



A Formalization of Java's Concurrent Access Modes

JOHN BENDER, University of California Los Angeles, USA

JENS PALSBERG, University of California Los Angeles, USA

Java's memory model was recently updated and expanded with new access modes. The accompanying documentation for these access modes is intended to make strong guarantees about program behavior that the Java compiler must enforce, yet the documentation is frequently unclear. This makes the intended program behavior ambiguous, impedes discussion of key design decisions, and makes it impossible to prove general properties about the semantics of the access modes.

In this paper we present the first formalization of Java's access modes. We have constructed an axiomatic model for all of the modes using the Herd modeling tool. This allows us to give precise answers to questions about the behavior of example programs, called litmus tests. We have validated our model using a large suite of litmus tests from existing research which helps to shed light on the relationship with other memory models. We have also modeled the semantics in Coq and proven several general theorems including a DRF guarantee, which says that if a program is properly synchronized then it will exhibit sequentially consistent behavior. Finally, we use our model to prove that the unusual design choice of a partial order among writes to the same location is unobservable in any program.

CCS Concepts: • **Theory of computation** → **Program semantics**; • **Software and its engineering** → **Concurrent programming languages**.

Additional Key Words and Phrases: Memory Models

ACM Reference Format:

John Bender and Jens Palsberg. 2019. A Formalization of Java's Concurrent Access Modes. *Proc. ACM Program. Lang.* 3, OOPSLA, Article 142 (October 2019), 28 pages. <https://doi.org/10.1145/3360568>

1 INTRODUCTION

The original Java memory model [Manson et al. 2005] was an early attempt to define the semantics of concurrent shared memory programs running on the Java platform, but the definitions were hard to understand and there was no easy way to check the behavior of example programs. It was also later discovered that it ruled out existing compiler optimizations which it claimed to support [Ševčík and Aspinall 2008]. Since then, researchers have made great advances in memory model design while studying other weak memory models like those for ARM [Alglave et al. 2008; Pulte et al. 2017], C11 [Batty et al. 2011; Kang et al. 2017; Lahav et al. 2017; Vafeiadis et al. 2015], Power [Alglave et al. 2014], and x86 [Owens et al. 2009].

Recently, the ninth version of the Java Development Kit updated and expanded Java's memory model using new "access modes". Though the design of the access modes is inspired by C11's memory orders [C11 2010], it differs in a few key ways. First, it sheds complicated legacy features like release sequences and release-consume accesses. Second, it includes a broad but simple mechanism to forbid so called "out of thin-air" behavior [Batty and Sewell 2014]. Finally, it makes no provision

Authors' addresses: John Bender, Computer Science, University of California Los Angeles, USA, johnbender@cs.ucla.edu; Jens Palsberg, Computer Science, University of California Los Angeles, USA, palsberg@ucla.edu.



This work is licensed under a Creative Commons Attribution 4.0 International License.

© 2019 Copyright held by the owner/author(s).

2475-1421/2019/10-ART142

<https://doi.org/10.1145/3360568>

for a total order on writes to the same location. Taken together this suggests new opportunities to use a simpler model, develop metatheory, and verify lock-free algorithms for the Java platform.

However, the documentation [JDK9 2017; Lea 2017, 2018] is frequently ambiguous. This makes it extremely difficult to provide definitive answers about program behavior and there is little hope of proving important properties about the semantics. Further, it impedes the discussion of key features of the model's design.

To address these issues, we present the first formalization of Java's access modes. Critically, our model is *precise* and *complete*, formalizing all of the main features of Java's access modes. Additionally, we have endeavoured to make the model as *readable* as possible, defining all modes and fences as small extensions to an intuitive notion of visibility. Specifically, we make the following contributions:

- an axiomatic model for all of Java's access modes, fences, and atomic read-writes;
- an instantiation of our model for the Herd tool, which allowed us to examine the outcomes for more than 80 test programs; and
- an instantiation of our model in Coq, which allowed us to prove three key theorems about the semantics that are novel for Java's access modes.

We have constructed our model using the cat language. This allowed us to leverage the Herd tool [Algave et al. 2014] to give definitive answers to questions about example programs, called litmus tests. Litmus tests are designed to highlight specific behaviors in memory models. Herd enumerates all possible executions for a litmus test and determines which are allowed according to the model. Then, if at least one execution is allowed, the behavior illustrated by the litmus test is allowed by the model.

We used Herd with more than 80 litmus tests drawn from prior research to help validate our model by comparing it to ARMv8, C11 and x86. Our goal for these comparisons is to show that there are no unexpected differences in behavior. For example, the conventional wisdom is that language memory models will exhibit more behaviors than architecture memory models due to aggressive compiler optimizations. Thus, Java's access modes should permit more behaviors than ARMv8. If that is not the case then it should be due to a deliberate design decision and not a bug in the definitions. For all 80 litmus tests our model behaves according to our expectations and the documented design of the access modes.

We have also formalized our model in Coq and we prove three key theorems: absence of causal cycles when all reads are "release" reads, sequentially consistent semantics under proper synchronization (DRF), and a guarantee that each stronger mode admits fewer executions [Vafeiadis et al. 2015]. These theorems further validate the definitions of our model, give more evidence that the model is complete with respect to the documentation, and clearly demonstrate that the semantics is suitable for formal reasoning.

Finally, the partial order on same-location writes in Java's access modes represents a significant departure from the conventions of existing memory models. We show that the impact of switching to a partial order is "unobservable" in any example program executed using our model and thereby show that Java's access modes can adopt a total order on writes to the same location. The simplicity of our model shines through in our proof, which makes it clear that our reasoning is not applicable to the more complex axiomatic models of RC11 and ARMv8.

The rest of this paper is laid out as follows. In Section 2 we use a simple example to illustrate the ambiguity of the documentation. In Section 3 we detail the features which distinguish Java's access modes from other memory models. In Section 4 we give a complete account of our axiomatic model. In Section 5 we detail the results of our comparisons with ARMv8, C11 and x86. In Section 6 we give a formalization of the model in Coq and use it to prove the theorems listed above. In Section

7 we demonstrate that the choice of a partial coherence order is not observable in our model. In Section 8 we discuss related work.

2 JAVA'S ACCESS MODES

Here we will introduce the basics of the Java access mode API and motivate our formalization by way of an example program. We will show how different interpretations of the documentation can cause one of three unwanted outcomes when using the access mode API.

Developers can make use of Java's Access Modes, hereafter "the JAM", through the `VarHandle` API included in the JDK version 9. There are four access modes: plain, opaque, release-acquire, and volatile. Regular reads and writes of shared variables are considered plain mode, and the `VarHandle` API allows reads and writes to be annotated with one of the other three modes. The specified intent of the JAM is that each mode provides progressively more guarantees about the behavior of its accesses at the expense of some performance.

$$\text{plain} \sqsubseteq \text{opaque} \sqsubseteq \text{release-acquire} \sqsubseteq \text{volatile}$$

Plain mode gives virtually no guarantees and the compiler is allowed free reign in optimizing such memory accesses. On the other hand volatile mode provides sequentially consistent (SC) semantics when it is used for all accesses, but requires a memory barrier for each access which can make execution slow ¹.

The `VarHandle` documentation [JDK9 2017; Lea 2017, 2018] is intended to answer questions about which access mode should be used for a read or write to maximize performance while providing enough guarantees to ensure program correctness. It cites three bad outcomes arising from the confusion around the old Java memory model as motivation for the new access modes: undersynchronized and broken code, oversynchronized and slow code, and platform specific code. However, the documentation is written in loose prose which can result in any one of the same three outcomes which the access modes were created to address.

To see this, consider the example execution of a message passing program in Figure 1. An execution graph like this catalogues the relationships and effects of each memory access from an example program. $M(x, n)$ represents a memory access of x with the value n where M can be write, W , a read, R and later an atomic read-write RW . Relationships between memory access are represented as directed edges. The dashed edge labeled with rf represent reads-from and the edges labeled with po represent the sequence of instructions as defined in the program. Note that a and b are concurrent with c and d .

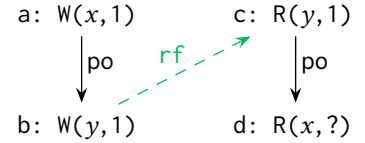


Fig. 1. Event Graph

Intuitively, if the read of y produces the value 1, it acts as a signal to the rest of the second thread that the write to x has completed. However, in the presence of weak memory accesses the sequence of execution may not follow program order. Thus, an important question is: which access modes should one use so that d reads the value 1 from a in keeping with program order, and not 0 from the initialization of x ? Depending on a reasonable interpretation of the documentation, our answer to this question may result in one of the aforementioned bad outcomes.

Undersynchronized, broken code. Consulting the documentation [Lea 2018] we find that opaque mode makes the following guarantee:

If Opaque (or any stronger) mode is used for all accesses to a variable, updates do not appear out of order.

¹This property of progressively greater guarantees is called monotonicity and we prove it as a theorem for our model in Section 6

Whether opaque mode is the correct mode for each access to ensure that *d* reads 1, depends heavily on the definition of “out of order”. Intuitively, we assume that *rf* implies that the read has observed the effects of the write. Then, if the “order” from the documentation for opaque mode is program order, we might reasonably conclude that updates (writes) are observed in program order. In which case, *d* observes the effects of *c* and *b* observes the effects of *a*. Then, through the *rf* edge we can conclude that *d* observes the effects of *a* and *d* will read 1.

Unfortunately this approach actually leaves the code undersynchronized and it will not work as described. It turns out, the documentation also says that opaque mode makes no guarantees about information that travels through reads:

...reading a value in Opaque mode need not tell you anything about values of any other variables.

Thus, whether or not writes are observed in program order, the read *c* may not tell us anything about what happened with *a* opening the way for *d* to read 0 from the initial write to *x*. These two pieces of documentation can lead to an inconsistent understanding of opaque mode. A reasonable reader may react defensively and look to the stronger modes to guarantee the correct program behavior.

Oversynchronized, slow code. A defensive approach to getting the right program behavior is to employ volatile mode. This would give all of the reads and writes in the program SC semantics, guaranteeing that everything is observed in program order. Thus, the *po* and *rf* edges would indeed result in *d* observing *a*, and *d* will read 1. However, if the compiler takes the approach outlined in the documentation it will insert a full fence for every SC access which will slow down execution.

Platform specific implementation. Alternately, we may note that x86 processors do not reorder writes with writes or reads with reads. If we assume opaque mode accesses will translate directly to hardware instructions, we can tag each access as opaque and use the same reasoning that we used for volatile mode. Unfortunately this reasoning is unsound for execution on ARM processors, which allow such reorderings, and the code would not be portable.

As it happens, the best solution is to use release-write for *b* and acquire-read for *c*. These modifications appear as the RA annotations for the write and read of *y* in Figure 2a. This will ensure the orderings, depicted as gray edges, for any access before *b* and any access after *c* with a minimum of synchronization. Thus, *a* is observed by *d* through those orderings and *d* will read the value 1. Critically, we can only demonstrate this after more precisely defining the guarantees provided by the opaque and release-acquire modes of the JAM. Our formalization brings this much needed clarity to the current specification.

3 DISTINGUISHING FEATURES OF THE JAM

While the JAM was inspired by the access modes of C11, it makes several departures from C11 and other memory models that are worthy of consideration and a challenge for formalization. First, it sheds legacy features like the release sequences and consume-reads of C11. Second, it contains a broad and simple definition of causal cycles for the purposes of ruling out thin-air reads. Finally it includes a non-total ordering on writes to a particular memory address, which is called the coherence order. Here we will discuss how each of these differences will impact any formalization effort for the JAM.

3.1 Letting Go of Release Sequences

Release sequences and consume-reads can be seen as specialized variants of the release-acquire memory order in C11 and release-acquire mode in the JAM. The idea is that, in certain cases, it's

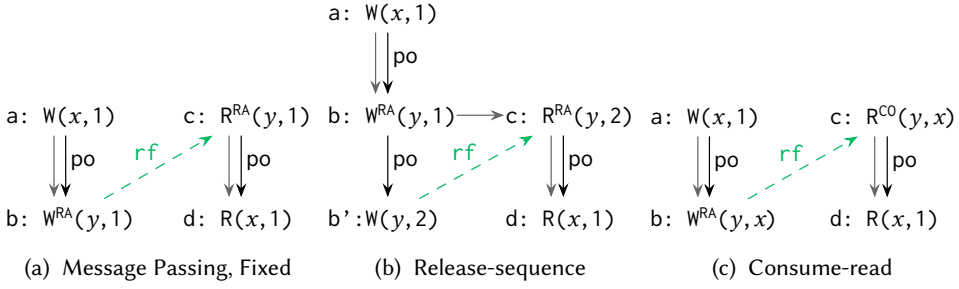


Fig. 2. Release sequences and Consume-reads

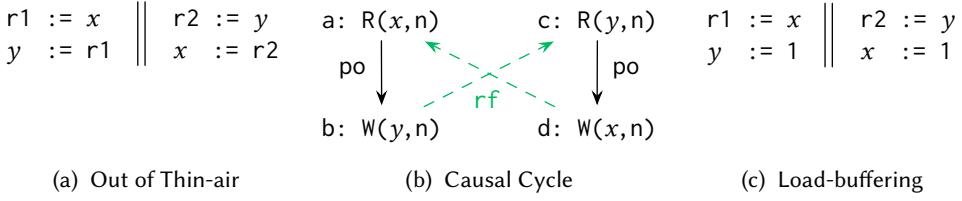


Fig. 3. Causal Cycles

possible to get the same guarantees as a release-write and acquire-read pair, like the one Figure 2a, with less synchronization. The result is faster execution in some contexts.

Release sequences can be used when a write that is after a release-write is read by an acquire-read. Consider the message passing variant in 2b where b is followed by another distinct write to y , b' . If c reads from b' then the release sequences of C11 would guarantee that d would see a through the orderings depicted with gray edges. Then d will read the value 1, just as it would if c read from b . Without the guarantees of release sequences, b' must also be a release-write to get the desired outcome.

Consume-reads are used with a release-write when the memory accesses which must be ordered after the consume-read are data dependent on the value of the read. Consider the message passing variant in 2c where c is annotated with CO and reads a pointer that determines the memory address that d reads from. Some architectures enforce the ordering of c and d in the presence of such data dependencies without the extra synchronization that would result from making c an acquire-read. This can speed up execution in those settings.

The JAM does not include the guarantees of release sequences or any way to annotate a read as a consume-read. According to the author of the specification [Lea 2019] this is a design choice made with the experience of C11 in mind and in favor of simplicity in the model². Consume-reads are deprecated in C++17 while release sequences are rarely used and can be emulated with fences where necessary. Removing these features gives us the opportunity to build a more readable and more easily testable model.

3.2 Acyclic Causality

In Figure 3a we have a classic example program which can exhibit a so called “thin-air” read under sufficiently weak memory models. The question is, at the end of execution can $x = y = 42$? Intuitively, assuming both x and y are initialized to 0, this program can’t generate 42 “out of thin-air”, but many axiomatic models (e.g. C11) do not exclude such executions from the set of all candidate executions.

²Notably, the next C11 standard will not provide ordering guarantees for the execution of Figure 2b [Boehm et al. 2018]

Observe that, in any execution that allows $x = y$, there must be a cycle in the program order and reads-from relations, as illustrated in Figure 3b. The JAM explicitly forbids such cycles, but at the cost of forbidding some behaviors which may be beneficial for performance.

For example, consider the classic load-buffering (LB) litmus test in Figure 3c, where the question is, can $r1 = r2 = 1$? If performance was the sole concern in the design of the JAM this behavior would be allowed because the reads can be reordered with the writes with the aim of improving performance. Unfortunately, this example also exhibits the same cycle in the program order and reads-from, so it is forbidden by the JAM.

The problem of differentiating these kinds of examples has been studied at great length by memory model researchers. Another approach is to forbid cycles in **rf** and a subset of program order based on a notion of dependency. Sadly, this too has very subtle issues, as outlined by [Batty and Sewell \[2014\]](#). The original, formal Java Memory Model [[Manson et al. 2005](#)] attempted to address the issue of causal cycles in its full generality. More recently the Promising semantics of [[Kang et al. 2017](#)] introduced a novel "promise" mechanism to model compiler optimizations for this purpose.

In all of these cases the complexity of the resulting models makes them hard to understand and hard to test. Instead, the JAM specification adopts a simple solution by forbidding cycles in the program order and reads-from. While this does forbid the behavior of the second example at the cost of some performance [[Ou and Demsky 2018](#)], it gives us yet another opportunity to build a simpler model.

3.3 Partial Coherence Order

The JAM specification makes no provision for a total ordering of writes to a given memory location which is a standard feature in other memory models. According to the author of the informal specification [[Lea 2019](#)] the motivation for this choice was to achieve the weakest possible definition for opaque mode. However, the consequences of a partial coherence order manifest in subtle ways.

For example, the standard definition of atomicity for read-writes relies on a total coherence order. Borrowing the definition from [Vafeiadis et al. \[2015\]](#), for a read-write, RW, the write it reads-from, W_1 , and the coherence order, $\xrightarrow{\text{co}}$, we have:

$$W_1 \xrightarrow{\text{rf}} RW \implies \neg \exists W_2, W_1 \xrightarrow{\text{co}} W_2 \xrightarrow{\text{co}} RW$$

Taken together with a total coherence order this means that the atomic pairs of writes and read-writes are always totally ordered. Since the JAM makes no such guarantees regarding normal writes there are ambiguities like the one in Figure 4. There, the write-read-write pairs are not ordered across threads since there is no global relationship between writes.

As a consequence, our model must include extra constraints for read-writes while being careful not to over-constrain normal writes which would otherwise be concurrent. We detail our approach in Section 4.7.

Separately, our drive for simplicity in the definitions of the model has yielded a key insight where the coherence order is concerned. While we have remained faithful to the documentation and personal communications with the author in modeling a partial coherence order, we will show that the effects of that choice are not observable under our model. Specifically, we show that it is impossible to construct a litmus test that behaves differently in the presence or absence of a total coherence order.

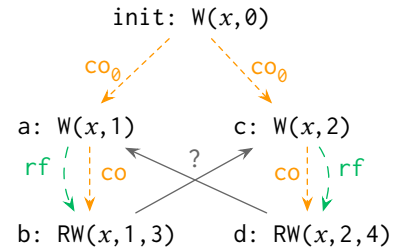


Fig. 4. Concurrent Read-writes

4 AXIOMATIC MODEL

In this section we will describe our axiomatic model of the JAM. To start, we will review the purpose and construction of memory models in the context of the JAM. We will then give an overview of Herd and the cat modeling language. We will also give a high level description of our modeling approach. Then we will provide the definitions and intuition for each part of the JAM: plain mode, opaque mode, fences, release-acquire mode, volatile mode, and atomic read-writes.

4.1 Modeling with Herd

The JAM defines the semantics of shared memory accesses that use the VarHandle API in concurrent Java programs. In particular it answers the question: what values can a read see during all of the possible executions of a program? Here we use a standard graph based representation for executions, where the nodes represent memory events like reads and writes and the edges represent relationships between those events. For example, recall the execution in Figure 1, there is a reads-from relationship between the write b and read c , ($b \xrightarrow{rf} c$) as well as program order relationships ($a \xrightarrow{po} b$, $c \xrightarrow{po} d$).

Thus, the values that a read can see are determined by the writes it can be paired with in the rf relation. Then, since a read can only be paired with one write per execution, the job of the memory model is to define the set of valid pairings by defining the set of valid executions.

To achieve this, the model uses the events and basic relations of an execution, like rf and po , to define other relations specific to the memory model. For example, in our model a reads-from relationship between a write and a read is one way to derive a visibility relationship. Then, to validate executions, the model includes constraints on an execution's relationships. For example, our model forbids causal cycles as described in Section 3.2. We will describe the definitions for both visibility and acyclic causality later.

Herd does the work of enumerating the set of possible executions for small test programs. For each execution it consults a specified memory model constructed with the cat modeling language to find out if the execution is valid.

To generate executions, Herd comes with built-in "architectures" that define the language (e.g. ARM assembly) used to construct litmus tests along with the events and relations that can appear in the enumerated executions. Herd architectures include the following built-in sets of events: writes (W), reads (R), initialization writes (IW), final writes (FW) and atomic read-writes (RMW). They may also include additional event sets. For example the C11 architecture includes the set of relaxed atomic accesses (RLX). Architectures also include the following built-in relations: reads-from (rf), program order (po), accesses to the same location (loc), and identity (id).

Note that a Herd architecture is different from a memory model as it does not validate executions. Rather, it provides the basic events and relations from which the memory model is constructed.

The cat modeling language allows the Herd user to define an axiomatic model in terms of operations over the event sets and relations of the architecture. It includes the composition ($;$), union (\cup), intersection ($\&$), compliment (\sim) and transitive closure ($+$) operators for relations. It is also possible to define new relations using (`let rel = ...`) and create filters over relations with (`let filter(rel) = ...`). For example, if we wished to define a relation between distinct writes to same location we might use: `let wloc = ~id & (W * W) & loc`. Finally, we can specify when an execution is valid by defining constraints on relations. For example, we can require that `wloc` be acyclic: `acyclic wloc`.

In the following subsections we will use cat to define each part of the JAM but the full model can be viewed in Appendix A.

4.2 Model Overview

At a high level, the definitions of our model focus on two key concepts. The first is an intuitive notion of visibility (**vo**), that represents when one memory access has "seen" the effects of another memory access. The second is an acyclicity requirement for the coherence order (**co**). Aside from the restriction on causal cycles, this is the only mechanism by which executions are forbidden. Thus, every unwanted execution must exhibit a cycle in the coherence order.

The core definitions of the model for plain and opaque mode provide mechanisms that generate coherence order edges based on visibility. For example, a read always chooses the last write it has seen and any such writes should be ordered in **co** accordingly. There are three more rules that generate coherence edges in our definition of plain and opaque mode following the work of Batty et al. [2011].

Every other part of the JAM, save for atomic read-writes, is modeled as a small extension to **vo**. In particular each mode builds on the visibility relationships of the previous weaker mode. For example volatile mode writes add both intra-thread and inter-thread visibility guarantees to those provided by release-writes.

As a result, we begin our exposition with plain and opaque mode and then extend the opaque mode definition of **vo** to model fences. We then leverage those visibility relationships to construct release-acquire mode and volatile mode, the latter extending the former. We conclude with an extension of the coherence order to model atomic read-writes.

4.3 Plain & Opaque Mode

In the JAM documentation, plain mode accesses are given virtually no guarantees when they occur in different threads without correct synchronization. On the other hand, opaque mode does provide some cross thread guarantees which form the basis for the rest of the memory model. There are some subtleties involved in the documented relationship between plain and opaque mode accesses so we will address them together. First, the main properties that opaque mode accesses guarantee according to the documentation are:

Bitwise Atomicity Reads will see the value of only one write. Opaque mode guarantees that reads will not see mixed bits from different writes.

Write Availability Writes can be read by later reads. The intent is to avoid a situation (e.g. in a spin-loop) where repeated reads never see a write in another thread because they are optimized by the compiler to execute only one time. When an optimization like this happens, the availability of the write for the read depends on when the read is executed [Corbet 2012].

Acyclic Causality A read should not influence its own value. As described in Section 3.2, this forbids counter-intuitive behavior like thin-air reads.

Coherence The order of writes should respect visibility and should agree with the way that reads observe their order. For example, one guarantee (of four we will define later) is that a read should be paired with the last write that it knows about and not an earlier one.

Our model of opaque mode focuses on acyclic causality and coherence. In keeping with other axiomatic models built for Herd [Alglave et al. 2008; Lahav et al. 2017; Pulte et al. 2017], our model pairs each read with a single write so all accesses are bitwise atomic. Additionally, Herd has limited support for loops so we do not account for the possibility of reads being optimized out of loops.

Plain Coherence. The JAM documentation does not include plain mode writes in its coherence guarantee. That is, following the documentation in this regard would mean that plain mode writes will not appear in **co** relationships. This follows the approach in the C11 documentation but it departs from formal models for C11.

To see why plain mode accesses should be included in the coherence ordering guarantees, note that according to the documentation plain mode accesses are safe to use within a critical section guarded by a lock according to the documentation. The idea is that code which is properly synchronized with locks will have single threaded semantics for the duration of the critical section. Thus, the model shouldn't require accesses to be annotated with anything stronger than plain mode.

Then, consider the message passing variant in Figure 5. Here, the purpose of this pattern is to signal the end of the critical section through atomic read-writes that unlock, b, and lock, c, the variable y . For any lock to work correctly, the accesses which are program order before the unlock, $a \xrightarrow{po} b$, must be visible to accesses after the lock, $c \xrightarrow{po} d$ (for now we leave the mechanism that enforces these orderings unspecified). In particular the effects of the write a should be visible to the write d , $a \xrightarrow{co?} d$.

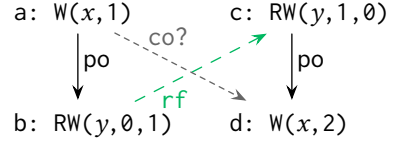


Fig. 5. Message Passing Coherence

However, even if the unlock and enforce program order of the plain writes, and thereby show that a is visible to d , we would not be able to derive $a \xrightarrow{co} d$ because the coherence rules do not apply to plain accesses. This stands in contrast to formal models of C11 [Lahav et al. 2017] where the coherence order and happens-before relations apply to plain accesses. As a result and in contrast to the documentation, our model includes coherence order guarantees for both plain and opaque mode accesses.

Figure 6 details our model of the JAM's plain and opaque modes. We include some additional sets O , RA , and V that represent opaque accesses, release-acquire accesses, and volatile accesses respectively. $po\text{-}loc$ and $invrf$ represent the program order restricted to the same location and inverted reads-from relations. Note, that we do not use Herd's built-in coherence relation but rather define our own without a total ordering.

We define the visibility order, vo , for plain and opaque mode using two relations, $po\text{-}loc$ and $rf+$. In the first case, given two memory accesses to the same location in the same thread the second should see the effects of the first. In the second case, the sequence of one or more reads guarantees that the final read is executing in a context where the effects of the write are visible. Note that a sequence of rf relationships may occur when the reads in the sequence are atomic read-writes. This definition for vo guarantees only the most basic forms of visibility which is important given the inclusion of plain mode accesses.

Then, to satisfy the coherence requirement we first ensure that the coherence order includes only distinct writes to the same location using the filter $wwco$, and then adopt the four behaviors originally identified for C11 by Batty et al. [2011]:

cows Two such writes ordered by visibility are similarly coherence ordered.

cows A read chooses the last write it has seen.

corw A write that follows a read of the same location in program order is coherence ordered after the write paired with the read.

```

let vo = rf+ | po-loc
let wwco(rel) = rel & ~id
               & loc & (W * W)

let cows = wwco(vo)
let cowl = wwco(vo; invrf)
let corw = wwco(vo; po-loc)
let corr = wwco(rf; po-loc; invrf)

let coinit = wwco((IW * W))
let cofw = wwco((W * FW))

let co = cows | cowl | corw
        | corr | cofw | coinit

acyclic co

let opq = O | RA | V
acyclic (po | rf) & opq

```

Fig. 6. Opaque Mode

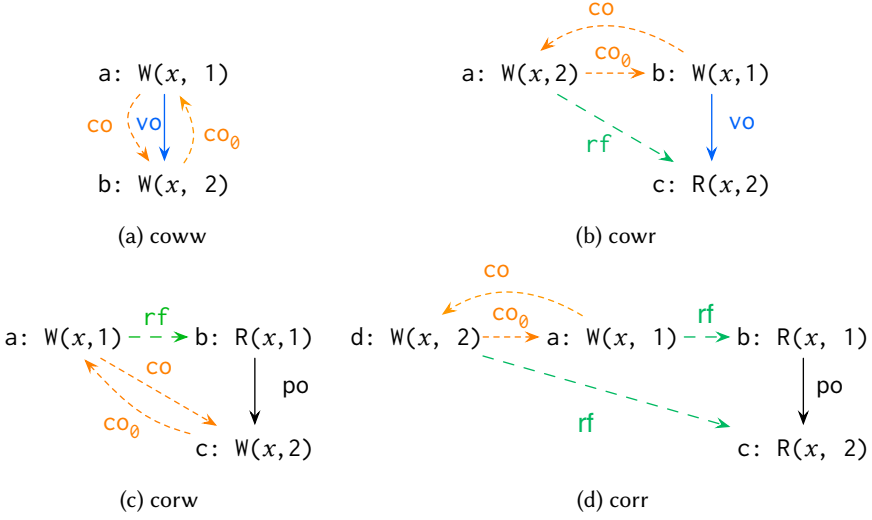


Fig. 7. Coherence

corr Program ordered reads to the same location order their writes in the same way.

We also order the initialization write before all writes of the same location, *coinit*, and we order all writes before the final write of the same location, *cofw*. In the rest of the paper these relationships will be distinguished as *co₀* for clarity but they are treated the same as any other *co* edge by the model.

These six rules make up the definition of the coherence order, *co*, until we discuss read-writes. Then, the primary mechanism by which the model forbids executions is through requiring the coherence order to be acyclic, *acyclic co*. In Figure 7, we give examples of how each coherence rule forbids an execution by showing a cycle in the coherence order. In every example, the final write of 1 is assumed to be after all other writes in the coherence order, *co₀*.

In Figure 7a, if the write *a* is visible to the write *b* then *a* is ordered before *b* and *a* can't be the last write. In Figure 7b, if the read *c* has seen the last write *b* then it should not be able to read an older write *a*. In Figure 7c, if we can follow the read of *a* to *b* to a later write, *c*, then *a* can't be the last write. In Figure 7d, given two reads in program order, *b* and *c*, if *b* has seen the last write then *c* should not be able to read an older write.

Finally, we must forbid causal cycles for opaque mode (and stronger) accesses. We define *opq* to be any opaque or stronger access since any guarantee provided by opaque should apply to stronger access modes. Then, we forbid cycles in the union of program order and reads-from, qualified for opaque mode accesses: *acyclic (po | rf) & opq*.

As we will see in the following sections, these base definitions allow us to model nearly every other component by extending *vo*. The lone exception is atomic read-writes, for which we extend the coherence order directly.

4.4 Fences

The JAM supports five types of fences: release, acquire, load-load, store-store, and full. Programs often include fences to enforce the visibility of effects between accesses before the fence to those after the fence. The full fence additionally provides inter-thread visibility guarantees.

We model the visibility guarantees of fences directly using two relations, *svo* and *spush* that represent specified relationships in the program following the style of Bender et al. [2015]; Cray and Sullivan [2015]. The guarantees of intra-thread ordering fences (release, acquire, load-load, and store-store) are modeled with specified visibility orders, *svo*. For example, where we might

normally order two read accesses with a load-load fence, instead we require a specified visibility order between them. Full fences, are modeled as specified push orders, `spush`, which include both inter-thread and intra-thread guarantees.

This approach allows us to extend visibility to account for synchronization in a way that is uniform for fences and their related access modes. Moreover, specified orders allow the compiler to make smart choices to ignore specified visibility based on any intra-thread ordering of instructions guaranteed by a target language or architecture as detailed in the work of Bender et al. [2015]; Sullivan [2015]. We will see an example of this when discussing atomic read-writes in Section 4.7. Further, specified orders have a natural implementation as fences which has been studied by Bender et al. [2015] and Sullivan [2015].

We will first give the definitions and some basic intuition for our approach and then a concrete example for each relation.

```
with to from linearisations(M \ IW, cofw | rf | into)
let spushto = to+ & (domain(spsh) * domain(spsh))
let rfso     = rf | svo | spsh | spushto;spsh
let vo       = rfso+ | po-loc
```

We extend the previous definition of `vo` by adding `svo`, `spush`, and `spushto`; `spush` to `rf` in a new relation `rfso` (reads-from, specified orders). Joining `svo` and `spush` with `rf` in `rfso` means that the intra-thread ordering guarantees of each type of specified order are included in `vo`. To capture the inter-thread visibility guarantees of full fences we use three auxiliary relationships, `to`, `into` and `spushto`.

Trace order, `to`, is a total ordering of all memory accesses (except for initial writes) constructed with Herd's built-in `linearisations` function and scoped to later definitions using `with ... from`. Trace order respects `cofw`, `rf`, and all intra-thread orderings, `into` (intra-thread trace order). `into`, not shown here, is the union of all specified visibility orders, push orders and later release-acquire/volatile accesses. Intuitively, each of the synchronization mechanisms captured by `into` guarantees that the associated memory accesses are executed in program order. Finally, `spushto` emulates the total ordering of full fences when combined with `spush` in `spushto`; `spush` assuring that, when there are two `spush` orders in the program the effects of the head of one is visible to the tail of the other depending on the trace order of the heads.

Specified Visibility Orders. To see how specified visibility orders create intra-thread visibility between accesses, consider the event graph in Figure 8 for the message passing example from Section 2.

By specifying that $a \xrightarrow{\text{svo}} b$ and $c \xrightarrow{\text{svo}} d$ we can show that d could not have read 0 from the initial write to x . First, we assume that the initial write to x is coherence order before a , $\text{init} \xrightarrow{\text{co}_0} a$. Then, assume that d has read from the initialization, $R(x, 0)$. We will show a contradiction. In the graph we have that, $a \xrightarrow{\text{svo}} b \xrightarrow{\text{rf}} c \xrightarrow{\text{svo}} d$. By the definition of `rfso` we have the three edges that combine to show $a \xrightarrow{\text{vo}} d$. Then, by `cowr` (Fig. 7b) we have that $a \xrightarrow{\text{co}} \text{init}$ which is a cycle in `co` and a contradiction.

Specified Push Orders. Specified push orders create visibility relationships in two ways. The first is an intra-thread visibility ordering between the two push ordered instructions like `svo`. This appears as `spush` in the definition of `rfso`.

The second is a cross-thread ordering that emulates the standard total ordering of two full fences. Given two push orders, the head of the first will be visible to tail of the second, or vice-versa based on the current ordering of the heads in `to`. The ordering of the heads is recorded as `spushto` which we connect with the additional visibility ordering of one tail using `spushto`; `spush`.

To see how push orders emulate full fences, consider the store-buffering litmus test in Figure 9. The question is, can both reads take their values from the initialization?

First, we assume the gray edge between the two fence instructions, $\text{fence} \rightarrow \text{fence}$, which represents one side of the total ordering provided by full fences. Then, we assume that d reads from the initial write to x and show a contradiction. Since fences also provide the intra-thread orderings $a \rightarrow \text{fence}$ and $\text{fence} \rightarrow d$, we can see that a is ordered before d . Then by cowr it must be that $a \xrightarrow{\text{co}} \text{init}$, which creates a cycle with $\text{init} \xrightarrow{\text{co}_0} a$ and a contradiction. On the other side of the ordering between fences a similar argument applies if b reads from initial write to y . Thus, b and c could not have read from the initial writes to x and y in the same execution.

Push orders are more direct but capture the same ordering. In Figure 10 both write-read pairs are push ordered.

First we assume one side of the total trace order between the two writes, $a \xrightarrow{\text{to}} c$. Again, we assume that d reads from the initial write to x and show a contradiction. By spushto ; spush , we have $a \xrightarrow{\text{vo}} d$ which means that d has seen a and the same reasoning from before applies with cowr .

Note, including spushto in vo in place of spushto ; spush would give the desired visibility between the heads and tails of two orders by transitivity, but it would also give an ordering to the heads of the push orders which does not otherwise exist.

4.5 Release-Acquire Mode

We can now make small extensions to our model to incorporate other access modes. First, we define release-acquire mode following the standard set by other models like C11 [C11 2010] and adopted by the JAM documentation.

```

let rel = W & (RA | V)
let acq = R & (RA | V)
let ra = po; [rel] | [acq]; po | rfso
let vo = ra+ | po-loc

```

We define release writes, rel , to be any write marked as release or volatile. We define acquire reads, acq to be any read marked as acquire or volatile. Again, any guarantee provided by release-acquire mode should hold for volatile mode.

We update the vo definition from opaque mode with fences by extending what was previously rfso and call it ra (release-acquire). We add visibility relationships from memory accesses for any location to a release write that is later in program order, $\text{po}; [\text{rel}]$. We also relationships from an acquire read to program order later opaque memory accesses for any location, $[\text{acq}]; \text{po}$. Note that the documentation makes clear that plain accesses should be ordered by release writes and acquire reads so long as the types are bitwise atomic.

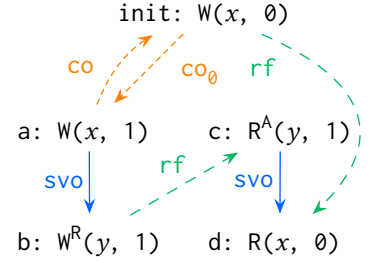


Fig. 8. Specified Visibility

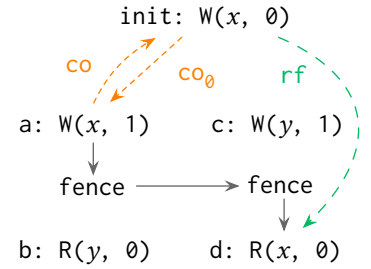


Fig. 9. Fences, One side

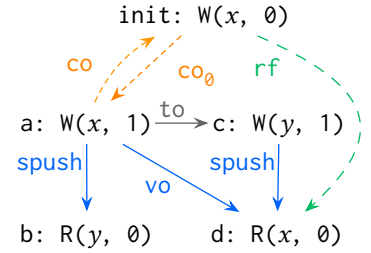


Fig. 10. Push order, One side

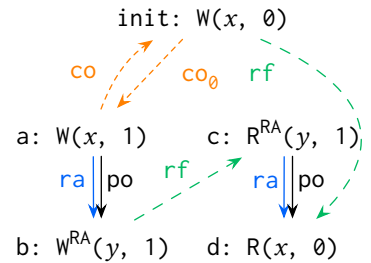


Fig. 11. Release-Acquire Visibility

To see, how the release-acquire extension to opaque mode works, consider the message passing example from the Section 2 updated to include ra in Figure 11. Note that, we use superscript RA for release writes, W^{RA} , and acquire reads, R^{RA} . Later we will use V for volatile accesses. Then, if we adopt a release write for b and an acquire read for c we can show that, $a \xrightarrow{vo} d$, and the reasoning is the same as for specified visibility orders in Section 4.4.

We note that our definitions suggest that if all reads were acquire reads then we could derive a visibility relationship between a read and itself in executions exhibiting causal cycles. In Section 6, we prove that visibility cycles are a contradiction in our model and that causal cycles are forbidden when all reads are acquire-reads.

4.6 Volatile Mode

Volatile mode is a further extension of the visibility order in release-acquire mode. The documentation provides a detailed account of how volatile accesses can be implemented in terms of full fences and release-acquire mode accesses. Volatile writes can be implemented as a release-write followed by a full fence and volatile reads can be implemented as a full fence followed by an acquire-read. Since specified push orders provide the same visibility guarantees as a full fence we can extend $spush$ with relationships from any earlier access to a volatile read and from a volatile write to any later access to achieve the semantics of the suggested implementation.

```
let vol = V
let volint = po; [vol & R] | [vol & W]; po
let push = spush | volint
let vvo = ra | pushto; push | push
let vo = vvo+ | po-loc
```

To that end, the volatile mode extension of our model updates vo by extending ra to vvo . Adding $volint$ (volatile internal) to $spush$ creates the necessary edges from any access to a volatile read that is later in program order, $po; [V \& R]$, and from an volatile write to any program order later opaque memory access, $[V \& W]; po$. These new edges preserve program order for any access before or after the volatile access when combined with the visibility definition of release-acquire and the extension of $spush$ to $push$ creates the necessary inter-thread visibility relationships.

Note that a simpler approach to achieving the inter-thread guarantees of volatile mode would be to translate the total trace ordering into visibility edges. That is, we could replace the definition of vol above with the following:

```
let vol = ra | spush | volint | to & (V * V)
```

This approach more directly encodes the cross thread ordering guarantee of sequential consistency, but this creates too many visibility relationships.

Consider a modified version of the classic IRIW litmus test in Figure 12a. The question is, can both d and f read from the initial write to x and y ? As outlined in the JAM documentation the acquire reads in this example are allowed to see the writes in different orders, so both can read from the initial writes for x and y . However, if volatile mode accesses are totally ordered as in the proposed definition, then we have either $a \xrightarrow{vo} b$ or $b \xrightarrow{vo} a$. In Figure 12a we have the first case. Then we have $a \xrightarrow{vo} d$ because $a \xrightarrow{vol} b \xrightarrow{rf} c \xrightarrow{ra} d$. Then by cowr, we have that $a \xrightarrow{co} init$, which is a contradiction. In the other case, $b \xrightarrow{vo} a$, we have a contradiction when f reads from the initial write to y .

By contrast, our approach has the same effect as push ordering accesses related by $volint$. As we will demonstrate later, using a matching litmus tests, this guarantees the correct behavior for the release-acquire variant of IRIW because it enforces no direct ordering between the two writes.

Importantly, it also gives the correct behavior when the reads are volatile, which should be SC semantics. Consider Figure 12b which has one of the two possible orderings given by $pushto; push$

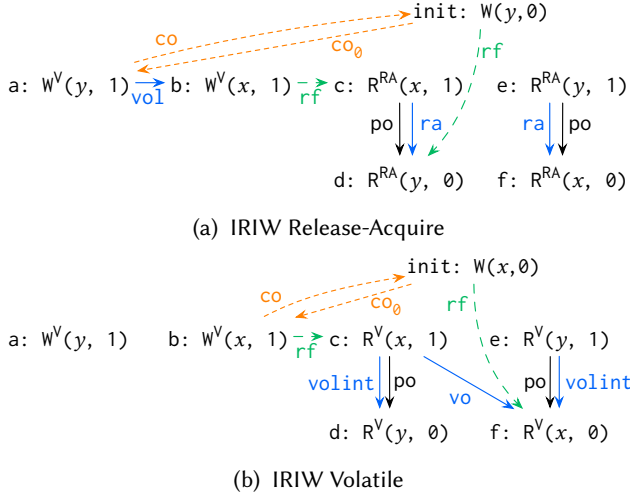


Fig. 12. IRIW Variants

when the reads are volatile. This correctly establishes $c \xrightarrow{vo} f$ creating a contradiction for the opposite thread's final read. As before the other direction of the total order forbids the other read.

4.7 Atomic Read-Writes

The behavior of atomic read-writes is the only part of the JAM that is not modeled by extending vo . Recall Figure 4 from Section 3. We must take care to ensure that we do not allow concurrent read-writes in the presence of a non-total coherence order. To achieve this we update the coherence order co with two additional rules:

```

let corwexcl = wwco((rf;[RW])^~1;co')
let corwtotl = wwco(((RW * W) | (W * RW)) & to)
let rec co = ... | corwexcl | corwtotl

```

The first, `corwexcl` ensures exclusivity in the relationship between the read-write and its paired write. Note that, for clarity, we separate out `corwexcl` even though it recursively refers to co in its definition. As illustrated in Figure 13, if there is a write co -after the one paired with the read-write, $a \xrightarrow{co} b$ then it is co -after the read-write, $c \xrightarrow{co} b$.

Recalling the example of Figure 4, this exclusivity is not enough to prevent concurrent read-read-write pairs to the same location. Since exclusivity uses the ordering of other writes with the paired write we could consider a total ordering just for paired writes.

```

let pairedw = domain(rf;[RW])
let corwtotl = wwco((pairedw * pairedw) & to)

```

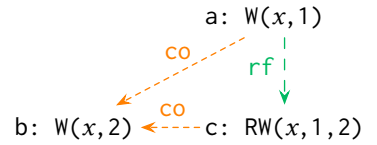


Fig. 13. Read-write Exclusivity

This suffices to forbid concurrent read-write chains to the same location. In Figure 4, by `corwtotl` we have $a \xrightarrow{co} c$ or $c \xrightarrow{co} a$. Then, in the first case, we have $b \xrightarrow{co} c$ by `corwexcl`. The other side is similar. However, this approach inadvertently orders writes that could be concurrent under a partial coherence order. Instead, as in the first definition, we choose a total ordering with the read-write itself.

Now any chain of read-writes will be exclusive without unnecessarily ordering regular writes. We will show that in Figure 4 either, $b \xrightarrow{co} c$ or $d \xrightarrow{co} a$. By `corwtotl` we have $d \xrightarrow{co} a$ or $a \xrightarrow{co} d$. In the first case we are done. In the second case, we consult `corwtotl` again and we have $b \xrightarrow{co} c$ or $c \xrightarrow{co} b$. In the first case we are done. In the second case, by assumption we have

$a \xrightarrow{\text{co}} d$ and $c \xrightarrow{\text{co}} b$. Then we can apply `corwexcl` to both to derive $c \xrightarrow{\text{co}} d$ and $d \xrightarrow{\text{co}} c$, which is a cycle and a contradiction. Importantly, this reasoning can be applied repeatedly to any chain of read writes to achieve exclusivity.

We note here, that JAM makes no intra-thread ordering guarantees for atomic read-writes even though they exist on some architectures like x86 [Owens et al. 2009].

5 VALIDATION

Our axiomatic model of the JAM is complete, covering all four modes, atomic read-writes, and all five fence types in less than 40 lines of definitions. Moreover each component is implemented as a modest extension to just two relations which makes it readable. Now we must demonstrate that the model is consistent with expectations about the behavior of the JAM as outlined in the documentation.

We have validated our formalization of the JAM by comparing litmus test outcomes for the JAM with the outcomes for ARMv8 [Pulte et al. 2017], RC11 [Lahav et al. 2017] and the model for x86 included with Herd. We use these comparisons to show that there are no unexpected differences in behavior between each pair of models. The comparisons are based on the relationships expressed in the documentation and an existing understanding of the relationships between language memory models and hardware memory models.

For example, we expect the JAM to permit more behaviors than ARMv8 with the exception of litmus tests like load-buffering (Fig. 3c). The reason is that ARMv8 is a compilation target for Java and we expect the compiled output of a Java program to exhibit fewer behaviors than the original program. Otherwise, any reasoning we may do for the Java program would be invalid for the compiled output. In the case of the load-buffering litmus test, the behavior is strictly forbidden by the JAM's definition of causal cycles but not on ARMv8.

We run the tests using the Herd tool [Alglave et al. 2014]. Herd exhaustively enumerates all possible executions of a litmus test and checks if each execution is allowed by the model being tested. We can then compare the three possible results for each model to see how often executions are allowed: Always, Sometimes, and Never. In terms of the behavior being tested Always and Sometimes mean that the behavior is allowed, while Never means that it is not allowed. Weaker memory models should allow more behaviors and see more Always and Sometimes results, while stronger memory models should allow fewer behaviors and see more Never results.

Our test suite is built with litmus tests taken from existing research for the three models we compare against. Importantly we did not modify any of these tests. However, Herd uses different built-in relations to refer to the access types of each model. For example, Herd provides the `L` built-in to refer to release-writes in ARMv8 litmus tests and the `REL` built-in to refer to release-writes in C11 litmus tests. Thus, for each comparison we include a mapping between the access types of the other model and the access modes of the JAM.

For each comparison we define the mapping between the built-in relations of the two models, we detail our expectations for the results, and then discuss the results of the comparison. As we will see our model behaves as expected in all cases.

5.1 Comparison with ARMv8

Herd provides several built-in relations for the ARMv8 litmus tests: `M` for normally memory operations, `L` for release writes, `A` for acquire reads, and `DMB.SY` for full fences. There are no `SC/Volatile` accesses.

Mapping. We map every memory access to opaque mode with `opq = M`. This is conservative with respect to our expectation that the JAM permits more behaviors than ARMv8. If we were to

Pulte et al. [2017]

| name | ARMv8 | JAM |
|-----------------------------------|------------------|--------------|
| WRC+addr | Never | Never |
| LB+data+data-wsi | Never | Never |
| W+RR | Never | Never |
| totalco | Never | Never |
| PPOCA | Sometimes | Sometimes |
| IRIW | Sometimes | Sometimes |
| IRIW+addr | Never | Sometimes |
| IRIW+poas+LL | Never | Sometimes |
| IRIW+poaps+LL | Never | Sometimes |
| MP+dmb.sy+addr-ws-rf-addr | Never | Sometimes |
| WW+RR+WW+RR+wsilp+poaa+wsilp+poaa | Never | Sometimes |
| LB | Sometimes | Never |

Fig. 14. ARMv8 Litmus Test Comparison

map some accesses to plain mode they would be allowed to exhibit causal cycles. Thus, mapping all memory accesses to opaque mode means the JAM model will permit fewer behaviors.

We map release-writes to release-acquire mode writes, $rel = L$ and acquire-reads to release-acquire mode reads, $acq = A$. Note, that M includes L and A so release writes and acquire reads will also be treated as opaque mode accesses. Finally we treat all accesses with a full fence between them in program order having a specified push order.

```

let opq = M
let rel = L
let acq = A
let spush = po; [DMB.SY]; po

```

Fig. 15. ARMv8 Mapping

Expectations. As stated, we expect the JAM to be weaker than ARMv8. That is, we expect that any time ARMv8 exhibits a behavior the JAM should too and there should be instances where the JAM exhibits behaviors that ARMv8 doesn't. The only exceptions are cases where the broad definition of causal cycles adopted by the JAM will rule out behaviors like the load-buffering example of Figure 3c.

Results. Aside from totalco and LB, all the tests come from the supplementary material accompanying the ARMv8 model of Pulte et al. [2017] which we use for comparison. Figure 14 shows the results of our comparison and they agree with our expectations: the JAM is at least as weak as ARMv8 except in the case of load-buffering (LB).

Many of the results owe to the fact that the JAM is not multi-copy atomic. That is, unlike ARMv8, different threads can see writes to different locations in different orders. So, IRIW+* and WRC+* are allowed for the JAM but not for ARMv8. The WW+RR+WW+RR+wsilp+poaa+wsilp+poaa litmus test is a variant of IRIW where all writes are release writes and all reads are acquire reads. The writes in this test use a weaker form of synchronization than the example in Figure 12a in Section 4.6 where the behavior is allowed by the JAM. As a result this behavior is allowed by the JAM. Finally, the LB test is identical to the example in Figure 3c from Section 3. This execution is allowed by the ARMv8 model but it is a cycle in $(po \mid rf) \ \& \ opq$ which is explicitly disallowed by the JAM.

5.2 Comparison with RC11

C11's atomic memory accesses can be annotated with memory orders. Herd provides the following built-in relations for the memory orders and accesses in the C11 litmus tests: M for plain memory access, RLX for relaxed memory order accesses, REL for release memory order accesses, ACQ for acquire memory order accesses, REL_ACQ for release-acquire memory order read-write accesses, SC for sequentially consistent memory order accesses, $F \ \& \ REL$ for release fences, $F \ \& \ AQR$ for acquire fences, $F \ \& \ SC$ for sequentially consistent (full) fences, and $F \setminus SC$ for all other fences.

```

let opq = RLX | ACQ | REL | ACQ_REL | SC
let rel = REL | ACQ_REL | SC
let acq = ACQ | ACQ_REL | SC
let vol = SC
let svo = po; [F & REL]; po; [W] | [R]; po; [F & ACQ]; po
let spush = po; [F & SC]; po

```

Mapping. Our mapping for the relations provided by Herd for C11 follows the informal relationship outlined in the documentation for the JAM [Lea 2018]. All plain accesses in C11 are treated as plain accesses in our mapping to the JAM. We map relaxed memory order accesses or stronger to opaque mode, $opq = RLX \mid \dots$, release memory order accesses or stronger to release mode, $rel = REL \mid \dots$, acquire memory order accesses or stronger to acquire mode, $acq = ACQ \mid \dots$, sequentially consistent memory order accesses to volatile mode, $vol = SC$. We also map release fences before writes, $po; [F \ \& \ REL]; po; [W]$, and acquire fences after reads, $[R]; po; [F \ \& \ ACQ]; po$, to specified visibility orders, svo . Finally, we map sequentially consistent fences, $po; [F \ \& \ SC]; po$, to push orders between program order earlier and program order later accesses.

Expectations. According to the documentation, the access modes of the JAM are inspired by the memory orders of C11 so this comparison is of particular importance. We expect the JAM to match C11 except in cases where release-sequences, consume-reads, or causal cycles are involved.

Results. In Figure 16 we have the results of the comparison with the RC11 model of [Lahav et al. 2017]. We include the tests from research by Wickerson and Batty [2015], [Vafeiadis et al. 2015], [Lahav et al. 2017]. In the case of [Wickerson and Batty 2015] and [Vafeiadis et al. 2015] we used the tests directly. In the case of [Lahav et al. 2017] we translated the tests from the paper to Herd litmus tests ourselves. The results largely agree with our expectations. The exceptions are places where the RC11 model breaks with the C11 specification. We will discuss each in turn.

The `cyc_na` test is the same as the load-buffering example from Section 3 but all the accesses are plain. The JAM allows this cycle in plain mode $po \mid rf$ because its acyclicity requirement only applies to opaque mode or stronger accesses. The RC11 model breaks with C11 by including an acyclicity requirement for $po \mid rf$ for all memory accesses.

The `lb` test is the same as `cyc_na` except that all of the accesses are relaxed memory order. In the mapping to the JAM this translates into opaque mode accesses which are not allowed to exhibit a cycle in $po \mid rf$. Thus both models forbid the load-buffering behavior in this case.

The `mp_relacq_rs` test leverages the release sequences of C11. Since the JAM does not include release sequences it does not forbid the behavior of this test.

Our test runs for `fig6` and `fig6_translated` timed out at 5 minutes. The reason for the timeouts is that the permutations of the trace order create a much larger set of possible executions for Herd to enumerate when compared with RC11. When the litmus tests reach a certain threshold of size and complexity Herd can't examine all of the possible executions within the time limit.

Otherwise, the behavior modeled by these tests highlights a quirk in C11's rules for SC accesses. Together, they demonstrate that strengthening the memory order of a particular relaxed store in `fig6` to an SC store in `fig6_translated` creates new behaviors. That is, the tests demonstrate that the memory orders of C11 are not monotonic. RC11 includes a fix proposed by [Vafeiadis et al. 2015] which forbids this behavior and thus the behavior is forbidden by RC11. In Section 6 we prove that the JAM's access modes are indeed monotonic in our model which means that stronger access modes exhibit fewer behaviors. Thus, the behavior in these tests is forbidden in our model and we expect that these behavior would never be allowed if Herd was able to examine all of the executions before the timeout.

| Vafeiadis et al. [2015] | | | Wickerson and Batty [2015] | | |
|-------------------------|------------------|------------------|----------------------------|------------------|------------------|
| name | RC11 | JAM | name | RC11 | JAM |
| a1 | Sometimes | Sometimes | cppmem_iriw_relacq | Sometimes | Sometimes |
| a1_reorder | Sometimes | Sometimes | cppmem_sc_atomics | Never | Never |
| a3 | Sometimes | Sometimes | iriw_sc | Never | Never |
| a3_reorder | Sometimes | Sometimes | mp_fences | Never | Never |
| a3v2 | Sometimes | Sometimes | mp_relacq | Never | Never |
| a4 | Never | Never | mp_relacq_rs | Never | Sometimes |
| a4_reorder | Sometimes | Sometimes | mp_relaxed | Sometimes | Sometimes |
| arfna | Never | Never | mp_sc | Never | Never |
| arfna_transformed | Never | Never | | | |
| b | Never | Never | Lahav et al. [2017] | | |
| b_reorder | Sometimes | Sometimes | name | RC11 | JAM |
| c | Never | Never | 2+2W | Never | Never |
| c_p | Never | Never | IRIW-acq-sc | Sometimes | Never |
| c_p_reorder | Never | Never | RWC+syncs | Never | Never |
| c_pq | Never | Never | W+RWC | Never | Never |
| c_pq_reorder | Never | Never | Z6.U | Sometimes | Never |
| c_q | Never | Never | | | |
| c_q_reorder | Never | Never | Herd X86 Tests | | |
| c_reorder | Never | Never | name | x86 | JAM |
| cyc | Never | Never | 4. SB | Sometimes | Sometimes |
| cyc_na | Never | Sometimes | 6. SB | Sometimes | timed out |
| fig1 | Always | Always | 6. SB+prefetch | Sometimes | timed out |
| fig6 | Never | timed out | CoRWR | Never | Never |
| fig6_translated | Never | timed out | iriw-internal | Sometimes | Sometimes |
| lb | Never | Never | iriw | Never | Sometimes |
| linearisation | Never | Never | podrw000 | Sometimes | Sometimes |
| linearisation2 | Never | Never | podrw001 | Sometimes | Sometimes |
| roachmotel | Never | Never | SB | Sometimes | Sometimes |
| roachmotel2 | Never | Never | SB+mfences | Never | Never |
| rseq_weak | Sometimes | Sometimes | SB+rfi-pos | Sometimes | Sometimes |
| rseq_weak2 | Always | Always | SB+SC | Sometimes | Sometimes |
| seq | Never | Never | X000 | Sometimes | Sometimes |
| seq2 | Never | Never | X001 | Sometimes | Sometimes |
| strengthen | Never | Never | X002 | Sometimes | Sometimes |
| strengthen2 | Never | Never | X003 | Sometimes | Sometimes |
| | | | X004 | Sometimes | Sometimes |
| New C11 Tests | | | X005 | Sometimes | Sometimes |
| name | RC11 | JAM | X006 | Sometimes | Sometimes |
| IRIW-sc-r1x-acq | Sometimes | Never | x86-2+2W | Never | Sometimes |

Fig. 16. RC11 & x86 Litmus Test Comparison

In the case of IRIW-acq-sc, Z6.U, and IRIW-sc-r1x-acq our model forbids the behavior in keeping with the C11 specification. We will consider each case and discuss why RC11 does not forbid each behavior.

The IRIW-acq-sc test appears in in Figure 17a. The reason this is forbidden in our model is that all accesses before a SC/Volatile read are ordered by `volint`. Then, because there are two such orders in the reading threads there is either a visibility order, $c \xrightarrow{vo} f$ or $d \xrightarrow{vo} e$. Thus, one of the two reads must see the non-initialization write using the same reasoning from Section 4.6 for Figure 12b. This effectively emulates the "leading" fence compilation scheme described in [Lahav et al. 2017], where a full fence is placed before the SC/Volatile accesses. The authors [Lahav 2019] pointed out that this scheme should forbid this behavior. By contrast RC11 relaxes the C11 model

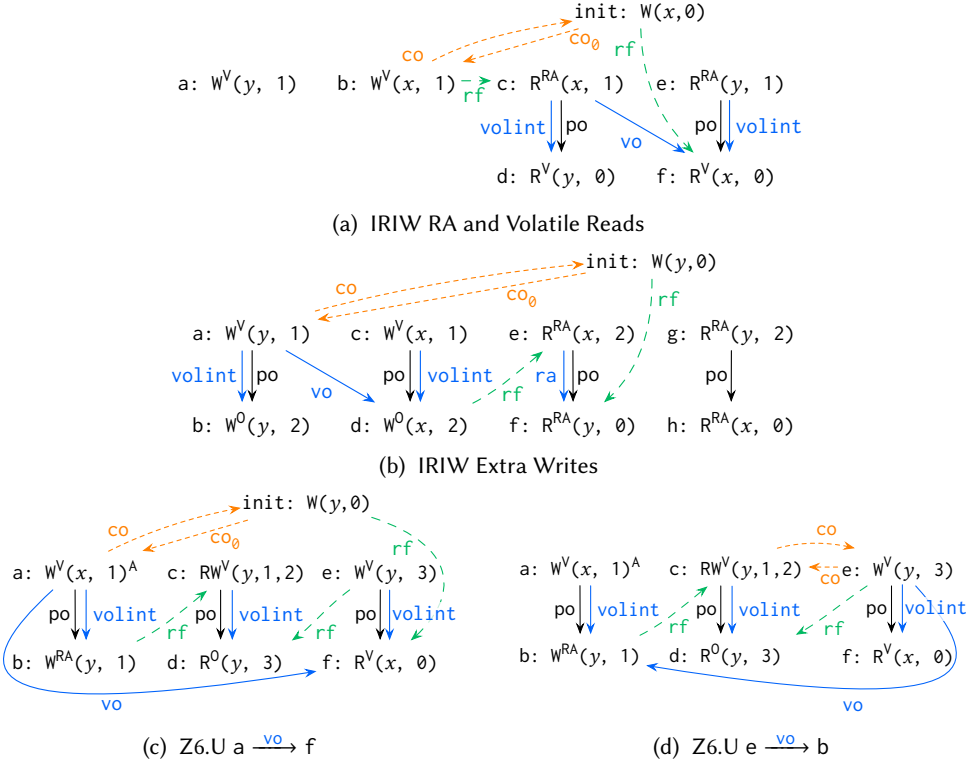


Fig. 17. Z6.U & IRIW RA with Volatile Reads Litmus Tests

to accommodate a leading or trailing fence compilation scheme. In this case, if the trailing fence scheme is used there's no ordering provided to any of the SC accesses and the execution is allowed.

Two possible executions for Z6.U test appear in Figures 17c and 17d. They represent each case of the visibility orders induced by $a \xrightarrow{volint} b$ and $e \xrightarrow{volint} f$. In 17c we have $a \xrightarrow{vo} f$, which means that f must have seen a and could not read from the initialization by $cowr$. In 17d we have $e \xrightarrow{vo} b$. We also have $b \xrightarrow{rf} c$. Then together we have $e \xrightarrow{vo} b \xrightarrow{rf} c$. Then because c and e are both writes to y we have $e \xrightarrow{co} c$ by $coww$. Separately, since $c \xrightarrow{volint} d$ we have $c \xrightarrow{vo} d$. Then since $e \xrightarrow{rf} d$ we have $c \xrightarrow{co} e$ by $cowr$. Thus we have a cycle in co and the execution is forbidden.

This effectively emulates the "trailing" fence compilation scheme described by Lahav et al. [2017], where a full fence is placed after SC/Volatile accesses. They point out that this scheme should forbid this behavior. Again, RC11 relaxes C11 to model both schemes and under a leading fence compilation scheme the behavior is allowed.

Finally, in Figure 17b we have our own variant of IRIW called IRIW-sc-rlx-acq based on the WW+RR+. . . test in the ARMv8 suite. In this case, if the compiler is inserting trailing fences after a and c , then either $a \xrightarrow{vo} d$ or $b \xrightarrow{vo} c$. Taking the first case we can construct $b \xrightarrow{vo+} f$ from $a \xrightarrow{vo} d \xrightarrow{rf} a \xrightarrow{ra} f$ and a cycle in co by $cowr$. The second case is similar but with the read of x in the 4th thread.

Again, the RC11 model does not forbid this execution. In our discussion with the authors [Lahav 2019], they suggested two possible interpretations for this behavior. The first is a leading fence compilation scheme which we have discussed above. The second is a merge of each pair of writes into a single write. In the second case we suggest that the Java compiler should not merge opaque mode (or stronger) accesses.

In summary, our model correctly forbids the behavior in each of these cases. The reason is that volatile reads would be preceded by a full fence and volatile writes would be followed by a full fence in the presence of mixed mode programs with SC accesses. Importantly, the tension between optimal compilation and a simple model exists here as it does with C11.

5.3 Comparison with x86

Herd provides the following built in relations for x86 litmus tests: M for memory accesses, SFENCE for fences for intra-thread ordering of writes with other accesses, LFENCE for ordering intra-thread ordering of reads with other accesses, and MFENCE as a full fence with cross thread ordering guarantees.

```
let opq = M
let svo = [W];po;[SFENCE];po | [R];po;[LFENCE];po
let spush = po;[MFENCE];po
```

Mapping. We map from all regular memory accesses to opaque mode `opq = M` in keeping with the opaque mode mapping from our comparison with ARMv8. We map SFENCE to specified visibility orders from writes to other memory accesses across the fence, `[W];po;[SFENCE];po`. We map LFENCE to specified visibility orders from reads to other memory accesses across the fence, `[R];po;[LFENCE];po`. Finally we map MFENCE as a specified push orders between any access before the fence to any access after the fence, `po;[MFENCE];po`.

Expectations. We expect the JAM to be weaker than x86 in all cases. The only weak behavior that x86 exhibits is reordering writes with reads (store-buffering) which the JAM also allows in opaque mode.

Results. Figure 16 shows the results of our comparison with the x86 model included with Herd. The litmus tests are also included with Herd for the model. The two tests which timed out, 6. SB+prefetch and 6. SB, are store buffering variants with additional threads and variables. As with the timed out tests in the comparison with RC11 this is due to the large number of executions created by the permutations of the trace order. Given that the JAM can reorder writes with reads we expect that the behavior of these tests is allowed. Otherwise the tests confirm our expectations.

6 METATHEORY

Here, we develop a metatheory for our model of the JAM. First, we detail the semantics and then sketch the proofs for the theorems listed in the introduction. We show that each mode, from Volatile to Opaque, admits strictly more executions, otherwise referred to as monotonicity. We also show that properly synchronized programs, i.e. race-free programs as defined by Boehm and Adve [2008], will only exhibit sequentially consistent behaviors as required by Java's data-race-free (DRF) guarantee. Finally, we show that acquire mode for all reads obviates the inclusion of the JAM's acyclic causality requirement.

These theorems further validate the definitions of our model, give more evidence that the model is complete with respect to the documentation, and clearly demonstrate that the semantics is suitable for formal reasoning. The semantics, lemmas and theorems have been mechanized in Coq³ and each theorem and lemma definition in this section includes a direct link to the relevant file and proof.

6.1 Semantics

The axioms of our model apply to memory event graphs, each of which represents one possible execution of a program. Thus far we have relied on Herd to generate executions for a program and then apply the axioms of our model to rule on whether an execution is allowed. Our formalization

³<https://bitbucket.org/ucla-pls/jam/>

in Coq replaces the Herd machinery with the history fragment of the semantics of [Crary and Sullivan \[2015\]](#) to model the set of possible executions. Here we will give enough detail about the history semantics to understand the content of our theorems and proofs. The full semantics is available in the Coq source⁴.

We use n and l to range over natural numbers and memory locations. We use i to range over unique memory access identifiers where previously we have used a, b, c , etc. We use m to represent one of the four access modes, P for plain, O for opaque, RA for release acquire and V for volatile. Memory accesses, a , can take the form of reads, l_m , and writes, $l_m := n$ with their accompanying modes as well as read-writes, $RW(l, n)$.

We use H to range over lists of memory events that represent executions. We call these lists histories and we use $H(h)$ to represent that the memory event h is in H . Memory events record the program's interactions with a history. For example, $H(\text{is}(i, l_m := n))$ records that the identifier i is a write to l of the value n and $H(\text{exec}(i))$ records that i has executed, subject to the axioms of the memory model. We use $i \xrightarrow{R}_H i$ to represent our model's relations as they apply to H .

We also require that histories satisfy some basic well-formedness conditions. For example, a read must take its value from a write to the same location which has been executed. Importantly, we require that the trace order respects any intra-thread ordering created by visibility generated from specified orders, release-acquire accesses, or volatile accesses. This is in keeping with our restriction on trace order as detailed in Section 4.4. We call this well-formedness condition executable. In what follows, the well-formedness conditions are unified as *trace coherence*. Otherwise, the only restrictions on histories come from the acyclicity conditions outlined in Section 4.3.

6.2 Theorems

To begin, we demonstrate the monotonicity of our access mode definitions. We define the reflexive ordering of the access modes as $P \sqsubseteq O \sqsubseteq RA \sqsubseteq V$ and extend it to accesses $l_{m_1} \sqsubseteq l_{m_2}$, $l_{m_1} := n_1 \sqsubseteq l_{m_2} := n_2$, $RW(l, n_1) \sqsubseteq RW(l, n_2)$ whenever $m_1 \sqsubseteq m_2$. As a technical matter we treat read-writes as always having the same order. We extend the order to histories by matching identifiers and ordering the accesses.

$$H_1 \sqsubseteq H_2 \triangleq \forall i \ a_1 \ a_2, H_1(\text{is}(i, a_1)) \wedge H_2(\text{is}(i, a_2)) \Rightarrow a_1 \sqsubseteq a_2$$

When the po, rf, and to relations of two histories H_1 and H_2 have the following relationships: $\xrightarrow{\text{po}}_{H_2} \subseteq \xrightarrow{\text{po}}_{H_1}$, $\xrightarrow{\text{to}}_{H_2} \subseteq \xrightarrow{\text{to}}_{H_1}$, $\xrightarrow{\text{rf}}_{H_2} \subseteq \xrightarrow{\text{rf}}_{H_1}$, then we say they *match*.

THEOREM 1 (MONOTONICITY).

[[coq/Monotonicity.v](#), [monotonicity](#)]

For two histories H_1 and H_2 , suppose that both match, both are trace coherent, and $H_2 \sqsubseteq H_1$. Further suppose that $\text{acyclic}(\xrightarrow{\text{co}}_{H_1})$ and that there are no specified visibility orders or push orders in H_2 , then $\text{acyclic}(\xrightarrow{\text{co}}_{H_2})$

We make two notes. First the absence of specified orders in H_2 is a technical convenience since specified order edges are not related to the strength of the access modes for reads and writes. Second, we focus on the acyclic coherence requirement because the match assumption means that it would be trivial to satisfy the acyclic causality requirement for H_2 supposing it is true of H_1 because po and rf have fewer edges in H_2 .

PROOF SKETCH. We assume $i \xrightarrow{\text{co}}_{H_2} i$ for some i and show that this must mean $i \xrightarrow{\text{co}}_{H_1} i$ which is a contradiction. This is straight forward by induction on $i \xrightarrow{\text{co}}_{H_2} i$, noting that each case of visibility in H_2 will exist in H_1 because of stronger access modes in H_1 . \square

For the last two theorems will first establish that the main component of visibility in our model, vvo, is irreflexive.

⁴<https://bitbucket.org/ucla-pls/jam/src/master/coq/Dynamic.v>

LEMMA 1 (IRREFLEXIVE VISIBILITY).

[[coq/Truncate.v](#), [vvop_irreflex](#)]

If H is trace coherent then, for all i , $\neg i \xrightarrow{vvo+}_H i$.

PROOF SKETCH. This follows from two facts. First, vvo derives from orderings that are always either program ordered from intra-thread synchronization (svo, spush, ra, volint) or trace ordered (pushto, rf). Second, both relations are total and to is consistent with the po edges for intra-thread visibility by the executable well-formedness condition. \square

Next we show that if a program is properly synchronized then it will only exhibit sequentially consistent behavior. We require the following standard definitions including the traditional notion of sequential consistency [[Shasha and Snir 1988](#)]:

$$\begin{aligned} i_1 \xrightarrow{fr}_H i_2 &\triangleq \exists i_3, i_3 \xrightarrow{rf}_H i_1 \wedge i_3 \xrightarrow{co}_H i_2 \\ i_1 \xrightarrow{com}_H i_2 &\triangleq i_1 \xrightarrow{co}_H i_2 \vee i_1 \xrightarrow{rf}_H i_2 \vee i_1 \xrightarrow{fr}_H i_2 \\ i_1 \xrightarrow{sc}_H i_2 &\triangleq i_1 \xrightarrow{po}_H i_2 \vee i_1 \xrightarrow{com}_H i_2 \end{aligned}$$

We also require a definition of proper synchronization in keeping with our focus on visibility.

$$i_1 \xrightarrow{sync}_H i_2 \triangleq \exists i_3, i_1 \xrightarrow{vvo+}_H i_3 \xrightarrow{po*}_H i_2 \wedge \forall i_4, i_3 \xrightarrow{po*}_H i_4 \Rightarrow i_3 \xrightarrow{vvo}_H i_4$$

The idea is that any conflicting access is visibility ordered by some mechanism, be it a specified order (fence) or a strong mode for i_3 . Then, following the definition of “type 2” data-races from [Boehm and Adve \[2008\]](#), we say that H is *race-free* when, for all conflicting accesses i_1 and i_2 , we have $i_1 \xrightarrow{sync}_H i_2$ or $i_2 \xrightarrow{sync}_H i_1$.

THEOREM 2 (DRF-SC).

[[coq/SC/DRF.v](#), [drf_sc](#)]

If H is trace coherent, race free and, $\text{acyclic}(\xrightarrow{co}_H)$, then $\text{acyclic}(\xrightarrow{sc}_H)$.

PROOF SKETCH. We assume $i \xrightarrow{sc+}_H i$ for some i and show a contradiction. By Lemma 1 it is enough to demonstrate a cycle in $vvo+$.

We can show that for any i_1 and i_2 , if we have $i_1 \xrightarrow{com}_H i_2$ then we have $i_1 \xrightarrow{sync}_H i_2$. Note that any accesses related by com are conflicting. Then we have either $i_1 \xrightarrow{sync}_H i_2$ or $i_2 \xrightarrow{sync}_H i_1$. In each case for com we can show that $i_2 \xrightarrow{sync}_H i_1$ creates a cycle in the coherence order so it must be $i_1 \xrightarrow{sync}_H i_2$.

Then, since po is irreflexive, we have that any sequence $i \xrightarrow{sc+}_H i$ must include at least one com edge. Then since com edges are also $sync$ edges, when we have $i_1 \xrightarrow{sc+}_H i_2$ we also have $i_1 \xrightarrow{po|com+}_H i_2$ and therefore $i_1 \xrightarrow{po|sync+}_H i_2$.

But then we can rearrange a cycle in $i \xrightarrow{po|sync+}_H i$ to be $i \xrightarrow{sync+}_H i$ by appending a leading po edge to the end. Observe that for any sequence $i_1 \xrightarrow{sync+}_H i_2$ we have $i_1 \xrightarrow{vvo+}_H i_2$, so we have $i \xrightarrow{vvo+}_H i$ as required. \square

Finally, we show that if all reads are acquire-reads then the JAM’s acyclic causality requirement is unnecessary. This theorem demonstrates the soundness of proposed compiler implementations for satisfying the acyclic causality requirement of the JAM [[Ou and Demsky 2018](#)].

THEOREM 3 (CAUSAL ACQUIRE-READS).

[[coq/ReleaseAcquire.v](#), [acq_causality](#)]

If H is trace coherent and all reads in H are acquire-reads, then $\text{acyclic}(\xrightarrow{po|rf}_H)$.

PROOF SKETCH. We assume $i \xrightarrow{po|rf+}_H i$ for some i and show a contradiction. By Lemma 1, it is enough to demonstrate a cycle in $vvo+$.

First, note that any sequence $i_1 \xrightarrow{rf}_H i_2 \xrightarrow{po}_H i_3$ implies $i_1 \xrightarrow{vvo+}_H i_3$ because rf implies vvo and because, by assumption, the read i_2 is an acquire read and i_3 is program order later, so $i_2 \xrightarrow{vvo}_H i_3$.

Let \xrightarrow{rfpoq}_H be a reads-from edge followed by an optional program order edge. Note that, by induction and the fact that rf and $rfpo$ imply vvo we can show that $i_1 \xrightarrow{rfpoq}_H i_2$ implies $i_1 \xrightarrow{vvo+}_H i_2$

Then, since po is irreflexive, we have that any sequence $i \xrightarrow{po \mid rf+}_H i$ must include at least one rf edge. Thus we can rearrange to get $i' \xrightarrow{rfpq+}_H i'$ and we have $i' \xrightarrow{vvo+}_H i'$ by the above, as required. \square

7 UNOBSERVABLE TOTAL COHERENCE ORDER

Here we will demonstrate that the effects of a total coherence order are unobservable in our model. This is a surprising result since a total coherence order is intuitively associated with maintaining per-location state and we would expect concurrent writes to have a material impact on the behavior of programs.

To start, recall that a feature of a memory model as exhibited by an example program is defined in binary terms. The feature is present when any executions are allowed and the feature is absent when no executions are allowed. Then we say that a feature is *observable* when we can construct any example program such that the feature's presence in the model changes whether any executions are allowed. So, to observe a total coherence order, we must construct a program that changes whether executions are allowed based on the presence of the total order. We will show that this is impossible under our model.

First, note that in our model there is no way to allow previously forbidden executions when adding coherence order edges. This means we can't observe the feature by going from some executions *with* a total order to no executions *without* a total order. Thus, any program that allows us to observe the total coherence order must forbid all executions in the presence of a total order. We will show that such a program will also forbid all executions when the total order is removed.

THEOREM 4. *It is impossible to construct an example program for the JAM, such that, for any two writes, if they are totally ordered there are no valid executions, and if they are not totally ordered there is at least one valid execution.*

PROOF. From a total order we have that all candidate executions can be divided between the two directions. For some writes to the same location a and b we have:

$$a \xrightarrow{co} b \vee \quad (1)$$

$$b \xrightarrow{co} a \quad (2)$$

We begin by eliminating all executions for the first direction, (1). This means we must construct a cycle in co , such that, for all executions including (1) we have :

$$a \xrightarrow{co} b \quad (1)$$

$$b \xrightarrow{co} \dots \xrightarrow{co} a \quad (3)$$

With this cycle, executions including (1) are forbidden by the acyclicity requirement for co . Observe that (3) must be present in all executions of the program. If it is not then we can divide the executions that do not have (3) between executions with (1) and executions with (2). For the executions without (3) and with (1) the execution is allowed and we have failed to forbid all executions. Now we must eliminate all executions for (2) with another cycle:

$$b \xrightarrow{co} a \quad (2)$$

$$a \xrightarrow{co} \dots \xrightarrow{co} b \quad (4)$$

Moreover (4) must persist for all executions for the same reason that (3) persisted. Then it must be that, for all executions we have:

$$b \xrightarrow{co} \dots \xrightarrow{co} a \quad (3)$$

$$a \xrightarrow{co} \dots \xrightarrow{co} b \quad (4)$$

Then, even if we remove the total order, (1) and (2), from the model, all executions of the program will still have a forbidden cycle by (3) and (4) and we can not demonstrate a single execution that is allowed. \square

Critically, this line of reasoning depends on the fact that we have a single acyclicity requirement in which `o` participates and both (3) and (4) form a cycle that violates that requirement. If that were not the case then (3) and (4) would not necessarily create a cycle or forbid any executions. As a result this reasoning does not apply to the ARMv8 and RC11 models as they have multiple acyclicity requirements involving the coherence order. In Appendix B we have included a litmus test for ARMv8 that shows it is possible to construct a program for that model that behaves differently with and without the total order.

Demonstrating that a total coherence order is not observable in our model means that the JAM may adopt a total coherence order without affecting the outcomes of the model. This would allow it to emulate other mainstream memory models in this respect and recover the intuitive notion of per-location state.

8 RELATED WORK

The JAM is intended as an update and expansion of Java's memory model. As a result, we have taken a fresh look at how to construct our model, but our work takes inspiration from a large body of research on such formalization efforts.

8.1 The Original Java Memory Model

The original Java Memory Model [Manson et al. 2005] included an early attempt to model standard compiler optimizations while ruling out thin-air reads. The specification in that work was part prose and part formal definitions and the "causality" mechanisms at the heart of the model made the definitions complex. Taken together these issues made the model unsuitable for the tasks which one normally formalizes a programming language, namely accurate discourse, metatheory, and algorithm verification. Later work by Ševčík and Aspinall [2008] manually examined a large suite of litmus tests to show that the model disallowed some standard compiler optimizations. Eventually the model was fully formalized in Coq by Huisman and Petri [2007] and Aspinall and Ševčík [2007], but, to the best of our knowledge, there is no way to easily test the behavior of example programs. By contrast we have constructed a readable and testable model for Java's new access modes.

8.2 C11, ARM, and x86

The C11 memory model, which served as the inspiration for Java's access modes, has seen extensive study. It was originally formalized by Batty et al. [2011] with later revisions to include read-modify-writes and fences [Sarkar et al. 2012]. The work of Vafeiadis et al. [2015], from which we draw the largest set of our C11 litmus tests, studied the soundness of common compiler optimizations under the model of C11 by Morisset et al. [2013]. Our monotonicity theorem is modeled after the same theorem from Vafeiadis et al. [2015]. Most recently, the axiomatic model of Lahav et al. [2017] incorporated the proposed fixes of Vafeiadis et al. [2015] and addressed the unsound compilation strategies which we discussed in Section 5. Also, further progress has been made on new models for C11 that more successfully support standard compiler optimizations without the problem of thin-air reads [Kang et al. 2017; Pichon-Pharabod and Sewell 2016]. However, like the original Java memory model, these models rely on complex formal constructs (promises and event structures respectively). Again, our semantics remains relatively simple thanks to the JAM's broad definition of causal cycles.

As outlined previously there are several important differences between C11 and Java. First, the partial coherence order required careful consideration. The consequences of this design choice

manifests most clearly where atomic read-writes are concerned. Second, the lack of legacy features, like C11's release sequences, and the simple mechanism that forbids causal cycles allowed us to build a simple model. In turn, the simplicity of the model makes it more readable than existing C11 models and it allowed us to argue forcefully that the model should adopt a total coherence order.

These differences appear most clearly in our litmus test comparison with RC11 model of [Lahav et al. \[2017\]](#). Of particular note are the `mp_relacq_rs` and `lb` tests. In the first case the JAM is weaker than C11 because it does not support release sequences, in the second case it is stronger because the C11 specification gives no concrete definition for how to rule out thin-air reads. Notably, the RC11 herd model also forbids causal cycles in `po | rf` so the load buffering behavior is forbidden. This method of preventing thin-air reads in RC11 is included to enable their proofs of soundness for compilation to POWER and not as a representation of the C11 specification.

Programs that mix SC/volatile access modes are the primary focus of [Lahav et al. \[2017\]](#). In particular the leading and trailing fence insertion approach to compilation was shown to be unsound for Power under models prior to RC11. We compared our model with RC11 in Section 5. Our semantics correctly forbids the behavior described in IRIW-sc-rlx-acq, IRIW-aqc-sc and Z6.U in keeping with the JAM documentation.

Hardware memory models have also seen extensive study. x86 was studied by [Owens et al. \[2009\]](#) and [Alglave et al. \[2014\]](#). We use the model included with Herd in our litmus test comparisons and we saw that x86 was stronger than the JAM, as expected. ARM processors have traditionally had a much weaker memory model when compared with x86. Recently the model for ARMv8 [\[Pulte et al. 2017\]](#) expanded the guarantees made by the architecture to include multi-copy-atomicity. We saw the effects of this in the behavior of the IRIW-* and WRC-* litmus tests from our comparison between the JAM and ARMv8. As expected the JAM is weaker than ARMv8 in every case except where cycles in `po | rf` are concerned.

8.3 The Relaxed Memory Calculus

Our mechanized semantics is based on the history fragment of the Relaxed Memory Calculus of [Crary and Sullivan \[2015\]](#). We use their concept of specified push orders and we draw inspiration for our definitions from their notion of visibility. Also, the proof of our theorems benefited greatly from the library of lemmas included with the RMC mechanization. However, their purpose was to model a weaker version of C11 in the interest of generality while, our goal is to model the JAM. Importantly, we do not employ the execution orders of RMC, our coherence definition is far more compact and we have added the `corr` rule to follow a more standard notion of coherence.

9 CONCLUSION

We have presented the first formal model for Java's access modes. Our model is precise, complete, and testable. We have validated our model against the expected behavior of example programs relative to three other mainstream memory models. We have further validated the model by using it to prove general theorems about the semantics of the Java access modes. Finally, we used our model to demonstrate that Java's access modes can adopt a total coherence order.

We believe our model is a strong foundation for a larger formalization effort around weak-memory concurrency on the Java platform. In particular we would like expand the model and use it to verify the lock-free algorithms of the standard library.

A FULL HERD MODEL

```
let opq = 0 | RA | V
let rel = W & RA
let acq = R & RA
let vol = V
```

```

(* release acquire ordering *)
let ra = po;[rel] | [acq];po

(* intra-thread volatile ordering *)
let volint = po;[vol & R] | [vol & W];po

(* intra-thread ordering constraints *)
let into = svo | spush | ra | volint

include "filters.cat"
let wwco(rel) = WW(rel) & loc & ~id
let coint = loc & (IW * (W \ IW))
let cofw = wwco((W * FW))

(* cross thread push ordering extended with volatile memory accesses, based on trace order*)
with to from linearisations(M \ IW, cofw | rf | into)
let push = spush | volint
let pushto = to+ & (domain(push) * domain(push))

(* extend ra visibility *)
let vvo = rf | svo | ra | push | pushto;push
let vo = vvo+ | po-loc

(* general definition from RC11, works for atomic rws and split instruction rws *)
let rmw-jom = [RMW] | rmw

(* coherence rules *)
let cown = wwco(vo)
let cownr = wwco(vo;invrf)
let corw = wwco(vo;po)
let corr = wwco(rf;po;invrf)
let cormwtotal = wwco(((range(rmw-jom) * _ ) | (_ * range(rmw-jom))) & to)
let cormwexcl = wwco((rf;rmw-jom)^-1;co-jom)
let rec co-jom = cown | cownr | corw | corr | cofw | coint | cormwtotal
| wwco((rf;rmw-jom)^-1;co-jom)

acyclic (po | rf) & opq
acyclic co-jom

```

B OBSERVABLE TOTAL COHERENCE FOR ARMV8

Here we demonstrate that it is possible to construct a program that is only forbidden due to the total coherence order of the ARMv8 memory model from [Pulte et al. \[2017\]](#). We start by noting that Herd model for ARMv8 has two acyclicity requirements that involve **co**, **acyclic po-loc | ca | rf as** external and irreflexive **ob as** external. Here, **ca** is defined as **ca = fr | co** and **ob** is defined as **ob = obs | ...** where **obs = ... | coe** and **coe** is coherence restricted to inter-thread relationships. Critically, they do not work together to form cycles, so we can use one with each side of the total order to demonstrate its observability in the model.

We have constructed a program (available as part of our [test suite](#)) which will exhibit a cycle in the first requirement for one side of the total order and a different, incompatible cycle, for the second requirement. We show this by running the program with Herd and allowing executions that violate these two requirements (we mark them with “flag” in the Herd parlance). The result in [Figure 18](#) is exactly two flagged executions. Each figure corresponds to one direction of the total

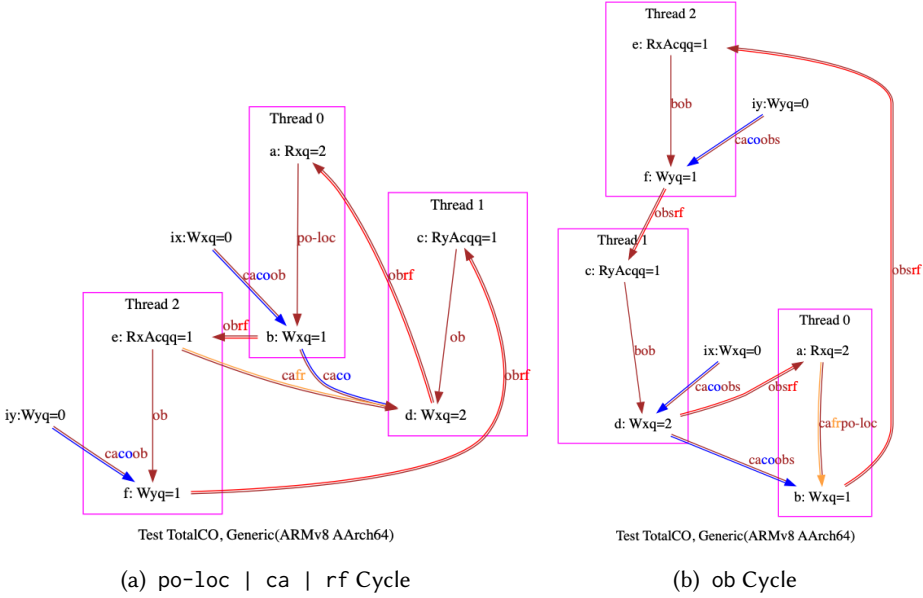


Fig. 18. Two Flagged Executions

ordering between b and d. We will inspect both to demonstrate that they are only forbidden as consequence of the total order, and thus if the total order was taken away they would both be allowed.

In Figure 18a note the following. All of the cycles for any kind of edge involve a $\text{po-loc} \rightarrow b$ which is not in ob. This means we can avoid a cycle in ob. Recall that if there were a cycle in ob then when we use ob for the other side of the total coherence order it would be a cycle regardless of the total coherence order. Also, note that if we take away $b \xrightarrow{\text{co}} d$ (blue) it will also remove $e \xrightarrow{\text{fr}} d$. Thus, if we take away $b \xrightarrow{\text{co}} d$ by removing the total coherence order of the model, there is no cycle left in the graph and the execution would be allowed.

In Figure 18b note the following. All of the cycles for any kind of edge involve $c \xrightarrow{\text{bob}} d$ which is not in po-loc | ca | rf. This means we can't establish a cycle in po-loc | ca | rf. Thus, if we take away $d \xrightarrow{\text{co}} b$ by removing the total coherence order of the model, there is no cycle in ob left in the graph and the execution would be allowed.

ACKNOWLEDGMENTS

We thank Doug Lea for discussions and encouragement, and we thank Shuyang Liu and Christian Kalhauge for helpful suggestions. This material is based upon work supported by the National Science Foundation under Grant No. 1815496. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author and do not necessarily reflect the views of the National Science Foundation.

REFERENCES

- 2010. Programming Languages-C++(final committee draft). C++ standards committee paper WG21/N3092= J16/10-0082.
- Jade Alglave, Anthony Fox, Samin Ishtiaq, Magnus O. Myreen, Susmit Sarkar, Peter Sewell, and Francesco Zappa Nardelli. 2008. The Semantics of Power and ARM Multiprocessor Machine Code. In *Proceedings of the 4th Workshop on Declarative Aspects of Multicore Programming (DAMP '09)*. ACM, New York, NY, USA, 13–24. <https://doi.org/10.1145/1481839.1481842>
- Jade Alglave, Luc Maranget, and Michael Tautschnig. 2014. Herding Cats: Modelling, Simulation, Testing, and Data Mining for Weak Memory. *ACM Trans. Program. Lang. Syst.* 36, 2, Article 7 (July 2014), 74 pages. <https://doi.org/10.1145/2627752>

- David Aspinall and Jaroslav Ševčík. 2007. Formalising Java's Data Race Free Guarantee. In *Proceedings of the 20th International Conference on Theorem Proving in Higher Order Logics (TPHOLs'07)*. Springer-Verlag, Berlin, Heidelberg, 22–37.
- Mark Batty, Scott Owens, Susmit Sarkar, Peter Sewell, and Tjark Weber. 2011. Mathematizing C++ Concurrency. *SIGPLAN Not.* 46, 1 (Jan. 2011), 55–66. <https://doi.org/10.1145/1925844.1926394>
- Mark Batty and Peter Sewell. 2014. The Thin-air Problem. <http://www.cl.cam.ac.uk/~pes20/cpp/notes42.html>. [Dec. 2018].
- John Bender, Mohsen Lesani, and Jens Palsberg. 2015. Declarative Fence Insertion. In *Proceedings of the 2015 ACM SIGPLAN International Conference on Object-Oriented Programming, Systems, Languages, and Applications (OOPSLA 2015)*. ACM, New York, NY, USA, 367–385. <https://doi.org/10.1145/2814270.2814318>
- Hans-J. Boehm and Sarita V. Adve. 2008. Foundations of the C++ Concurrency Memory Model. In *Proceedings of the 29th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI '08)*. ACM, New York, NY, USA, 68–78. <https://doi.org/10.1145/1375581.1375591>
- Hans J Boehm, Olivier Giroux, and Viktor Vafeiadis. 2018. P0982R1: Weaken Release Sequences. <http://www.openstd.org/jtc1/sc22/wg21/docs/papers/2018/p0982r1.html>. [Aug. 2019].
- Jon Corbet. 2012. ACCESS_ONCE(). <https://lwn.net/Articles/508991/>. [Dec 2018].
- Karl Crary and Michael Sullivan. 2015. A Calculus for Relaxed Memory. In *Proceedings of POPL'15, ACM Symposium on Principles of Programming Languages*.
- Marieke Huisman and Gustavo Petri. 2007. The Java memory model: a formal explanation. *VAMP 7* (2007), 81–96.
- JDK9. 2017. VarHandle API Docs. <https://docs.oracle.com/javase/9/docs/api/java/lang/invokedynamic/VarHandle.html>. [Mar. 2019].
- Jecheon Kang, Chung-Kil Hur, Ori Lahav, Viktor Vafeiadis, and Derek Dreyer. 2017. A Promising Semantics for Relaxed-memory Concurrency. In *Proceedings of POPL'17, SIGPLAN-SIGACT Symposium on Principles of Programming Languages*.
- Ori Lahav. 2019. personal communication.
- Ori Lahav, Viktor Vafeiadis, Jecheon 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 (PLDI 2017)*. ACM, New York, NY, USA, 618–632. <https://doi.org/10.1145/3062341.3062352>
- Doug Lea. 2017. JEP 193. <http://openjdk.java.net/jeps/193>. [Mar. 2019].
- Doug Lea. 2018. Using JDK 9 Memory Order Modes. <http://gee.cs.oswego.edu/dl/html/j9mm.html>. [Online, March 2019].
- Doug Lea. 2019. personal communication.
- Jeremy Manson, William Pugh, and Sarita V. Adve. 2005. The Java Memory Model. In *Proceedings of the 32nd ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL '05)*. ACM, New York, NY, USA, 378–391.
- Robin Morisset, Pankaj Pawan, and Francesco Zappa Nardelli. 2013. Compiler testing via a theory of sound optimisations in the C11/C++ 11 memory model. In *ACM SIGPLAN Notices*, Vol. 48. ACM, 187–196.
- Peizhao Ou and Brian Demsky. 2018. Towards Understanding the Costs of Avoiding Out-of-thin-air Results. *Proc. ACM Program. Lang.* 2, OOPSLA, Article 136 (Oct. 2018), 29 pages. <https://doi.org/10.1145/3276506>
- Scott Owens, Susmit Sarkar, and Peter Sewell. 2009. A Better x86 Memory Model: X86-TSO. In *Proceedings of the 22nd International Conference on Theorem Proving in Higher Order Logics (TPHOLs '09)*. Springer-Verlag, Berlin, Heidelberg, 391–407. https://doi.org/10.1007/978-3-642-03359-9_27
- Jean Pichon-Pharabod and Peter Sewell. 2016. A Concurrency Semantics for Relaxed Atomics That Permits Optimisation and Avoids Thin-air Executions. *SIGPLAN Not.* 51, 1 (Jan. 2016), 622–633. <https://doi.org/10.1145/2914770.2837616>
- Christopher Pulte, Shaked Flur, Will Deacon, Jon French, Susmit Sarkar, and Peter Sewell. 2017. Simplifying ARM Concurrency: Multicopy-atomic Axiomatic and Operational Models for ARMv8. *Proc. ACM Program. Lang.* 2, POPL, Article 19 (Dec. 2017), 29 pages. <https://doi.org/10.1145/3158107>
- Susmit Sarkar, Kayvan Memarian, Scott Owens, Mark Batty, Peter Sewell, Luc Maranget, Jade Alglave, and Derek Williams. 2012. Synchronising c/c++ and power. *Acm Sigplan Notices* 47, 6 (2012), 311–322.
- Jaroslav Ševčík and David Aspinall. 2008. On Validity of Program Transformations in the Java Memory Model. In *ECOOP 2008 – Object-Oriented Programming: 22nd European Conference Paphos, Cyprus, July 7-11, 2008 Proceedings*, Jan Vitek (Ed.). Springer Berlin Heidelberg, Berlin, Heidelberg, 27–51. https://doi.org/10.1007/978-3-540-70592-5_3
- Dennis Shasha and Marc Snir. 1988. Efficient and Correct Execution of Parallel Programs That Share Memory. *ACM Trans. Program. Lang. Syst.* 10, 2 (April 1988), 282–312. <https://doi.org/10.1145/42190.42277>
- Michael J. Sullivan. 2015. Low-level Concurrent Programming Using the Relaxed Memory Calculus. <https://www.msully.net/stuff/thesprop.pdf>. PhD. Dissertation [Online, accessed July 2017].
- Viktor Vafeiadis, Thibaut Balabonski, Soham Chakraborty, Robin Morisset, and Francesco Zappa Nardelli. 2015. Common Compiler Optimisations Are Invalid in the C11 Memory Model and What We Can Do About It. In *Proceedings of the 42nd Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL '15)*. ACM, New York, NY, USA, 209–220. <https://doi.org/10.1145/2676726.2676995>
- John Wickerson and Mark Batty. 2015. Taming the complexities of the C11 and OpenCL memory models. *arXiv preprint arXiv:1503.07073* (2015).