

Automatically Comparing Memory Consistency Models

John Wickerson

Imperial College London, UK
j.wickerson@imperial.ac.uk

Tyler Sorensen

Imperial College London, UK
t.sorensen15@imperial.ac.uk

Mark Batty

University of Kent, UK
m.j.batty@kent.ac.uk

George A. Constantinides

Imperial College London, UK
g.constantinides@imperial.ac.uk



Abstract

A *memory consistency model* (MCM) is the part of a programming language or computer architecture specification that defines which values can legally be read from shared memory locations. Because MCMs take into account various optimisations employed by architectures and compilers, they are often complex and counterintuitive, which makes them challenging to design and to understand.

We identify four tasks involved in designing and understanding MCMs: generating conformance tests, distinguishing two MCMs, checking compiler optimisations, and checking compiler mappings. We show that all four tasks are instances of a general constraint-satisfaction problem to which the solution is either a program or a pair of programs. Although this problem is intractable for automatic solvers when phrased over programs directly, we show how to solve analogous constraints over program *executions*, and then construct programs that satisfy the original constraints.

Our technique, which is implemented in the Alloy modelling framework, is illustrated on several software- and architecture-level MCMs, both axiomatically and operationally defined. We automatically recreate several known results, often in a simpler form, including: distinctions between variants of the C11 MCM; a failure of the ‘SC-DRF guarantee’ in an early C11 draft; that x86 is ‘multi-copy atomic’ and Power is not; bugs in common C11 compiler optimisations; and bugs in a compiler mapping from OpenCL to AMD-style GPUs. We also use our technique to develop and validate a new MCM for NVIDIA GPUs that supports a natural mapping from OpenCL.

Categories and Subject Descriptors C.1.4 [Processor Architectures]: Parallel processors; D.3.4 [Programming Languages]: Compilers; F.3.2 [Logics and Meanings of Programs]: Semantics of Programming Languages

Keywords C/C++, constraint solving, graphics processor (GPU), model checking, OpenCL, program synthesis, shared memory concurrency, weak memory models

1. Introduction

In the specification of a concurrent programming language or a parallel architecture, the *memory consistency model* (MCM) defines which values can legally be read from shared memory locations [4]. MCMs have to be general enough to enable portability, but specific enough to enable efficient implementations. They must also admit optimisations employed by architectures (such as store buffering and instruction reordering [30]) and by compilers (such as common subexpression elimination and constant propagation [7]). This profusion of design goals has led to MCMs for languages (such as C11 [38] and OpenCL 2.0 [42]), for CPU architectures (such as x86, ARM, and IBM Power), and for GPU architectures (such as AMD and NVIDIA), that are complicated and counterintuitive. In particular, all of these MCMs permit executions that are not *sequentially consistent* (SC), which means that they do not correspond to a simple interleaving of concurrent instructions [45]. As a result, designing and reasoning about MCMs is extremely challenging.

Responding to this challenge, researchers have built numerous automatic tools (see §8). These typically address the question of whether a program P , executed under an MCM M , can reach the final state σ . Put another way: can the *litmus test* (P, σ) pass under M ? While useful, there are several other questions whose answers are valuable for MCM reasoning and development. Four that have appeared frequently in the literature are:

- Q1** Which programs can be run to test whether a compiler or machine conforms to a given MCM? [11, 24]
- Q2** Is one MCM more permissive than another? That is, is there a litmus test that can pass under one but must fail under the other? [11, 14, 43, 48, 50, 58, 63]
- Q3** Can ‘strengthening’ a program (syntactically) ever enable additional behaviours? For instance, can we take a litmus test that must fail, impose additional sequencing or dependencies between its instructions (or, in the C11 case, give an atomic operation a stronger ‘memory order’ [38 (§7.17.3)]), and thereby allow it to pass? [21, 22, 55, 77, 79, 80]
- Q4** Is a given software/architecture/compiler mapping correct? Or is there a litmus test that must fail under the software-level MCM, but which, when compiled, can pass under the architecture-level MCM? [15, 16, 46, 47, 81, 82]

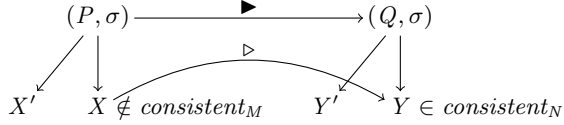


Figure 1. Diagram for explaining Key Idea 2

1.1 Key Idea 1: Generalising the Question

Our first key idea is the observation that all four of the questions listed above can be answered, sometimes positively and sometimes negatively, by exhibiting programs P and Q and state σ such that the litmus test (P, σ) *must fail* under a given MCM M , P and Q are related by a given binary relation \triangleright on programs, and (Q, σ) *can pass* under a given MCM N . That is, each corresponds to finding inhabitants of the following set.

Definition 1 (General problem). $g(M, N, \triangleright) \stackrel{\text{def}}{=} \{(P, Q, \sigma) \mid \sigma \notin \text{obs}_M(P) \wedge P \triangleright Q \wedge \sigma \in \text{obs}_N(Q)\}$

where obs returns the set of final states that can be observed after executing a given program under a given MCM.¹

Our four questions correspond to the following specialisations of g 's parameters:

- Q1** $g(M, 0, id)$ consists of tests that check conformance to M , where 0 is the MCM that allows all executions and id is the identity relation;
- Q2** $g(M, N, id)$ consists of tests that cannot pass under M but can under N ;
- Q3** $g(M, M, \text{'is weaker than'})$ consists of monotonicity violations in M , given a relation that holds when one program is (syntactically) 'weaker than' another; and
- Q4** $g(M, N, \text{'compiles to'})$ consists of bugs in a compiler mapping from M to N , given a relation that holds when one program 'compiles to' another.

1.2 Key Idea 2: Constraining Executions, not Programs

Having captured in Def. 1 our four questions, we now ask whether answers can be generated using an automatic constraint solver. Unfortunately, Def. 1 is problematic because showing that a litmus test must fail under M requires universal quantification over executions. That is, we must show that there *exists* a litmus test (P, σ) such that *all* possible executions of P under M do not reach σ . As Milicevic et al. observe, 'this " $\exists \forall$ " quantifier pattern is a particularly difficult instance of higher-order quantification' [52], and our pilot experiments confirmed it to be intractable in practice.

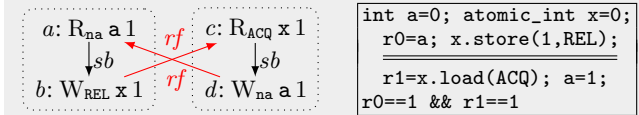
Our second 'key idea', then, involves rephrasing the constraints to be not over *programs*, but over *program executions*, where they become much cheaper to solve (see Def. 2, below), and then to recover litmus tests from these executions. We explain how this is possible with reference to Fig. 1.

We start by finding individual executions X and Y , such that X is inconsistent under M , Y is consistent under N , and $X \triangleright Y$. From X , we construct a litmus test (P, σ) that can behave like X , and from Y , we construct a litmus test (Q, σ) that can behave like Y . The \triangleright relation between X and Y ensures that these litmus tests will have the same final state σ and that $P \triangleright Q$. We now have a litmus test that *can fail* under M and another that *can pass* under N . But g requires a litmus test that *must fail* under M , and perhaps

¹Elements of $g(M, N, \triangleright)$ can also be seen as counterexamples to \triangleright implying observational refinement [33]: $P \triangleright Q \not\Rightarrow \text{obs}_N(Q) \subseteq \text{obs}_M(P)$.

(P, σ) can still pass by taking a different execution X' . Example 1 illustrates how such a situation can arise.

Example 1. The diagram below (left) depicts a C11 execution. Dotted rectangles group events into threads; W and R denote write and read events; REL and ACQ denote atomic 'release' and 'acquire' accesses; na denotes a non-atomic memory access; rf is the 'reads from' relation; and sb means 'sequenced before'. We extract a litmus test from this executions (below right), by mapping write events to store instructions, reads to loads, and having the final state enforce the desired rf relation.



The execution is deemed inconsistent in C11. This is because the successful release/acquire synchronisation implies that b 'happens-before' c , hence that a happens-before d , and hence that it is impossible for a to read from d . As a result, it is tempting to use the litmus test derived from this execution for conformance testing (**Q1**). In fact, the litmus test is not useful for this purpose because it has a data race: the non-atomic store to x goes ahead even if the release/acquire synchronisation on x fails. Racy C11 programs have undefined semantics, so non-conformance cannot be deduced from the test passing.

To guard against situations like the one above, we require that X is also 'dead'. Semantically, X is dead if whenever X is inconsistent, and (P, σ) is a 'minimal' litmus test constructed from X , then no execution of P that leads to σ is allowed (which implies, in particular, that P is race-free). This ensures that P not only *can fail*, but *must fail*, as required. We obtain semantic deadness via a syntactic approximation ($dead_M$), which simply involves a few MCM-specific constraints on the shape of X . For instance, we require that whenever two non-atomic accesses are prevented from racing by release/acquire synchronisation (as a and d are in Example 1), one of the accesses must have a control dependency on the acquire event. That is, if we add a control dependency edge from c to d , then this execution would be in $dead_M$. This ensures semantic deadness because the resultant litmus test, now having `if (r1==1) a=1` rather than just `a=1`, is now race-free. It therefore becomes a useful conformance test.

Formally, we reduce our general constraint-solving problem to finding inhabitants of the following set.

Definition 2 (General problem over executions). $\tilde{g}(M, N, \triangleright) \stackrel{\text{def}}{=} \{(X, Y) \in \mathbb{X}^2 \mid X \notin \text{consistent}_M \wedge X \in \text{dead}_M \wedge X \triangleright Y \wedge Y \in \text{consistent}_N\}$.

Analogies to our four specific problems, **Q1–Q4**, can be obtained by specialising \tilde{g} 's parameters like we did in §1.1.

We remark that deadness ensures the *soundness* of our solving strategy, but because we obtain semantic deadness via a syntactic approximation, it may spoil its *completeness*. Nonetheless, although our technique is incomplete, we demonstrate in the following subsection that it is *useful*.

1.3 Applications

We have implemented our technique in the Alloy modelling framework [39]. An Alloy model comprises a set of classes plus a set of constraints that relate objects and fields in those classes. If further provided with upper bounds on the number of objects in each class, Alloy can compile the constraints down to a SAT query, and then invoke a SAT solver to search for a satisfying instance.

We have applied our technique to a range of MCMs, including both software-level and architecture-level MCMs, both CPU and GPU varieties, and both operationally-defined and axiomatically-defined. Our results fall into two categories: automatic recreations of results that have previously been manually generated, and new results.

Recreated results We have rediscovered litmus tests that witness:

- the impact of three proposed changes to the C11 axioms (§4.1, §4.2, §4.3) – and our distinguishing litmus tests are substantially simpler than the originals in two cases;
- a violation of the supposedly guaranteed sequentially-consistent semantics for data-race-free programs (the ‘SC-DRF guarantee’ [5]) in a early draft of the C11 standard (§4.5) – similar to that reported by Batty et al. [16];
- that x86 is ‘multi-copy atomic’ [23, 72] but Power is not (§4.4);
- the C11 MCM behaving non-monotonically, by allowing tests to pass only if sequencing is added or a memory order is strengthened (§5) – and in the second case, our litmus test is simpler than that found manually by Vafeiadis et al. [77]; and
- two bugs in a published compiler mapping from OpenCL to AMD-style GPUs [62] (§6.1) – one of which is substantially simpler than the original found by Wickerson et al. [82].

New results Our main new result (§6.2) concerns the mapping of OpenCL to PTX, an assembly-like language for NVIDIA GPUs [61]. We first use **Q4** to show that a ‘natural’ OpenCL/PTX compiler mapping is unsound for an existing formalisation of the PTX MCM by Alglave et al. [9], but sound for a stronger PTX MCM that we propose. We then use **Q2** to generate litmus tests that distinguish the two PTX MCMs, which we use to validate our stronger MCM experimentally against actual NVIDIA GPUs.

In summary, we make the following contributions.

1. We show that four frequently-asked questions about MCMs can be viewed as instances of one general formula (g , Def. 1).
2. We rephrase the formula to constrain executions rather than litmus tests (\hat{g} , Def. 2), so that it can be tractably explored using a constraint-solving tool.
3. We implement our approach in Alloy, and use it to automatically reproduce several results obtained manually in previous work, often finding simpler examples.
4. We present a new, experimentally-validated MCM for PTX, and an OpenCL/PTX compiler mapping, and use Alloy to validate the mapping against the PTX MCM.

Our supplementary material [1] contains our Alloy models and our PTX testing results.

2. Executions

The semantics of a program is a set of *executions*. This section describes our formalisation of executions, both in general (§2.1) and specifically for C11, OpenCL, and PTX (§2.2), and then explains how MCMs decide which executions are allowed (§2.3).

Program executions are composed of *events*, each representing the execution of a program instruction. Most existing MCM frameworks embed several pieces of information within each event, such as the location it accesses, the thread it belongs to, and the value it reads and/or writes (e.g., [5, 11, 16, 55, 68]). In our work, events are *pure*, in Needham’s sense [57]: they are given meaning simply by their membership of, for instance, the ‘read events’ set, or the ‘same location’ relation. This formulation brings three benefits. First, it means that we can easily build a hierarchy of executions to

E	$\subseteq \mathbb{E}$	all events in the execution
R, W	$\subseteq \mathbb{E}$	events that read (resp. write) a location
F	$\subseteq \mathbb{E}$	fence events
nal	$\subseteq \mathbb{E}$	events that access a non-atomic location
sb	$\subseteq \mathbb{E}^2$	sequenced-before
ad, cd, dd	$\subseteq \mathbb{E}^2$	address/control/data dependency
$sthd, sloc$	$\subseteq \mathbb{E}^2$	same thread, same location
rf, co	$\subseteq \mathbb{E}^2$	reads-from, coherence order
rfe	$\subseteq \mathbb{E}^2$	$\stackrel{\text{def}}{=} rf - sthd$

- 1 $R \cup W \cup F \cup nal \subseteq E$
- 2 $(R \cup W) \not\cap F$
- 3 $sb \subseteq sthd$
- 4 $strict-po(sb)$
- 5 $ad \subseteq [R]; sb; [R \cup W]$
- 6 $dd \subseteq [R]; sb; [W]$
- 7 $cd \subseteq [R]; sb$
- 8 $equiv(sthd, E)$
- 9 $equiv(sloc, R \cup W)$
- 10 $nal; sloc = nal$
- 11 $rf \subseteq (W \times R) \cap sloc$
- 12 $rf; rf^{-1} \subseteq id$
- 13 $strict-po(co)$
- 14 $co \cup co^{-1} = (W - nal)^2 \cap sloc - id$

Figure 2. Basic executions, \mathbb{X}

unify several levels of abstraction. For instance, starting with a top-level ‘execution’ class, we can obtain a class of ‘C11 executions’ just by adding extra fields such as ‘the set of events with acquire semantics’; we need not modify the type of events. Second, it means that the same events can appear (with different meanings) in two executions, thus reducing the total number of events needed in our search space. Third, we avoid the need to define (and therefore, in a bounded search query, set the number of) locations, threads, and values.

Notation Our MCMs are written in Alloy’s modelling language, but in this paper we opt for more conventional set-theoretic notation. For a binary relation r , \bar{r} is its complement, r^{-1} is its inverse, $r^?$ is its reflexive closure, r^+ is its transitive closure, and r^* is its reflexive, transitive closure. The *strict-po* predicate holds of binary relations that are acyclic and transitive, and *equiv*(r, s) holds when r is a subset of s^2 , reflexive over s , symmetric, and transitive. We compose an m -ary relation r_1 with an n -ary relation r_2 (where $m, n \geq 1$) via $r_1; r_2 \stackrel{\text{def}}{=} \{(x_1, \dots, x_{m-1}, z_1, \dots, z_{n-1}) \mid \exists y. (x_1, \dots, x_{m-1}, y) \in r_1 \wedge (y, z_1, \dots, z_{n-1}) \in r_2\}$, and we lift a set to a subset of the identity relation via $[s] \stackrel{\text{def}}{=} \{(e, e) \mid e \in s\}$. We write $imm(r)$ for $r - (r; r^+)$, and $s_1 \not\cap s_2$ for $s_1 \cap s_2 = \emptyset$.

2.1 Basic Executions

Let \mathbb{E} be a set of events.

Definition 3 (Basic executions). The set \mathbb{X} of basic executions comprises structures with fourteen fields, governed by well-formedness constraints (Fig. 2). We write f_X for the field f of execution X , omitting the subscript when it is clear from the context. The constraints can be understood as follows. The subsets R , W , F , and nal are all drawn from the events E that appear in the execution **1**. In particular, compound read-modify-write (RMW) events belong to *both* R and W . Fences are distinct from reads and writes **2**. *Sequenced before* is an intra-thread strict partial order **3** **4**. (We allow sb within a thread to be partial because in C-like languages, the evaluation order of certain components, such as the operands of the $+$ -operator, is not specified.) Address dependencies are either read-to-read or read-to-write **5**, data dependencies are read-to-write **6**, and control dependencies originate from reads **7**. The *sthd* relation forms an equivalence among all events **8**, while *sloc* forms an equivalence among reads and writes **9**. We allow a distinction between ‘atomic’ and ‘non-atomic’ locations; MCMs that do not make this distinction (such as architecture-level MCMs)

$A \subseteq \mathbb{E}$ atomic events
 $acq, rel \subseteq \mathbb{E}$ events that have acquire (resp. release) semantics
 $sc \subseteq \mathbb{E}$ events that have SC semantics

$$\begin{aligned}
 & \textcircled{15} \quad acq \cup rel \cup sc \cup (R \cap W) \cup F \subseteq A \subseteq E \\
 & \textcircled{16} \quad R \cap sc \subseteq acq \subseteq R \cup F \quad \textcircled{17} \quad W \cap sc \subseteq rel \subseteq W \cup F \\
 & \textcircled{18} \quad F \cap sc \subseteq acq \cap rel \quad \textcircled{19} \quad R - A \subseteq nal \subseteq E - A
 \end{aligned}$$

Figure 3. C11 executions, \mathbb{X}_{C11} (extending \mathbb{X})

$dv \subseteq \mathbb{E}$ events that have whole-device scope
 $swg \subseteq \mathbb{E}^2$ same workgroup

$$\textcircled{20} \quad sthd \subseteq swg \quad \textcircled{21} \quad equiv(swg, E) \quad \textcircled{22} \quad dv \subseteq A$$

Figure 4. OpenCL executions, \mathbb{X}_{OpenCL} (extending \mathbb{X}_{C11})

simply constrain the set of non-atomic locations to be empty. The *nal* set consists only of complete *sloc*-classes $\textcircled{10}$. A relation *rf* is a well-formed *reads-from* if it relates writes to reads at the same location $\textcircled{11}$ and is injective $\textcircled{12}$. The inter-thread reads-from (*rfe*) is derived from *rf*. A relation *co* is a well-formed *coherence order* if when restricted to writes on a single atomic location it forms a strict total order; that is, it is acyclic and transitive $\textcircled{13}$, and it relates every pair of distinct writes to the same atomic location $\textcircled{14}$.

Remark 4. We emphasise that elements of \mathbb{E} have no intrinsic structure, only identity. Nevertheless, when drawing executions, we tag events with their type: R for elements of $R - W$, W for elements of $W - R$, and C (for ‘compound’) for elements of $R \cap W$. We indicate *sthd* equivalence classes with dotted rectangles, and *sloc* equivalence classes using named representatives, e.g. *x* and *y*. We also tag events with the values read/written, but this is purely for readability.

Remark 5 (Initial writes). Like some prior MCM formalisations [50, 69], but unlike most (e.g. [10, 11, 16, 43]), our executions exclude initial writes. We found that Alloy’s performance degrades rapidly as the upper bound on $|\mathbb{E}|$ increases; by avoiding initial writes, this bound can be lowered. Removing initial writes makes rf^{-1} a *partial* function, which complicates the definition of *from-read*, as described below.

Definition 6. *From-read* relates each read to all of the writes that are *co*-later than the write that the read observed [6]:

$$fr \stackrel{\text{def}}{=} ((rf^{-1}; co) \cup fr_{\text{init}}) - id$$

where $fr_{\text{init}} \stackrel{\text{def}}{=} ([R] - (rf^{-1}; rf)); sloc; [W]$. In the absence of initial writes, fr_{init} handles reads that observe the initial value.

2.2 Language-Specific Executions

We can obtain executions for various languages as subclasses of \mathbb{X} .

Definition 7 (C11 executions). C11 executions (\mathbb{X}_{C11}) are structures that inherit all the fields and well-formedness conditions from basic executions, and add those listed in Fig. 3. The new fields originate from the ‘memory orders’ that are attached to atomic operations in C11 [38 (§7.17.3)]. Acquire events, release events, SC events, RMWs, and fences are all atomic $\textcircled{15}$. Atomic events that are neither acquires nor releases correspond to C11’s ‘relaxed’ memory order. Acquire semantics is given to *all* SC reads and *only* reads and fences $\textcircled{16}$. Release semantics is given to *all* SC writes and *only* writes and fences $\textcircled{17}$. SC fences have both acquire and release semantics $\textcircled{18}$. Non-atomic reads access only non-atomic locations, and atomic operations never access non-atomic locations $\textcircled{19}$.

$dv \subseteq \mathbb{E}$ events that have whole-device scope
 $swg \subseteq \mathbb{E}^2$ same workgroup (‘co-operative thread array’)

$$\begin{aligned}
 & \textcircled{23} \quad sthd \subseteq swg \quad \textcircled{24} \quad equiv(swg, E) \quad \textcircled{25} \quad dv \subseteq F \quad \textcircled{26} \quad nal = \emptyset \\
 & \textcircled{27} \quad sthd - id \subseteq sb \cup sb^{-1} \quad \textcircled{28} \quad R \not\cap W
 \end{aligned}$$

Figure 5. PTX executions, \mathbb{X}_{PTX} (extending \mathbb{X})

Definition 8 (OpenCL executions). OpenCL [42] extends C11 with hierarchical models of execution and memory that reflect the GPU architectures it was primarily designed to target. Accordingly, OpenCL executions (\mathbb{X}_{OpenCL} , Fig. 4) extend C11 executions first by partitioning threads into one or more *workgroups* via the *swg* equivalence $\textcircled{20}$ $\textcircled{21}$, and second by allowing some atomic operations to be tagged as ‘device scope’ $\textcircled{22}$, which ensures they are visible throughout the device. Other atomics (i.e., with only ‘workgroup scope’) can efficiently synchronise threads within a workgroup but are unsuitable for inter-workgroup synchronisation. We do not consider OpenCL’s local memory in this work, and we restrict our attention to the single-device case.

Definition 9 (PTX executions). Like OpenCL executions, PTX executions (\mathbb{X}_{PTX} , Fig. 5) gather threads into groups $\textcircled{23}$ $\textcircled{24}$. Some PTX fences (*membar.g1*) have whole-device scope $\textcircled{25}$, and others (*membar.cta*) have only workgroup scope. There are no non-atomic locations $\textcircled{26}$, sequencing is total within each thread $\textcircled{27}$, and we do not consider RMWs $\textcircled{28}$.

Remark 10. When drawing C11 executions, we identify the sets A , acq , rel and sc by tagging events in $E - A$ with *na*, those in $A - acq - rel$ with *RLX*, those in $acq - rel - sc$ with *ACQ*, those in $rel - acq - sc$ with *REL*, those in $acq \cap rel - sc$ with *AR*, and those in sc with *SC*. In OpenCL or PTX executions, we tag events in $A - dv$ with *WG*, and those in dv with *DV*.

2.3 Consistent, Race-Free, and Allowed Executions

Each MCM M defines sets $consistent_M$ and $racefree_M$ of executions. (For architecture-level MCMs, which do not define races, the latter contains all executions.) The sets work together to define the executions allowed under M , as follows.

Definition 11 (Allowed executions). The executions of a program P that are allowed under an MCM M are

$$\llbracket P \rrbracket_M \stackrel{\text{def}}{=} \begin{cases} \llbracket P \rrbracket_0 \cap consistent_M & \text{if } \llbracket P \rrbracket_0 \cap consistent_M \subseteq racefree_M \\ \top & \text{otherwise} \end{cases}$$

where \top stands for an appropriate universal set. Here, $\llbracket P \rrbracket_0$ is the set of P ’s *candidate executions*. These can be thought of as the executions allowed under an MCM that imposes no constraints, and are discussed separately (see Def. 16). The equation above says that the allowed executions of P are its consistent candidates, unless a consistent candidate is *racy*, in which case *any* execution is allowed. (This is sometimes called ‘catch-fire’ semantics [19].)

The consistency and race-freedom axioms for C11 and OpenCL (minus ‘consume’ atomics) are defined in Figure 6 and explained below. We have included some recently-proposed simplifications. In particular, the simpler release sequence (proposed by Vafeiadis et al. [77]) makes deadness easier to define (§3.4), and omitting the total order ‘*S*’ over SC events (as proposed by Batty et al. [14]) avoids having to iterate over all total orders when showing an execution to be inconsistent.

Happens before (*hb*) edges are created by sequencing and by a release *synchronising with* (*sw*) an acquire in a different thread. Synchronisation begins with a release write (or a release fence that

$$\begin{array}{l}
\text{consistent}_{\text{C11}} \\
\textcircled{29} \text{acyclic}((hb \cap sloc) \cup rf \cup co \cup fr) \\
\textcircled{30} rf; [nal] \subseteq imm([W]; (hb \cap sloc)) \\
\textcircled{31} \text{acyclic}((Fsb^?; (co \cup fr \cup hb); sbF^?) \cap sc^2(\overline{incl})) \\
\\
\text{racefree}_{\text{C11}} \\
\textcircled{32} cnf - A^2 - sthd \subseteq hb \cup hb^{-1} \\
\textcircled{33} cnf \cap sthd \subseteq sb \cup sb^{-1} \quad \textcircled{34} cnf - incl \subseteq hb \cup hb^{-1} \\
\\
Fsb \stackrel{\text{def}}{=} [F]; sb \quad sbF \stackrel{\text{def}}{=} sb; [F] \quad rs \stackrel{\text{def}}{=} (sb \cap sloc)^*; rf^* \\
sw \stackrel{\text{def}}{=} [rel]; Fsb^?; [W \cap A]; rs; rf; [R \cap A]; sbF^?; [acq] \\
\overline{incl} \stackrel{\text{def}}{=} dv^2 \cup swg, \quad hb \stackrel{\text{def}}{=} ((sw(\overline{incl}) - sthd) \cup sb)^+ \\
cnf \stackrel{\text{def}}{=} ((W \times W) \cup (R \times W) \cup (W \times R)) \cap sloc - id
\end{array}$$

Figure 6. The C11 and OpenCL MCMs

precedes an atomic write) and ends with an acquire read (or an acquire fence that follows an atomic read) and the read observes either that write or a member of the write’s *release sequence* (rs). An event’s release sequence comprises all the writes to the same location that are sequenced after the event, as well as the RMWs (which may be in another thread) that can be reached from one of those writes via a chain of rf edges [77 (§4.3)]. In OpenCL, synchronisation only occurs between events with *inclusive scopes* ($incl$), which means that if the events are in different workgroups they must both be annotated with ‘device’ scope. Happens-before edges between events accessing the same location, together with rf , co , and fr edges, must not form cycles $\textcircled{29}$ [77 (§5.3)]. A read of a non-atomic location must observe a write that is still *visible* (vis) $\textcircled{30}$. The *sequential consistency* (SC) axiom $\textcircled{31}$ requires, essentially, that the co , hb and fr edges between sc events do not form cycles [14 (§3.2)] An execution has a *data race* unless every pair of conflicting (cnf) events in different threads, not both atomic, is in hb $\textcircled{32}$. It has an *unsequenced race* unless every pair of conflicting events in the same thread is in sb $\textcircled{33}$. An OpenCL execution has a *heterogeneous race* [34] unless every pair of conflicting events with non-inclusive scopes is in hb $\textcircled{34}$.

3. Obtaining Litmus Tests

If we find executions that solve \tilde{g} (Def. 2), we need to ‘lift’ these executions up to the level of programs in order to obtain a solution to the original problem, g (Def. 1).

We begin this section by defining a language for these generated programs (§3.1). The language is designed to be sufficiently expressive that for any discovered execution X , there exists a program in the language that can behave like X . In particular, the program must be able to create arbitrary sequencing patterns and dependencies. Beyond this, we keep the language as small as possible to keep code generation simple. The language is also generic, so that it can be used to generate both assembly and high-level language tests.

We go on to define a function that obtains litmus tests from executions (§3.2), and show that, under an additional assumption about executions, the function is total (§3.3). We then define our deadness constraint on executions (§3.4), and conclude with a discussion of some of the practical aspects of generating litmus tests (§3.5).

$$\begin{array}{l}
\mathbb{A}_s ::= s \mid \mathbb{A}_s + 0 \times \mathbb{R} \\
\mathbb{I} ::= \text{st}(\mathbb{A}_L, \mathbb{A}_V) \mid \text{ld}(\mathbb{R}, \mathbb{A}_L) \mid \text{cas}(\mathbb{A}_L, \mathbb{V}, \mathbb{A}_V) \mid \text{fe} \\
\mathbb{C} ::= \mathbb{I}^\ell \mid \mathbb{C} +^\ell \mathbb{C} \mid \mathbb{C};^\ell \mathbb{C} \mid \text{if}^\ell \mathbb{R} = \mathbb{V} \text{ then } \mathbb{C} \\
\mathbb{P} ::= \mathbb{C} \parallel \dots \parallel \mathbb{C}
\end{array}$$

Figure 7. A programming language

3.1 Programming Language

The syntax of our generic programming language is defined in Fig. 7. We postulate a set \mathbb{V} of values (containing zero), a set \mathbb{R} of registers, and a set \mathbb{L} of (shared) locations, which is subdivided into atomic (\mathbb{L}_a) and non-atomic (\mathbb{L}_{na}) locations. Every register and location is implicitly initialised to zero. Let \mathbb{A}_s denote a language of *expressions* over a set s . Observe that since the only construction is the addition of a register multiplied by zero, these expressions only evaluate to elements of s . We use them to create artificial dependencies (i.e. *syntactic* but not *semantic* dependencies). The set \mathbb{I} of *instructions* includes stores, loads, compare-and-swaps (CAS’s), and fences. The $\text{cas}(x, v, v')$ instruction compares x to v (and if the comparison succeeds, sets x to v'), returning the observed value of x in either case. Observe that, artificial dependencies notwithstanding, stores and CAS’s write only constant values to locations. Let \mathbb{C} denote a set of *components*, each formed from sequenced ($;$) or unsequenced ($+$) composition, or one-armed conditionals that test for a register having a constant value. We have no need for loops because the executions we generate are finite. Each component has a globally-unique label, ℓ . Let \mathbb{P} be the set of *programs*, each a parallel composition of components. The top-level components in a program are called *threads*.

Some languages attach extra information to instructions, such as their atomicity and memory order (in C11), their memory scope (in OpenCL), or whether they are ‘locked’ (in x86). Accordingly, when using the generic language in Fig. 7 to represent one of these languages, we introduce an additional function that stores extra information for each instruction label.

We ensure that the programs we generate are *well-formed* – this is necessary for ensuring that they provide valid solutions to g .

Definition 12 (Well-formed programs). A program is well-formed if: (1) different stores/CAS’s to the same location store different values, (2) each register is written at most once, (3) stores/CAS’s never store zero, (4) if-statements never test for zero, and (5) whenever an if-statement tests register r , there is a load/CAS into r sequenced earlier in the same thread.

3.2 From Executions to Litmus Tests

Our strategy for solving g , as outlined in §1.2, is summarised by the following ‘proof rule’:

$$\begin{array}{l}
\text{Step 1:} \quad (X, Y) \in \tilde{g}(M, N, \triangleright) \\
\text{Step 2:} \quad (P, \sigma) \in \text{lit}^{\min}(X) \quad (Q, \sigma') \in \text{lit}(Y) \quad \sigma = \sigma' \quad P \blacktriangleright Q \\
\hline
\text{Result:} \quad (P, Q, \sigma) \in g(M, N, \triangleright) \quad (1)
\end{array}$$

The purpose of this subsection is to define the lit and lit^{\min} functions. We begin by defining a more general constraint, lit' .

Definition 13. The predicate $\text{lit}'(X, P, \sigma, \text{disabled}, \text{failures})$ serves to connect executions X with litmus tests (P, σ) . The *disabled* argument is a subset of P ’s components, which we interpret as those that do not actually execute when creating the execution X (because they are guarded by an if-statement whose test failed). The *failures* argument is a subset of P ’s CAS instructions, which we interpret as those that fail to carry out the ‘swap’ when creating execution X . (Characterising executions by

which instructions fail and which are disabled is only sensible because of our restriction to loop-free programs.) The predicate $\mathit{lit}'(X, P, \sigma, \mathit{disabled}, \mathit{failures})$ holds whenever there exists a bijection μ between P 's enabled instructions and X 's events such that the following conditions all hold (in which we abbreviate $\mu(i)$ as μ_i):

Conditionals For an if-statement with condition ' $r = v$ ', the body is enabled iff both the if-statement is enabled and $\sigma(r) = v$.

Disabled loads If a load/CAS into register r is disabled, then $\sigma(r) = 0$.

Unguarded components Any component not guarded by an if-statement is enabled.

CAS failures For every enabled CAS i in P : i is in $\mathit{failures}$ iff there is no j with $(\mu_j, \mu_i) \in \mathit{rf}$ that writes the value i expects.

Instruction types For every enabled instruction i in P : i is a store iff $\mu_i \in W - R$; i is a load or a failed CAS iff $\mu_i \in R - W$; i is a successful CAS iff $\mu_i \in R \cap W$; and i is a fence iff $\mu_i \in F$.

Non-atomic locations If i is an enabled non-fence instruction in P then $\mu_i \in \mathit{nal}$ iff i 's location is in \mathbb{L}_{na} .

Threading, locations, and dependencies For every enabled instructions i and j in P : $(\mu_i, \mu_j) \in \mathit{sthd}$ iff i and j are in the same thread in P ; $(\mu_i, \mu_j) \in \mathit{sloc}$ iff i and j both access the same location; $(\mu_i, \mu_j) \in \mathit{cd}$ iff there is an enabled if-statement, say T , such that i is sequenced before T , i writes to the register T tests, and j is in T 's body; $(\mu_i, \mu_j) \in \mathit{ad}$ iff the expression for j 's location depends (syntactically) on the register written by i ; and $(\mu_i, \mu_j) \in \mathit{dd}$ iff j writes an expression that depends (syntactically) on the register i writes.

Sequenced composition For every enabled $;$ -operator in P : if i and j are enabled instructions in the left and right operands (respectively), then $(\mu_i, \mu_j) \in \mathit{sb}$.

Unsequenced composition For every enabled $+$ -operator in P : if i and j are enabled instructions in the left and right operands (respectively), then $\{(\mu_i, \mu_j), (\mu_j, \mu_i)\} \not\in \mathit{sb}$.

Final registers For every enabled load/CAS j in P on location x into register r : either (1) there exists an enabled store/CAS i that writes v to x , $(\mu_i, \mu_j) \in \mathit{rf}$ and $\sigma(r) = v$, or (2) $\sigma(r) = 0$ and μ_j is not in the range of rf .

Final locations For every location x : either (1) there is an enabled store/CAS i that writes v to x , μ_i has no successor in co , and $\sigma(x) = v$, or (2) x is never written and $\sigma(x) = 0$.

Definition 14 (Obtaining litmus tests). Let $\mathit{lit}(X)$ be the set of all litmus tests (P, σ) for which $\mathit{lit}'(X, P, \sigma, \mathit{disabled}, \mathit{failures})$ holds for some instantiation of $\mathit{disabled}$ and $\mathit{failures}$.

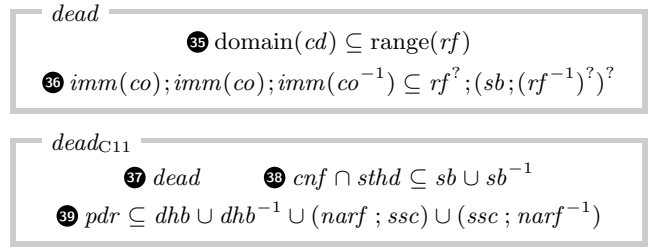
Definition 15 (Obtaining minimal litmus tests). Let $\mathit{lit}^{\min}(X)$ be the set of all litmus tests (P, σ) for which $\mathit{lit}'(X, P, \sigma, \emptyset, \emptyset)$ holds. That is, when P behaves like X , all of its instructions are executed and all of its CAS's succeed. We say that the litmus test (P, σ) is *minimal* for X , or, dually, X is a *maximal* execution of P .

Definition 16 (Candidate executions). We define P 's candidate executions by inverting lit : $\llbracket P \rrbracket_0 \stackrel{\text{def}}{=} \{X \mid \exists \sigma. (P, \sigma) \in \mathit{lit}(X)\}$.

Definition 17 (Observable final states). We can now formally define the notion of observation we employed in Def. 1: $\mathit{obs}_M(P)$ is the set of final states σ for which $(P, \sigma) \in \mathit{lit}(X)$ for some $X \in \llbracket P \rrbracket_M$.

3.3 Totality of lit^{\min}

We now explain why the lit^{\min} function is currently not total (that is: there exist well-formed executions X for which $\mathit{lit}^{\min}(X)$ is empty), and how we can impose an additional restriction on executions to make it become total.



$$\begin{aligned} \mathit{pdr} &\stackrel{\text{def}}{=} \mathit{cnf} - A^2 & \mathit{cde} &\stackrel{\text{def}}{=} (\mathit{rfe} \cup \mathit{cd})^*; \mathit{cd} & \mathit{drs} &\stackrel{\text{def}}{=} \mathit{rs} - ([R]; \overline{\mathit{cde}}) \\ \mathit{dsw} &\stackrel{\text{def}}{=} \mathit{sw} \cap (((F\mathit{sb}^?; [\mathit{rel}]; \mathit{drs}^?) - (\mathit{cd}^{-1}; \overline{\mathit{cde}})); \mathit{rf}) \\ \mathit{dhb} &\stackrel{\text{def}}{=} \mathit{sb}^?; (\mathit{dsw}; \mathit{cd})^* & \mathit{ssc} &\stackrel{\text{def}}{=} \mathit{id} \cap \mathit{cde} & \mathit{narf} &\stackrel{\text{def}}{=} \mathit{rf} \cap \mathit{nal}^2 - \mathit{hb} \end{aligned}$$

Figure 8. The *dead* constraint and its specialisation for C11

The problem is: our programming language cannot express all of the sequencing patterns that an execution can capture in a valid sb relation. For instance, it is not possible to write a program that generates exactly the following sb edges:



This is because, armed only with sequenced ($;$) and unsequenced ($+$) composition (see Fig. 7), it is only possible to produce sb relations that are in the set of *series-parallel partial orders* [53].² Helpfully, series-parallel partial orders are characterised exactly by a simple check: that they do not contain the ‘N’-shaped sub-graph shown in (2) [78]. Accordingly, we impose one further well-formedness constraint on executions:

$$\begin{aligned} \nexists a, b, c, d \in E. \\ \{(a, c), (a, d), (b, c)\} \in \mathit{sb} \wedge \{(a, b), (b, c), (c, d)\} \not\in \mathit{sb}^?. \end{aligned}$$

Our lit^{\min} function now becomes total (and hence lit too). Indeed, it is straightforward to determinise the constraints listed in Def. 13 into an algorithm for *constructing* litmus tests from executions (and we have implemented this algorithm, in Java).

3.4 Dead Executions

When searching for inconsistent executions, we restrict our search to those that are also *dead*.

Definition 18 (Semantic deadness). The set of executions that are (semantically) dead under MCM M is given by $\mathit{semdead}_M \stackrel{\text{def}}{=} \{X \in \mathbb{X} \mid X \notin \mathit{consistent}_M \Rightarrow \forall P, \sigma, X', \sigma'. ((P, \sigma) \in \mathit{lit}^{\min}(X) \wedge (P, \sigma') \in \mathit{lit}(X') \wedge X' \in \mathit{consistent}_M) \Rightarrow (X' \in \mathit{racefree}_M \wedge \sigma' \neq \sigma)\}$.

That is, for any minimal litmus test (P, σ) for X , no consistent candidate execution of P is racy or reaches σ . In other words: (P, σ) must fail under M .

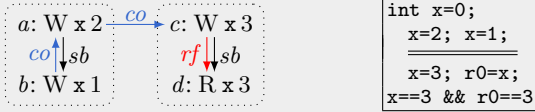
We employ syntactic approximations of semantic deadness. Figure 8 shows how this approximation is defined for any architecture-level MCM that enforces coherence (*dead*), and how it is strengthened to handle the C11 MCM (*dead*_{C11}).

At the architecture level, we need not worry about races, so ensuring deadness is straightforward. In what follows, let X be an

²That said, if we extended our language to support fork/join parallelism, then the execution in (2) would become possible: $\mathit{t1}=\mathit{fork}(b); a; \mathit{t2}=\mathit{fork}(c); \mathit{join}(\mathit{t1}); d; \mathit{join}(\mathit{t2})$.

inconsistent execution, (P, σ) be a minimal litmus test of X , and X' be another candidate execution of P . First, we require every event that is the source of a control dependency to read from a non-initial write $\textcircled{35}$. This ensures that P need not contain ‘if $r = 0$ ’. Such programs are problematic because we cannot tell whether the condition holds because r was set to zero or because r was never assigned. Second, we require that every co edge (except the last) is justified by an sb edge $\textcircled{36}$. This condition ensures that X' cannot be made consistent simply by inverting one or more of X 's co edges. The construction $imm(co)$ obtains event pairs (e_1, e_2) that are consecutive in co ; composing this with ‘ $imm(co)$; $imm(co^{-1})$ ’ restricts our attention to those pairs for which it is possible to take a further co step from e_2 . The last co edge does not need justifying with an sb edge because it is directly observable in σ .

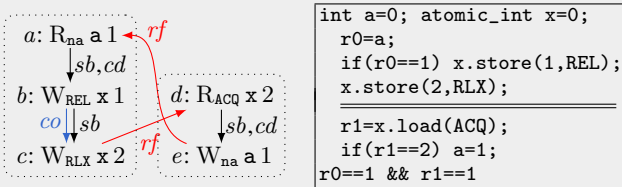
Example 2. The basic execution below (left) is inconsistent in any MCM that imposes coherence (because co is contradicting sb), but the litmus test obtained from it (below right) does not necessarily fail because its final state can be obtained via a consistent execution that simply reverses the co edge from (b, a) to (a, b) .



At the software level, we must also worry about races. First, we forbid all unsequenced races $\textcircled{38}$, because if X does not have an unsequenced race, then neither will X' , because unsequenced races are not affected by runtime synchronisation behaviour. Second, condition $\textcircled{39}$ is concerned with *potential data races* (pdr): events that conflict and are not both atomic. Although X is inconsistent (which renders any races in X irrelevant) we worry that X' might be consistent and racy. So, we require pdr -linked events to be also in the *dependable happens-before* (dhb) relation, or to exhibit a *self-satisfying cycle* (ssc), both of which we explain below.

Dependable happens-before This is a restriction of ordinary happens-before (hb). Essentially: if e_1 and e_2 are in dhb in X , and they map to instructions i_1 and i_2 respectively in P , then if i_1 and i_2 are represented by events in X' , those events are guaranteed to be related by happens-before.

Example 3. The execution below (left) has fixed the shortcoming in Example 1 by adding control dependencies, but it has introduced the C11 *release sequence* as a further complexity.



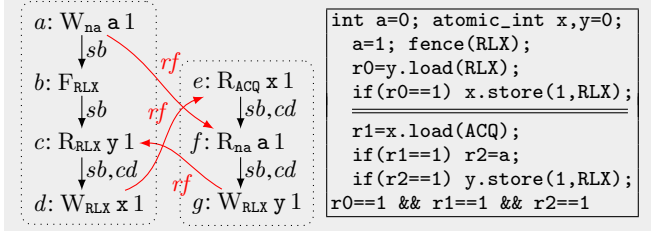
Although event d synchronises with b , it actually obtains its value from c , in b 's release sequence. The execution is not semantically dead because its litmus test (above right) is racy: if $r0$ is not assigned 1, then the release store is not executed; this means that $r1$ can read 2 without synchronisation having occurred; this leads to a race on a . By subtracting $(cd^{-1}; cde)$ in the definition of dsw , we ensure that whatever controls b also controls c , and this rules out undesirable executions like the one above.

Moreover, the $([R]; cde)$ in the definition of drs ensures that if b is an RMW, it controls the execution of c . The effect is that c is not executed if the CAS corresponding to the RMW fails.

Self-satisfying cycles An event occurs in a self-satisfying cycle (ssc) if it is connected to itself via a chain of cd and rf edges that ends with a control dependency. The if-statements that create these dependencies are constructed such that their bodies are only executed if the desired rf edges are present (cf. Def. 13).

A potential race between e_1 and e_2 is deemed dead if both access a non-atomic location, e_1 is observed by e_2 but does not happen before it, and e_2 is in a self-satisfying cycle. The reasoning is as follows. First, note that the execution is inconsistent, because reads of non-atomic locations cannot observe writes that do not happen before them (c.f. $\textcircled{40}$). Second, the self-satisfying cycle ensures that e_2 reads from e_1 in every candidate execution that includes those events. Therefore, every candidate execution is inconsistent and any data races can be safely ignored. We illustrate this reasoning in the following example.

Example 4. The execution below (left) is inconsistent because f reads a non-atomic location from a write (a) that does not happen before f . It is dead because f is in a self-satisfying cycle. Its litmus test (below right) is not racy because the conditionals ensure that the right-hand thread's load of a is only executed if it obtains the value 1, which means that it reads from the left-hand thread's store to a , which means that the execution is inconsistent and hence that any races can be ignored.



Checking deadness The constraints that define $dead_M$ are quite subtle, particularly for complex MCMs like C11. Fortunately, we can use Alloy to check that these constraints imply semantic deadness, by seeking elements of $dead_M - semdead_M$. That is: we search for an execution X that is deemed dead, but which gives rise to a litmus test (P, σ) that has, among its candidates, a consistent execution X' that is either racy or reaches σ . We were able to check that $dead_M \subseteq semdead_M$ holds for all executions with no more than 8 events; beyond this, Alloy timed out (after four hours).

Remark 19. We are using Alloy here to search all candidate executions of a program, yet in §1.2 we argued that it is impractical to phrase constraints over programs because to do so would necessitate an expensive search over all candidate executions of a program. We emphasise that this is not a contradiction. The formula to which we objected had a problematic $\exists \forall$ pattern: ‘show that *there exists* a program such that *for all* of its candidate executions, a certain property holds’, whereas satisfying $dead_M - semdead_M$ requires only existential quantification.

3.5 Practical Considerations

We have found that Step 2 of our proof rule (1) can still be hard for Alloy to solve, often leading to timeouts, particularly when X and Y are large. The search space is constrained quite tightly by the values of X and Y , but there are still many variables involved. One of the degrees of flexibility in choosing P and Q is in the handling of if-statements. For instance, both $\textcircled{40}$ and $\textcircled{41}$ below give rise to the same executions because of our well-formedness restrictions on programs (Def. 12):

$\textcircled{40}$ if b then $(C_1 ; C_2)$ $\textcircled{41}$ (if b then C_1) ; (if b then C_2).

$$\text{consistent}_{SC} \quad \textcircled{44} \text{acyclic}(rf \cup co \cup fr \cup sb)$$

Figure 10. The SC MCM (following Shasha et al. [70])

$$\begin{aligned} \text{strengthen}(X, Y) \quad & \textcircled{45} E_X = E_Y \quad \textcircled{46} R_X = R_Y \quad \textcircled{47} W_X = W_Y \\ & \textcircled{48} F_X = F_Y \quad \textcircled{49} \text{nal}_X = \text{nal}_Y \quad \textcircled{50} sb_X \subseteq sb_Y \\ & \textcircled{51} ad_X \subseteq ad_Y \quad \textcircled{52} cd_X \subseteq cd_Y \quad \textcircled{53} dd_X \subseteq dd_Y \\ & \textcircled{54} sthd_X \subseteq sthd_Y \quad \textcircled{55} sloc_X = sloc_Y \quad \textcircled{56} rf_X = rf_Y \\ & \textcircled{57} co_X = co_Y \end{aligned}$$

Figure 11. ‘Strengthening’ an execution

the intuition that if two reads, say r_1 and r_2 (with $(r_1, r_2) \in ppo$), observe two writes, say w_1 and w_2 respectively, then any write *co*-later than w_2 must follow w_1 in the single copy of memory. The *wo* relation is the union of two cases: in the first, both reads observe write events, and in the second, one reads from the initial value (reusing our definition of fr_{init} from Def. 6). Note that MCA is parameterised by the given model’s *ppo*.⁵

With this formal definition of MCA, we can seek executions allowed by a given MCM but disallowed by MCA. We tested x86 and Power, and as expected, Power does not guarantee MCA (Alloy finds a counterexample similar to (3)) but x86 does.

4.5 Seeking SC-DRF Violations

We used Alloy to seek violations of the SC-DRF guarantee [5] in an early draft of the C11 MCM [37]. The SC-DRF guarantee, in the C11 context, states that if a program is race-free and contains no non-SC atomic operations, then it has only SC semantics.

We sought an execution X that is dead (so that its corresponding litmus test is race-free), uses no non-SC atomics, and is consistent under the (draft) C11 MCM but inconsistent under SC:

$$A_X \subseteq sc_X \wedge X \in \text{dead}_{C11} \cap \text{consistent}_{C11\text{-Draft}} - \text{consistent}_{SC}$$

where the SC MCM is characterised in Fig. 10 and C11-Draft has all the axioms listed by Batty et al. [14 (Def. 11)], minus their ‘S4’ axiom. Alloy found an example similar to that reported by Batty et al. [16 (§4, *Sequential consistency for SC atomics*)]. The non-SC execution is consistent under C11-Draft because no rule forbids SC reads to observe initial writes. The ‘S4’ axiom was added [16 (§2.7)] to fix exactly this issue.

5. Application: Checking Monotonicity (Q3)

In this section, we rediscover two monotonicity violations in the C11 MCM, whereby a new behaviour is enabled either by sequentialisation (§5.1) or by strengthening a memory order (§5.2).

Checking monotonicity requires the relation defined in Fig. 11, which holds when one execution is ‘stronger’ than another. There is almost an isomorphism between X and Y , except that Y may have extra *sb* and dependency edges $\textcircled{50}$ $\textcircled{51}$ $\textcircled{52}$ $\textcircled{53}$, and it may interleave multiple threads together $\textcircled{54}$.

⁵ We have an alternative formulation for happens-before-based models.

$$\begin{aligned} \text{compile}(\pi, X, Y) \quad & \textcircled{58} [E_X] = \pi; \pi^{-1} \quad \textcircled{59} [E_Y] \subseteq \pi^{-1}; \pi \\ & \textcircled{60} \pi^{-1}; sb_X^?; \pi = sb_Y^? \quad \textcircled{61} ad_X = \pi; ad_Y; \pi^{-1} \\ & \textcircled{62} dd_X = \pi; dd_Y; \pi^{-1} \quad \textcircled{63} cd_X; \pi = \pi; cd_Y \\ & \textcircled{64} sloc_X = \pi; sloc_Y; \pi^{-1} \quad \textcircled{65} \pi^{-1}; sthd_X; \pi = sthd_Y \\ & \textcircled{66} rf_X = \pi; rf_Y; \pi^{-1} \quad \textcircled{67} co_X = \pi; co_Y; \pi^{-1} \end{aligned}$$

Figure 12. Compiling an execution

5.1 Monotonicity of C11 w.r.t. Sequencing

One way to extend the *strengthen* relation to the C11 setting is to add the following constraints:

$$A_X = A_Y \quad acq_X = acq_Y \quad rel_X = rel_Y \quad sc_X = sc_Y.$$

With this notion of strengthening, Alloy found a pair of 6-event executions that witness a monotonicity violation in C11. They are almost identical to those given by Vafeiadis et al. [77 (Fig. 1)], though slightly less elegant: Alloy chose one of the writes to be SC when a relaxed write would have sufficed.

5.2 Monotonicity of C11 w.r.t. Memory Orders

Another way to extend the *strengthen* relation is to add

$$\begin{aligned} A_X = A_Y \quad acq_X \subseteq acq_Y \quad rel_X \subseteq rel_Y \quad sc_X \subseteq sc_Y \\ sb_X = sb_Y \quad ad_X = ad_Y \quad cd_X = cd_Y \quad dd_X = dd_Y \\ sthd_X = sthd_Y \end{aligned}$$

which allows memory orders to be strengthened but forbids changes to sequencing or threading. Using this notion of strengthening, Alloy found a 7-event monotonicity violation in C11: the execution previously given in Example 4. That execution is inconsistent, but if the relaxed fence is strengthened to a release fence, it becomes consistent. This example is similar to one due to Vafeiadis et al. [77 (§3, *Strengthening is Unsound*)], but it is simpler: where theirs requires 8 events, 4 locations, and 3 threads, our Alloy-generated example requires only 7 events, 3 locations, and 2 threads.

6. Application: Checking Compiler Mappings (Q4)

We now report on our use of Alloy to investigate OpenCL compiler mappings for AMD-style GPUs (§6.1) and NVIDIA GPUs (§6.2), and a C11 compiler mapping for Power (§6.3).

Checking compiler mappings requires a relation that holds when one execution ‘compiles’ to another. Of course, real compilers do not compile on a ‘per-execution’ level, but from a pair of executions related in this way, we can obtain a pair of programs that are related by the compiler mapping at the level of program code. Compiler mappings are more complicated than the strengthenings considered in §5, because they may introduce additional events, such as fences. For this reason, we introduce an additional relation, π , to connect each original event to its compiled event(s). (In §5, π is the identity relation.)

We are able to handle ‘straight line’ mappings like C11/x86 [16], OpenCL/AMD (§6.1) and OpenCL/PTX (§6.2). We can also handle more complicated mappings like C11/Power [15, 66], which includes RMWs and introduces control dependencies, but our examples are limited to the execution level, because at the code level, C11 RMWs map to Power loops, which our loop-free output programming language cannot express.

Execution X maps to execution Y if there exists a π for which $\text{compile}(\pi, X, Y)$ holds, as defined in Fig. 12. Here, π is an injective, surjective, multivalued function from X ’s events to Y ’s

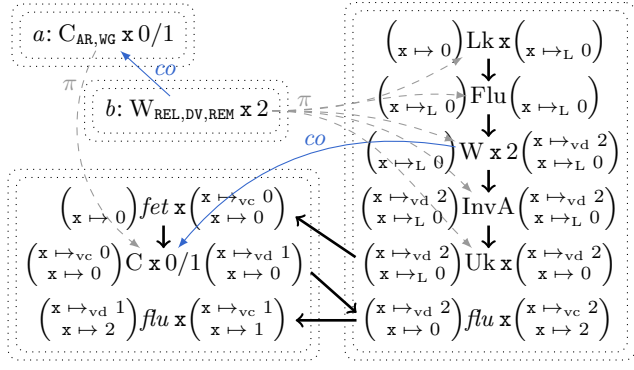


Figure 13. RMW atomicity bug in the OpenCL/AMD mapping

events 58–59, that preserves sequencing 60, dependencies 61–62–63, locations 64, threads 65, the reads-from relation 66, and the coherence order 67. The location-preservation and thread-preservation constraints have different shapes because, for instance, when a mapping sends source event e to a set $\{e'_i\}_i$ of target events, it must ensure that *one* e'_i has e 's location (the other events will be fences), but that *every* e'_i is in the same thread.

6.1 Compiling OpenCL to AMD-Style GPUs

Orr et al. [62] describe a compiler mapping from OpenCL to an AMD-style GPU architecture. Actually, they support OpenCL extended with ‘remote-scope promotion’, in which a workgroup-scoped event a can synchronise with an event b in a different workgroup if b is annotated as ‘remote’. Wickerson et al. [82] report on two bugs in this scheme: a failure to implement message-passing, and a failure of RMW atomicity.

We have formalised Orr et al.’s architecture-level MCM and compiler mapping in Alloy [1], following the formalisation by Wickerson et al., and then used Alloy to search for bugs – essentially automating the task that Wickerson et al. conducted manually. The architecture-level MCM is *operational*, which means that consistent executions are obtained constructively, not merely characterised by axioms. This means that the MCM is more complex to express in Alloy. Essentially, the MCM involves a single global memory in which entries can be temporarily locked, several processing elements partitioned into compute units, and a write-back write-allocate cache per compute unit.

Fig. 13 depicts the RMW atomicity bug discovered by Alloy. The top left of the figure shows a 2-event execution that is inconsistent (and dead) in OpenCL, and the right shows a corresponding 9-event execution that is observable on hardware. (Wickerson et al.’s bug is similar, but requires an extra ‘fetch’ event.) All events are grouped by thread (inner dotted rectangles) and by workgroup (outer dotted rectangles). Events in the architecture-level execution are totally ordered (thick black arrows). We track the local and global state before and after each event, writing $\begin{pmatrix} \sigma_l \\ \sigma_g \end{pmatrix}$ for local state σ_l and global state σ_g . The global state comprises global memory entries, some of which are locked (L). The local state comprises the compute unit’s cache entries, each either valid (v) or invalid (i), and either clean (c) or dirty (d). The π edges show how the software-level events are mapped to architecture-level events. The RMW at workgroup scope (a) maps to a single RMW, while the remote (REM) write to x (b) is implemented by first flushing all dirty local cache entries (Flu), then doing the write (W), then invalidating all entries in all caches (InvA), all while preventing concurrent access to x in the global memory (Lk, Uk). Architecture-level executions also include the actions of the ‘environment’ fetching (*fet*) and flushing (*flu*) entries to and from global memory, as well as

$$\begin{array}{l} \text{consistentPTX1} \\ \textcircled{68} \text{ } \textit{acyclic}((sb \cap \textit{sloc} - R^2) \cup \textit{rf} \cup \textit{co} \cup \textit{fr}) \\ \textcircled{69} \text{ } \textit{acyclic}(hb_{wg}) \quad \textcircled{70} \text{ } \textit{acyclic}(hb_{dv}) \end{array}$$

$$\begin{array}{l} \textit{rmo} \stackrel{\text{def}}{=} \textit{rfe} \cup \textit{co} \cup \textit{fr} \quad \textit{f}_{dv} \stackrel{\text{def}}{=} sb; [F \cap dv]; sb \quad \textit{f}_{any} \stackrel{\text{def}}{=} sb; [F]; sb \\ \textit{hb}_{wg} \stackrel{\text{def}}{=} (\textit{rmo} \cup \textit{f}_{any}) \cap \textit{swg} \quad \textit{hb}_{dv} \stackrel{\text{def}}{=} \textit{rmo} \cup \textit{f}_{dv} \end{array}$$

Figure 14. The PTX1 MCM

71	store(na RLX, s)	\rightsquigarrow st.cg	
72	store(REL, s)	\rightsquigarrow F _s ; st.cg	
73	store(SC, s)	\rightsquigarrow F _s ; st.cg; F _s	where
74	load(na, s)	\rightsquigarrow ld.cg	F _{wg} $\stackrel{\text{def}}{=}$ membar.cta
75	load(RLX ACQ, s)	\rightsquigarrow ld.cg; F _s	F _{dv} $\stackrel{\text{def}}{=}$ membar.gl
76	load(SC, s)	\rightsquigarrow F _s ; ld.cg; F _s	
77	fence(s)	\rightsquigarrow F _s	

Table 1. OpenCL/PTX compiler mapping, program code level

derived *rf* and *co* relations. The undesirable but observable execution, in which $x = 1$ in the final state, arises because the mapping fails to propagate $x = 2$ to the global memory before releasing the lock on x . The fix is to flush immediately after the write.

To make Alloy’s search for OpenCL/AMD compiling bugs tractable, we found it necessary to make several simplifications: we deleted the QuickRelease buffers [29]; we allowed multiple locations to be fetched and flushed in a single action (which reduces the total number of actions required); we hard-coded the system to have exactly two workgroups with one thread each; and we maximised sharing between global and local memory entry objects. These changes are not necessarily semantics-preserving, but any bogus solutions found using the simplified MCM can be simply tested manually against the full one.

6.2 Compiling OpenCL to PTX

In this subsection, we develop and check a compiler mapping from OpenCL to PTX. First, using Q4, we show that a natural mapping is invalid for an existing formalisation of the PTX MCM (PTX1) [9], but valid for a stronger model that we develop (PTX2). Then, we use Q2 to generate litmus tests that distinguish PTX2 from PTX1, which we use to confirm that our PTX2 remains empirically sound for NVIDIA GPUs.

Figure 14 defines the PTX1 MCM. The model enforces coherence (but not between reads, à la Sparc RMO [8]) 68, allows any fence to restore SC within a workgroup 69, and allows device-scoped fences to restore SC throughout the device 70. We omit dependencies because they are difficult to test: the PTX compiler often removes our artificially-inserted dependencies.

Table 1 defines the OpenCL/PTX mapping we use. OpenCL stores are mapped to PTX stores (st.cg) 71–72–73, and OpenCL loads are mapped to PTX loads (ld.cg) 74–75, with fences (membar) placed before and/or after the memory access depending on whether the OpenCL instruction is non-atomic, relaxed, acquire, release, or SC. PTX fences are required even for relaxed OpenCL loads because PTX allows consecutive loads from the same address to be re-ordered (see 68). The PTX fences are scoped to match the scope s of the OpenCL instruction. In line with prior work on the PTX MCM [9], we exclude RMWs, local memory, and multi-GPU interactions.

Figure 15 shows how the code-level mappings in (the first two rows of) Tab. 1 are encoded at the execution level. OpenCL workgroups correspond to PTX workgroups 76. An OpenCL write with-

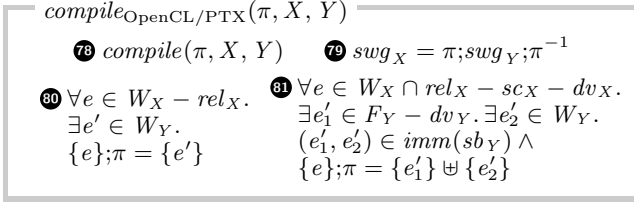


Figure 15. OpenCL/PTX mapping, execution level (excerpt)

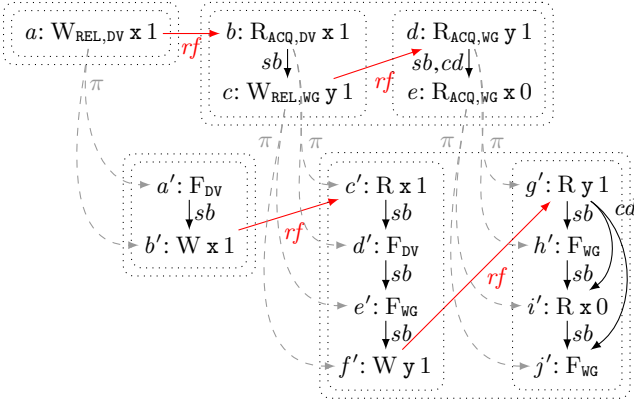


Figure 16. A WRC bug in the OpenCL/PTX1 mapping

out release semantics (e) maps to a single PTX write (e') $\textcircled{80}$, while an OpenCL workgroup-scoped release write (e) maps to a workgroup-scoped fence (e'_1) sequenced before a PTX write (e'_2) $\textcircled{81}$. The other rows are handled similarly.

We used Alloy to check this mapping against PTX1, and found an execution that is disallowed by OpenCL (Fig. 16, top) but allowed, after compilation, by PTX1 (Fig. 16, bottom). Note that the outer dotted rectangles group threads by workgroup. The crux here is *cumulative* synchronisation across scopes [36 (§1.7.1)]. The left thread synchronises at workgroup scope with the middle thread (via a and b), which synchronises at device scope with the right thread (via c and d). If PTX1 supported cumulative synchronisation across scopes, the left and right threads would now be synchronised, and the execution above, in which e observes a stale y , would be forbidden, just as it is in OpenCL.

We could fix the mapping for PTX1 by upgrading all the PTX fences in Tab. 1 to device scope. However, because widely-scoped fences incur high performance overhead on NVIDIA GPUs [71], we opt instead to strengthen PTX1 to support our mapping, by enforcing cumulative synchronisation across scopes (obtaining PTX2). This entails changing the definition of hb_{dv} (Fig. 14) to

$$hb_{dv} \stackrel{\text{def}}{=} rmo \cup (hb_{wg}^*; f_{dv}; hb_{wg}^*)$$

so that device-scoped fences (f_{dv}) can restore SC throughout the device even if they are preceded or followed by workgroup-scoped synchronisation (hb_{wg}). After this change, Alloy finds no bugs in our mapping with up to 5 software-level events. It times out (after four hours) when checking larger executions.

It remains to show that PTX2 is sound w.r.t. empirical GPU testing data. To do this, we use **Q2** to find litmus tests (P, σ) that are allowed under PTX1 but disallowed under PTX2. We wish to find not just a single litmus test (as in §4), but as many as possible, so we run Alloy iteratively, each time disallowing the execution shape found previously, until it is unable to find more. We then

Task	$ \mathbb{E} $	t_{enc}/s	t_{sol}/s		
1 Q2 C11-SRA vs. C11 (§4.1)	6	0.7	0.6	G	✓
2 Q2 C11-SwRf vs. C11 (§4.2)	7	0.8	625	G	✗
3 Q2 C11-SwRf vs. C11 (§4.2)	12	2	214	P	✓
4 Q2 C11 vs. C11-Orig (§4.3)	5	0.4	0.3	G	✓
5 Q2 MCA vs. x86 (§4.4)	9	0.8	607	P	✗
6 Q2 MCA vs. Power (§4.4)	6	2	0.06	G	✓
7 Q2 SC vs. C11-Draft (§4.5)	4	0.4	0.04	G	✓
8 Q2 PTX2 vs. PTX1 (§6.2)	7	0.7	4	G	✓
9 Q3 C11 (sequencing) (§5.1)	5	0.5	163	M	✗
10 Q3 C11 (sequencing) (§5.1)	6	0.7	5	P	✓
11 Q3 C11 (mem. orders) (§5.2)	7	0.9	51	G	✓
12 Q4 C11 / x86	5 + 5	0.7	13029	P	✗
13 Q4 C11 / Power (§6.3)	6 + 8	8	91	P	✓
14 Q4 OpenCL / AMD (§6.1)	2 + 9	6	1355	G	✓
15 Q4 OpenCL / AMD (§6.1)	4 + 10	16	4743	P	✓
16 Q4 OpenCL / PTX1 (§6.2)	5 + 8	2	11	P	✓
17 Q4 OpenCL / PTX2 (§6.2)	5 + 15	4	9719	P	✗

Table 2. Summary of tasks. We record the number ($|\mathbb{E}|$) of events in the search space, the time taken (in seconds) to encode (t_{enc}) and then solve (t_{sol}) the SAT query, whether the fastest solver was MiniSat (M), Glucose (G), or Plingeling (P), and whether a solution was found (✓) or not (✗). For **Q4** tasks, $m + n$ means that \mathbb{E} is partitioned into m software-level and n architecture-level events.

check testing results (or if results do not exist for a given test, we run new tests) to confirm that (P, σ) cannot pass on actual GPUs.

Using this method, Alloy finds all 14 distinguishing 7-event tests (e.g. WRC [50]), plus 12 of the distinguishing 8-event tests (e.g. IRIW [19]) before timing out. We are able to query Alglave et al.’s experimental results [9] for 22 of these 26 tests. The rest are single-address tests (which arise because PTX does not guarantee coherence in general). These we ran on an NVIDIA GTX Titan using the GPU-litmus tool [9]. We found no behaviours that are allowed by PTX1, disallowed by PTX2, and empirically observable on a GPU [1]. Thus, subject to the available testing data, strengthening PTX1 to PTX2 is sound, and thus the natural OpenCL/PTX compiler mapping can be used.

6.3 Compiling C11 to Power

Work in progress by Lahav et al. [44] has uncovered a bug in the C11/Power mapping previously thought to have been proven sound by Batty et al. [15].⁶ Before becoming aware of their work, we had used Alloy to verify the soundness of the mapping for up to 6 software-level events and up to 6 architecture-level events. Incrementing these bounds any further resulted in intractable solving times, which explains why the bug, which requires 6 software-level events and 13 architecture-level events, had not previously been found by Alloy. In order to recreate Lahav et al.’s result, we modified the C11/Power mapping so as not to place fences at the start or the end of a thread. Removing these redundant fences allows the bug to be expressed using just 8 architecture-level events, and found automatically by Alloy in a reasonable time (see Tab. 2 in the next section).

7. Performance Evaluation

In this section, we report on an empirical investigation of how our design decisions affect Alloy’s SAT-solving performance.

⁶Concurrent work by Manerkar et al. [51] has shown the C11/ARMv7 mapping to be similarly flawed.

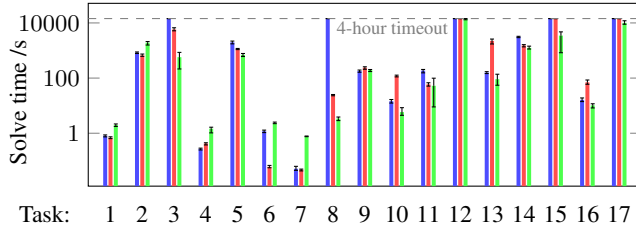


Figure 17. Comparing MiniSat (blue), Glucose (red) and Plingeling (green)

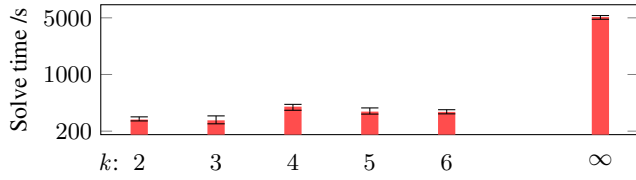


Figure 18. Performance of Task 13 with fixpoints unrolled k times

Table 2 summarises the tasks on which we have evaluated our technique. We used a machine with four 16-core 2.1 GHz AMD Opteron processors and 128 GB of RAM.

7.1 Choice of SAT Solver

Figure 17 summarises the performance of three SAT solvers on our tasks: Glucose (version 2.1) [13], MiniSat (version 2.2) [26], and Plingeling [17]. Each bar shows the mean solve time over 4 runs, plus the minimum and maximum time.⁷ Plingeling is able to complete all tasks, whereas Glucose and MiniSat time out on three and five of them respectively. On the tasks all three solvers complete, MiniSat takes an average of 9 minutes, Glucose takes 8, and Plingeling takes 6. Plingeling’s relatively high startup overhead is evident on the quicker tasks.

7.2 Combined vs. Separate Event Sets

For **Q4** tasks (§6), we can choose to draw software-level and architecture-level events either from a single set \mathbb{E} or from disjoint partitions \mathbb{E}_{sw} and \mathbb{E}_{arch} . The former approach is attractive because it gives the user one fewer parameter to control, and it minimises the total number of events required to find solutions because a single element of \mathbb{E} can represent both a software-level event and an architecture-level event. However, as shown in Tab. 2, we choose the latter approach. To see why, consider Task 17. To validate the OpenCL/PTX mapping for OpenCL executions of up to 5 events, we must consider PTX executions of up to 15 events (since each OpenCL event can map to up to three PTX events). We found that setting $|\mathbb{E}_{\text{sw}}| = 5$ and $|\mathbb{E}_{\text{arch}}| = 15$ led to a tractable constraint-solving problem for Alloy, but that merely setting $|\mathbb{E}| = 15$, which involves fewer events in total but includes OpenCL executions with more than 5 events, rendered the problem insoluble.

7.3 Recursive Definitions vs. Fixed Unrolling

Alglave et al.’s formalisation of the Power MCM in the `.cat` format involves several recursively-defined relations, but Alloy does not support definitions of the form ‘`let rec r = f(r)`’. Accordingly, in our formalisation of the Power MCM [1], we expand the recursive construction explicitly, by requiring the existence of an r satisfying $f(r) \subseteq r$ and $\forall r'. f(r') \subseteq r' \Rightarrow r \subseteq r'$. The lat-

⁷A more thorough comparison of SAT solvers would control for the order of clauses, which greatly influences performance [59].

ter constraint involves universal quantification over relations and hence requires the higher-order AlloyStar solver [52]. A first-order alternative is simply to unroll the recursive definitions a few times; that is, to set $r \stackrel{\text{def}}{=} f^k(\emptyset)$ for a fixed k . We found, for the small search scopes involved in our work, that $k \geq 2$ is sufficient for avoiding false positives when checking a compiler mapping from C11 (Task 13). Figure 18 shows that the proper fixpoint construction (i.e., $k = \infty$) is much more expensive than a fixed unrolling.

8. Related Work

Existing tools for MCM reasoning typically take a concurrent program as input, and produce all of the executions allowed under a given MCM. Some rely on SAT solvers to discover executions [18, 25, 75, 83], while others enumerate them explicitly [10, 12, 16, 60, 67, 68]. Our work tackles, in a sense, the ‘inverse’ problem: we start with executions that witness interesting MCM properties, and go on to produce programs that can give rise to these executions.

Other works tackling Q1–Q4 Regarding **Q1**, Darbari et al. [24] have automatically generated conformance tests for HSA [35] from an Event-B specification [2] of its MCM. Alglave et al.’s `diy` tool [11] generates conformance tests for a range of MCMs based on *critical cycles* [70]. We find that our Alloy-based approach has several advantages over `diy` when generating conformance tests. First, we can straightforwardly handle custom MCM constructs (e.g., C11 memory orders) while `diy` currently does not. Second, we generate only tests needed for conformance, while `diy` generates many more. Third, `diy` can also miss some tests required to distinguish two MCMs (such as the single-address tests we saw in §6.2), if not guided carefully through user-supplied critical cycles.

Regarding **Q2**, there is a long history of unifying frameworks for comparing MCMs [3, 23, 28, 31, 32, 70], and of manual proofs that different formulations of the same MCM are equivalent [50, 63]. On the automation side, Mador-Haim et al. [48] have, like ourselves, used a SAT solver as part of a tool for generating litmus tests that distinguish MCMs. However, where we use the solver to *generate* tests, they just use it to *check* the behaviour of pre-generated tests. Given that generating all 6-instruction litmus tests takes them ‘a few minutes’ (in 2010), we expect that their approach of explicit enumeration would not scale to find the 12-instruction litmus tests that are sometimes necessary to distinguish MCMs (see §4.2), and which Alloy is able to find in just 4 minutes.

Prior work has proved (manually) the validity of compiler optimisations in non-SC MCMs [79, 80]. Since optimisations should not introduce new behaviours, this problem is related to monotonicity (**Q3**). On the automation side, the Traver tool [21] uses an automated theorem prover to verify/falsify a given program transformation against a non-SC MCM. Unlike our work, Traver does not support multi-threaded optimisations such as linearisation. Chakraborty et al. have built a tool [22] that automatically verifies that LLVM optimisations preserve C11 semantics. Like Traver, their tool only checks specific *instances* of an optimisation, while our approach is able to check optimisations themselves. Morisset et al. [55] use random testing to validate optimisations (albeit those not involving atomic operations) against the C11 MCM; Vafeiadis et al. [77] then show (manually) that several of these optimisations are invalid when atomic operations are involved. Our work, in turn, shows how several of Vafeiadis et al.’s results can be recreated automatically, often in a simpler form (§5).

Regarding **Q4**, prior work has proved (manually) the correctness of both compiler mappings [15, 16, 82] and full compilers [81] in a non-SC context. These proofs all involve intensive manual effort, in contrast to our lightweight automatic checking, though our checking can of course only guarantee the absence of bugs within Alloy’s search scope. On the automation side, Lustig et

al. have built tools for finding and verifying semantics-preserving translations, but where we focus on language/architecture translation, they focus on architecture/microarchitecture [46] and architecture/architecture [47] translation. Very recent work by Tripel et al. [76] has produced a framework that can check language/architecture mappings; this works by enumerating standard litmus tests and then simulating each one against both the language-level MCM and, after compilation, the architecture-level MCM.

Reflections on Alloy Alloy is a mature, open-source, and widely-used modelling application, and its trio of features – a modular and object-oriented modelling language, an automatic constraint solver, and a graphical visualiser – makes it ideal for MCM development. Although Tab. 2 shows several lengthy solving times, those figures are obtained once the search space has been set as large as computational feasibility allows. Given smaller search spaces, as would be appropriate during MCM prototyping, Alloy is suitable for interactive use.

When does Alloy’s failure to distinguish two MCMs become a proof that they are equivalent? Mador-Haim et al. [49] prove that 6 events are enough, but their result applies only to multi-copy atomic, architecture-level MCMs (see §4.4). Momtahan [54] gives a result for general Alloy models, but imposes strong restrictions on quantifiers that our models do not meet.

Ivy [64] defines a relational modelling language similar to Alloy’s. Where Alloy ensures the decidability of instance-finding by restricting to a finite search scope (which limits its usefulness for verification), Ivy instead restricts formulas to the form $\exists \bar{x}. \forall \bar{y}. \phi$ (for quantifier-free ϕ). If our models can be rephrased to fit into Ivy’s restricted language, there is the potential not just to ‘debug’ MCM properties, but to *verify* them. Another language related to Alloy, Rosette [74], is used in very recent work by Bornholt et al. [20] to solve the problem of synthesising a MCM from a set of desired litmus test outcomes.

We find that Alloy has several advantages over other frameworks that have been used to reason about MCMs, such as Isabelle (e.g., [16]), Lem [56] (e.g., [15]), Coq (e.g., [77]), and .cat [12] (e.g., [14]). A key advantage is that the entire memory modelling process can be conducted within Alloy: the Alloy modelling language can express programming languages, compiler mappings, MCMs, and properties to test, the Alloy Analyzer can discover solutions, and the Alloy Visualizer can display solutions using theming customised for the model at hand. Alglave et al.’s .cat framework, like Alloy, allows MCM axioms to be expressed in the concise propositional relation calculus [73], but Alloy also supports the more powerful predicate calculus as a fallback. As such, Alloy is expressive enough to capture both axiomatic and operational MCMs, while remaining sufficiently restrictive that fully-automatic property checking is computationally tractable.

9. Conclusion

By solving relational constraints between executions and then lifting solutions to litmus tests, Alloy can generate conformance tests, compare MCMs (against each other and against general properties like SC-DRF and multi-copy atomicity), and check monotonicity and compiler mappings. As such, we believe that Alloy should become an essential component of the memory modeller’s toolbox. Indeed, we are already working with two large processor vendors to apply our technique to their recent and upcoming architectures and languages. Other future work includes applying our technique to more recent MCMs that are defined in a non-axiomatic style [27, 40, 41, 65].

Although Alloy’s lightweight, automatic approach cannot give the same universal assurance as fully mechanised theorems, we have found it invaluable in practice, because even (and perhaps

especially) in the complex and counterintuitive world of non-SC MCMs, Jackson’s maxim [39] holds true: *Most bugs have small counterexamples*.

Acknowledgements

We thank Alastair Donaldson, Carsten Fuhs, Daniel Jackson, Luc Maranget, and Matthew Parkinson for their valuable feedback. We acknowledge the support of a Research Fellowship from the Royal Academy of Engineering and the Lloyd’s Register Foundation (Batty), a Research Chair from the Royal Academy of Engineering and Imagination Technologies (Constantinides), an EPSRC Impact Acceleration Award, EPSRC grants EP/I020357/1, EP/K034448/1, and EP/K015168/1, and a gift from Intel Corporation.

References

- [1] Supplementary material for this paper is available in the ACM digital library, and in the following GitHub repository. URL: <https://johnwickerson.github.io/memalloy>.
- [2] Jean-Raymond Abrial, Michael Butler, Stefan Hallerstede, Thai Son Hoang, Farhad Mehta, and Laurent Voisin. “Rodin: an open toolset for modelling and reasoning in Event-B”. In: *Int. J. Softw. Tools Technol. Transfer* 12 (2010).
- [3] Sarita V. Adve. “Designing Memory Consistency Models For Shared-Memory Multiprocessors”. PhD thesis. University of Wisconsin-Madison, 1993.
- [4] Sarita V. Adve and Kourosh Gharachorloo. “Shared memory consistency models: A tutorial”. In: *IEEE Computer* 29.12 (1996).
- [5] Sarita V. Adve and Mark D. Hill. “Weak Ordering - A New Definition”. In: *Int. Symp. on Computer Architecture (ISCA)*. 1990.
- [6] Mustaque Ahamad, Rida A. Bazzi, Ranjit John, Prince Kohli, and Gil Neiger. “The Power of Processor Consistency”. In: *ACM Symp. on Parallelism in Algorithms and Architectures (SPAA)*. 1993.
- [7] Alfred V. Aho, Monica S. Lam, Ravi Sethi, and Jeffrey D. Ullman. *Compilers: Principles, Techniques, & Tools*. Second edition. Addison-Wesley, 2006.
- [8] Jade Alglave. “A formal hierarchy of weak memory models”. In: *Formal Methods in System Design* 41 (2012).
- [9] Jade Alglave, Mark Batty, Alastair F. Donaldson, Ganesh Gopalakrishnan, Jeroen Ketema, Daniel Poetzl, Tyler Sorensen, and John Wickerson. “GPU concurrency: weak behaviours and programming assumptions”. In: *Int. Conf. on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*. 2015.
- [10] Jade Alglave, Daniel Kroening, and Michael Tautschnig. “Partial Orders for Efficient Bounded Model Checking of Concurrent Software”. In: *Computer Aided Verification (CAV)*. 2013.
- [11] Jade Alglave, Luc Maranget, Susmit Sarkar, and Peter Sewell. “Fences in Weak Memory Models”. In: *Computer Aided Verification (CAV)*. 2010.
- [12] Jade Alglave, Luc Maranget, and Michael Tautschnig. “Herding cats: modelling, simulation, testing, and data-mining for weak memory”. In: *ACM Trans. on Programming Languages and Systems (TOPLAS)* 36.2 (2014).
- [13] Gilles Audemard and Laurent Simon. “Predicting Learnt Clauses Quality in Modern SAT Solvers”. In: *Int. Joint Conf. on Artificial Intelligence (IJCAI)*. 2009.
- [14] Mark Batty, Alastair F. Donaldson, and John Wickerson. “Overhauling SC atomics in C11 and OpenCL”. In: *ACM*

- Symp. on Principles of Programming Languages (POPL)*. 2016.
- [15] Mark Batty, Kayvan Memarian, Scott Owens, Susmit Sarkar, and Peter Sewell. “Clarifying and Compiling C/C++ Concurrency: from C++11 to POWER”. In: *ACM Symp. on Principles of Programming Languages (POPL)*. 2012.
- [16] Mark Batty, Scott Owens, Susmit Sarkar, Peter Sewell, and Tjark Weber. “Mathematizing C++ Concurrency”. In: *ACM Symp. on Principles of Programming Languages (POPL)*. 2011.
- [17] Armin Biere. *Lingeling, Plingeling, PicoSAT and PrecoSAT at SAT Race 2010*. Tech. rep. 10/1. Institute for Formal Models and Verification, Johannes Kepler University, 2010.
- [18] Jasmin Christian Blanchette, Tjark Weber, Mark Batty, Scott Owens, and Susmit Sarkar. “Nitpicking C++ Concurrency”. In: *Int. Symp. on Principles and Practice of Declarative Programming (PPDP)*. 2011.
- [19] Hans-J. Boehm and Sarita V. Adve. “Foundations of the C++ Concurrency Memory Model”. In: *ACM Conf. on Programming Language Design and Implementation (PLDI)*. 2008.
- [20] James Bornholt and Emina Torlak. *Synthesizing Memory Models from Litmus Tests*. Tech. rep. UW-CSE-16-10-01. University of Washington, 2016.
- [21] Sebastian Burckhardt, Madanlal Musuvathi, and Vasu Singh. “Verifying Local Transformations on Relaxed Memory Models”. In: *Int. Conf. on Compiler Construction (CC)*. 2010.
- [22] Soham Chakraborty and Viktor Vafeiadis. “Validating Optimizations of Concurrent C/C++ Programs”. In: *Int. Symp. on Code Generation and Optimization (CGO)*. 2016.
- [23] William W. Collier. *Reasoning about Parallel Architectures*. Prentice Hall, 1992.
- [24] Ashish Darbari, Iain Singleton, Michael Butler, and John Colley. “Formal Modelling, Testing and Verification of HSA Memory Models using Event-B”. Draft. 2016. URL: <http://arxiv.org/pdf/1605.04744v1.pdf>.
- [25] Brian Densky and Patrick Lam. “SATCheck: SAT-Directed Stateless Model Checking for SC and TSO”. In: *ACM Int. Conf. on Object-Oriented Programming, Systems, Languages, and Applications (OOPSLA)*. 2015.
- [26] Niklas Eén and Niklas Sörensson. “An Extensible SAT-solver”. In: *Theory and Applications of Satisfiability Testing (SAT)*. 2003.
- [27] Shaked Flur, Kathryn E. Gray, Christopher Pulte, Susmit Sarkar, Ali Sezgin, Luc Maranget, Will Deacon, and Peter Sewell. “Modelling the ARMv8 Architecture, Operationally: Concurrency and ISA”. In: *ACM Symp. on Principles of Programming Languages (POPL)*. 2016.
- [28] Kouros Gharachorloo. “Memory Consistency Models for Shared-Memory Multiprocessors”. PhD thesis. Stanford University, 1995.
- [29] Blake A. Hechtman, Shuai Che, Derek R. Hower, Yingying Tian, Bradford M. Beckmann, Mark D. Hill, Steven K. Reinhardt, and David A. Wood. “QuickRelease: A Throughput-oriented Approach to Release Consistency on GPUs”. In: *IEEE Int. Symp. on High Performance Computer Architecture (HPCA)*. 2014.
- [30] John L. Hennessy and David A. Patterson. *Computer Architecture: A Quantitative Approach*. Fifth edition. Morgan Kaufmann, 2012.
- [31] Lisa Higham, LillAnne Jackson, and Jalal Kawash. “Specifying Memory Consistency of Write Buffer Multiprocessors”. In: *ACM Trans. on Programming Languages and Systems (TOPLAS)* 25.1 (2007).
- [32] Lisa Higham, Jalal Kawash, and Nathaly Verwaal. “Defining and Comparing Memory Consistency Models”. In: *Int. Conf. on Parallel and Distributed Computing Systems (PDCS)*. 1997.
- [33] C. A. R. Hoare. “Proof of Correctness of Data Representations”. In: *Acta Informatica* 1 (1972).
- [34] Derek R. Hower, Blake A. Hechtman, Bradford M. Beckmann, Benedict R. Gaster, Mark D. Hill, Steven K. Reinhardt, and David A. Wood. “Heterogeneous-race-free Memory Models”. In: *Int. Conf. on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*. 2014.
- [35] HSA Foundation. *HSA Platform System Architecture Specification*. Version 1.0, 2015. URL: <http://www.hsafoundation.com/standards/>.
- [36] IBM. *Power ISA (Version 2.06B)*. 2010.
- [37] ISO/IEC. *Programming languages – C++*. Draft N3092, 2010. URL: <http://www.open-std.org/jtc1/sc22/wg21/docs/papers/2010/n3092.pdf>.
- [38] ISO/IEC. *Programming languages – C*. International standard 9899:2011, 2011.
- [39] Daniel Jackson. *Software Abstractions – Logic, Language, and Analysis*. Revised edition. MIT Press, 2012.
- [40] Alan Jeffrey and James Riely. “On Thin Air Reads: Towards an Event Structures Model of Relaxed Memory”. In: *ACM/IEEE Symp. on Logic in Computer Science (LICS)*. 2016.
- [41] Jeehoon Kang, Chung-Kil Hur, Ori Lahav, Viktor Vafeiadis, and Derek Dreyer. “A Promising Semantics for Relaxed-Memory Concurrency”. In: *ACM Symp. on Principles of Programming Languages (POPL)*. 2017.
- [42] Khronos Group. *The OpenCL Specification*. Version 2.0, 2013. URL: <https://www.khronos.org/registry/cl/>.
- [43] Ori Lahav, Nick Giannarakis, and Viktor Vafeiadis. “Taming Release-Acquire Consistency”. In: *ACM Symp. on Principles of Programming Languages (POPL)*. 2016.
- [44] Ori Lahav, Viktor Vafeiadis, Jeehoon Kang, Chung-Kil Hur, and Derek Dreyer. *Repairing Sequential Consistency in C/C++11*. Tech. rep. MPI-SWS-2016-011. MPI-SWS, 2016.
- [45] Leslie Lamport. “How to Make a Multiprocessor Computer That Correctly Executes Multiprocess Programs”. In: *IEEE Transactions on Computers* C-28.9 (1979).
- [46] Daniel Lustig, Michael Pellauer, and Margaret Martonosi. “PipeCheck: Specifying and Verifying Microarchitectural Enforcement of Memory Consistency Models”. In: *Int. Symp. on Microarchitecture (MICRO)*. 2014.
- [47] Daniel Lustig, Caroline Trippel, Michael Pellauer, and Margaret Martonosi. “ArMOR: Defending Against Memory Consistency Model Mismatches in Heterogeneous Architectures”. In: *Int. Symp. on Computer Architecture (ISCA)*. 2015.
- [48] Sela Mador-Haim, Rajeev Alur, and Milo M. K. Martin. “Generating Litmus Tests for Contrasting Memory Consistency Models”. In: *Computer Aided Verification (CAV)*. 2010.
- [49] Sela Mador-Haim, Rajeev Alur, and Milo M. K. Martin. “Litmus Tests for Comparing Memory Consistency Models: How Long Do They Need to Be?” In: *Design Automation Conference (DAC)*. 2011.
- [50] Sela Mador-Haim, Luc Maranget, Susmit Sarkar, Kayvan Memarian, Jade Alglave, Scott Owens, Rajeev Alur, Milo M. K. Martin, Peter Sewell, and Derek Williams. “An Axiomatic Memory Model for POWER Multiprocessors”. In: *Computer Aided Verification (CAV)*. 2012.

- [51] Yatin A. Manerkar, Caroline Trippel, Daniel Lustig, Michael Pellauer, and Margaret Martonosi. “Counterexamples and Proof Loophole for the C/C++ to POWER and ARMv7 Trailing-Sync Compiler Mappings”. 2016. URL: <http://arxiv.org/pdf/1611.01507v1.pdf>.
- [52] Aleksandar Milicevic, Joseph P. Near, Eunsuk Kang, and Daniel Jackson. “Alloy*: A General-Purpose Higher-Order Relational Constraint Solver”. In: *Int. Conf. on Software Engineering (ICSE)*. 2015.
- [53] Rolf H. Möhring. “Computationally tractable classes of ordered sets”. In: *Algorithms and Order*. Ed. by Ivan Rival. Springer, 1989.
- [54] Lee Momtahan. “Towards a Small Model Theorem for Data Independent Systems in Alloy”. In: *Electronic Notes in Theoretical Computer Science* 128 (2005).
- [55] Robin Morisset, Pankaj Pawan, and Francesco Zappa Nardelli. “Compiler Testing via a Theory of Sound Optimisations in the C11/C++11 Memory Model”. In: *ACM Conf. on Programming Language Design and Implementation (PLDI)*. 2013.
- [56] Dominic P. Mulligan, Scott Owens, Kathryn E. Gray, Tom Ridge, and Peter Sewell. “Lem: reusable engineering of real-world semantics”. In: *ACM Int. Conf. on Functional Programming (ICFP)*. 2014.
- [57] Roger M. Needham. “Names”. In: *Distributed Systems*. Ed. by Sape Mullender. ACM Press, 1989.
- [58] Kyndylan Nienhuis, Kayvan Memarian, and Peter Sewell. “An Operational Semantics for C/C++11 Concurrency”. In: *ACM Int. Conf. on Object-Oriented Programming, Systems, Languages, and Applications (OOPSLA)*. 2016.
- [59] Mladen Nikolić. “Statistical Methodology for Comparison of SAT Solvers”. In: *Theory and Applications of Satisfiability Testing (SAT)*. 2010.
- [60] Brian Norris and Brian Densky. “CDSChecker: Checking Concurrent Data Structures Written with C/C++ Atomics”. In: *ACM Int. Conf. on Object-Oriented Programming, Systems, Languages, and Applications (OOPSLA)*. 2013.
- [61] NVIDIA. *Parallel Thread Execution ISA, version 4.3*. 2015. URL: http://docs.nvidia.com/cuda/pdf/ptx_isa_4.3.pdf.
- [62] Marc S. Orr, Shuai Che, Ayse Yilmazer, Bradford M. Beckmann, Mark D. Hill, and David A. Wood. “Synchronization Using Remote-Scope Promotion”. In: *Int. Conf. on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*. 2015.
- [63] Scott Owens, Susmit Sarkar, and Peter Sewell. “A Better x86 Memory Model: x86-TSO”. In: *Theorem Proving in Higher Order Logics (TPHOLs)*. 2009.
- [64] Oded Padon, Kenneth L. McMillan, Aurojit Panda, Mooly Sagiv, and Sharon Shoham. “Ivy: Safety Verification by Interactive Generalization”. In: *ACM Conf. on Programming Language Design and Implementation (PLDI)*. 2016.
- [65] Jean Pichon-Pharabod and Peter Sewell. “A Concurrency Semantics for Relaxed Atomics that Permits Optimisation and Avoids Thin-Air Executions”. In: *ACM Symp. on Principles of Programming Languages (POPL)*. 2016.
- [66] Susmit Sarkar, Kayvan Memarian, Scott Owens, Mark Batty, Peter Sewell, Luc Maranget, Jade Alglave, and Derek Williams. “Synchronising C/C++ and POWER”. In: *ACM Conf. on Programming Language Design and Implementation (PLDI)*. 2012.
- [67] Susmit Sarkar, Peter Sewell, Jade Alglave, Luc Maranget, and Derek Williams. “Understanding POWER Multiprocessors”. In: *ACM Conf. on Programming Language Design and Implementation (PLDI)*. 2011.
- [68] Susmit Sarkar, Peter Sewell, Francesco Zappa Nardelli, Scott Owens, Tom Ridge, Thomas Braibant, Magnus O. Myreen, and Jade Alglave. “The Semantics of x86-CC Multiprocessor Machine Code”. In: *ACM Symp. on Principles of Programming Languages (POPL)*. 2009.
- [69] Ali Sezgin. “Formalization and Verification of Shared Memory”. PhD thesis. University of Utah, 2004.
- [70] Dennis Shasha and Marc Snir. “Efficient and Correct Execution of Parallel Programs that Share Memory”. In: *ACM Trans. on Programming Languages and Systems (TOPLAS)* 10.2 (1988).
- [71] Tyler Sorensen and Alastair F. Donaldson. “Exposing Errors Related to Weak Memory in GPU Applications”. In: *ACM Conf. on Programming Language Design and Implementation (PLDI)*. 2016.
- [72] Daniel J. Sorin, Mark D. Hill, and David A. Wood. *A Primer on Memory Consistency and Cache Coherence*. Ed. by Mark D. Hill. Vol. 16. Synthesis Lectures on Computer Architecture. Morgan & Claypool, 2011.
- [73] Alfred Tarski. “On the Calculus of Relations”. In: *Journal of Symbolic Logic* 6.3 (1941), pp. 73–89.
- [74] Emina Torlak and Rastislav Bodik. “Growing Solver-Aided Languages with Rosette”. In: *Onward!* 2013.
- [75] Emina Torlak, Mandana Vaziri, and Julian Dolby. “MemSAT: Checking Axiomatic Specifications of Memory Models”. In: *ACM Conf. on Programming Language Design and Implementation (PLDI)*. 2010.
- [76] Caroline Trippel, Yatin A. Manerkar, Daniel Lustig, Michael Pellauer, and Margaret Martonosi. “Exploring the Trisection of Software, Hardware, and ISA in Memory Model Design”. Draft. 2016. URL: <http://arxiv.org/pdf/1608.07547v1.pdf>.
- [77] Viktor Vafeiadis, Thibaut Balabonski, Soham Chakraborty, Robin Morisset, and Francesco Zappa Nardelli. “Common compiler optimisations are invalid in the C11 memory model and what we can do about it”. In: *ACM Symp. on Principles of Programming Languages (POPL)*. 2015.
- [78] Jacobo Valdes, Robert E. Tarjan, and Eugene L. Lawler. “The recognition of Series Parallel digraphs”. In: *ACM Symp. on Theory of Computing (STOC)*. 1979.
- [79] Jaroslav Ševčík. “Safe Optimisations for Shared-Memory Concurrent Programs”. In: *ACM Conf. on Programming Language Design and Implementation (PLDI)*. 2011.
- [80] Jaroslav Ševčík and David Aspinall. “On Validity of Program Transformations in the Java Memory Model”. In: *Europ. Conf. on Object-Oriented Programming (ECOOP)*. 2008.
- [81] Jaroslav Ševčík, Viktor Vafeiadis, Francesco Zappa Nardelli, Suresh Jagannathan, and Peter Sewell. “Relaxed-Memory Concurrency and Verified Compilation”. In: *ACM Symp. on Principles of Programming Languages (POPL)*. 2011.
- [82] John Wickerson, Mark Batty, Alastair F. Donaldson, and Bradford M. Beckmann. “Remote-scope promotion: clarified, rectified, and verified”. In: *ACM Int. Conf. on Object-Oriented Programming, Systems, Languages, and Applications (OOPSLA)*. 2015.
- [83] Yue Yang, Ganesh Gopalakrishnan, Gary Lindstrom, and Konrad Slind. “Nemos: A Framework for Axiomatic and Executable Specifications of Memory Consistency Models”. In: *Int. Parallel and Distributed Processing Symp. (IPDPS)*. 2004.