, 1 \rangle\}$ <div style="text-align: center;"> </div> <p style="text-align: center;">(a)</p>	$I \triangleq (x_s=1 \Rightarrow x_v=1) \wedge (z_s=1 \Rightarrow z_v=1)$ $\wedge (z_p=1 \Rightarrow z_s=1) \wedge (y_v=1 \Rightarrow x_v=1)$ $\wedge (z_v=1 \Rightarrow x_s=1) \wedge (z_p=1 \Rightarrow x_p=1)$ $\mathcal{R}_1 \triangleq \{I, P_1\} \quad \mathcal{R}_2 \triangleq \{I, Q_1, Q_2, Q_3\}$ $\mathcal{G}_1 \triangleq \{\langle I, x_v, 1 \rangle, \langle P_1, y_v, 1 \rangle\} \quad \mathcal{G}_2 \triangleq \{\langle Q_3, z_v, 1 \rangle\}$ <div style="text-align: center;"> </div> <p style="text-align: center;">(b)</p>
--	--

Fig. 4. Proof sketches of Example 1 (a) and Example 2 (b)

given by  $I$ . The first three conjuncts capture the chronological order between the different versions of  $x$  and  $z$ . The penultimate conjunct states that when  $z_v=1$  (i.e. once the second thread executes  $z:=1$  after having read 1 for  $y_v$  in a), then  $x_s=1$  (i.e. **flush**  $x'$  must have already executed). The last conjunct is that of main interest, corresponding to the desired property in Fig. 1e. Note that the assertions at each program point are stable, and that  $\langle \mathcal{R}_1; \mathcal{G}_1 \rangle$  and  $\langle \mathcal{R}_2; \mathcal{G}_2 \rangle$  are non-interfering.

**Example 2.** We proceed with the example in Fig. 4b which is an adaptation of Fig. 4a with the **flush** moved to the right thread after reading from  $y$ . As such, the invariant  $I$  is similar to that of Fig. 4a, with the main difference lying in the fourth conjunct. In particular, when the left thread executes  $y:=1$  yielding  $y_v=1$ , the earlier  $x:=1$  has already executed, i.e.  $x_v=1$ . However, due to the absence of an intervening **flush** between the two writes, unlike in Fig. 4a we cannot assert  $x_s=1$ .

**Example 3** (Atomic persists). Our next example in Fig. 5a is inspired by the *persistent transactions* of Raad et al. [2019], where the authors showed how to ensure multiple stores on different locations (appear to) persist atomically before subsequent stores. Let us write  $\tau_1$  and  $\tau_2$  for the left and right threads, respectively. As shown in Fig. 5a,  $\tau_1$  writes 1 to  $x$  and  $y$ , and  $\tau_2$  writes 1 to  $z$  only if  $x$  contains 1, i.e. only if  $\tau_1$  has already executed  $x:=1$ . Our goal in this example is to ensure that the writes on  $x$  and  $y$  (by  $\tau_1$ ) both persist before the write on  $z$  (by  $\tau_2$ ). That is, we must ensure that the **flush** instructions of  $\tau_1$  are executed before  $\tau_2$  executes  $z:=1$ .

To this end, since  $\tau_2$  writes to  $z$  only after reading 1 from  $x$  (written by  $\tau_1$ ), we use a *lock* to control the accesses on  $x$ . More concretely,  $\tau_1$  acquires the *lx* lock on  $x$  at the beginning and releases it only after executing its **flush** instructions. Similarly,  $\tau_2$  acquires the *lx* lock prior to accessing  $x$ . As such, if  $\tau_2$  reads 1 for  $x$  (i.e. observes  $x:=1$  by  $\tau_1$ ) and executes  $z:=1$ , then it must have acquired *lx* after it was released by  $\tau_1$ , i.e. after  $\tau_1$  executed its **flush** instructions, as required.

The *lx* lock is acquired by calling **lock**, implemented as a spin lock in Fig. 5a: the implementation of **lock**( $lx, a, v$ ) loops until  $lx$  is free (i.e.  $lx=0$ ), at which point it acquires it by atomically setting it to the non-zero value  $v$ . In order to distinguish which thread currently holds the lock,  $\tau_1$  and  $\tau_2$  respectively write the distinct values 1 and 2 to  $lx$  upon acquiring it. Lastly, the  $a$  argument denotes a thread-local register used as the loop flag by **lock**, and is passed on to **lock** by the calling thread.

We present a proof sketch of this program in Fig. 5a, where for brevity we have omitted the RG contexts. As before, our goal is to establish the  $I$  invariant, with its first, second and fourth

$I \triangleq$ $(x_s=1 \Rightarrow x_v=1) \wedge (y_s=1 \Rightarrow y_v=1) \wedge (lx_v \in \{0, 2\} \wedge x_v=1 \Rightarrow x_s=y_s=1)$ $\wedge (z_s=1 \Rightarrow z_v=1) \wedge (z_v=1 \Rightarrow x_s=y_s=1) \wedge (z_p=1 \Rightarrow x_p=y_p=1)$ <div style="display: flex; justify-content: space-between; align-items: center;"> <div style="width: 45%;"> <pre> {I} lock(lx, a, 1); {I ∧ lx_v=1} x := 1; {I ∧ lx_v=1 ∧ x_v=1} y := 1; {I ∧ lx_v=1 ∧ x_v=y_v=1} flush x; {I ∧ lx_v=1 ∧ x_s=y_s=1} flush y; {I ∧ lx_v=1 ∧ x_s=y_s=1} lx := 0; {I} </pre> </div> <div style="width: 45%; border-left: 1px solid black; padding-left: 10px;"> <pre> {I} lock(lx, b, 2); {I ∧ lx_v=2} c := x; {I ∧ lx_v=2 ∧ c=x_v} if (c = 1) {I ∧ lx_v=2 ∧ x_s=y_s=1} z := 1; {I ∧ lx_v=2} flush z; {I ∧ lx_v=2} {I ∧ lx_v=2} lx := 0; {I} </pre> </div> </div> <hr style="border: 0.5px dashed black;"/> <div style="display: flex; justify-content: space-between; align-items: center;"> <div style="width: 45%;"> <pre> {I} lock(lx, a, v) ≜ a := 1; {I ∧ a=1} while(a ≠ 0) {I}   a := CAS(lx, 0, v) {I ∧ (a=0 ⇒ lx_v=v)} {I ∧ lx_v=v} </pre> <p style="text-align: center;">(a)</p> </div> </div>	$I \triangleq y_p=1 \Rightarrow x_p=1 \quad I_v \triangleq y_v=1 \Rightarrow x_v=1$ $x \equiv v \stackrel{\text{def}}{\Leftrightarrow} x_v=x_s=x_p=v$ $x \in \{0, 1\} \stackrel{\text{def}}{\Leftrightarrow} x \equiv 0 \vee x \equiv 1$ <div style="display: flex; justify-content: space-between; align-items: center; margin-top: 10px;"> <div style="width: 45%;"> <pre> {P': I ∧ x ∈ {0, 1} ∧ y ≡ 0} C ≜ x := 1; {I ∧ x_v=1 ∧ x_p ∈ {0, 1} ∧ y ≡ 0} flush x; {I ∧ x_v=x_s=1 ∧ x_p ∈ {0, 1} ∧ y ≡ 0} y := 1; {Q': I ∧ x_v=y_v=1 ∧ x_p, y_p ∈ {0, 1}} </pre> </div> <div style="width: 45%; border-left: 1px solid black; padding-left: 10px;"> <pre> {R: I ∧ I_v ∧ x ∈ {0, 1} ∧ y ∈ {0, 1}} C_rec ≜ a := y; {R ∧ y ≡ a} if (a = 0) {R ∧ y ≡ 0} {P'} C {Q'} </pre> </div> </div> <hr style="border: 0.5px dashed black;"/> <div style="display: flex; justify-content: space-between; align-items: center;"> <div style="width: 45%;"> <pre> C_rec ⊢ {P: x ≡ 0 ∧ y ≡ 0} C {Q: x_v=1} </pre> <p style="text-align: center;">(b)</p> </div> </div>
---	---

Fig. 5. Proof sketches of Example 3 (a) and Example 4 (b)

conjuncts describing the chronological order on different versions of  $x$ ,  $y$  and  $z$ . The third conjunct states that once  $\tau_1$  has finished executing, i.e. it has released the lock ( $lx \in \{0, 2\}$ ) having written 1 to  $x$  ( $x_v=1$ ), its flush instructions have also executed ( $x_s=y_s=1$ ). Similarly, the penultimate conjunct states that once  $\tau_2$  has written 1 to  $z$  ( $z_v=1$ ) having read 1 for  $x$ , i.e. after acquiring  $lx$  once it is released by  $\tau_1$  (see above), then the  $\tau_1$  **flush** instructions must have already executed ( $x_s=y_s=1$ ). The last conjunct describes the desired persist ordering: if the write on  $z$  has persisted ( $z_p=1$ ), then so must have both writes on  $x$  and  $y$  ( $x_p=y_p=1$ ), as required.

**Example 4** (Recovery). In Fig. 5b we show how to use the (REC) rule to reason about recovery. The original program  $C$  is similar to that in Fig. 1b: first 1 is written to  $x$  and persisted by calling **flush**, and then 1 is written to  $y$ . As in Fig. 1b, the intervening **flush** ensures that  $x := 1$  persists before  $y := 1$  and thus  $y_p=1 \Rightarrow x_p=1$  as described by  $I$ , assuming initially  $x \equiv 0$  (i.e.  $x_v=x_s=x_p=0$ ) and  $y \equiv 0$ .

As  $I \triangleq y_p=1 \Rightarrow x_p=1$  is invariant throughout the execution ( $I$  only concerns persistent versions), the recovery  $C_{\text{rec}}$  treats the write on  $y$  as a flag to ascertain if  $x := 1$  has persisted. That is, if  $y$  holds 1 upon recovery, then  $x$  must also hold 1 and thus  $C_{\text{rec}}$  simply returns; otherwise,  $C_{\text{rec}}$  reruns  $C$ .

We present a proof sketch of  $C$  and  $C_{\text{rec}}$  in Fig. 5b. We first establish  $\langle T; T \rangle \vdash \{P'\} C \{Q'\}$  and  $\langle T; T \rangle \vdash \{R\} C_{\text{rec}} \{Q'\}$ , and then use (REC) to prove  $C_{\text{rec}} \vdash \{P\} C \{Q\}$ . Note that as mandated by the (REC) premise, we must show  $P \Rightarrow P'$ ,  $Q' \Rightarrow Q$ , and  $Q' \Rightarrow R[x_p/x_s][x_p/x_v]$  which follow immediately.

## 5 THE Ix86<sub>sim</sub> OPERATIONAL SEMANTICS

We present the Ix86<sub>sim</sub> operational semantics discussed in §2.3. Recall that in Ix86<sub>sim</sub> a memory operation on  $x$  manipulates  $x_v$ ; when  $x \in X$ , **flush**  $x$  copies  $x'_v$  to  $x'_s$  for each  $x' \in X$ ;  $x_v$  may be

$$\begin{array}{l}
(a := e) \triangleq a := e \quad (a := x) \triangleq a := x_v \quad (x := e) \triangleq x_v := e \quad (a := \text{CAS}(x, e_1, e_2)) \triangleq a := \text{CAS}(x_v, e_1, e_2) \\
(a := \text{FAA}(x, e)) \triangleq a := \text{FAA}(x_v, e) \quad (\text{sfence}) \triangleq \text{skip} \quad (\text{mfence}) \triangleq \text{mfence} \quad (\text{flush } x) \triangleq \langle \text{persist } X \rangle \\
(\text{if } (e) \text{ then } c_1 \text{ else } c_2) \triangleq \text{if } (e) \text{ then } \langle c_1 \rangle \text{ else } \langle c_2 \rangle \quad (\text{while}(e) \text{ c}) \triangleq \text{while}(e) \langle c \rangle \\
\langle c_1; c_2 \rangle \triangleq \langle c_1 \rangle; \langle c_2 \rangle \quad \llbracket c_1 \parallel \dots \parallel c_n \rrbracket \triangleq \langle c_1 \rangle \parallel \dots \parallel \langle c_n \rangle \parallel c_s \parallel c_p \\
\text{with } c_s \triangleq \text{while}(\ast) \langle \text{pick } x; x_s := x_v \rangle \quad c_p \triangleq \text{while}(\ast) \langle \overline{x_p} := \overline{x_s}; \rangle
\end{array}$$

Fig. 6. Ix86<sub>sim</sub> program translation where we assume  $x \in X$ 

copied to  $x_s$  non-deterministically; and  $\overline{x_s}$  may be non-deterministically copied to  $\overline{x_p}$ . To model this, we define a *translation function* that transforms Px86<sub>sim</sub> programs to access the instrumented memory.

**Translation.** Our translation function,  $\llbracket \cdot \rrbracket$ , is defined in Fig. 6 and uses the auxiliary function,  $\langle \cdot \rangle$ , to translate sequential programs. As discussed, memory accesses on  $x$  are simply translated to access  $x_v$ ; conditionals, loops and sequential composition are translated inductively; and **flush**  $x$  is translated to **persist**  $X$  (when  $x \in X$ ) *atomically*, as indicated by  $\langle \cdot \rangle$ . That is, in one computation step **persist**  $X$  reads the value of  $x_v$  and subsequently copies to  $x_s$  for each  $x \in X$ , as we describe shortly.

The **mfence** instructions are left unchanged by the translation, while **sfence** instructions are translated as **skip**. This is because as discussed in §3, **sfence** behaves as a no-op in the absence of **flush/flush**<sub>opt</sub>. As such, since our translation eliminates **flush** and our language excludes **flush**<sub>opt</sub>, **sfence** is simply translated to **skip**. Lastly, the translation of a concurrent program is obtained from the point-wise translation of each thread, run in parallel with  $c_s$  and  $c_p$ . Intuitively,  $c_s$  models the non-deterministic propagation of  $x_v$  to  $x_s$  for an arbitrary  $x$ , carried out atomically. Analogously,  $c_p$  models the non-deterministic propagation of  $\overline{x_s}$  to  $\overline{x_p}$ , for all locations. We write  $\tau_s$  and  $\tau_p$  for the threads executing  $c_s$  and  $c_p$ , respectively. Given a program  $C$ , we write  $C^{\text{SP}}$  for  $C \parallel c_s \parallel c_p$ .

We next describe the Ix86<sub>sim</sub> operational semantics by separating the transitions of its *program* and *storage* subsystems. The former describe the steps in program execution, e.g. how a conditional branch is triggered. The latter describe how the storage subsystem (the non-volatile memory and thread-local buffers in Fig. 2a) determine the execution steps. The Ix86<sub>sim</sub> operational semantics is then defined by combining the transitions of its program and storage subsystems.

**Program Transitions.** The Ix86<sub>sim</sub> program transitions are given at the top of Fig. 7 and are defined via the transitions of their constituent threads. Thread transitions are of the form:  $c, S \xrightarrow{\tau:l} c', S'$ , where  $c, c' \in \text{SCOM}$  denote translated sequential programs, and  $S, S' \in \text{STACK}$  denote *stacks* mapping registers to values. The  $\tau:l$  marks the transition by recording the executing thread  $\tau$ , and the transition *label*  $l$ . A label may be  $\epsilon$  for silent transitions of no-ops;  $(R, x, v)$  for reading  $v$  from  $x$ ;  $(W, x, v)$  for writing  $v$  to  $x$ ;  $(U, x, v, v')$  for a successful RMW (update) modifying the value of  $x$  to  $v'$  when its value matches  $v$ ; **MF** for executing an **mfence**;  $(\text{FL}, X)$  for persisting the  $X$  cache line;  $(S, x_s, v)$  for the atomic propagation of  $v$  to  $x_s$ ; and  $(P, \overline{x_p^i}, v^i)$  for the atomic propagation of  $v^i$  to  $\overline{x_p^i}$ .

Given an expression  $e$ , we write  $S(e)$  for the value to which  $e$  evaluates under  $S$ ; this definition is standard and omitted. Most program transitions are standard. The (P-MF) transition describes executing an **mfence**. The (P-CAS0) transition describes the unsuccessful execution of **CAS**( $x, e_1, e_2$ ); i.e. when the value read ( $v$ ) is different from  $S(e_1)$ . The (P-CAS1) transition dually describes the successful execution of **CAS**. Note that in the failure case no update takes place and the transition is labelled with a read, and not an update as in the success case. The (P-FAA) transition behaves analogously. When executing **persist**  $X$  (i.e. a translated **flush**  $x$  with  $x \in X$ ), the volatile-to-synchronous propagation of  $X$  is modelled by the  $(\text{FL}, X)$  transition in



---

**Program transitions:**  $\text{COM} \times \text{STACK} \xrightarrow{\text{TID:LAB} \cup \{\epsilon\}} \text{COM} \times \text{STACK}$        $S \in \text{STACK} \triangleq \text{REG} \xrightarrow{\text{fin}} \text{VAL}$   
 $\text{LAB} \triangleq \left\{ (\mathbb{R}, x_v, v), (\mathbb{W}, x_v, v), (\mathbb{U}, x_v, v, v'), \text{MF}, (\text{FL}, X), (\mathbb{S}, x_s, v), (\mathbb{P}, x_p^i, v^i) \mid x, x^i \in \text{LOC} \wedge X \in \text{CL} \wedge v, v', v^i \in \text{VAL} \right\}$

$$\frac{S' = S[a \mapsto S(e)]}{a := e, S \xrightarrow{\tau:\epsilon} \text{skip}, S'} \text{ (P-ASSIGN)} \quad \frac{S' = S[a \mapsto v]}{a := x_v, S \xrightarrow{\tau:(\mathbb{R}, x_v, v)} \text{skip}, S'} \text{ (P-READ)} \quad \frac{}{\text{mfence}, S \xrightarrow{\tau:\text{MF}} \text{skip}, S} \text{ (P-MF)}$$

$$\frac{S(e) = v}{x_v := e, S \xrightarrow{\tau:(\mathbb{W}, x_v, v)} \text{skip}, S} \text{ (P-WRITE)} \quad \frac{S(e_1) \neq v \quad S' = S[a \mapsto v]}{a := \text{CAS}(x_v, e_1, e_2), S \xrightarrow{\tau:(\mathbb{R}, x_v, v)} \text{skip}, S'} \text{ (P-CAS0)}$$

$$\frac{S(e) = v \quad S' = S[a \mapsto v_0]}{a := \text{FAA}(x_v, e), S \xrightarrow{\tau:(\mathbb{U}, x_v, v_0, v_0 + v)} \text{skip}, S'} \text{ (P-FAA)} \quad \frac{S(e_1) = v_1 \quad S(e_2) = v_2 \quad S' = S[a \mapsto v_1]}{a := \text{CAS}(x_v, e_1, e_2), S \xrightarrow{\tau:(\mathbb{U}, x_v, v_1, v_2)} \text{skip}, S'} \text{ (P-CAS1)}$$

$$\frac{}{\langle \text{persist } X \rangle, S \xrightarrow{\tau:(\text{FL}, X)} \text{skip}, S} \text{ (P-ATOMFL)} \quad \frac{S(e) \neq 0 \Rightarrow c = c_1 \quad S(e) = 0 \Rightarrow c = c_2}{\text{if } (e) \text{ then } c_1 \text{ else } c_2, S \xrightarrow{\tau:\epsilon} c, S} \text{ (P-IF)}$$

$$\frac{}{\text{while}(e) \ c, S \xrightarrow{\tau:\epsilon} \text{if } (e) \text{ then } c; (\text{while}(e) \ c) \text{ else skip}, S} \text{ (P-WHILE)} \quad \frac{}{\text{skip}; c, S \xrightarrow{\tau:\epsilon} c, S} \text{ (P-SEQ1)}$$

$$\frac{}{\langle \text{pick } x; x_s := x_v \rangle, S \xrightarrow{\tau_s:(\mathbb{S}, x_s, v)} \text{skip}, S} \text{ (P-ATOMS)} \quad \frac{c_1, S \xrightarrow{\tau:l} c'_1, S'}{c_1; c_2, S \xrightarrow{\tau:l} c'_1; c_2, S'} \text{ (P-SEQ2)}$$

$$\frac{}{\langle x_p^i := x_s^i \rangle, S \xrightarrow{\tau_p:(\mathbb{P}, x_p^i, v^i)} \text{skip}, S} \text{ (P-ATOMP)} \quad \frac{C(\tau), S \xrightarrow{\tau:l} c, S'}{C, S \xrightarrow{\tau:l} C[\tau \mapsto c], S'} \text{ (P-PAR)}$$


---

**Storage transitions:**  $\text{IMEM} \times \text{IBMAP} \xrightarrow{\text{TID:LAB} \cup \{\epsilon\}} \text{IMEM} \times \text{IBMAP}$        $B \in \text{IBMAP} \triangleq \text{TID} \xrightarrow{\text{fin}} \text{IBUFF}$   
 $b \in \text{IBUFF} \triangleq \text{SEQ} \langle (\text{ILOc} \times \text{VAL}) \cup \mathcal{P}(\text{LOC}) \rangle$

$$\frac{B(\tau) = b}{\text{IM}, B \xrightarrow{\tau:(\mathbb{W}, x_v, v)} \text{IM}, B[\tau \mapsto b.\langle x_v, v \rangle]} \text{ (M-WRITE)} \quad \frac{B(\tau) = \epsilon \quad \text{IM}(x_v) = v_1}{\text{IM}, B \xrightarrow{\tau:(\mathbb{U}, x_v, v_1, v_2)} \text{IM}[x_v \mapsto v_2], B} \text{ (M-RMW)}$$

$$\frac{\text{rd}(\text{IM}, B(\tau), x_v) = v}{\text{IM}, B \xrightarrow{\tau:(\mathbb{R}, x_v, v)} \text{IM}, B} \text{ (M-READ)} \quad \frac{B(\tau) = \epsilon}{\text{IM}, B \xrightarrow{\tau:\text{MF}} \text{IM}, B} \text{ (M-MF)} \quad \frac{\text{IM}(x_v) = v}{\text{IM}, B \xrightarrow{\tau_s:(\mathbb{S}, x_s, v)} \text{IM}[x_s \mapsto v], B} \text{ (M-ATOMS)}$$

$$\frac{\forall i. \text{IM}(x_s^i) = v^i}{\text{IM}, B \xrightarrow{\tau_p:(\mathbb{P}, x_p^i, v^i)} \text{IM}[x_p^i \mapsto v^i], B} \text{ (M-ATOMP)} \quad \frac{B(\tau) = b}{\text{IM}, B \xrightarrow{\tau:(\text{FL}, X)} \text{IM}, B[\tau \mapsto b.X]} \text{ (M-ATOMFL)}$$

$$\frac{B(\tau) = \langle x_v, v \rangle.b}{\text{IM}, B \xrightarrow{\tau:\epsilon} \text{IM}[x_v \mapsto v], B[\tau \mapsto b]} \text{ (M-PROPW)} \quad \frac{B(\tau) = X.b \quad X = \overline{x^i} \quad \forall i. \text{IM}(x_v^i) = v^i}{\text{IM}, B \xrightarrow{\tau:\epsilon} \text{IM}[x_s^i \mapsto v^i], B[\tau \mapsto b]} \text{ (M-PROPFL)}$$

with  $\text{rd}(\text{IM}, b, x_v) \triangleq \begin{cases} v & \text{if } \exists b_1, b_2. b = b_1.\langle x_v, v \rangle.b_2 \wedge \forall v'. \langle x_v, v' \rangle \notin b_2 \\ \text{IM}(x_v) & \text{otherwise} \end{cases}$

---

Fig. 7. The  $\text{lx86}_{\text{sim}}$  program transitions (above); the  $\text{lx86}_{\text{sim}}$  storage transitions (below)

(P-ATOMFL). The volatile-to-synchronous propagation of  $c_s$  is modelled by  $(S, x_s, v)$  in (P-ATOMS); *mutatis mutandis* for (P-ATOMP).

$$\begin{array}{c}
\text{COM} \vdash \text{CONF} \xrightarrow{\text{TID} \times \text{LAB} \cup \{\epsilon, \zeta\}} \text{CONF} \\
\frac{\text{C}, \text{S} \xrightarrow{\tau:\epsilon} \text{C}', \text{S}'}{\Delta \vdash \text{C}, \text{S}, \text{IM}, \text{B} \xrightarrow{\tau:\epsilon} \text{C}', \text{S}', \text{IM}, \text{B}} \text{(SILENTP)} \\
\frac{\text{IM}, \text{B} \xrightarrow{\tau:\epsilon} \text{IM}', \text{B}'}{\Delta \vdash \text{C}, \text{S}, \text{IM}, \text{B} \xrightarrow{\tau:\epsilon} \text{C}, \text{S}, \text{IM}', \text{B}'} \text{(SILENTS)} \\
\text{CONF} \triangleq \text{COM} \times \text{STACK} \times \text{IMEM} \times \text{IBMAP} \\
\frac{\text{C}, \text{S} \xrightarrow{\tau:l} \text{C}', \text{S}' \quad \text{IM}, \text{B} \xrightarrow{\tau:l} \text{IM}', \text{B}'}{\Delta \vdash \text{C}, \text{S}, \text{IM}, \text{B} \xrightarrow{\tau:l} \text{C}', \text{S}', \text{IM}', \text{B}'} \text{(STEP)} \\
\frac{\text{B}_0 \triangleq \lambda \tau. \epsilon \quad \text{IM}' = \text{IM}[\overline{x_s \mapsto \text{IM}(x_p)}][\overline{x_v \mapsto \text{IM}(x_p)}]}{\text{C}_{\text{rec}} \vdash \text{C}, \text{S}, \text{IM}, \text{B} \xrightarrow{\zeta} \text{C}_{\text{rec}}, \text{S}_0, \text{IM}', \text{B}_0} \text{(CRASH)}
\end{array}$$

Fig. 8. The  $\text{Ix86}_{\text{sim}}$  operational semantics

**Storage Transitions.** The  $\text{Ix86}_{\text{sim}}$  storage transitions are given at the bottom of Fig. 7 and are of the form:  $\text{IM}, \text{B} \xrightarrow{\tau:l} \text{IM}', \text{B}'$ , where  $\text{IM}, \text{IM}'$  denote the instrumented memory, and  $\text{B}, \text{B}'$  denote the *buffer map*, associating each thread with its *buffer*. Each buffer entry may be of the form: (1)  $\langle x_v, v \rangle$ , denoting a pending write of value  $v$  on  $x_v$ ; or (2)  $X \subseteq \text{Loc}$ , denoting a pending **flush** on cache line  $X$ . When a thread writes  $v$  to  $x_v$ , this is recorded in its buffer as the  $\langle x_v, v \rangle$  entry, as described by (M-WRITE). Similarly, when a thread makes a (FL,  $X$ ) transition (i.e. executes **persist**  $X$  translated from **flush**  $x$  with  $x \in X$ ), this is recorded in its buffer as  $X$ , as shown in (M-ATOMFL). Recall that when a thread reads from  $x_v$ , it first consults its own buffer, followed by the memory (if no write to  $x_v$  is found in the buffer). This lookup chain is captured by  $\text{rd}(\text{IM}, \text{b}, x_v)$  in the premise of (M-READ).

The (M-MF) rule ensures that an **mfence** proceeds only when the buffer of the executing thread is empty, as stipulated by the  $\text{B}(\tau)=\epsilon$  premise. In the (M-RMW) rule, when executing an RMW instruction on  $x_v$ , a similar lookup chain is followed to determine the value of  $x_v$ , as with a read. To ensure their atomicity, RMW instructions may only proceed when the buffer of the executing thread is drained. Moreover, the resulting update is committed directly to the memory, bypassing the thread buffer. This is to ensure that the resulting update is immediately visible to other threads.

As with (P-ATOMS), (M-ATOMS) describes the non-deterministic copying of  $x_v$  to  $x_s$ , with the result written directly to memory, provided that the buffer of the executing thread (i.e.  $\tau_s$ ) is empty; (P-ATOMP) behaves analogously. Lastly, (M-PROPW) describes the debuffering of a pending write in a thread-local buffer and propagating it to memory. Similarly, (M-PROPFL) describes the debuffering of a pending flush on  $X$ , where the value of  $x_v^i$  (in  $\text{IM}(x_v^i)$ ) is propagated to  $\text{IM}(x_s^i)$ , for each  $x^i \in X$ .

**Combined Transitions.** The  $\text{Ix86}_{\text{sim}}$  operational semantics in Fig. 8 is defined by combining the program and storage transitions, under a *recovery program*,  $\text{C}_{\text{rec}}$ , run after a crash. The (SILENTP) rule describes the case when the program subsystem takes a silent step and thus the storage subsystem is left unchanged; similarly for (SILENTS). The (STEP) rule describes the case when the program and storage subsystems both take the *same* transition (with the same label) and thus the transition effect is that of their combined effects. Lastly, the (CRASH) rule describes the case when the program crashes: the registers (in  $\text{S}$ ) and the thread-local buffers are lost (as they are volatile), and are thus reset; the memory is left largely unchanged (as it is non-volatile); and the execution is restarted with the recovery program. Note that upon a crash the persistent versions in the resulting memory ( $\text{IM}'$ ) remain unchanged, while the synchronous and volatile versions are lost and are simply overwritten by the persistent versions. This is because if  $x_v \neq x_p$  (resp.  $x_s \neq x_p$ ) upon a crash, then intuitively the write responsible for the current value of  $x_v$  (resp.  $x_s$ ) has not yet reached the persistent memory and is thus lost after recovery.

**$Ix86_{sim}$  Subsumes  $Px86_{sim}$ .** In [Thm. 1](#) below we show that our  $Ix86_{sim}$  model is weaker than  $Px86_{sim}$  and subsumes all its behaviours. This then allows us to establish the soundness of POG over the simpler  $Ix86_{sim}$  model. To prove that  $Ix86_{sim}$  subsumes  $Px86_{sim}$ , we show that for all non-instrumented memories  $M$  produced by a  $Px86_{sim}$  trace, there exists an instrumented memory  $IM'$  produced by an  $Ix86_{sim}$  trace such that  $M$  and  $IM$  agree on the (persisted) values of all locations.

**Theorem 1.** *For all memories  $M$  produced by a  $Px86_{sim}$  trace, there exists an instrumented memory  $IM'$  produced by an  $Ix86_{sim}$  trace such that:  $\forall x. M(x)=IM(x_p)$ .*

## 6 POG SOUNDNESS

We prove that the POG proof rules in [Fig. 3](#) are sound. We proceed with several auxiliary definitions.

**Assertion Semantics.** We define the set of *states* as  $STATE \triangleq IMEM \times STACK$ , and use  $\sigma$  as a metavariable for states. Assertions are interpreted as sets of states as expected: the instrumented memory provides the interpretation of the  $x_v, x_s, x_p$  constants, and the stack provides the interpretation of the  $a$  constants. In what follows we write  $\llbracket P \rrbracket$  for  $\{\sigma \mid \sigma \models P\}$ .

**$\mathcal{R}/\mathcal{G}$  Semantics.**  $\mathcal{R}/\mathcal{G}$  components are interpreted as relations on states. The guarantee  $\mathcal{G}$  is interpreted point-wise as the smallest preorder that admits all constituent guarded assignments. That is, if  $\langle x_v, e, P \rangle \in \mathcal{G}$  and  $(IM, S) \in \llbracket P \rrbracket$ , then  $((IM, S), (IM', S)) \in \llbracket \mathcal{G} \rrbracket^g$ , when  $IM' = IM[x_v \mapsto S(e)]$ :

$$\llbracket \emptyset \rrbracket^g \triangleq \text{id} \quad \llbracket \mathcal{G} \cup \{\langle x_v, e, P \rangle\} \rrbracket^g \triangleq (\llbracket \mathcal{G} \rrbracket^g \cup \{((IM, S), (IM[x_v \mapsto S(e)], S)) \mid (IM, S) \models P\})^*$$

Dually, the rely  $\mathcal{R}$  is interpreted point-wise as the largest relation on states that preserves the stability of all constituent assertions. That is, if  $P \in \mathcal{R}$ ,  $\sigma \in \llbracket P \rrbracket$  and  $(\sigma, \sigma') \in \llbracket \mathcal{R} \rrbracket^r$ , then  $\sigma' \in \llbracket P \rrbracket$ :

$$\llbracket \emptyset \rrbracket^r \triangleq STATE \times STATE \quad \llbracket \mathcal{R} \cup \{P\} \rrbracket^r \triangleq \llbracket \mathcal{R} \rrbracket^r \cap \{(\sigma, \sigma') \mid \sigma \models P \Rightarrow \sigma' \models P\}$$

**Configuration Safety.** We define the semantics of POG triples via an auxiliary predicate,  $\text{safe}_n(C^{sp}, \sigma, R, G, Q, I)$ , stating that executing the translated program  $C^{sp}$  over the  $\sigma$  state is safe with respect to the interpreted  $R, G$  relations, post-states  $Q$  and invariant  $I$  for up to  $n$  steps.

**Definition 3** (Configuration safety). For all  $C, IM, S$  and  $R, G \subseteq STATE \times STATE$ ,  $Q \subseteq STATE$ :  $\text{safe}_0(C^{sp}, (IM, S), R, G, Q, I)$  always holds; and  $\text{safe}_{n+1}(C^{sp}, (IM, S), R, G, Q, I)$  holds iff:

- (1)  $(IM, S) \in I$  and
- (2) if  $C = C_{\text{skip}}$ , then  $(IM, S) \in Q$ , where  $C_{\text{skip}} \triangleq \lambda \tau. \text{skip}$
- (3) for all  $\sigma$ : if  $((IM, S), \sigma) \in R \cup A_{\text{sp}}$ , then  $\text{safe}_n(C^{sp}, \sigma, R, G, Q, I)$
- (4) for all  $\tau, l, IM_1, IM_2, B_1, B_2, C', IM', S', l \neq \perp$ :

$$\text{if } C, IM_1, B_1, S \xrightarrow{\tau, l} C', IM_2, B_2, S' \wedge IM = \Downarrow(IM_1, B_1, \tau) \wedge IM' = \Downarrow(IM_2, B_2, \tau)$$

$$\text{then } ((IM, S), (IM', S)) \in G \cup A_{\text{sp}} \wedge \text{safe}_n((C')^{sp}, (IM', S'), R, G, Q, I)$$

where  $\Downarrow(IM, B, \tau) = IM' \stackrel{\text{def}}{\Leftrightarrow} (IM, B) \xrightarrow{\tau, \epsilon}^* IM', B' \wedge B'(\tau) = \epsilon$

$$\text{and } A_{\text{sp}} \triangleq \left\{ ((IM, S), (IM[x_p \mapsto IM(x_s)], S)), ((IM, S), (IM[y_s \mapsto y_p], S)) \mid y \in \text{Loc} \right\}^*$$

Recall that a translated program (via  $\llbracket \cdot \rrbracket$ ) is of the form  $C^{sp} \triangleq C \parallel c_s \parallel c_p$ . A configuration is always safe for zero steps. For  $n+1$  steps, a configuration is safe if: (1)  $(IM, S)$  is a state in the invariant  $I$ ; (2) whenever the program has finished execution (i.e.  $C = C_{\text{skip}}$ ), then  $(IM, S)$  must be a post-state in  $Q$ ; (3) whenever the environment changes the state according to the rely (in  $R$ ) or performs a (volatile-to-synchronous or synchronous-to-persistent) propagation (in  $A_{\text{sp}}$ ), then the resulting configuration remains safe for a further  $n$  steps; and (4) whenever the thread performs an  $Ix86_{sim}$  transition, its changes to the state are either those permitted by the guarantee (in  $G$ ) or those of propagation (in  $A_{\text{sp}}$ ), and the new configuration remains safe for  $n$  more steps.

Note that the  $\text{Ix86}_{\text{sim}}$  transitions (Fig. 8) are over an instrumented memory and a buffer map. As such, the memory  $\text{IM}$  in the pre-state of the thread  $\tau$  describes several storage pairs of the form  $(\text{IM}_1, \text{B}_1)$  such that  $\Downarrow(\text{IM}_1, \text{B}_1, \tau) = \text{IM}$ . That is, once the pending entries of  $\tau$  in  $\text{B}_1(\tau)$  are propagated to  $\text{IM}_1$  (via  $\xrightarrow{\tau:\epsilon}$  storage transitions in rule (M-PROP $\text{W}$ ) of Fig. 8), the resulting memory corresponds to  $\text{IM}$ . Intuitively,  $\text{IM}$  denotes the *view* of  $\tau$  of the storage subsystem, in that  $\tau$  observes its pending writes and flushes in  $\text{B}_1(\tau)$ , even though they have not yet reached the memory.

We next define *valid POG judgements*, and show that POG triples (Fig. 3) yield valid judgements. Note that the invariant of a triple,  $Q_p$ , is obtained from the postcondition  $Q$  by erasing the values of registers as well as the volatile/synchronous versions, and replacing them with arbitrary values.

**Definition 4.** A judgement  $\langle \mathcal{R}; \mathcal{G} \rangle \Vdash \{P\} C \{Q\}$  is valid iff for all  $n, \sigma \in \llbracket P \rrbracket$ ,  $\text{safe}_n(\llbracket C \rrbracket, \sigma, \llbracket \mathcal{R} \rrbracket^r, \llbracket \mathcal{G} \rrbracket^g, \llbracket Q \rrbracket, \llbracket Q_p \rrbracket)$  holds with  $Q_p \triangleq \exists \bar{v}_a, \bar{v}_v, \bar{v}_s. Q[\bar{v}_a/\bar{a}][\bar{v}_v/\bar{x}_v][\bar{v}_s/\bar{x}_s]$ . A judgement  $\text{C}_{\text{rec}} \Vdash \{P\} C \{Q\}$  is valid iff for all  $(\text{IM}, \text{S}) \in \llbracket P \rrbracket$ ,  $\text{IM}', \text{S}'$ , if  $\text{C}_{\text{rec}} \vdash C, \text{S}, \text{IM}, \text{B}_0 \Rightarrow^* \text{C}_{\text{skip}}, \text{S}', \text{IM}', \text{B}_0$ , then  $(\text{IM}', \text{S}') \in \llbracket Q \rrbracket$ .

**Theorem 2** (POG soundness). For all POG triples  $\langle \mathcal{R}; \mathcal{G} \rangle \Vdash \{P\} C \{Q\}$  and recovery triples  $\text{C}_{\text{rec}} \vdash \{P\} C \{Q\}$  in Fig. 3, the POG judgements  $\langle \mathcal{R}; \mathcal{G} \rangle \Vdash \{P\} C \{Q\}$  and  $\text{C}_{\text{rec}} \Vdash \{P\} C \{Q\}$  are valid.

Finally, we establish the following *adequacy* theorem showing that POG is sound with respect to  $\text{Px86}_{\text{sim}}$ . Note that it is sufficient to consider *closed* triples with RG contexts  $\langle \top; \top \rangle$ , i.e. whole programs (not run in parallel with another program) under an empty environment.

**Theorem 3** (Adequacy). For all  $\langle \top; \top \rangle \Vdash \{P\} C \{Q\}$ ,  $(\text{IM}, \text{S}) \models P$ , and  $\text{M}, \text{M}', \text{S}'$ , if  $\text{IM}$  and  $\text{M}$  agree (i.e.  $\forall x. \text{M}(x) = \text{IM}(x_v) = \text{IM}(x_s) = \text{IM}(x_p)$ ) and starting from  $(\text{M}, \text{S})$  with empty (thread-local and persistent) buffers, program  $C$  runs to completion under  $\text{Px86}_{\text{sim}}$  and yields  $(\text{M}', \text{S}')$  and empty buffers, then there exists  $\text{IM}'$  such that  $\text{IM}'$  and  $\text{M}'$  agree and  $(\text{IM}', \text{S}') \models Q$ .

## 7 A TWO-STEP TRANSFORMATION FOR ELIMINATING $\text{flush}_{\text{opt}}$ INSTRUCTIONS

We describe the most common use-case of  $\text{flush}_{\text{opt}}$  instructions, *epoch persistency*, where the use of  $\text{flush}_{\text{opt}}$  (rather than  $\text{flush}$ ) may prove advantageous for performance. We observe that programs using  $\text{flush}_{\text{opt}}$  for epoch persistency follow a certain pattern, and describe how we transform such programs to ones that use  $\text{flush}$  instead, without altering their persistency behaviour.

**Epoch Persistency using  $\text{flush}_{\text{opt}}$ .** Recall from §2 that  $\text{flush}_{\text{opt}}$  instructions provide weaker ordering constraints and can be reordered e.g. with respect to writes on different cache lines. As such, in order to mitigate the weak ordering constraints on  $\text{flush}_{\text{opt}}$  and to ensure their execution by a particular program point, they are typically followed by an **sfence**/**mfence**/RMW in program text; see e.g. Fig. 1d. Indeed, Fig. 1d is an example of *epoch persistency*. More concretely, by combining  $\text{flush}_{\text{opt}}$  and **sfence**/**mfence**/RMW instructions, one can divide an execution into distinct *epochs*, where the writes in each epoch may persist in an arbitrary order, while the writes of earlier (in program order) epochs persist before those in later epochs. This is illustrated in Fig. 9a, where  $x := 1$  and  $y := 1$  both persist before  $z := 1$ , whereas  $x := 1$  and  $y := 1$  themselves may persist in either order; i.e. the write on  $x, y$  persist in the first epoch before the write on  $z$  in the second epoch.

Let  $L$  and  $L'$  denote locations to be persisted in two consecutive epochs; one can then enforce epoch persistency by following the pattern below in three steps:

- (1) executing the writes on  $L$  and their corresponding  $\text{flush}_{\text{opt}}$  for each cache line, provided that each  $\text{flush}_{\text{opt}}$  on a cache line  $X$  follows the writes in  $L$  on  $X$ ;
- (2) executing an *epoch barrier*, namely an **sfence**, **mfence** or RMW; and (EPOCH)
- (3) executing the writes on  $L'$ .

$c_1 \triangleq x := 1;$ $\mathbf{flush}_{\text{opt}} x;$ $y := 1;$ $\mathbf{flush}_{\text{opt}} y;$ $\mathbf{sfence};$ $z := 1;$ (a)	$c_2 \triangleq x := 1;$ $y := 1;$ $\mathbf{flush}_{\text{opt}} x;$ $\mathbf{flush}_{\text{opt}} y;$ $\mathbf{sfence};$ $z := 1;$ (b)	$c_3 \triangleq x := 1;$ $y := 1;$ $\mathbf{flush} x;$ $\mathbf{flush} y;$ $\mathbf{sfence};$ $z := 1;$ (c)	$c_i \parallel \begin{array}{l} a := z; \\ \mathbf{if} (a=1) \\ w := 1; \end{array}$ (d)
$\not\downarrow: z=1 \Rightarrow x=y=1$	$\not\downarrow: z=1 \Rightarrow x=y=1$	$\not\downarrow: z=1 \Rightarrow x=y=1$	$\not\downarrow: w=1 \Rightarrow x=y=1$

Fig. 9. Epoch persistency using  $\mathbf{flush}_{\text{opt}}$  (a); the program obtained from  $c_1$  by reordering  $\mathbf{flush}_{\text{opt}}x$  (b); the program obtained from  $c_2$  by converting  $\mathbf{flush}_{\text{opt}}$  to  $\mathbf{flush}$  (c); a common concurrent use-case of  $\mathbf{flush}_{\text{opt}}$  (d), where  $c_i \in \{c_1, c_2, c_3\}$ . The  $c_1, c_2, c_3$  programs all have the same persistency behaviour (observing the same values upon recovery), and can be used interchangeably in (d) without changing the persistency behaviour.

The combination of  $\mathbf{flush}_{\text{opt}}$  instructions in (1) and the epoch barrier in (2) ensures that the writes on  $L$  in (1) persist before those on  $L'$  in (3) (e.g.  $L=\{x, y\}$  and  $L'=\{z\}$  in Fig. 9a, assuming that  $x$  and  $y$  are in distinct cache lines). Note that in step (1) it is sufficient to execute one  $\mathbf{flush}_{\text{opt}}$  per cache line as  $\mathbf{flush}_{\text{opt}}$  persists all (earlier) pending writes on the same cache line. For instance, if  $y := 1$  in Fig. 9a is replaced with  $x' := 1$  (where  $x$  and  $x'$  are on the same cache line), step (1) may simply comprise  $x := 1; x' := 1; \mathbf{flush}_{\text{opt}} x$  (without a separate  $\mathbf{flush}_{\text{opt}} x'$ ). This epoch persistency pattern is used by e.g. Raad et al. [2020] to implement several persistent libraries.

Note that replacing each  $\mathbf{flush}_{\text{opt}}$  in Fig. 9a with  $\mathbf{flush}$  also achieves epoch persistency, albeit at a finer granularity (one write per epoch). This is because each  $\mathbf{flush}$  is ordered with respect to all writes and thus introduces a new epoch. Using  $\mathbf{flush}_{\text{opt}}$  thus admits more than one write per epoch, and may improve performance as it allows the writes in the same epoch to persist in any order.

**Two-Step Transformation.** Note that as  $\mathbf{flush}_{\text{opt}}$  instructions are not ordered with respect to writes on different cache lines, each  $\mathbf{flush}_{\text{opt}}$  in step (1) can be reordered after the writes on different cache lines *without changing the persistency behaviour*. For instance,  $\mathbf{flush}_{\text{opt}} x$  in  $c_1$  of Fig. 9a can be reordered past the  $y := 1$  write to obtain  $c_2$  in Fig. 9b, where  $c_1$  and  $c_2$  have equivalent persistency behaviours. As such, step (1) of the pattern above can further be split as: (1.a) executing the writes on  $L$ ; (1.b) executing  $\mathbf{flush}_{\text{opt}}$  for each  $L$  cache line.

Indeed, this intuition informs the first step of our transformation towards eliminating  $\mathbf{flush}_{\text{opt}}$ : given a program of the form  $C \triangleq c_a; \mathbf{flush}_{\text{opt}} x; c_b; c$  with  $x \in X$ , if  $c$  is the *first* epoch barrier following  $\mathbf{flush}_{\text{opt}} x$  (i.e.  $c_b$  contains no  $\mathbf{sfence}/\mathbf{mfence}/\mathbf{RMW}$ ), then one can rewrite  $C$  as  $C' \triangleq c_a; c_b; \mathbf{flush}_{\text{opt}} x; c$  without changing its persistency behaviour (i.e.  $C$  and  $C'$  yield the same values for all locations upon recovery), *provided that*  $c_b$  contains no writes on  $X$  and no read instructions. Note that  $c_b \triangleq y := 1$  in Fig. 9a simply meets the stipulated provisos. We shortly elaborate on these provisos on  $c_b$  in §7.1, and note that it is always possible to achieve epoch persistency in such a way that fulfils these conditions. Note that one can repeatedly apply this transformation to push down all  $\mathbf{flush}_{\text{opt}}$  in step (1) just before the epoch barrier in (2), thus splitting (1) as (1.a) and (1.b).

Next, we observe that once all  $\mathbf{flush}_{\text{opt}}$  have been pushed down before the epoch barrier, one can almost always replace each  $\mathbf{flush}_{\text{opt}}$  with a corresponding  $\mathbf{flush}$  without altering its persistency behaviour. This is thanks to the strong ordering between each  $\mathbf{flush}_{\text{opt}}$  and the subsequent epoch barrier. For instance,  $c_2$  in Fig. 9b can be transformed to  $c_3$  in Fig. 9c, while leaving its persistency behaviour unchanged. This rewriting constitutes the second and last step of our transformation for eliminating  $\mathbf{flush}_{\text{opt}}$  instructions. In §7.2 we elaborate on the only scenarios under which this rewriting may alter the persistency behaviour, and note that such scenarios do not arise realistically.

$x := 1;$ $\mathbf{flush}_{\text{opt}} x;$ $x' := 1;$ $\mathbf{sfence};$ $z := 1;$ (a)	$x := 1;$ $x' := 1;$ $\mathbf{flush}_{\text{opt}} x;$ $\mathbf{sfence};$ $z := 1;$ (b)	$x := 1;$ $\mathbf{flush}_{\text{opt}} x;$ $\mathbf{sfence};$ $x' := 1;$ $z := 1;$ (c)
$\checkmark: z=1 \Rightarrow x=1$ $z=1 \wedge x'=0 \checkmark$	$\checkmark: z=1 \Rightarrow x=x'=1$ $z=1 \wedge x'=0 \times$	$\checkmark: z=1 \Rightarrow x=1$ $z=1 \wedge x'=0 \checkmark$

Fig. 10. An epoch persistency example not following the (EPOCH) pattern (a); the program obtained from (a) by reordering  $\mathbf{flush}_{\text{opt}} x$  after  $x' := 1$  (where  $x, x' \in X$ ), thus altering its persistency behaviour (b); a program with equivalent persistency behaviour to (a) that follows the (EPOCH) pattern (c).

Finally, note that this two-step transformation can be used to eliminate  $\mathbf{flush}_{\text{opt}}$  instructions in concurrent programs by applying the transformation to each thread containing  $\mathbf{flush}_{\text{opt}}$ . This is illustrated in Fig. 9d, where e.g.  $c_1$  in the left thread can be first transformed to  $c_2$  and subsequently to  $c_3$ , while preserving the persistency behaviour of the concurrent program at each step.

### 7.1 Caveats of Reordering $\mathbf{flush}_{\text{opt}}$ after Later Instructions (Transformation Step 1)

Recall that in the first step of our transformation we push down a  $\mathbf{flush}_{\text{opt}} x$  instruction with  $x \in X$  just before the epoch barrier, provided that there are *no writes on  $X$*  and *no reads* between  $\mathbf{flush}_{\text{opt}} x$  and the epoch barrier, as otherwise this transformation alters the persistency behaviour.

In the case of writes on  $X$ , consider the example in Fig. 10a where  $x, x' \in X$  and the  $x' := 1$  write is between  $\mathbf{flush}_{\text{opt}} x$  and  $\mathbf{sfence}$ . Fig. 10b depicts the program obtained from Fig. 10a by reordering  $\mathbf{flush}_{\text{opt}} x$  after  $x' := 1$ . Recall that executing  $\mathbf{flush}_{\text{opt}} x$  persists *all earlier* writes on  $X$ ; as such  $\mathbf{flush}_{\text{opt}} x$  in Fig. 10b is guaranteed to persist both  $x := 1$  and  $x' := 1$  (i.e.  $z=1 \Rightarrow x=x'=1$  upon recovery), while  $\mathbf{flush}_{\text{opt}} x$  in Fig. 10a is guaranteed to persist only  $x := 1$  (i.e.  $z=1 \Rightarrow x=1$  upon recovery). The two programs may therefore have different persistency behaviours: it is possible to observe  $z=1 \wedge x'=0$  after recovery in Fig. 10a but not in Fig. 10b.

Note that Fig. 10a does not adhere to (EPOCH): either  $x := 1$  and  $x' := 1$  are to persist in the same epoch and the program should have been rewritten as in Fig. 10b, or they are to persist in separate epochs in which case the program could have been rewritten as in Fig. 10c. That is, the intended persistency behaviour of Fig. 10a is ambiguous, and it is better practice to rewrite it as either Fig. 10b or Fig. 10c, both of which adhere to (EPOCH). We thus argue that it is possible to achieve epoch persistency *without* an intervening write on  $X$  between  $\mathbf{flush}_{\text{opt}} x$  and its epoch barrier. As such, it is possible to apply our first transformation step while preserving the persistency behaviour. Observe that the first transformation step in both Fig. 10b and Fig. 10c is idempotent ( $\mathbf{flush}_{\text{opt}} x$  is already before  $\mathbf{sfence}$ ) and thus trivially preserves the persistency behaviour.

In the case of reads, consider the example in Fig. 11a where the  $a := x$  read is between  $\mathbf{flush}_{\text{opt}} x$  and  $\mathbf{sfence}$ . Fig. 11b depicts the program obtained from Fig. 11a by reordering  $\mathbf{flush}_{\text{opt}} x$  after  $a := x$ . Note that the right thread executes  $y := 1$  only when  $a := x$  reads 2 from  $x$ , which in turn implies that  $x := 2$  in the left thread is store-ordered after  $x := 1$  in the right (otherwise  $a := x$  would read 1 from  $x$ ). That is,  $y=1$  implies the following store order in both Fig. 11a and Fig. 11b:  $x := 1 \rightarrow x := 2 \rightarrow a := x$ . In Fig. 11b we further have  $a := x \rightarrow \mathbf{flush}_{\text{opt}} x$ . Put together, this ensures  $x := 1 \rightarrow x := 2 \rightarrow \mathbf{flush}_{\text{opt}} x$  when  $y=1$  in Fig. 11b. Consequently, executing  $\mathbf{flush}_{\text{opt}} x$  first persists  $x := 1$  and then persists  $x := 2$ , as executing  $\mathbf{flush}_{\text{opt}} x$  persists all pending writes on  $x$  in the store order. As such,  $z=y=1 \Rightarrow x=2$  upon recovery in Fig. 11b. By contrast, in Fig. 11a the  $a := x \rightarrow \mathbf{flush}_{\text{opt}} x$  order does not hold, and thus  $\mathbf{flush}_{\text{opt}} x$  is only guaranteed to persist  $x := 1$ , yielding  $z=y=1 \Rightarrow x \in \{1, 2\}$  upon recovery. The two programs may thus have different persistency behaviours: it is possible to observe  $z=y=x=1$  after recovery in Fig. 11a but not in Fig. 11b.

$x := 2;$	$x := 1;$ $\mathbf{flush}_{\text{opt}} x;$ $a := x;$ $\mathbf{if} (a=2) \ y := 1;$ $\mathbf{sfence};$ $z := 1;$	$x := 2;$	$x := 1;$ $a := x;$ $\mathbf{flush}_{\text{opt}} x;$ $\mathbf{if} (a=2) \ y := 1;$ $\mathbf{sfence};$ $z := 1;$	$x := 2;$	$x := 1;$ $\mathbf{flush}_{\text{opt}} x;$ $\mathbf{sfence};$ $a := x;$ $\mathbf{if} (a=2) \ y := 1;$ $z := 1;$
	(a)		(b)		(c)
$\zeta: z=y=1 \Rightarrow x \in \{1, 2\} \quad z=y=x=1 \checkmark$		$\zeta: z=y=1 \Rightarrow x=2 \quad z=y=x=1 \times$		$\zeta: z=y=1 \Rightarrow x \in \{1, 2\} \quad z=y=x=1 \checkmark$	

Fig. 11. An epoch persistency example not following the (EPOCH) pattern (a); the program obtained from (a) by reordering  $\mathbf{flush}_{\text{opt}} x$  after a read ( $a := x$ ), thus altering its persistency behaviour (b); a program with equivalent persistency behaviour to (a) that follows the (EPOCH) pattern (c).

Note that once again Fig. 11a does not adhere to (EPOCH) and its intended persistency behaviour is ambiguous. More concretely, either: (1)  $x, y$  are on different cache lines, in which case the absence of a corresponding  $\mathbf{flush}_{\text{opt}} y$  implies that  $y := 1$  is to persist in the epoch after  $x := 1$ , and is thus better practice to rewrite the program as in Fig. 11c; or (2)  $x, y$  are on the same cache line but  $y := 1$  is to persist in the epoch after  $x := 1$  and thus as in the previous case it is clearer to rewrite the program as in Fig. 11c; or (3)  $x, y$  are on the same cache line and  $y := 1$  is to persist in the same epoch as  $x := 1$ , in which case  $\mathbf{flush}_{\text{opt}} x$  must follow  $y := 1$ . That is, in all three cases it is possible to rewrite Fig. 11c such that adheres to (EPOCH) and avoids the intervening read between  $\mathbf{flush}_{\text{opt}}$  and the epoch barrier. We thus argue that it is possible to achieve epoch persistency *without* an intervening read between a  $\mathbf{flush}_{\text{opt}}$  and the epoch barrier, and thus it is often possible to apply our first transformation step while preserving the persistency behaviour. Finally, we prove that the first step of our transformation is sound in that it does not alter the persistency behaviour.

**Theorem 4** (Step 1 soundness). *Given  $C$  with  $C(\tau) = c_a; \mathbf{flush}_{\text{opt}} x; c_b; c$  and  $x \in X$ , if  $c$  is the first epoch barrier after  $\mathbf{flush}_{\text{opt}} x$  (i.e.  $c_b$  contains no  $\mathbf{sfence}/\mathbf{mfence}/\mathbf{RMW}$ ) and  $c_b$  contains no writes on  $X$  and no reads, then  $C$  and  $C'$  have equivalent persistency behaviours, where  $C' \triangleq C[\tau \mapsto c_a; c_b; \mathbf{flush}_{\text{opt}} x; c]$ .*

## 7.2 Caveats of Converting $\mathbf{flush}_{\text{opt}}$ to $\mathbf{flush}$ (Transformation Step 2)

Recall that once all  $\mathbf{flush}_{\text{opt}}$  instruction have been pushed down just before the epoch barrier, our second transformation step replaces each  $\mathbf{flush}_{\text{opt}}$  with a corresponding  $\mathbf{flush}$ . However, in the case of a *blind persist* this transformation may alter the persistency behaviour of the program.

A blind persist denotes a scenario where a persist operation on  $x$  is executed without a previous access (read/write) on  $x$ . An example of this is illustrated in Fig. 12a, where  $\mathbf{flush}_{\text{opt}} x$  is issued without any prior access on  $x$  (assuming  $x, y$  are on different cache lines). Fig. 12b depicts the program obtained from Fig. 12a by replacing  $\mathbf{flush}_{\text{opt}}$  with  $\mathbf{flush}$ . Note that the left thread executes  $w := 1$  only when  $a=2$ , which in turn implies that  $y := 1$  in the left thread is store-ordered before  $y := 2$  in the right (otherwise  $a=2$  would not be possible). That is,  $w=1$  implies the following store order in both Fig. 12a and Fig. 12b:  $x := 1 \rightarrow y := 1 \rightarrow y := 2$ . As  $\mathbf{flush}$  instructions are ordered with respect to *all* writes, in Fig. 12b we further have  $y := 2 \rightarrow \mathbf{flush} x$ . Put together, this ensures  $x := 1 \rightarrow \mathbf{flush} x$  when  $w=1$  in Fig. 12b, and thus executing  $\mathbf{flush} x$  persists  $x := 1$ ; i.e.  $w=z=1 \Rightarrow x=1$  upon recovery in Fig. 12b. By contrast, as  $\mathbf{flush}_{\text{opt}}$  instructions may be reordered with respect to writes on different cache lines, in Fig. 12a the  $y := 2 \rightarrow \mathbf{flush}_{\text{opt}} x$  order does not hold, and thus  $\mathbf{flush}_{\text{opt}} x$  may not persist  $x := 1$ , yielding  $w=z=1 \Rightarrow x \in \{0, 1\}$  upon recovery. The two programs may thus have different persistency behaviours: it is possible to observe  $w=z=1 \wedge x=0$  after recovery in Fig. 12a but not in Fig. 12b.

$c_1 \triangleq x := 1;$ $y := 1;$ $a := y;$ <b>if</b> ( $a=2$ ) $w := 1;$	$y := 2;$ $\mathbf{flush}_{\text{opt}} x;$ <b>sfence</b> ; $z := 1;$	$c_1$	$y := 2;$ <b>if</b> ( $c=1$ ) $\mathbf{flush}_{(\text{opt})} x;$ <b>sfence</b> ; $z := 1;$	$c_1$	$y := 2;$ $x := 2;$ $\mathbf{flush}_{(\text{opt})} x;$ <b>sfence</b> ; $z := 1;$
(a)	(b)	(c)	(d)	(e)	(f)
$\zeta: w=z=1 \Rightarrow x \in \{0, 1\}$ $w=z=1 \wedge x=0$ ✓	$\zeta: w=z=1 \Rightarrow x=1$ $w=z=1 \wedge x=0$ ✗	$\zeta: w=z=1 \Rightarrow x=1$ $w=z=1 \wedge x=0$ ✗	$\zeta: w=z=1 \Rightarrow x=2$ $w=z=1 \wedge x=0$ ✗	$\zeta: w=z=1 \Rightarrow x=1$ $w=z=1 \wedge x=0$ ✗	$\zeta: w=z=1 \Rightarrow x=2$ $w=z=1 \wedge x=0$ ✗

Fig. 12. A blind persist example with  $x \in X$ ,  $y \notin X$  (a); the program obtained from (a) by replacing  $\mathbf{flush}_{\text{opt}}$  with  $\mathbf{flush}$ , altering its persistency (b); more realistic examples analogous to (a) with non-blind persists, where  $\mathbf{flush}/\mathbf{flush}_{\text{opt}}$  can be used interchangeably (denoted by  $\mathbf{flush}_{(\text{opt})}$ ) without altering the persistency (c, d).

Note that blind persists are uncommon: persist instructions are expensive and are not typically issued without ascertaining that there is a corresponding write pending to be persisted. More concretely, prior to issuing a  $\mathbf{flush}_{\text{opt}}/\mathbf{flush}$  on  $x$ , the existence of pending writes on  $x$  is usually ascertained by either reading from  $x$ , or by having written to  $x$  earlier (in the same thread), as shown in Fig. 12c and Fig. 12d, respectively, where  $\mathbf{flush}_{(\text{opt})}$  denotes either  $\mathbf{flush}_{\text{opt}}$  or  $\mathbf{flush}$ . For instance, recall that in epoch persistency (EPOCH) each persist is non-blind as it is preceded by a write on the same cache line. In more realistic scenarios such as (EPOCH), Figs. 12c and 12d,  $\mathbf{flush}_{\text{opt}}$  instructions can thus be replaced with corresponding  $\mathbf{flush}$  without altering the persistency behaviour.

Finally, we prove that the second transformation step is sound (see the accompanying technical appendix for the definition of blind persists).

**Theorem 5** (Step 2 soundness). *Given  $C$  with  $C(\tau) = c_a; \mathbf{flush}_{\text{opt}} x; c_b$ , if  $\mathbf{flush}_{\text{opt}}$  is not blind and  $c$  is an epoch barrier, then  $C$  and  $C' \triangleq C[\tau \mapsto c'; \mathbf{flush} x; c]$  have equivalent persistency behaviours.*

## 8 CONCLUSIONS, RELATED AND FUTURE WORK

We presented POG, the *first* program logic for reasoning about persistency behaviours under the  $\text{Px86}_{\text{sim}}$  fragment that excludes  $\mathbf{flush}_{\text{opt}}$  instructions. We used POG to verify several representative examples. To establish the soundness of POG, we developed an intermediate operational model,  $\text{Ix86}_{\text{sim}}$ , which simplifies  $\text{Px86}_{\text{sim}}$  by forgoing its persistent buffer and modelling its persist orderings by tracking three different versions for each location. We demonstrated that  $\text{Ix86}_{\text{sim}}$  subsumes  $\text{Px86}_{\text{sim}}$  and emulates all its valid behaviours. We then proved that POG is sound with respect to  $\text{Ix86}_{\text{sim}}$  and thus also with respect to  $\text{Px86}_{\text{sim}}$ . As we note below (see future work),  $\text{Ix86}_{\text{sim}}$  is a valuable contribution in its own right, as it facilitates automated verification techniques for persistency behaviours. Finally, in order to extend the reasoning principles of POG to the full  $\text{Px86}_{\text{sim}}$  that also contains  $\mathbf{flush}_{\text{opt}}$  instructions, we presented a two-step transformation mechanism that allows us, in most cases, to rewrite a program using  $\mathbf{flush}_{\text{opt}}$  instructions, to an equivalent one that uses  $\mathbf{flush}$  instructions instead without altering its persistency behaviour. As such, to reason about a program  $C$  that uses  $\mathbf{flush}_{\text{opt}}$ , one can first use our transformation to rewrite  $C$  to a program  $C'$  with equivalent persistency behaviour, and then use POG to reason about  $C'$ .

We based POG on the OGRA program logic [Lahav and Vafeiadis 2015]. As such, since OGRA is proved sound for the release-acquire (RA) consistency model and RA is a weaker model than x86-TSO, POG is also sound for RA. Consequently, POG is incomplete for reasoning about x86-TSO in that it cannot be used to prove the absence of behaviours that are admissible under RA but not x86-TSO. For instance, given the program  $C \triangleq (x := 1; y := 1) \parallel (y := 2; x := 2)$ , we cannot use POG to prove that the final states of  $C$  exclude those in which  $x=1 \wedge y=2$  holds.



However, we note that almost all existing program logics in the literature, including those for the strong sequential consistency model, are incomplete (e.g. [Dinsdale-Young et al. 2010; Jones 1983; Jung et al. 2015; Kaiser et al. 2017; Lahav and Vafeiadis 2015; Nanevski et al. 2014; Raad et al. 2015; Svendsen et al. 2018; Turon et al. 2014; Vafeiadis and Narayan 2013]) as their main aim to enable simple high-level reasoning principles that apply to common patterns, albeit at the cost of compromising completeness for certain cases. Nevertheless, as shown in the literature, incompleteness can be remedied to some extent by including ghost state [Jacobs and Piessens 2011]. Specifically, it is possible to prove the example above using ghost state, provided that we also show that introducing such state is sound under  $\text{Ix86}_{\text{sim}}/\text{x86-TSO}$ , which we leave for future work.

**Related work.** Although existing literature on NVM has grown rapidly in recent years, formally verifying programs and algorithms that operate on NVM has largely remained unexplored. Friedman et al. [2018] developed several persistent queue implementations using the Intel-x86 persist instructions (e.g. **flush**); Zuriel et al. [2019] also developed two persistent set implementations using Intel-x86 persist instructions. Both [Friedman et al. 2018; Zuriel et al. 2019] argue that their implementations are correct by providing an informal argument at the level of program traces. Derrick et al. [2019] provided a formal correctness proof of the queue implementation by Friedman et al. [2018]; this proof is also at the level of program traces. Moreover, all three of [Derrick et al. 2019; Friedman et al. 2018; Zuriel et al. 2019] assume that the underlying memory model is sequential consistency (SC) [Lamport 1979], rather than Intel x86-TSO. Raad et al. [2019] recently developed a persistent transactional library on top of the ARM architecture; they later adapted this implementation to the  $\text{Px86}_{\text{sim}}$  architecture. In both cases they provide a formal proof of their implementation correctness on top of the corresponding architecture. Nevertheless, these proofs are low-level in that they operate at the level of execution traces, rather than the program syntax.

To our knowledge, no existing work provides a syntactic proof system for verifying the persistency guarantees of concurrent programs, especially in the presence of relaxed (out-of-order) or asynchronous persists. The most closely related work to ours are those of [Chen et al. 2015; Ntzik et al. 2015]. Chen et al. [2015] present crash Hoare logic (CHL) for reasoning about the crashing behaviour of the FSCQ file system. CHL is more restricted than POG in two ways. First, CHL can only be used for *sequential* programs and does not support concurrent reasoning. Second, in contrast to POG's support for explicit **flush** instructions, CHL does not support the Unix explicit persist instruction **fsync**. Ntzik et al. [2015] extend the Views framework [Dinsdale-Young et al. 2013] to support fault conditions. As an extension of Views, this work supports concurrency; however, their support for persistent reasoning is rather limited: (1) it assumes that the underlying memory model is SC and does not account for weak concurrent behaviours; (2) it does not distinguish between stores and persists, and thus assumes that all stores persist synchronously and in the store order; and consequently (3) does not support any explicit persist instructions such as **flush**.

**Future work.** We plan to build on our work here in several ways. First, we will use POG to verify existing implementations of persistent libraries and data structures such as [Friedman et al. 2018; Intel 2015; Zuriel et al. 2019]. Second, we will build automated techniques such as model checking (MC) for verifying persistency. We plan to do this by extending existing MC algorithms that already support x86-TSO (e.g. [Abdulla et al. 2015; Kokologiannakis et al. 2019a,b]) with the atomic propagation constructs of  $\text{Ix86}_{\text{sim}}$ . This will allow us to leverage cutting-edge MC tools to verify persistency with minimal implementation overhead. Lastly, building on the ideas underpinning POG, we will devise a similar program logic for reasoning about persistency under the ARMv8 architecture [Raad et al. 2019].

## ACKNOWLEDGMENTS

We thank the OOPSLA 2020 reviewers for their valuable feedback. Azalea Raad was supported in part by a European Research Council (ERC) Consolidator Grant for the project “RustBelt”, under the European Union Horizon 2020 Framework Programme (grant agreement number 683289). Ori Lahav was supported by the Israel Science Foundation (grant number 5166651), by Len Blavatnik and the Blavatnik Family foundation, and by the Alon Young Faculty Fellowship.

## REFERENCES

- Parosh Aziz Abdulla, Stavros Aronis, Mohamed Faouzi Atig, Bengt Jonsson, Carl Leonardsson, and Konstantinos Sagonas. 2015. Stateless Model Checking for TSO and PSO. In *Tools and Algorithms for the Construction and Analysis of Systems*, Christel Baier and Cesare Tinelli (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 353–367.
- Arm. 2018. ARM Architecture Reference Manual ARMv8, for ARMv8-A architecture profile (DDI 0487D.a). [https://static.docs.arm.com/ddi0487/da/DDI0487D\\_a\\_armv8\\_arm.pdf](https://static.docs.arm.com/ddi0487/da/DDI0487D_a_armv8_arm.pdf)
- Haogang Chen, Daniel Ziegler, Tej Chajed, Adam Chlipala, M. Frans Kaashoek, and Nikolai Zeldovich. 2015. Using Crash Hoare Logic for Certifying the FSCQ File System. In *Proceedings of the 25th Symposium on Operating Systems Principles* (Monterey, California) (*SOSP '15*). ACM, New York, NY, USA, 18–37. <https://doi.org/10.1145/2815400.2815402>
- Edmund Clarke, Daniel Kroening, and Flavio Lerda. 2004. A Tool for Checking ANSI-C Programs. In *Tools and Algorithms for the Construction and Analysis of Systems*, Kurt Jensen and Andreas Podelski (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 168–176.
- Jeremy Condit, Edmund B. Nightingale, Christopher Frost, Engin Ipek, Benjamin Lee, Doug Burger, and Derrick Coetzee. 2009. Better I/O Through Byte-addressable, Persistent Memory. In *Proceedings of the ACM SIGOPS 22Nd Symposium on Operating Systems Principles* (Big Sky, Montana, USA) (*SOSP '09*). ACM, New York, NY, USA, 133–146. <https://doi.org/10.1145/1629575.1629589>
- John Derrick, Simon Doherty, Brijesh Dongol, Gerhard Schellhorn, and Heike Wehrheim. 2019. Verifying Correctness of Persistent Concurrent Data Structures. In *Formal Methods – The Next 30 Years*, Maurice H. ter Beek, Annabelle McIver, and José N. Oliveira (Eds.). Springer International Publishing, Cham, 179–195.
- Thomas Dinsdale-Young, Lars Birkedal, Philippa Gardner, Matthew Parkinson, and Hongseok Yang. 2013. Views: Compositional Reasoning for Concurrent Programs. In *Proceedings of the 40th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages* (Rome, Italy) (*POPL '13*). ACM, New York, NY, USA, 287–300. <https://doi.org/10.1145/2429069.2429104>
- Thomas Dinsdale-Young, Mike Dodds, Philippa Gardner, Matthew J. Parkinson, and Viktor Vafeiadis. 2010. Concurrent Abstract Predicates. In *ECOOP 2010 – Object-Oriented Programming*, Theo D’Hondt (Ed.). Springer Berlin Heidelberg, Berlin, Heidelberg, 504–528.
- Michal Friedman, Maurice Herlihy, Virendra Marathe, and Erez Petrank. 2018. A Persistent Lock-free Queue for Non-volatile Memory. In *Proceedings of the 23rd ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming* (Vienna, Austria) (*PPoPP '18*). ACM, New York, NY, USA, 28–40. <https://doi.org/10.1145/3178487.3178490>
- Vaibhav Gogte, Stephan Diestelhorst, William Wang, Satish Narayanasamy, Peter M. Chen, and Thomas F. Wenisch. 2018. Persistency for Synchronization-free Regions. In *Proceedings of the 39th ACM SIGPLAN Conference on Programming Language Design and Implementation* (Philadelphia, PA, USA) (*PLDI 2018*). ACM, New York, NY, USA, 46–61. <https://doi.org/10.1145/3192366.3192367>
- V. Gogte, W. Wang, S. Diestelhorst, P. M. Chen, S. Narayanasamy, and T. F. Wenisch. 2020. Relaxed Persist Ordering Using Strand Persistency. In *2020 ACM/IEEE 47th Annual International Symposium on Computer Architecture (ISCA)*. 652–665.
- Shiyong Huang and Jeff Huang. 2016. Maximal Causality Reduction for TSO and PSO. In *Proceedings of the 2016 ACM SIGPLAN International Conference on Object-Oriented Programming, Systems, Languages, and Applications* (Amsterdam, Netherlands) (*OOPSLA 2016*). ACM, New York, NY, USA, 447–461. <https://doi.org/10.1145/2983990.2984025>
- Intel. 2015. Persistent Memory Programming. <http://pmem.io/>
- Intel. 2019. 3D XPoint. <https://www.intel.com/content/www/us/en/architecture-and-technology/intel-optane-technology.html>
- Intel. 2019. Intel 64 and IA-32 Architectures Software Developer’s Manual (Combined Volumes). <https://software.intel.com/sites/default/files/managed/39/c5/325462-sdm-vol-1-2abcd-3abcd.pdf> Order Number: 325462-069US.
- Joseph Izraelevitz, Hammurabi Mendes, and Michael L. Scott. 2016. Linearizability of Persistent Memory Objects Under a Full-System-Crash Failure Model. In *Distributed Computing*, Cyril Gavoille and David Ilcinkas (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 313–327.
- Joseph Izraelevitz, Jian Yang, Lu Zhang, Juno Kim, Xiao Liu, Amirsaman Memaripour, Yun Joon Soh, Zixuan Wang, Yi Xu, Subramanya R. Dulloor, Jishen Zhao, and Steven Swanson. 2019. Basic Performance Measurements of the Intel Optane

- DC Persistent Memory Module. arXiv:1903.05714 [cs.DC]
- Bart Jacobs and Frank Piessens. 2011. Expressive Modular Fine-Grained Concurrency Specification. *SIGPLAN Not.* 46, 1 (Jan. 2011), 271–282. <https://doi.org/10.1145/1925844.1926417>
- C. B. Jones. 1983. Tentative Steps Toward a Development Method for Interfering Programs. *ACM Trans. Program. Lang. Syst.* 5, 4 (Oct. 1983), 596–619. <https://doi.org/10.1145/69575.69577>
- Arpit Joshi, Vijay Nagarajan, Marcelo Cintra, and Stratis Viglas. 2015. Efficient Persist Barriers for Multicores. In *Proceedings of the 48th International Symposium on Microarchitecture (Waikiki, Hawaii) (MICRO-48)*. ACM, New York, NY, USA, 660–671. <https://doi.org/10.1145/2830772.2830805>
- Ralf Jung, David Swasey, Filip Sieczkowski, Kasper Svendsen, Aaron Turon, Lars Birkedal, and Derek Dreyer. 2015. Iris: Monoids and Invariants as an Orthogonal Basis for Concurrent Reasoning. In *Proceedings of the 42nd Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (Mumbai, India) (POPL '15)*. Association for Computing Machinery, New York, NY, USA, 637–650. <https://doi.org/10.1145/2676726.2676980>
- J. Kaiser, Hoang-Hai Dang, D. Dreyer, O. Lahav, and Viktor Vafeiadis. 2017. Strong Logic for Weak Memory: Reasoning About Release-Acquire Consistency in Iris. In *ECOOP*.
- T. Kawahara, K. Ito, R. Takemura, and H. Ohno. 2012. Spin-transfer torque RAM technology: Review and prospect. *Microelectronics Reliability* 52, 4 (2012), 613 – 627. <https://doi.org/10.1016/j.microrel.2011.09.028> Advances in non-volatile memory technology.
- Michalis Kokologiannakis, Azalea Raad, and Viktor Vafeiadis. 2019a. Effective Lock Handling in Stateless Model Checking. *Proc. ACM Program. Lang.* 3, OOPSLA, Article 173 (Oct. 2019), 26 pages. <https://doi.org/10.1145/3360599>
- Michalis Kokologiannakis, Azalea Raad, and Viktor Vafeiadis. 2019b. Model Checking for Weakly Consistent Libraries. In *Proceedings of the 40th ACM SIGPLAN Conference on Programming Language Design and Implementation (Phoenix, AZ, USA) (PLDI 2019)*. ACM, New York, NY, USA, 96–110. <https://doi.org/10.1145/3314221.3314609>
- Aasheesh Kolli, Vaibhav Gogte, Ali Saidi, Stephan Diestelhorst, Peter M. Chen, Satish Narayanasamy, and Thomas F. Wenisch. 2017. Language-level Persistency. In *Proceedings of the 44th Annual International Symposium on Computer Architecture (Toronto, ON, Canada) (ISCA '17)*. ACM, New York, NY, USA, 481–493. <https://doi.org/10.1145/3079856.3080229>
- Aasheesh Kolli, Jeff Rosen, Stephan Diestelhorst, Ali Saidi, Steven Pelley, Sihang Liu, Peter M. Chen, and Thomas F. Wenisch. 2016. Delegated Persist Ordering. In *The 49th Annual IEEE/ACM International Symposium on Microarchitecture (Taipei, Taiwan) (MICRO-49)*. IEEE Press, Piscataway, NJ, USA, Article 58, 13 pages. <http://dl.acm.org/citation.cfm?id=3195638.3195709>
- Ori Lahav and Viktor Vafeiadis. 2015. Owicki-Gries Reasoning for Weak Memory Models. In *Automata, Languages, and Programming*, Magnús M. Halldórsson, Kazuo Iwama, Naoki Kobayashi, and Bettina Speckmann (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 311–323.
- Ori Lahav, Viktor Vafeiadis, Jeehoon Kang, Chung-Kil Hur, and Derek Dreyer. 2017. Repairing Sequential Consistency in C/C++11. In *Proceedings of the 38th ACM SIGPLAN Conference on Programming Language Design and Implementation (Barcelona, Spain) (PLDI 2017)*. ACM, New York, NY, USA, 618–632. <https://doi.org/10.1145/3062341.3062352>
- Leslie Lamport. 1979. How to Make a Multiprocessor Computer That Correctly Executes Multiprocess Programs. *IEEE Trans. Computers* 28, 9 (Sept. 1979), 690–691. <https://doi.org/10.1109/TC.1979.1675439>
- Benjamin C. Lee, Engin Ipek, Onur Mutlu, and Doug Burger. 2009. Architecting Phase Change Memory As a Scalable Dram Alternative. In *Proceedings of the 36th Annual International Symposium on Computer Architecture (Austin, TX, USA) (ISCA '09)*. ACM, New York, NY, USA, 2–13. <https://doi.org/10.1145/1555754.1555758>
- Aleksandar Nanevski, Ruy Ley-Wild, Ilya Sergey, and Germán Andrés Delbianco. 2014. Communicating State Transition Systems for Fine-Grained Concurrent Resources. In *Programming Languages and Systems*, Zhong Shao (Ed.). Springer Berlin Heidelberg, Berlin, Heidelberg, 290–310.
- Faisal Nawab, Joseph Izraelevitz, Terence Kelly, Charles B. Morrey, Dhruva R. Chakrabarti, and Michael James Scott. 2017. Dali: A Periodically Persistent Hash Map. In *DISC*.
- Gian Ntzik, Pedro da Rocha Pinto, and Philippa Gardner. 2015. Fault-Tolerant Resource Reasoning. In *Programming Languages and Systems*, Xinyu Feng and Sungwoo Park (Eds.). Springer International Publishing, Cham, 169–188.
- Susan Owicki and David Gries. 1976. An axiomatic proof technique for parallel programs I. *Acta Informatica* 6, 4 (01 Dec 1976), 319–340. <https://doi.org/10.1007/BF00268134>
- Steven Pelley, Peter M. Chen, and Thomas F. Wenisch. 2014. Memory Persistency. In *Proceeding of the 41st Annual International Symposium on Computer Architecture (Minneapolis, Minnesota, USA) (ISCA '14)*. IEEE Press, Piscataway, NJ, USA, 265–276. <http://dl.acm.org/citation.cfm?id=2665671.2665712>
- Azalea Raad and Viktor Vafeiadis. 2018. Persistence Semantics for Weak Memory: Integrating Epoch Persistency with the TSO Memory Model. *Proc. ACM Program. Lang.* 2, OOPSLA, Article 137 (Oct. 2018), 27 pages. <https://doi.org/10.1145/3276507>
- Azalea Raad, Jules Villard, and Philippa Gardner. 2015. CoLoSL: Concurrent Local Subjective Logic. In *Proceedings of the 24th European Symposium on Programming (ESOP'15) (Lecture Notes in Computer Science)*, Jan Vitek (Ed.), Vol. 9032. Springer, 710–735.

- Azalea Raad, John Wickerson, Gil Neiger, and Viktor Vafeiadis. 2020. Persistency Semantics of the Intel-x86 Architecture. *Proc. ACM Program. Lang.* 4, POPL, Article 11 (Jan. 2020), 31 pages. <https://doi.org/10.1145/3371079>
- Azalea Raad, John Wickerson, and Viktor Vafeiadis. 2019. Weak Persistency Semantics from the Ground Up: Formalising the Persistency Semantics of ARMv8 and Transactional Models. *Proc. ACM Program. Lang.* 3, OOPSLA, Article 135 (Oct. 2019), 27 pages. <https://doi.org/10.1145/3360561>
- Peter Sewell, Susmit Sarkar, Scott Owens, Francesco Zappa Nardelli, and Magnus O. Myreen. 2010. X86-TSO: A Rigorous and Usable Programmer’s Model for x86 Multiprocessors. *Commun. ACM* 53, 7 (July 2010), 89–97. <https://doi.org/10.1145/1785414.1785443>
- Filip Sieczkowski, Kasper Svendsen, Lars Birkedal, and Jean Pichon-Pharabod. 2015. A Separation Logic for Fictional Sequential Consistency. In *Programming Languages and Systems*, Jan Vitek (Ed.). Springer Berlin Heidelberg, Berlin, Heidelberg, 736–761.
- D. B. Strukov, G. S. Snider, D. R. Stewart, and R. S. Williams. 2008. The missing memristor found. *Nature* 453 (2008), 80 – 83.
- Kasper Svendsen, Jean Pichon-Pharabod, Marko Doko, Ori Lahav, and Viktor Vafeiadis. 2018. A Separation Logic for a Promising Semantics. In *Programming Languages and Systems*, Amal Ahmed (Ed.). Springer International Publishing, Cham, 357–384.
- Aaron Turon, Viktor Vafeiadis, and Derek Dreyer. 2014. GPS: Navigating Weak Memory with Ghosts, Protocols, and Separation. In *Proceedings of the 2014 ACM International Conference on Object Oriented Programming Systems Languages & Applications* (Portland, Oregon, USA) (OOPSLA ’14). ACM, New York, NY, USA, 691–707. <https://doi.org/10.1145/2660193.2660243>
- Viktor Vafeiadis and Chinmay Narayan. 2013. Relaxed Separation Logic: A Program Logic for C11 Concurrency. In *Proceedings of the 2013 ACM SIGPLAN International Conference on Object Oriented Programming Systems Languages & Applications*. 867–884.
- Yoav Zuriel, Michal Friedman, Gali Sheffi, Nachshon Cohen, and Erez Petrank. 2019. Efficient Lock-free Durable Sets. *Proc. ACM Program. Lang.* 3, OOPSLA, Article 128 (Oct. 2019), 26 pages. <https://doi.org/10.1145/3360554>