

A model of speculative evaluation

CRAIG DISSELKOEN, University of California San Diego and Mozilla Research Internship

RADHA JAGADEESAN, DePaul University

ALAN JEFFREY, Mozilla Research

JAMES RIELY, DePaul University

This paper studies information flow caused by speculation mechanisms in hardware and software. The Spectre attack shows that there are practical information flow attacks which use an interaction of dynamic security checks, speculative evaluation and cache timing. Previous formal models of program execution have not been designed to model speculative evaluation, and so do not capture attacks such as Spectre. In this paper, we propose a model based on pomsets which is designed to model speculative evaluation. The model provides a compositional semantics for a simple shared-memory concurrent language, which captures features such as data and control dependencies, relaxed memory and transactions. We provide models for existing information flow attacks based on speculative evaluation and transactions, and new information flow attacks on compiler optimizations. The new attacks are experimentally validated against gcc and clang. A simple temporal logic provides reasoning principals, including composition and proofs using invariants.

1 INTRODUCTION

This paper studies information flow caused by speculation mechanisms in hardware and software.

Information flow from high (hi) to low (lo) provides a formal foundation for end-to-end security. Informally, a program is secure if there is no observable dependency of lo-observables on hi-inputs. The precise formalization of this intuitive idea has been the topic of extensive research (e.g., see [?] for a detailed survey till 2006), and can be generalized to account for a variety of language features and observables such as non-determinism [?], concurrency [?], reactivity [?], and probability [?]. The static and dynamic enforcement of these definitions in general purpose languages [?] has also been studied extensively and has influenced language design and implementation.

A key parameter in the definitions cited above is the notion of the *observational power* of the attacker model. Whereas the classical input-output behavior is often an adequate foundation, it has long been known [??] that side-channels that leak information arise from other observables such as execution time and power consumption.

The focus of this paper is the development of a compositional semantic model of executions of programs to explicate side channel attacks that arise from speculation. In particular, we model conditionals such that it is possible for the behaviour of a conditional ($\text{if } (M) \{ C \} \text{ else } \{ D \}$) to depend on the behaviour of D even when the condition M is true.

Our study addresses several sources of speculation in the concurrent execution of programs.

- Pipelined microprocessors use heuristics to predict the path of a program, for example by guessing memory dependencies, or the outcome of conditionals. Execution proceeds along

Authors' addresses: Craig Disselkoen, University of California San Diego, Mozilla Research Internship, cdisselk@cs.ucsd.edu; Radha Jagadeesan, DePaul University, rjagadeesan@cs.depaul.edu; Alan Jeffrey, Mozilla Research, ajeffrey@mozilla.com; James Riely, DePaul University, jriely@cs.depaul.edu.

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

© 2018 Copyright held by the owner/author(s).

XXXX-XXXX/2018/7-ART

<https://doi.org/10.1145/nnnnnnnn.nnnnnnn>

the predicted path until it is possible to validate the prediction; the execution is committed if the prediction is correct; otherwise it is rolled back. The speculation does not affect the observable input-output behavior of the program, thus ensuring correctness with respect to the usual intended semantics of the program. However, since there are clear timing differences between the cases where the prediction is and is not successful, it raises the concern that the predictive mechanisms themselves could cause information leaks via timing side-channels; a fear that is realized with devastating impact by the Spectre family of attacks [?].

- Several modern microprocessors [?] support transactions to aid in the design of correct and efficient concurrent programs. Transactions are executed optimistically; they are committed if there are no conflicts, and aborted otherwise. All memory effects of an aborted transaction are rolled back; so there is no way for a concurrent observer to detect an aborted transaction. However, the thread of the aborted transaction, via abort-handler code, gets notified of the cancellation of the transaction.

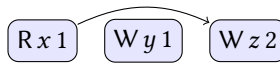
Thus, when a transaction of a thread aborts, the thread can deduce plausible information about specific memory accesses of a concurrently executing transaction. In combination with the techniques used in Spectre, this potential has been exploited recently to accelerate and scale up the efficacy of the attacker in the Spectre family of attacks [?].

- Compiler optimizations may depend on both branches of a conditional, for example one that replaces $(\text{if } (M) \{ C \} \text{ else } \{ C \})$ by C . In particular, $(\text{if } (r) \{ x := 1 \} \text{ else } \{ x := 1 \})$ can be optimized to $(x := 1)$, which does not have a control dependency from r to the assignment to x . In contrast, the program $(\text{if } (r) \{ x := 1 \} \text{ else } \{ x := 2 \})$ cannot be optimized, so the control dependency cannot be removed.

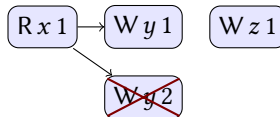
This would be fine if there were no program constructs which could observe control dependencies, but unfortunately the relaxed memory models for Java [?], C++ [?] and LLVM [?] allow for such observations. This means there is the possibility for information flows caused by optimizing compilers, which we investigate in this paper.

Whereas information flows caused by speculation in hardware, or by transactional memory are known [? ?], these attacks on compiler optimizations, and the relaxed memory models that justify them, are new. In this paper we provide both a formal model for such attacks, and experimental evidence about their practicality.

Information flow attacks motivate the main technical development of this paper. Our model is based on *partially ordered multisets* [? ?] (“pomsets”), whose labels are given by read and write actions. These can be visualized as a graph where the edges indicate dependencies, for example $(r := x; y := 1; z := r + 1)$ has an execution modelled by the pomset:



The edge from $(R x 1)$ to $(W z 2)$ indicates a data dependency. The novel aspect of the model is that events have *preconditions* which may be false. These are used in giving the semantics of conditionals, for example $(\text{if } (x) \{ y := 1; z := 1 \} \text{ else } \{ y := 2; z := 1 \})$ has an execution:



The edges from $(R x 1)$ to $(W y 1)$ and $(W y 2)$ indicate control dependencies. The presence of a crossed out $(W y 2)$ indicates an event with an unsatisfiable precondition,

The novel contributions of this paper are:

- a model of program execution that includes speculation (§2),
- examples showing how the model can be applied, including information flow attacks on hardware, optimizing compilers, and transactional memory (§3), and
- experimental evidence about how practical it is to mount the new class of attacks (§4).

2 MODEL

The model used in this paper is one of sets of pomsets with event labels of the form $(\phi \mid a)$, where ϕ is the event's precondition (such as $M = v$) and a is the event's action (such as $\mathsf{W} x v$). For example the semantics of the program $(x := M)$ includes the case where M is v , which is written to x , and is captured by the one-event pomset:

$$M = v \mid \mathsf{W} x v$$

We make few requirements of the logic of preconditions, save that it has includes equalities between expressions, is closed under substitution, and supports a notion of implication.

For example, the set of pomsets $\llbracket r := y; x := r + 1 \rrbracket$ contains:

$$\mathsf{R} y 1 \rightarrow \mathsf{W} x 2$$

The semantics is defined compositionally. First, $\llbracket x := r + 1 \rrbracket$ contains the pomset:

$$r = 1 \mid \mathsf{W} x 2$$

Next, we perform the substitution of r with 1 in every precondition, to get that $\llbracket x := r + 1 \rrbracket[1/r]$ contains the pomset:

$$1 = 1 \mid \mathsf{W} x 2$$

and since $(1 = 1)$ is a tautology, we elide it:

$$\mathsf{W} x 2$$

This substitution is performed in defining $\llbracket r := y; x := r + 1 \rrbracket$, which contains the pomset:

$$\mathsf{R} y 1 \rightarrow \mathsf{W} x 2$$

as required. There is an ordering $(\mathsf{R} y 1) < (\mathsf{W} x 2)$ because the precondition $(r = 1)$ depends on r . If the precondition was independent of r then there would be no ordering, for example $\llbracket r := y; x := r + 1 - r \rrbracket$ contains the pomset:

$$\mathsf{R} y 1 \quad \mathsf{W} x 1$$

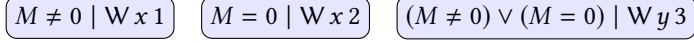
since the precondition $(r + 1 - r = 1)$ is independent of r .

The main novelty of our semantics, since it is designed to model speculative evaluation, is in modelling conditionals. In most sets-of-pomsets semantics, a pomset in $\llbracket \text{if}(M) \{ C \} \text{ else } \{ D \} \rrbracket$ would either be given by a pomset in $\llbracket C \rrbracket$ or a pomset in $\llbracket D \rrbracket$. To model speculative evaluation, we need to allow a pomset in $\llbracket \text{if}(M) \{ C \} \text{ else } \{ D \} \rrbracket$ to be given by both a pomset in $\llbracket C \rrbracket$ and a pomset in $\llbracket D \rrbracket$. For example, $\llbracket \text{if}(M) \{ x := 1 \} \text{ else } \{ x := 2 \} \rrbracket$ contains the pomset:

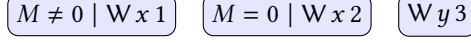
$$M \neq 0 \mid \mathsf{W} x 1 \quad M = 0 \mid \mathsf{W} x 2$$

that is we have recorded behaviour from both branches of execution. Moreover, an action which is performed on both sides of the conditional can be merged, producing only one event in the resulting

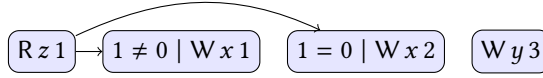
pomset. The precondition of the merged event is the disjunction of the preconditions of the original events. For example $\llbracket \text{if } (M) \{ x := 1; y := 3 \} \text{ else } \{ x := 2; y := 3 \} \rrbracket$ contains the pomset:



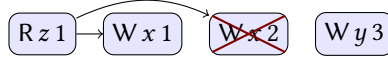
and since $(M \neq 0) \vee (M = 0)$ is a tautology, this is:



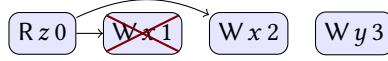
Combining this model of conditionals with the model of memory using substitutions gives that $\llbracket \text{if } (z) \{ x := 1; y := 3 \} \text{ else } \{ x := 2; y := 3 \} \rrbracket$ contains the pomset:



and since $(1 \neq 0)$ is a tautology and $(1 = 0)$ is unsatisfiable, this is:



Similarly, $\llbracket \text{if } (z) \{ x := 1; y := 3 \} \text{ else } \{ x := 2; y := 3 \} \rrbracket$ contains the pomset:



Note that this semantics captures control dependencies such as $(R z 0) < (W x 1)$, independencies such as $(R z 0) \not< (W y 3)$, and failed speculations such as the crossed out $(W x 1)$.

In summary, the features we need of the underlying data model are:

- *actions*, which may read or write to memory locations, and
- *preconditions*, which form a logic closed under substitution.

from which we can define the operations used in defining the semantics of programs, which include:

- *prefixing* $a \rightarrow \mathcal{P}$, which adds an event with precondition ϕ and action a to pomsets in \mathcal{P} ,
- *substitution* $\mathcal{P}[M/x]$, which performs a substitution on every precondition in \mathcal{P} , and
- *concurrency* $\mathcal{P}_1 \parallel \mathcal{P}_2$, which unions pomsets from \mathcal{P}_1 and \mathcal{P}_2 , allowing events to be merged.

We make data models precise in §2.1, define pomsets in §2.2, and operations on sets of pomsets in §2.3, which are used to give a compositional semantics for a simple imperative language.

2.1 Data models

A *data model* consists of:

- a set of *memory locations* \mathcal{X} , ranged over by x and y ,
- a set of *registers* \mathcal{R} , ranged over by r and s ,
- a set of *values* \mathcal{V} , ranged over by v and w ,
- a set of *expressions* \mathcal{E} , ranged over by M and N ,
- a set of *logical formulae* Φ , ranged over by ϕ and ψ , and
- a set of *actions* \mathcal{A} , ranged over by a and b ,

such that:

- values include at least the constants 0 and 1,
- expressions include at least registers and values,
- expressions are closed under substitutions of the form $M[N/r]$,
- formulae include at least true, false, and equalities of the form $(M = N)$ and $(x = N)$,

- formulae are closed under negation, conjunction, disjunction,
- formulae are closed under substitutions of the form $\phi[x/r]$ or $\phi[N/x]$,
- there is a relation \models between formulae, and
- there are partial functions R and $W : \mathcal{A} \rightarrow (X \times \mathcal{V})$,

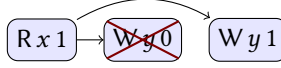
We shall say a reads v from x whenever $R(a) = (x, v)$, and a writes v to x whenever $W(a) = (x, v)$, and We shall say ϕ implies ψ whenever $\phi \models \psi$, ϕ is a *tautology* whenever $\text{true} \models \phi$, ϕ is *unsatisfiable* whenever $\phi \models \text{false}$, and ϕ is *independent of x* whenever $\phi \models \phi[v/x] \models \phi$ for every v . In examples, the actions are of the form $(R\ x\ v)$, which reads v from x , and $(W\ x\ v)$, which writes v to x . Going forward, we assume a given data model, though some examples in §3 make use of particular \mathcal{A} .

2.2 3-valued pomsets

Recall the definition of a pomset from [?]:

Definition 2.1. A pomset (E, \leq, λ) with alphabet Σ is a partial order (E, \leq) together with $\lambda : E \rightarrow \Sigma$.

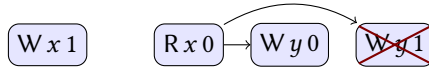
Going forward, we fix the alphabet $\Sigma = (\Phi \times \mathcal{A})$. We will write $(\phi \mid a)$ for the pair (ϕ, a) , elide ϕ when ϕ is a tautology, and write ⌫ when ϕ is unsatisfiable. We lift terminology from logical formulae and actions to events, for example if $\lambda(e) = (\phi \mid a)$ and then we say e is unsatisfiable whenever ϕ is unsatisfiable, e writes v to x whenever a writes v to x , and so forth. We visualize a pomset as a graph where the nodes are drawn from E , each node e is labelled with $\lambda(e)$, and an edge $d \rightarrow e$ corresponds to an ordering $d \leq e$. For example:



is a visualization of the pomset where:

$$0 \leq 1 \quad 0 \leq 2 \quad \lambda(0) = (\text{true}, R\ x\ 1) \quad \lambda(1) = (\text{false}, W\ y\ 0) \quad \lambda(2) = (\text{true}, W\ y\ 1)$$

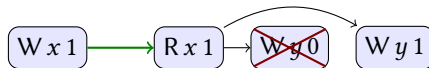
We are building a compositional semantics of shared memory concurrency, which means we require a notion of when a read has a matching write. This is a property we require of closed programs, but *not* of open programs. For example a program whose semantics includes:



may be put in parallel with another program which writes 0 to x . If the program is closed with respect to x though, such an execution cannot exist, so we need each read of x to have a matching write. This is captured by defining when e reads x from d [?]. A preliminary definition (which, as we shall see, needs to be strengthened) is:

- $d < e$,
- if e is satisfiable, then d is a tautology,
- d writes v to x , and e reads v from x , and
- there is no $d < c < e$ such that c writes to x .

In diagrams, for readability we often highlight the reads-from edges, for example:



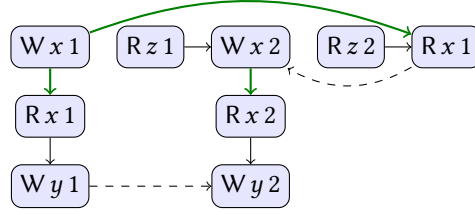
Unfortunately by itself, this is not enough. The problem is the final clause saying that there does not exist an x -blocking event c between d and e . Unfortunately, concurrency can turn events that were

- if e is satisfiable, then d is a tautology,
- d writes v to x , and e reads v from x , and
- if c writes to x then either $d \not\prec c$ or $c \not\prec e$.

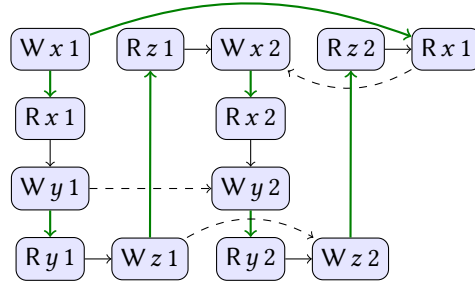
Definition 2.5. A 3-valued pomset is x -closed if for every $e \in E$:

- e is independent of x , and
- if e reads from x , then there is a d such that e can read x from d .

In our previous example, in order for $(R\ x\ 1)$ to read from $(W\ x\ 1)$, we either need $(W\ x\ 1) \not\prec (W\ x\ 2)$ or $(W\ x\ 2) \not\prec (W\ x\ 1)$, for example:



then putting this in parallel as before results in:



but this is *not* a valid 3-valued pomset, since $(W\ x\ 2) < (R\ x\ 1)$ but also $(W\ x\ 2) \not\prec (R\ x\ 1)$, which is a contradiction.

2.3 Sets of 3-valued pomsets

Our model of programs is going to be sets of 3-valued pomsets. In this section we define the operations on pomsets which are used in giving the semantics. These are operations such as prefixing, parallel composition, restriction, and so on; they are familiar from models of concurrency such as [?], but adapted here to the setting of speculative evaluation.

Definition 2.6. Let $P_0 \in (\mathcal{P}_1 \parallel \mathcal{P}_2)$ whenever there are $P_1 \in \mathcal{P}_1$ and $P_2 \in \mathcal{P}_2$ such that:

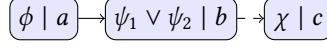
- $E_0 = E_1 \cup E_2$,
- if $e \leq_1 d$ or $e \leq_2 d$ then $e \leq_0 d$,
- if $e \not\prec_1 d$ or $e \not\prec_2 d$ then $e \not\prec_0 d$,
- if $\lambda_0(e) = (\phi_0 \mid a)$ then either:
 - $\lambda_1(e) = (\phi_1 \mid a)$ and $\lambda_2(e) = (\phi_2 \mid a)$ and ϕ_0 implies $\phi_1 \vee \phi_2$,
 - $\lambda_1(e) = (\phi_1 \mid a)$ and $e \notin E_2$ and ϕ_0 implies ϕ_1 , or
 - $\lambda_2(e) = (\phi_2 \mid a)$ and $e \notin E_1$ and ϕ_0 implies ϕ_2 .

We use $\mathcal{P}_1 \parallel \mathcal{P}_2$ in defining the semantics of conditionals and concurrency. It contains the union of pomsets from \mathcal{P}_1 and \mathcal{P}_2 , allowing overlap as long as they agree on actions. For example, if \mathcal{P}_1

and \mathcal{P}_2 contain:



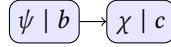
then $\mathcal{P}_1 \parallel \mathcal{P}_2$ contains:



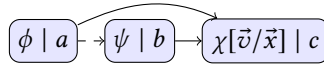
Definition 2.7. Let $a \rightarrow \mathcal{P}$ be the set \mathcal{P}' where $P' \in \mathcal{P}'$ whenever there is $P \in \mathcal{P}$ such that:

- $E' = E \cup \{c\}$,
- if $d \leq e$ then $d \leq' e$,
- if $d \not\leq e$ then $d \not\leq' e$,
- $\lambda'(c) = (\phi, a)$, and
- if $\lambda(e) = (\psi | b)$ then:
 - $\lambda'(e) = (\psi' | b)$,
 - ψ' implies $\begin{cases} \psi[v/x] & \text{if } a \text{ reads } v \text{ from } x \text{ and } c <' e \\ \psi[v/x] \wedge \psi & \text{if } a \text{ reads } v \text{ from } x \\ \psi & \text{otherwise} \end{cases}$
 - if a and b both write to y , then $c \not\leq' e$.

Prefixing is used to define the semantics of reads and writes, and adds a new event c with action a and a precondition which implies ϕ . It performs a substitution on the preconditions of any subsequent events. The tricky part of the definition is the last two clauses, which places requirements on $c < e$ and $c \not\leq e$. In normal prefixing, these would always be required, but we only require it when either the precondition of e depends on a variable read by a , or if c and e write to the same variable. For example, if \mathcal{P} contains:



where a reads v from x and writes w to y , b writes to y , and ψ is independent of x , then $(\phi | a) \rightarrow \mathcal{P}$ contains:



Definition 2.8. Let $\mathcal{P}[M/x]$ be the set \mathcal{P}' where $P' \in \mathcal{P}'$ whenever there is $P \in \mathcal{P}$ such that:

- $E' = E$,
- if $d \leq e$ then $d \leq' e$, and
- if $d \not\leq e$ then $d \not\leq' e$, and
- if $\lambda(e) = (\psi | a)$ then $\lambda'(e) = (\psi[M/x] | a)$.

and similarly for $\mathcal{P}[x/r]$.

Definition 2.9. Let $(\phi | \mathcal{P})$ be the subset of \mathcal{P} such that $P \in \mathcal{P}$ whenever:

- if $\lambda(e) = (\psi | a)$ then ϕ implies ψ .

Definition 2.10. Let $(\forall x . \mathcal{P})$ be the subset of \mathcal{P} such that $P \in \mathcal{P}$ whenever P is x -closed.

We can use the operations defined above to give the semantics of a simple concurrent imperative programming language.

Definition 2.11.

$$\begin{aligned}
\llbracket \text{skip} \rrbracket &= \{\emptyset\} \\
\llbracket x := M; C \rrbracket &= \bigcup_v (M = v \mid (\mathsf{W} x v) \rightarrow \llbracket C \rrbracket[M/x]) \\
\llbracket r := x; C \rrbracket &= \llbracket C \rrbracket[x/r] \cup \bigcup_v (\mathsf{R} x v \rightarrow \llbracket C \rrbracket[x/r]) \\
\llbracket \text{if } (M) \{ C \} \text{ else } \{ D \} \rrbracket &= (M \neq 0 \mid \llbracket C \rrbracket) \parallel (M = 0 \mid \llbracket D \rrbracket) \\
\llbracket C \parallel D \rrbracket &= \llbracket C \rrbracket \parallel \llbracket D \rrbracket \\
\llbracket \text{var } x; C \rrbracket &= \nu x. \llbracket C \rrbracket
\end{aligned}$$

A write generates a write event that may be visible to other threads. A read may see a thread-local value, or it may generate a read event that must be justified by another thread. In the latter case, occurrences of r are replaced with x (rather than v) to ensure that dependencies are tracked properly.

We have completed the formal definition of our model of speculative evaluation, and now turn to examples of this model in use.

3 EXAMPLES

3.1 Sequential memory accesses

In the semantics of memory, there are two very different ways memory can be accessed: sequentially or concurrently. These are modelled differently, since hardware and compilers give very different guarantees about their behaviour. In this section, we discuss the sequential semantics, and leave the concurrent semantics to §3.2.

Consider the program $(x := 0; y := x + 1)$. One execution of this program is where the write to y uses the sequential value of x , which is 0:

$$\boxed{\mathsf{W} x 0} \quad \boxed{\mathsf{W} y 1}$$

To see how this execution is modelled, we first expand out the syntax sugar to get the program $(x := 0; r := x; y := r + 1; \text{skip})$. Now $\llbracket \text{skip} \rrbracket$ is just $\{\emptyset\}$, and $\llbracket y := r + 1; \text{skip} \rrbracket$ includes:

$$(r + 1 = 1 \mid (\mathsf{W} y 1) \rightarrow \llbracket \text{skip} \rrbracket[1/y])$$

which contains the pomset:

$$\boxed{r + 1 = 1 \mid \mathsf{W} y 1}$$

expressing that this program can write 1 to y , as long as the precondition $(r + 1 = 1)$ is satisfied. Now $\llbracket r := x; y := r + 1; \text{skip} \rrbracket$ has two cases, the sequential case (which does not introduce a read action) and the concurrent case (which does). For the moment, we are interested in the sequential case, which is:

$$\llbracket y := r + 1; \text{skip} \rrbracket[x/r]$$

which contains the pomset:

$$\boxed{x + 1 = 1 \mid \mathsf{W} y 1}$$

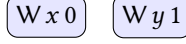
In this pomset, the precondition is $(x + 1 = 1)$, which specifies a property of the thread-local value of x . Finally $\llbracket x := 0; r := x; y := r + 1; \text{skip} \rrbracket$ includes:

$$(0 = 0 \mid (\mathsf{W} x 0) \rightarrow \llbracket r := x; y := r + 1; \text{skip} \rrbracket[0/x])$$

which contains the pomset:



all of whose preconditions are tautologies, so this has the expected behaviour:

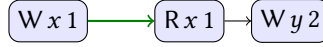


Note that there is no requirement of order between $(W x 0)$ and $(W y 1)$.

This example demonstrates how preconditions capture the sequential semantics of memory. In an execution containing an event with label $(\phi \mid a)$, one way the precondition ϕ can be discharged is by an assignment $x := M$, which performs a substitution $[M/x]$. This is a variant of the Hoare semantics for assignment, where if C has precondition ϕ then $x := M; C$ has precondition $\phi[M/x]$.

3.2 Concurrent memory accesses

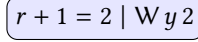
We now turn to the case of concurrent accesses to memory. Consider a concurrent version of the program from §3.1: $(x := 1 \parallel y := x + 1)$. One execution of this program is where the write to y performs a concurrent read of x :



To see how this execution is modelled, we first expand out the syntax sugar to get the program $(x := 1; \text{skip} \parallel r := x; y := r + 1; \text{skip})$. As before, $\llbracket y := r + 1; \text{skip} \rrbracket$ includes:

$$(r + 1 = 2 \mid (W y 2) \rightarrow \llbracket \text{skip} \rrbracket[2/y])$$

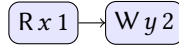
which contains the pomset:



As before, $\llbracket r := x; y := r + 1; \text{skip} \rrbracket$ has two cases. We are now interested in the concurrent case, which includes:

$$(R x 1) \rightarrow \llbracket y := r + 1; \text{skip} \rrbracket[x/r]$$

which contains the pomset:



Note that $(R x 1)$ reads 1 from x , and while $(x + 1 = 2)[1/x]$ is a tautology, $(x + 1 = 2)$ is not, and so there is a dependency $(R x 1) < (W y 2)$.

Now, $\llbracket x := 1; \text{skip} \rrbracket$ includes the pomset:



and so $\llbracket x := 1; \text{skip} \parallel r := x; y := r + 1; \text{skip} \rrbracket$ includes:



as expected, including a reads-from dependency $(W x 1) < (R x 1)$.

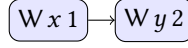
This example demonstrates how read and write events capture the concurrent semantics of memory. In an execution containing an event with label $(R x v)$, if the execution is x -closed, then there must be an event it reads from, for example one labelled $(W x v)$.

3.3 Independent writes

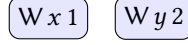
Consider an example with two independent writes ($x := 1$; $y := 2$). This has semantics including:

$$(W x 1) \rightarrow (W y 2) \rightarrow \{\emptyset\}$$

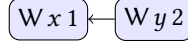
One of the executions this contains is:



but it also contains:



and:

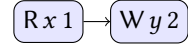


since there is no requirement that $(W x 1) < (W y 2)$.

Thus, the semantics of $(x := 1$; $y := 2)$ is the same as the semantics of $(y := 2$; $x := 1)$.

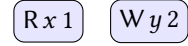
3.4 Independent reads and writes

Whereas write prefixing introduces weak dependencies on events which write to the same location, read prefixing introduces strong dependencies on preconditions which depend on the location being read. For example in §3.2 we saw that the program $(y := x + 1)$ includes the pomset:



but since $(x + 1 = 2)$ depends on x , we have the requirement that $(R x 1) \leq (W y 2)$.

This is in contrast to the program $(r := x$; $y := r + 2 - r)$. Since $(x + 2 - x = 2)$ is independent of x (at least for integer arithmetic) this contains:



and so the semantics of $(r := x$; $y := r + 2 - r)$ is the same as the semantics of $(y := 2$; $r := x)$.

Note this this example shows that we are not just dealing with a syntactic notion of dependency, which is common in hardware models of memory. In syntactic dependency, since r occurs free in $(y := r + 2 - r)$, there would be a dependency between $(r := x)$ and $(y := r + 2 - r)$. In contrast, this model is based on logical implication, which can be interpreted semantically.

3.5 Control dependencies

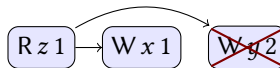
Conditionals introduce control dependencies, for example consider the program:

$$r := z; \text{ if } (r) \{ x := 1 \} \text{ else } \{ y := 2 \}$$

This includes executions in which the false branch is taken:



and ones where the true branch is taken:



In both cases, we record the actions in the branch that was not taken. This is a novel feature of this model, and is intended to capture speculative evaluation. In §3.7 we will show how this model

captures Spectre-like information flow attacks, once the attacker is provided with the ability to observe such speculations.

To see how these executions are modelled, consider the semantics of $\llbracket x := 1; \text{skip} \rrbracket$, which contains any pomset of the form:

$$\phi \mid W x 1$$

in particular it contains:

$$r \neq 0 \mid W x 1$$

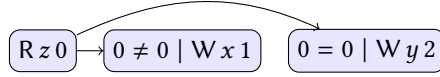
Similarly $\llbracket y := 2; \text{skip} \rrbracket$ contains:

$$r = 0 \mid W y 2$$

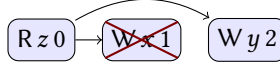
and so $\llbracket \text{if } (r) \{ x := 1; \text{skip} \} \text{ else } \{ y := 2; \text{skip} \} \rrbracket$ contains:

$$r \neq 0 \mid W x 1 \quad r = 0 \mid W y 2$$

Now, the semantics of concurrent read performs substitutions, for example:



which gives the required pomset:



Note that the precondition $r = 0$ is dependent on r , and so there is a dependency $(R z 0) < (W y 2)$, modelling the control dependency introduced by the conditional.

3.6 Control independencies

In most models of control dependencies, the dependency relation is syntactic, based on whether the action occurs inside syntactically inside a conditional. In contrast, the notion in this model is semantic: if an action can occur on both sides of a conditional, there is no control dependency. Consider a variant of the example from §3.5:

$$r := z; \text{if } (r) \{ x := 1 \} \text{ else } \{ x := 1 \}$$

This has the expected execution in which the control dependencies exist:



but it also has an execution in which the two writes of 1 to x are merged, resulting in no dependency:

$$R z 0 \quad W x 1$$

To see how this arises, consider the definition of $\llbracket \text{if } (r) \{ x := 1; \text{skip} \} \text{ else } \{ x := 1; \text{skip} \} \rrbracket$:

$$\mathcal{P}_1 \parallel \mathcal{P}_2 \quad \text{where} \quad \mathcal{P}_1 = (r \neq 0 \mid \llbracket x := 1; \text{skip} \rrbracket) \quad \text{and} \quad \mathcal{P}_2 = (r = 0 \mid \llbracket x := 1; \text{skip} \rrbracket)$$

Now, one pomset in \mathcal{P}_1 is:

$$r \neq 0 \mid W x 1$$

that is P_1 where:

$$E_1 = \{e\} \quad \lambda_1(e) = (r \neq 0, W x 1)$$

and similarly, one pomset in \mathcal{P}_2 is:

$$r = 0 \mid W x 1$$

that is P_2 where:

$$E_2 = \{e\} \quad \lambda_2(e) = (r = 0, W x 1)$$

Crucially, in the definition of $\mathcal{P}_1 \parallel \mathcal{P}_2$ there is *no* requirement that E_1 and E_2 are disjoint, and in this case they overlap at e . As a result, one pomset in $\mathcal{P}_1 \sqcup \mathcal{P}_2$ is P_0 where:

$$E_0 = \{e\} \quad \lambda_0(e) = (r \neq 0 \vee r = 0, W x 1)$$

that is:

$$W x 1$$

Note that this pomset has no precondition dependent on r , since $(r \neq 0 \vee r = 0)$ does not depend on r , which is why we end up with an execution without a control dependency:

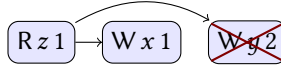
$$R z 0 \quad W x 1$$

This semantics captures compiler optimizations which may, for example merge code executed on both branches of a conditional, or hoist constant assignments out of loops.

We can now see the counterintuitive behavior of conditionals in the presence of control dependencies. There are programs such as `if (z) { x := 1 } else { x := 1 }` executions in which $(W x 1)$ is independent of $(R z 1)$:

$$R z 1 \quad W x 1$$

while programs such as `if (z) { x := 1 } else { y := 2 }` only have executions in which $(W x 1)$ is dependent on $(R z 1)$:



so these programs have different dependency relations, depending on conditional branches that were not taken. In §3.9 we shall see that this has security implications, since relaxed memory can observe dependency. The attack is similar to Spectre, so we shall take a detour to see how Spectre can be modeled in this setting.

3.7 Spectre

We give a simplified model of Spectre attacks, ignoring the details of timing. In this model, we extend programs with the ability to tell whether a memory location has been touched (in practice this is implemented using timing attacks on the cache). For example, we can model Spectre by:

```

var a; if (canRead(SECRET)) { a[SECRET] := 1 }
else if (touched a[0]) { x := 0 }
else if (touched a[1]) { x := 1 }
  
```

This is a low-security program, which is attempting to discover the value of a high-security variable SECRET. The low-security program is allowed to attempt to escalate its privileges by checking that it is allowed to read a high-security variable:

```

if (canRead(SECRET)) { ... code allowed to read SECRET ... } else { ... }
  
```

In this case, $\text{canRead}(\text{SECRET})$ is false, so the fallback code is executed. Unfortunately, the escalated code is speculatively evaluated, which allows information to leak by testing for which memory locations have been touched.

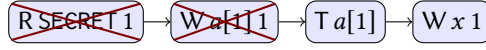
We model the touched test by introducing a new action (Tx) and defining:

$$\llbracket \text{if}(\text{touched } x) \{ C \} \text{ else } \{ D \} \rrbracket = ((Tx) \rightarrow \llbracket C \rrbracket) \cup \llbracket D \rrbracket$$

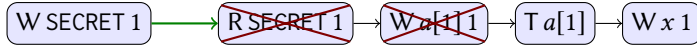
The additional requirement we need to add for x -closure is:

- if $\lambda(e) = (\phi \mid Tx)$ then there is $d < e$ where d reads or writes x .

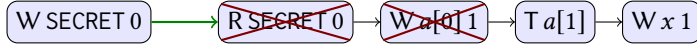
Note that there is no requirement that d be satisfiable, and indeed Spectre has the execution:



Putting this in parallel with a high-security write to SECRET gives:



but due the requirement of a -closure we do *not* have:



Thus, the attacker has managed to leak the value of a high-security location to a low-security one: if $(Wx\ 1)$ is observed, the SECRET must have been 1.

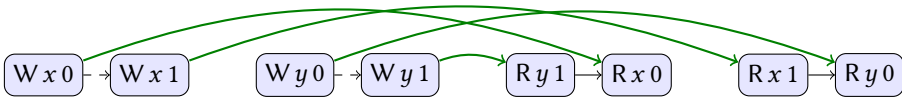
This shows how our model of speculation can express (very abstract, untimed) Spectre attacks.

3.8 Relaxed memory

In §3.9 we present an information flow attack on relaxed memory, similar to Spectre in that it relies on speculative evaluation. Unlike Spectre it does not depend on timing attacks, but instead is based on the sensitivity of relaxed memory to data dependencies. For this reason, we present a simple model of relaxed memory, which is strong enough to capture this attack. The model includes concurrent memory accesses, which can introduce concurrent reads-from. Since we are allowing events to be partially ordered, this gives a simple model of relaxed memory, for example an independent read independent write (IRIW) example is:

$$x := 0; x := x + 1 \parallel y := 0; y := y + 1 \parallel \text{if}(x) \{ r := y \} \parallel \text{if}(y) \{ s := x \}$$

which includes the execution:



This model does not introduce thin-air reads (TAR), for example the TAR pit is $(x := y \parallel y := x)$ but an attempt to produce a value from thin air fails, for the usual reason of producing a cycle in \leq , as shown on the left below:



This cycle can be broken by removing a dependency, for example $(x := y \parallel r := x; y := r + 1 - r)$ has the execution on the right above. Note that $(Rx\ 1) \not\leq (Wy\ 1)$, so this does not introduce a cycle.

Although it is not the primary focus of this paper, our model may be an attractive model of relaxed memory. Many prior models either permit thin-air executions that our model forbids or forbid desirable executions that our model permits.

In §5, we develop a logic which allows us to prove that our semantics forbids thin air examples that are permitted by prior speculative models [? ? ?].

Our model passes all of the causality test cases [?]. Significantly, this includes test case 9, which is forbidden by [?], one of the few models that disallows the thin air example from §5. We present this test case in the appendix, where we also discuss the thread inlining examples from [?].

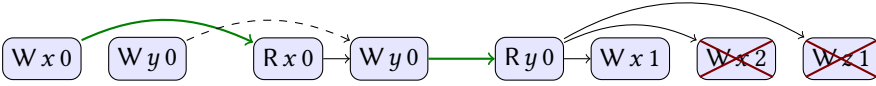
[?] showed that the thin-air problem has no per-candidate-execution solution for C++. This result does not apply to our model, as the semantics of a conditional can depend on the semantics of both branches.

3.9 Information flow attacks on relaxed memory

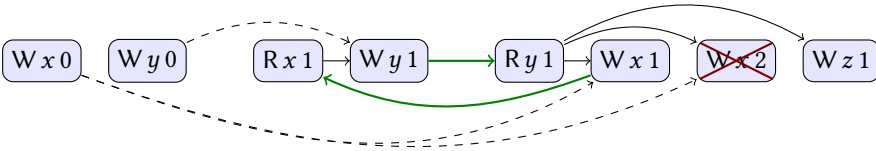
Consider an attacker program, again using security checks to try to learn a SECRET. Whereas SPECTRE uses hardware capabilities, which have to be modeled by adding extra capabilities to the language, this new attacker works by exploiting relaxed memory which can result in unexpected information flows. The attacker program is:

```
var x:=0; var y:=0;
y:=x || if (y==0) { x:=1 }
      else if (canRead(SECRET)) { x:= SECRET }
      else { x:=1; z:=1 }
```

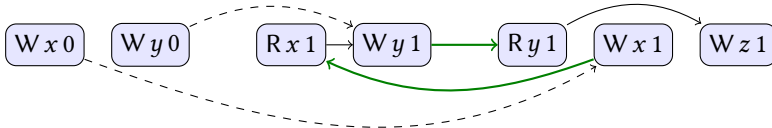
In the case where SECRET is 2, this has many executions, one of which is:



but there are no executions which exhibit (W z 1), since any attempt to do so produces a cycle:



In the case where SECRET is 1, there is an execution:



Note that in this case, there is no dependency from (R y 1) to (W x 1), which is what makes this execution possible. Thus, if the attacker sees an execution with (W z 1), they can conclude that SECRET is 1, which is an information flow attack.

This attack is not just an artifact of the model, since the same behavior can be exhibited by compiler optimizations. Consider the program fragment:

```
if (y = 0) { x := 1 } else if (canRead(SECRET)) { x := SECRET } else { x := 1; z := 1 }
```

Now, in the case where SECRET is a constant 1, the compiler can inline it:

```
if (y = 0) { x := 1 } else if (canRead(SECRET)) { x := 1 } else { x := 1; z := 1 }
```

and lift the assignment to x out of the if statement:

```
x := 1; if (y = 0) { } else if (canRead(SECRET)) { } else { z := 1 }
```

After these optimizations, a sequentially consistent execution exhibits (Wz1). We discuss the practicality of this attack further in §4.

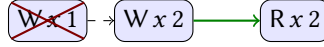
This approach can be generalized to detect information flows in arbitrary code. If we replace the code fragment above with:

```
if (y = 0) { x := P(0) } else { x := P(1); z := 1 }
```

then (Wz1) is possible only if P is independent of its input. Thus, the conditional is able to capture multiple executions, as in [?].

3.10 Dead store elimination

A common compiler optimization is *dead store elimination*, in which writes are omitted if they will be overwritten by a subsequent write later in the same thread. We can model eliminated writes by ones with an unsatisfiable precondition. For example, one execution of $(x := 1; x := 2) \parallel (r := x)$ is:



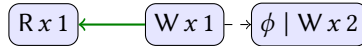
Recall that for any satisfiable e , if e reads x from y then d is a tautology. This means that, although we can eliminate (Wx1) we cannot eliminate (Wx2).

One heuristic that a compiler might adopt is to only eliminate writes that are guaranteed to be followed by another write to the same variable. This can be formalized by saying that d is eliminatable if there is a $e \not\vdash d$ such that e is a tautology and d writes to every location e writes to. A model of dead store elimination is one where, in every pomset, every eliminatable event is unsatisfiable. This simple model includes the examples above.

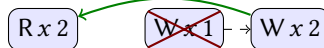
Note that if dead store elimination is *always* performed, then there is an information flow attack similar to the one in §3.9. Consider the program:

```
r := x || x := 1;
if (canRead(SECRET)) { if (SECRET != 0) { x := 2 } }
else { x := 2 }
```

In the case that SECRET is 0, there is an execution:



where ϕ is $(\neg \text{canRead}(\text{SECRET}))$, which is not a tautology, and so the (Wx1) event is not eliminated. In the case that SECRET is not 0, the matching execution is:



Now the (Wx2) event is a guaranteed write, so the (Wx1) is eliminated, and so cannot be read. In the case that the attacker can rely on dead store elimination taking place, this is an information flow: if the attacker observes x to be 1, then they know SECRET is 0. We return to this attack in §4.

3.11 Release/acquire synchronization

In relaxed memory models, synchronization actions act as memory fences: that is, they are a barrier to reordering memory accesses. In this section, we present a simple model of release/acquire fencing. In §3.12, we show that this can be scaled up to a model of transactional memory.

We assume there are sets Rel and $\text{Acq} \subseteq \mathcal{A}$. We say that a is a *release action* if $a \in \text{Rel}$ and a is an *acquire action* if $a \in \text{Acq}$. In a pomset, a release event is one labelled with a release action, and an acquire event is one labelled by an acquire action. The semantics of fences are given by adding an extra constraint to the definition of $(\phi \mid a) \rightarrow \mathcal{P}$ (recalling that c is the a -labelled event being introduced):

- $c \leq e$ whenever c is an acquire event or e is a release event.

This constraint ensures that events are ordered before a release and after an acquire.

In examples, we will use releasing writes and acquiring reads:

- $(\text{Rel } x \ v)$, a release action that writes v to x , and
- $(\text{Acq } x \ v)$, an acquire action that reads v from x .

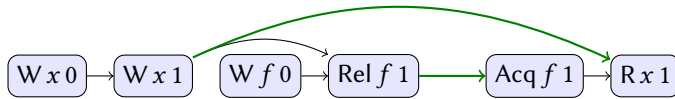
The semantics of programs with releasing write and acquiring read are the same as for regular write and read and, but with $\text{Rel } x \ v$ replacing $\text{W } x \ v$ and $\text{Acq } x \ v$ replacing $\text{R } x \ v$.

$$\begin{aligned} \llbracket \text{rel } x := M; C \rrbracket &= \bigcup_v (M = v \mid \text{Rel } x \ v) \rightarrow \llbracket C \rrbracket [M/x] \\ \llbracket \text{acq } r := x; C \rrbracket &= \llbracket C \rrbracket [x/r] \cup \bigcup_v (\text{true} \mid \text{Acq } x \ v) \rightarrow \llbracket C \rrbracket [x/r] \end{aligned}$$

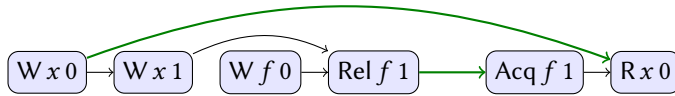
For example, consider the program:

$(x := 0; f := 0; x := 1; \text{rel } f := 1;) \parallel (\text{acq } r := f; s := x;)$

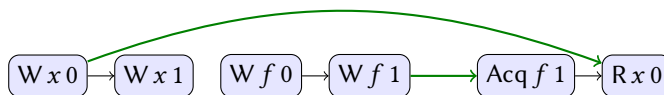
This has an execution:



but *not*:



since $(W x 0) < (W x 1) < (R x 0)$, so this pomset does not satisfy the requirements to be a x -closed. If we replace the release with a plain write, then the outcome $(\text{Acq } f \ 1)$ and $(R x \ 0)$ is possible:



since no order is required between $(W x \ 1)$ and $(W f \ 1)$. Symmetrically, if we replace the acquire of the original program with a plain read, then the outcome $(R f \ 1)$ and $(R x \ 0)$ is possible.

3.12 Transactions

We present a model of serializable transactions that is sufficient to capture PRIME+ABORT attacks [?]. We make several simplifying assumptions: transactions are serializable, strongly isolated [?], and only abort due to cache conflicts. To model the latter, we assume that the set of locations \mathcal{X} is partitioned into *cache sets*.

The action $(B\ v) \in \text{Acq}$ represents the begin of a transaction with id v and $(C\ v) \in \text{Rel}$ represents the corresponding commit, with semantics:

$$\begin{aligned} \llbracket r := \text{begin } v; D \rrbracket &= (\text{true} \mid B\ v) \rightarrow \llbracket C \rrbracket[v/r] \\ \llbracket r := \text{commit } v; D \rrbracket &= (\text{true} \mid C\ v) \rightarrow \llbracket D \rrbracket[1/r] \cup (\text{false} \mid C\ v) \rightarrow \llbracket D \rrbracket[0/r] \end{aligned}$$

At top level, we require that pomsets be *serializable*, as defined below.

Definition 3.1. We say that event c *matches* b if $\lambda(c) = (C\ v)$ and $\lambda(b) = (B\ v)$. We say that a begin event *aborts* if every matching commit is unsatisfiable. We say that begin event b *begins* e if $b \leq e$ and there is no intervening matching commit; in this case e *belongs to* b . We say that commit event c *commits* e if $e \leq c$ and there is no intervening matching begin. A pomset is *serializable* if:

- (1) no two begins have the same id,
- (2) every commit follows the matching begin,
- (3) \leq totally orders tautological begins and commits,
- (4) if b begins e , but not d , then $d \leq e$ implies $d \leq b$ and $e \not\leq d$ implies $b \not\leq d$
- (5) if c ends e , but not d , then $e \leq d$ implies $c \leq d$ and $d \not\leq e$ implies $d \not\leq c$,
- (6) if e and d belong to b and read the same location, then both read the same value,
- (7) if e belongs to b , then e implies some matching c that ends e , and
- (8) if b aborts then there must be some e that belongs to b and some tautologous $d \leq e$ that does not belong to b such that e and d touch locations in the same cache set.

In discussion, we identify transactions by their unique begin event. A transaction that does not abort is *successful*. Conditions 1-5 ensure serializability of committed transactions. Conditions 4-6 also ensure strong isolation for non-transactional events [?]. Condition 7 ensures that all events in aborted transactions are unsatisfiable. Condition 8 requires that aborts only occur due to caching conflicts—this is similar to the treatment of the touch operation in §3.7.

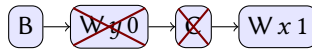
To model PRIME+ABORT, we require an honest agent whose cache-set access depends upon a secret. If $a[0]$ and $a[1]$ belong to separate cache sets, then such an honest agent is:

$a[\text{SECRET}] := 1$

The attack relies on discovery of some y which belongs to the cache-set of $a[1]$. Then the program

`begin; y:=0; r:=commit; if (r==0) {x:=1;}`

can write 1 to x only if the SECRET is 1, in which case the following execution is possible.



4 EXPERIMENTS

One theme of this paper is that optimizations not typically part of formal abstractions can result in information-flow leaks. This is typified by the Spectre attack, which leverages speculative execution, a hardware optimization. §3.9 and §3.10 presented other attacks along the same line, which leverage relaxed memory models and dead store elimination respectively. These attacks may result, not from hardware optimizations, but from common *compiler* optimizations. These attacks also, unlike Spectre, do not rely on timing side channels, or indeed timers of any kind, bypassing many common Spectre mitigations [?].

In this section we present concrete implementations of the attacks outlined in §3.9 and §3.10, in both cases leveraging compiler optimizations to construct an information flow attack. The attacker model for these attacks (detailed in §4.1) is currently unrealistic in a real-world sense, rendering these attacks proof-of-concepts rather than immediately exploitable vulnerabilities. However, we believe the novelty of their general mechanisms may lead to interesting discussion; and with much more development, these attacks may evolve into genuine threats against real-world targets such as JIT compilers. We demonstrate the efficacy of both of our concrete proof-of-concept attacks against the clang and gcc C compilers.

All of our experiments are performed on a **{describe machine}** with clang version **{clang version}** and gcc version **{gcc version}**.

4.1 Attacker model

In our attacker model, we assume that there is a SECRET which an attacker wishes to learn; for instance, SECRET may be a cryptographic key hardcoded into the application. This SECRET is known at compile time, but may not be accessed except behind a security check. Since the attacker is running with low security privileges, the security check always fails, so the attacker can only access SECRET in dead code. The attacker has no capabilities other than writing and executing code – in particular the attacker may not disassemble the compiler or libraries to learn the SECRET directly; may not examine the internal state of the compiler; may not access timers of any kind; and may not leverage hardware side channels. The attacker’s goal is to learn the value of the SECRET.

As a hypothetical concrete example, suppose there is a library which contains a hardcoded SECRET such as an API or signing key, which cannot be accessed directly, only through a function guarded by a security check:

```
private static uint SECRET = 0x1234;
public uint get_secret() {
    if (canRead(SECRET)) { return SECRET; }
    else { return 0; }
}
```

As noted above, this is not necessarily a realistic attacker model, since in most cases secrets are only known at run time rather than compile time, which means that the attacks presented in this section are more of theoretical interest than practical concern. However, the mechanism of the attack is novel and could potentially be applied in other contexts. For instance, many real-world contexts allow attackers (untrusted or third-party entities) to write code in a scripting language which is then compiled alongside and integrated into a larger application, often using a just-in-time (JIT) compiler. JavaScript code from third-party websites running in a browser is a common example of this. Our attack gives an attacker similar capabilities against a compiler, except it considers the simpler setting of using C code against a C compiler. One could imagine a similar attack using JavaScript against browser JIT compilers, where the compiler may have access to interesting secrets in the browser itself, and may be able to optimize based on those secrets. We plan to explore JavaScript attacks of this type as future work.

4.2 Load-store reordering attack

We begin by examining the attack in §3.9 in more detail, subject to the attacker model given above. In particular, we show that by exploiting compiler optimizations which perform load-store reordering, an attacker can learn the value of a compile-time SECRET despite only being allowed to use it inside dead code, that is, code that can never be executed at runtime. This attack was tested and works against gcc version **{gcc version}**.

SECRET == 0	SECRET == 1
mov f(%rip), %eax	mov f(%rip), %eax
mov \$1, x(%rip)	mov y(%rip), %eax
test %eax, %eax	mov \$1, x(%rip)
je label1	test %eax, %eax
mov %0, x(%rip)	sete %eax
label1:	ret
mov y(%rip), %eax	
test %eax, %eax	
sete %eax	
ret	

Fig. 1. (Simplified) x86 assembly output from gcc for the main thread of the load-store reordering attack. In particular, note that the order between `mov $1, x(%rip)` and `mov y(%rip), %eax` is different in the two cases. The call to `canRead(SECRET)` has been inlined; we implemented `canRead(x)` as `return f`; where `volatile bool f = false`; . Thus, gcc preserves the read of `f` even when its value is unused, as in the case on the right.

The form of the attack presented in §3.9 works in theory, but in practice, just because a compiler is *allowed* to perform a load-store reordering doesn't mean that it *will*. We found that gcc and clang chose to read `y` into a register first (before writing to `x`), regardless of the value of `SECRET`. However, we did find a related pattern in which gcc will emit a different ordering of the read of `y` and the write of `x` depending on the value of a `SECRET`:

```
var x:=0; var y:=0;
y:=x || x:=1;
    if (canRead(SECRET)) { x:=SECRET; }
    if (y > 0) { return 0 }
    else { return 1 }
```

Figure 1 shows the assembly output of gcc in the cases where `SECRET` is 0 and 1 respectively. In the case that `SECRET` is 1, gcc removes the `if` statement entirely, and moves the read of `y` above the write of `x`. However, when `SECRET` is 0, the `if` statement must remain intact, and gcc does not move the read of `y`. This means that if `SECRET` is 1, the second thread will always read `y == 0` and always return 1. However, if `SECRET` is 0, it is possible that the first thread may observe `x == 1` and write `y := 1` in time for the second thread to observe `y == 1` and thus return 0. In this way, we leverage compiler load-store reordering to learn the value of a compile-time `SECRET`.

We extend this attack to leak a secret consisting of an arbitrary number `N` of bits. To do this, we simply compile `N` copies of the test function, each performing a boolean test on a single bit of the secret. The function used for reading the `k`th bit is as follows (for `N <= 64`):

```
var x:=0; var y:=0;
y:=x || x:=1;
    if (canRead(SECRET)) { x:=(SECRET & (1 << k)) ? 1 : 0; }
    if (y > 0) { return 0 }
    else { return 1 }
```

Following the same analysis as above, this function will always return 1 if the appropriate bit of `SECRET` is 1, but may return 0 if the appropriate bit of `SECRET` is 0. The extension of the attack

Redundancy	Bandwidth (bits/s)	Bitwise Acc	Per-run Acc
1	3.17 million	90.13%	0.0%
2	1.62 million	96.77%	0.7%
3	1.07 million	98.84%	3.9%
4	812 thousand	99.55%	13.5%
5	652 thousand	99.83%	34.0%
7	466 thousand	99.97%	71.8%
10	322 thousand	99.998%	96.6%
15	216 thousand	100.00%	100.0%

Fig. 2. Performance results for the load-store reordering attack when leaking a 2048-bit secret. ‘Redundancy’ is the number of redundant runs performed for error correction; more redundant runs improves accuracy but reduces bandwidth. ‘Bandwidth’ is the number of bits leaked per second after accounting for any error correction. ‘Bitwise Accuracy’ is the percentage of bits that were correct, while ‘Per-run Accuracy’ is the percentage of full 2048-bit secrets that were correct in all bit positions. **{Note: numbers are not final (collected on Craig’s machine inside a VM), but give an idea of where we stand.}**

to the general case with truly arbitrary N is straightforward; SECRET becomes an array of 64-bit values, and we use $k / 64$ and $1 \ll (k \& 63)$ as the array index and bitmask respectively.

We make three additional tweaks to improve the reliability so that the attacker can confidently infer the value of SECRET based on the observed return values of the function. First, rather than performing $y := x$ only once in the first thread, we perform $y := x$ continuously in a loop. This maximizes the probability that, once $x := 1$ occurs in the second thread, y will be immediately assigned 1 by the first thread and the second thread will be able to read $y == 1$.

Second, we wish to lengthen the timing window between $x := 1$ and the read of y in the second thread (in the case where the appropriate bit of SECRET is 0 and the read of y remains below $x := 1$). However, we wish to do this in a way that does not block the reordering of the read of y upwards in the case where the appropriate bit of SECRET is 1. We do this by inserting many copies of the line

```
if (canRead(SECRET)) { x := (SECRET & (1 << k)) ? 1 : 0; }
```

instead of just one. In the case where the appropriate bit of SECRET is 0, this results in many calls to `canRead(SECRET)` and many conditional jumps, which in practice creates a timing window for the first thread to perform $y := x$. However, in the case where the appropriate bit of SECRET is 1, all of these inserted lines can be removed just as a single copy could be. In practice, we found that inserting too many copies of the line prevents gcc from reordering the read of y above the write to x as desired; inserting 30 copies was sufficient to create a timing window while still allowing the desired reordering.

Finally, we redundantly execute the entire attack several times, noting the return value of the function in each case. We note that if *any* of the redundant runs produces a return value of 0 for a particular bit position, we can be certain that the corresponding bit of SECRET *must* be 0, as it implies the read of y was not reordered upwards in that particular function. On the other hand, the more runs that produce a return value of 1 for a particular bit position, the more certain we can be that the read of y was reordered above the $x := 1$ assignment, and the appropriate bit of SECRET is 1.

Figure 2 gives the performance results for this attack against gcc version **{ver}**. The attack can sustain hundreds of thousands of bits per second leaked with near-perfect accuracy, or millions of bits per second with error rates of a few percent. Note that this bandwidth assumes that all copies of the attack function are already compiled; the cost of compilation is not included here.

4.3 Dead store elimination attack

In this section we return to the attack in §3.10 based on dead store elimination. We show that in our attacker model (given in §4.1), the attacker is able to exploit dead store elimination to again learn the value of a compile-time SECRET despite only being allowed to use it inside dead code, that is, code that can never be executed at runtime. This attack is even more efficient than the attack on load-store reordering, and further, we were able to demonstrate its effectiveness against both gcc and clang.

We start from the simple form of the attack presented in §3.10, and extend it to leak a secret consisting of an arbitrary number N of bits. As we did in the load-store reordering attack, we again compile N copies of the test function, each performing a boolean test on a single bit of the secret. The function used for reading the k th bit is as follows (for $N \leq 64$):

```
var x:=0;
  r:=x || x:=1;
      if (canRead(SECRET)) { if (SECRET & (1 << k) !=0) { x:=2 } }
      else { x:=2 }
```

Then, we test each function in turn, each time noting the value of r observed by the ‘listening’ thread. If the appropriate bit of SECRET is 1, the $x := 2$ assignment is guaranteed to happen, so the compiler can eliminate the $x := 1$ assignment as a dead store and we will observe $r == 2$; however, if the appropriate bit of SECRET is 0, the $x := 1$ assignment cannot be eliminated, and we will observe $r == 1$ with some probability. The extension of the attack to the general case with truly arbitrary N is straightforward and proceeds exactly as it did for the attack on load-store reordering.

We make three additional tweaks to improve the reliability so that the attacker can confidently infer the value of SECRET based on the observed values of r . These three tweaks strongly resemble the reliability tweaks we made to the load-store reordering attack and differ only in a few details.

First, rather than simply observing x with $r := x$ in the ‘listening’ thread, we continuously load x in a loop until a nonzero value is observed — i.e., we perform $\text{do } \{r := x\} \text{ while } (r == 0)$. This remedies the case where $r := x$ could observe a value of x from ‘before’ either of the two possible writes performed by the other thread. **do we need to explicitly say that we ensure x is initialized to 0 — and coherently seen as such by both threads — before starting the attack?**

Second, we insert additional time-consuming computation immediately following the $x := 1$ operation in the ‘main’ thread. This lengthens the timing window in which x has the value 1, increasing the likelihood that the ‘listening’ thread will be able to observe $x == 1$ (unless the $x := 1$ write was eliminated, of course). Inserting this computation can be done without interfering with the dead store elimination process itself, so that the compiler will continue to eliminate the $x := 1$ write if and only if the appropriate bit of SECRET was 1. For gcc, we have a fair amount of freedom with the time-consuming computation — for instance, we can use an arbitrarily long loop. In fact, we can perform a further optimization by monitoring the value of the variable r (written to by the listening thread) and breaking out of the loop early if we see that the listening thread has already observed $x == 1$. However, with clang, we cannot use a loop at all — the time-consuming computation must be branch-free, and furthermore must not consist of too many instructions. This is because clang’s dead store elimination pass operates only within basic blocks, and uses a heuristic to stop scanning the basic block early if it is too large. Nonetheless, we find that even with these restrictions, we are able to construct a reliable and fast attack against both clang and gcc.

Finally, we redundantly execute the entire attack several times, noting the final value of r (the first observed nonzero value of x) in each case. We note that if *any* of the redundant runs produces

$r == 1$ for a particular bit position, we can be certain that the corresponding bit of SECRET *must* be 0, as it implies that the $x := 1$ write was not eliminated in that particular function. On the other hand, the more runs that observe $r == 2$ in a particular bit position despite our other reliability-increasing measures taken above, the more certain we can be that the $x := 1$ write was eliminated in that function, and the appropriate bit of SECRET is 1.

Unlike the load-store reordering attack, our implementation of the dead store elimination attack has two important “knobs” which trade off reliability vs. performance, rather than only one. First, we have the length of time which the writing thread attempts to “stall” immediately after the $x := 1$ write. Second, we have the number of entire redundant runs of the attack that are performed before the attacker reaches her conclusion. Increased reliability can be achieved by adjusting either of these knobs, and they each have (different) effects on the overall performance of the attack. After exploring the parameter space, we found that 3 redundant runs is sufficient to provide near-100% accuracy while allowing us to maximize the speed of the attack. Specifically, on our machine, our attack on gcc reaches speeds of {exact gcc leak speed} bits leaked per second ({exact gcc raw leak speed} ‘raw’ bits leaked per second, that is, before error correction) with {exact gcc accuracy}, while our attack on clang reaches speeds of {exact clang leak speed} bits leaked per second ({exact clang raw leak speed} ‘raw’ bits leaked per second) with {exact clang accuracy}. In particular, this means our attack can leak a 2048-bit cryptographic key in under {exact speed} ms, on either gcc or clang, with probability {exact probability} that there are exactly zero bit errors in the leaked key, or probability {exact probability} that there is at most one bit error in the leaked key.

5 LOGIC

In this section, we develop sufficient logical infrastructure to prove that our semantics disallows thin air executions. We present a variant of the TAR-pit example from §3.8 which poses difficulties under many speculative semantics.

We adapt past linear temporal logic (PLTL) [?] to pomsets by dropping the previous instant operator and adopting strict versions of the temporal operators. The atoms of our logic are write and read events. Given an pomset P and event e , define:

$$\begin{aligned}
 P, e &\models W x v, \text{ if } \lambda(e) = (\text{true}, W x v) \\
 P, e &\models R x v, \text{ if } \lambda(e) = (\text{true}, R x v) \\
 P, e &\models \phi \wedge \psi, \text{ if } P, e \models \phi \text{ and } P, e \models \psi \\
 P, e &\models \text{true} \\
 P, e &\models \neg\phi, \text{ if } P, e \not\models \phi \\
 P, e &\models \Box^{-1}\phi, \text{ if } (\forall d \leq e, d \neq e) P, d \models \phi
 \end{aligned}$$

Define $P \models \phi$ if $(\forall e \in E) P, e \models \phi$ and $\mathcal{P} \models \phi$ if $(\forall P \in \mathcal{P}) P \models \phi$.

Let $\Diamond^{-1}\phi$ be defined as $\neg(\Box^{-1}\neg\phi)$. In addition, let false, \vee and \Rightarrow be defined in the standard way.

The past operators do not include the current instant, and thus they do *not* satisfy the rule $\Box^{-1}\phi \Rightarrow \Diamond^{-1}\phi$. However, they do satisfy:

$$\begin{aligned}
 (\phi \Rightarrow \Diamond^{-1}\phi) &\Rightarrow \neg\phi && \text{(Coinduction)} \\
 (\Box^{-1}\phi \Rightarrow \phi) &\Rightarrow \phi && \text{(Induction)}
 \end{aligned}$$

Note that $P \models \phi \wedge \Box^{-1}\phi$ whenever $P \models \phi$.

We now present two proof rules. The first rule captures the semantics of local variables. Define $\text{closed}(x) = (R x v \Rightarrow \Diamond^{-1}W x v)$. Although this definition does not mention intervening writes, it

is sufficient for our example. It is straightforward to establish that following rule is sound:

$$\frac{\phi \text{ is independent of } x \quad P \models \text{closed}(x) \Rightarrow \phi}{\forall x. P \models \phi} \quad (\text{Closing } x)$$

The second rule describes composition, in the style of Abadi and Lamport [?]. To simplify the presentation, we present the special case with a single invariant. In order to state the theorem, we generalize the satisfaction relation to include environment assumptions. Let $\text{Models}(\phi) = \{\mathcal{P} \mid \mathcal{P} \models \phi\}$ be the set of pomsets that satisfy ϕ . We say that ϕ is prefix closed if $\text{Models}(\phi)$ is prefix-closed¹. Define $\phi, \mathcal{P} \models \psi$ if $\text{Models}(\phi) \parallel \mathcal{P} \models \psi$.

PROPOSITION 5.1. *Let ϕ be prefix-closed. Let $\mathcal{P}_1, \mathcal{P}_2$ be augmentation-closed². Then:*

$$\frac{\phi, \mathcal{P}_1 \models \phi \quad \phi, \mathcal{P}_2 \models \phi}{\mathcal{P}_1 \parallel \mathcal{P}_2 \models \phi} \quad (\text{Composition})$$

PROOF SKETCH. We will show that all prefixes in the prefix closures of $\mathcal{P}_1 \parallel \mathcal{P}_2$ satisfy the required property. Proof proceeds by induction on prefixes of $P \in \mathcal{P}_1 \parallel \mathcal{P}_2$.

The case for empty prefix follows from assumption that ϕ is prefix closed.

For the inductive case, consider $P \in \mathcal{P}_1 \parallel \mathcal{P}_2$ where $P_i \in \mathcal{P}_i$. Since \mathcal{P}_1 and \mathcal{P}_2 are augmentation closed, we can assume that the restriction of P to the events of P_i coincides with P_i , for $i = 1, 2$. Consider a prefix P' derived by removing a maximal element e from P . Suppose e comes from P_1 (the other case is symmetric). Since P_2 is a prefix of P' and $P' \models \phi$ by induction hypothesis, we deduce that $P_2 \models \phi$. Since $P_1 \in \mathcal{P}_1$, by assumption $\phi, \mathcal{P}_1 \models \phi$ we deduce that $P \models \phi$. \square

We now turn to the example program, which is a variant of [? , Figure 8]:

var $x := 0$; var $y := 0$; var $z := 0$; ($y := x \parallel \text{if } (z) \{ x := 1 \} \text{ else } \{ x := y; a := y \} \parallel z := 1$)

This program is allowed to write 1 to a under many speculative memory models [? ? ?], even though the read of 1 from y in the else branch of the second thread arises out of thin air. In contrast, we prove the formula $\neg \Diamond^{-1}(\text{W } a \text{ } 1)$ holds for the models of this program in our semantics. We start with the following invariant, which holds for each of the three threads, and thus, by composition, for the aggregate program:

$$[\Diamond^{-1}(\text{W } y \text{ } 1) \Rightarrow \Diamond^{-1}(\text{R } x \text{ } 1)] \wedge [\Diamond^{-1}(\text{W } a \text{ } 1) \Rightarrow (\Diamond^{-1}(\text{R } y \text{ } 1) \wedge \Box^{-1}(\text{W } x \text{ } 1 \Rightarrow \Diamond^{-1}(\text{R } y \text{ } 1)))]$$

Closing y , we have, $\Diamond^{-1}(\text{R } y \text{ } 1) \Rightarrow \Diamond^{-1}(\text{W } y \text{ } 1)$ which we substitute into the left conjunct to get:

$$\Diamond^{-1}(\text{R } y \text{ } 1) \Rightarrow \Diamond^{-1}(\text{R } x \text{ } 1)$$

which in turn we substitute into the right conjunct to get:

$$\Diamond^{-1}(\text{W } a \text{ } 1) \Rightarrow (\Diamond^{-1}(\text{R } x \text{ } 1) \wedge \Box^{-1}(\text{W } x \text{ } 1 \Rightarrow \Diamond^{-1}(\text{R } x \text{ } 1)))$$

Closing x , we can replace $\Diamond^{-1}(\text{R } x \text{ } 1)$ with $\Diamond^{-1}(\text{W } x \text{ } 1)$:

$$\Diamond^{-1}(\text{W } a \text{ } 1) \Rightarrow (\Diamond^{-1}(\text{W } x \text{ } 1) \wedge \Box^{-1}(\text{W } x \text{ } 1 \Rightarrow \Diamond^{-1}(\text{W } x \text{ } 1)))$$

Applying coinduction to the right conjunct, we have:

$$\Diamond^{-1}(\text{W } a \text{ } 1) \Rightarrow (\Diamond^{-1}(\text{W } x \text{ } 1) \wedge \Box^{-1}(\neg \text{W } x \text{ } 1))$$

Simplifying, we have, as required:

$$\Diamond^{-1}(\text{W } a \text{ } 1) \Rightarrow \text{false}$$

¹ P' is a prefix of P if $E' \subseteq E$, $e \in E'$ and $d \leq e$ imply $d \in E'$, and $(\lambda', \leq', \dagger')$ coincide with (λ, \leq, \dagger) for elements of E' .

² P' is an augmentation of P if $E' = E$, $e \leq d$ implies $e \leq' d$, $e \dagger d$ implies $e \dagger' d$, and if $\lambda(e) = (\psi \mid b)$ then $\lambda'(e) = (\psi' \mid b)$ where ψ' implies ψ .

6 CONCLUSIONS AND FUTURE WORK

In this paper, we have presented a model of speculative evaluation, and showed that it captures non-trivial properties of speculations produced by hardware, compiler optimizations, and transactions. These properties include information flow attacks: in the case of hardware and transactions this is modelling known attacks [? ?], but in the case of compiler optimizations the attacks are new, and were discovered as a direct result of developing the model. We have experimentally validated that the attacks can be carried out against gcc and clang, though only against secrets known at compile time.

The model of relaxed memory used in this paper is deliberately simplified, compared for example to C11 [? ?]. In particular our model of reads-from is strong, and could be weakened by replacing the requirement $d < e$ in Definition 2.4 by $e \not\prec d$. It remains to be seen how this impacts the model, in particular the logical formulation of x -closure in §5 as $((R\ x\ v) \Rightarrow \Diamond^{-1}(W\ x\ v))$ would no longer be sound. The model is also not considering coherence, though we speculate it can be added by requiring that for each x , \prec form a total order when restricted to events that write to x .

The design space for transactions is very rich [?]. We have only presented one design choice, and it remains to be seen how other design choices could be adopted. For example we have chosen to pun between commits which are aborted due to transaction failure, and ones which are aborted for other reasons such as failed speculation.

One interesting feature of this model is that (in the language of [?]) it is a *per-candidate execution model*, in that the correctness of an execution only requires looking at that one execution, not at others. This is explicit in memory models such as [? ?] in which “alternative futures” are explored, in a style reminiscent of Abramsky’s bisimulation as a testing equivalence [?]. Models of information flow are similar, in that they require comparing different runs to test for the presence of dependencies [?]. In contrast, the model presented here explicitly captures dependency in the pomset order, and models multiple runs by giving the semantics of `if` in terms of a concurrent semantics of both branches. In the parlance of information flow [?], the humble conditional suffices to construct a composition operator to detect information flow in the presence of speculation.

A MEMORY MODEL EXAMPLES

[?] developed a set of twenty causality test cases in the process of revising the Java Memory Model (JMM) [?]. Using hand calculation, we have confirmed that our model gives the desired result for all twenty cases, unrolling loops as necessary. Our model also gives the desired results for all of the examples in [?, §4] and all but one in [?, §5.3]: redundant-write-after-read-elimination—this counterexample applies to any sensible non-coherent semantics. Our model agrees with the JMM on the “surprising and controversial behaviors” of [?, §8], and thus fails to validate thread inlining.

In this section, we discuss three of the causality test cases and the thread inlining from [?]. In presenting the examples, we unroll loops, correct typos and simplify the code.

A.1 Causality test case 8

Test case 8 asks whether:

$$\text{var } x := 0; \text{ var } y := 0; (\text{if } (x < 2) \{ y := 1 \} \parallel x := y)$$

may read 1 for both x and y . This behavior is allowed, since “interthread analysis could determine that x and y are always either 0 or 1.” This breaks the dependency between the read of x and the write to y in the first thread, allowing the write to be moved earlier.

A.2 Causality test case 9

Test case 8 asks whether:

$$\text{var } x := 0; \text{ var } y := 0; (\text{if } (x < 2) \{ y := 1 \} \parallel x := y \parallel y := 2;)$$

may read 1 for both x and y . This behavior is also allowed. This is “similar to test case 8, except that x is not always 0 or 1. However, a compiler might determine that the read of x by thread 1 will never see the write by thread 3 (perhaps because thread 3 will be scheduled after thread 1)”

Note that [?] fails on this case.

A.3 Causality test case 14

Test case 14 asks whether:

$$\text{var } a := 0; \text{ var } b := 0; \text{ var } y := 0; (\text{if } (a) \{ b := 1 \} \text{ else } \{ y := 1 \} \parallel \text{while } (y + b == 0) \{ \text{skip} \} a := 1)$$

may read 1 for a and b , yet 0 for y . Here a and b are regular variables and y is volatile, which is equivalent to release/acquire in this example. This behavior is also disallowed, since “in all sequentially consistent executions, [the read of a gets 0] and the program is correctly synchronized. Since the program is correctly synchronized in all SC executions, no non-sc behaviors are allowed.”

Unrolling the loop once, we have

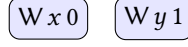
$$\text{var } a := 0; \text{ var } b := 0; \text{ var } y := 0; (\text{if } (a) \{ b := 1 \} \text{ else } \{ y := 1 \} \parallel \text{if } (y \vee b) \{ a := 1 \})$$

A.4 Thread inlining

One property one could ask of a model of shared memory is thread inlining: any execution of $\llbracket P; Q \rrbracket$ is an execution of $\llbracket P \parallel Q \rrbracket$. This is *not* a goal of our model, and indeed is not satisfied, due to the different semantics of concurrent and sequential memory accesses. We demonstrate this by considering an example from the Java Memory Model [?], which shows that Java does not satisfy thread inlining either.

The lack of thread inlining is related to the different dependency relations introduced by sequential and concurrent access. Recall from §3.1 that the program $(x := 0; y := x + 1;)$ has

execution:



but that $(x := 1; \parallel y := x+1;)$ has:

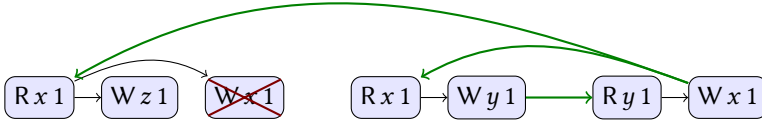


That is, in the sequential case there is no dependency from the write of x to the write of y , but in the concurrent case there is such a dependency.

This can be used to construct a counter-example to thread inlining, based on [?, Ex 11]:

$x := 0; \text{ if } (x == 1) \{ z := 1; \} \text{ else } \{ x := 1; \} \parallel y := x; \parallel x := y;$

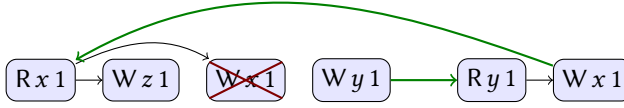
This has no execution containing $(Wz\ 1)$. Any attempt to build such an execution results in a cycle:



Inlining the thread $(y := x)$ gives [?, Ex 12]:

$x := 0; \text{ if } (x == 1) \{ z := 1; \} \text{ else } \{ x := 1; \} y := x; \parallel x := y;$

with execution:



To see why this execution exists, consider the program fragment:

$\text{if } (x == 1) \{ z := 1; \} \text{ else } \{ x := 1; \} y := x;$

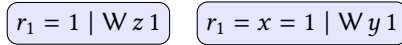
Removing the syntax sugar, this is:

```

r1 := x; if (r1 == 1) {
  z := 1; r2 := x; y := r2; skip
} else {
  x := 1; r3 := x; y := r3; skip
}

```

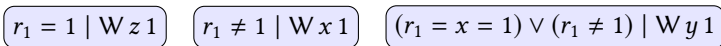
Now, $\llbracket z := 1; r_2 := x; y := r_2; \text{skip} \rrbracket$ includes pomset:



and $\llbracket x := 1; r_3 := x; y := r_3; \text{skip} \rrbracket$ includes pomset:



so $\llbracket \text{if } (r_1 = 1) \{ z := 1; r_2 := x; y := r_2; \text{skip} \} \text{ else } \{ x := 1; r_3 := x; y := r_3; \text{skip} \} \rrbracket$ includes:



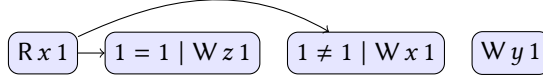
which means $\llbracket \text{if } (r_1 = 1) \{ z := 1; r_2 := x; y := r_2; \text{skip} \} \text{ else } \{ x := 1; r_3 := x; y := r_3; \text{skip} \} \rrbracket [x/r_1]$ includes:

$$(x = 1 \mid Wz1) \quad (x \neq 1 \mid Wx1) \quad ((x = x = 1) \vee (x \neq 1)) \mid Wy1$$

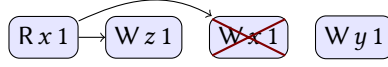
Now $(x = x = 1) \vee (x \neq 1)$ is a tautology, so this is just:

$$(x = 1 \mid Wz1) \quad (x \neq 1 \mid Wx1) \quad Wy1$$

and so $\llbracket r_1 := x; \text{if } (r_1 = 1) \{ z := 1; r_2 := x; y := r_2; \text{skip} \} \text{ else } \{ x := 1; r_3 := x; y := r_3; \text{skip} \} \rrbracket$ includes:



which simplifies to:



as required. The rest of the example is straightforward, and shows that our semantics agrees with the JMM in not supporting thread inlining.