

# AUTOMATICALLY ELIMINATING SPECULATIVE LEAKS WITH BLADE

ANONYMOUS AUTHOR(S)

**Abstract** We introduce BLADE, a new approach to automatically and efficiently synthesizing provably correct repairs for transient execution vulnerabilities like Spectre. BLADE is built on the insight that to stop speculative execution attacks, it suffices to *cut* the dataflow from expressions that speculatively introduce secrets (*sources*) to those that leak them through the cache (*sinks*), rather than prohibiting speculation altogether. We formalize this insight in a *static type system* that (1) types each expression as either *transient*, *i.e.*, possibly containing speculative secrets or as being *stable*, and (2) prohibits speculative leaks by requiring that all *sink* expressions are stable. We introduce `protect`, a new abstract primitive for fine grained speculation control that can be implemented via existing architectural mechanisms, and show how our type system can automatically synthesize a *minimal* number of `protect` calls needed to ensure the program is secure. We evaluate BLADE by using it to repair several verified, yet vulnerable WebAssembly implementations of cryptographic primitives. BLADE can fix existing programs that leak via speculation *automatically*, without user intervention, and *efficiently* using two orders of magnitude fewer fences than would be added by existing compilers, thereby ensuring security with minimal performance overhead.

## 1 Introduction

Implementing secure cryptographic algorithms is hard. The code must not only be functionally correct and memory safe, it must avoid divulging secrets indirectly through side channels like control-flow, memory-access patterns, or execution time. Consequently, much recent work focuses on how to ensure implementations do not leak secrets *e.g.*, via type systems [12, 39], verification [4], and program transformations [6].

Unfortunately, these efforts are foiled by speculative execution. Even if secrets are closely controlled via guards and access checks, the processor can simply ignore those checks when executing speculatively. An attacker can exploit this to leak secrets in turn.

In principle, memory fences block speculation, and hence, offer a way to recover the original security guarantees. In practice, however, fences pose a confounding dilemma. Programmers can either rely on heuristic approaches for inserting fences [37], but then forgo guarantees about the absence of side-channels. Alternatively, they can recover security guarantees by conservatively inserting fences after every load, but endure the huge performance costs.

In this paper, we introduce BLADE, a new approach to automatically, provably and efficiently eliminate speculation-based leakage. BLADE is based on the key insight that to prevent leaking data via speculative execution, it is unnecessary to stop *all* speculation as done by traditional memory fences. Instead, it suffices to *cut* the data flow from expressions (*sources*) that speculatively introduce secrets to those that leak them through the cache (*sinks*). We develop this insight into an automatic enforcement algorithm via four contributions.

**1. A Semantics for Speculation.** Our first contribution is a formal operational semantics for a simple While language that precisely captures how speculation can occur and what an attacker can observe via speculation (§ 3). To prevent leakage, we propose and formalize the semantics of an abstract primitive called `protect` that does not exist in today’s hardware but captures the essence of several primitives proposed in recent work [2, 32]. Furthermore, this primitive can be implemented in software *e.g.*, via *speculative load hardening* [30]. Crucially, and in contrast to a regular fence which stops *all* speculation, `protect` only stops speculation for a given *variable*. For example  $x := \text{protect}(e)$  ensures that  $e$ ’s value is only assigned to  $x$  *after*  $e$  has been assigned its *stable*, non-speculative value.

**2. A Type System for Speculation.** Our second contribution is an approach to conservatively approximating the dynamic semantics of speculation via a *static type system* that types each expression as either *transient* (**T**), *i.e.*, expressions that may contain speculative secrets, or *stable* (**S**), *i.e.*, those that cannot (§ 4.1). Our system prohibits speculative leaks by requiring that all *sink* expressions that can influence intrinsic attacker visible behavior (*e.g.*, cache addresses) are typed as *stable*. We connect the static and dynamic semantics by proving that well-typed programs are indeed secure, *i.e.*, satisfy a correctness condition called *speculative non-interference* [17] which states that the program does not leak under speculative execution more than it would under sequential execution.

**3. Automatic Protection.** Existing programs that are free of `protect` statements are likely insecure under speculation and will be rejected by our type system. Thus, our third contribution is an algorithm that automatically synthesizes a *minimal* number of `protect` statements to ensure that the program satisfies speculative non-interference. To this end, we extend the type checker to construct a *def-use graph* that captures the data-flow between program expressions. A *cut-set* in the graph is a set of variables whose removal eliminates

```

1 1 void SHA2_update_last(int *input_len, ...)
2 2 {
3 3     if (! valid(input_len)) { ... }
4 4     int len = *input_len;
5 5     int *dst3 = ... + len;
6 6     _mm_lfence();
7 7     int *dst3_safe = protect(.. + len);
8 8     ...
9 9     *dst3_safe = pad;
10 10    ...
11 11 }

```

**Figure 1.** Code fragment from the HACL\* SHA2 implementation, containing a potential speculative execution vulnerability that leaks *explicitly* through the cache by writing memory at a secret-tainted address (line 9). A naive patch is shown in **red**, the patch computed by BLADE is shown in **green**.

all paths from secret-sources to observable-sinks. We show that inserting a `protect` statement for each variable in a cut-set suffices to yield a program that is well-typed, and hence, secure with respect to speculation (§5.3). Happily, finding such cuts is an instance of the classic max-flow/min-cut problem, so existing polynomial time algorithms let us efficiently synthesize `protect` statements that resolve the dilemma of enforcing security with minimal performance overhead.

**4. Evaluation.** Our final contribution is an implementation of our method in a tool called BLADE, and an evaluation using BLADE to repair verified yet vulnerable (to transient execution attacks) programs: the WebAssembly implementations of the signal messaging Protocol and its respective cryptographic libraries [29], and a number of verified cryptographic algorithms from [38] (§ 6). Our evaluation shows that BLADE can automatically compute fixes for existing programs. Compared to an existing fully automatic protection as implemented in existing compilers (notably Clang), BLADE inserts two orders of magnitude fewer fences and thus imposes negligible performance overhead.

## 2 Overview

In this section, we present two potential speculative execution vulnerabilities in HACL\*— a verified cryptographic library — that were discovered by BLADE and discuss how BLADE repairs the vulnerabilities by inserting `protect` statements. We then show how BLADE computes the repairs via our minimal fence inference algorithm and finally how BLADE proves that the repairs are indeed correct, via our transient-flow type system.

### 2.1 Two Speculation Bugs and Their Fixes

Figure 1 shows a code fragment from a function in the implementation of the SHA2 hash in HACL\*. Though BLADE operates on WebAssembly, we present equivalent simplified C code for readability. The function takes as input a pointer `input_len`, validates the input (line 3), loads from memory the public length of the hash (line 4), calculates a target address `dst3` (line 5), and finally pads the buffer pointed to by `dst3` (line 9).

**1. Leaking Through a Memory Write.** During normal, sequential execution this code is not a problem: the function validates the input to prevent classic buffer overflows vulnerabilities. However, an attacker can exploit the function to leak sensitive data during speculation. To do this, the attacker first has to modify the value that the pointer `input_len` holds during speculation. Since `input_len` is a function parameter, this can be achieved *e.g.*, by calling the function repeatedly with legitimate addresses, training the branch predictor to predict the next input to be valid. After (mis)training the branch predictor, the attacker manipulates `input_len` to point to an address containing secret data (*e.g.*, the secret key used by the hash function) and calls the function again, this time with an invalid pointer. As a result of the mistraining, the branch predictor causes the processor to skip validation and erroneously load the secret into `len`, which in turn, is used to calculate pointer `dst3`. The buffer pointed to by `dst3` is then written in line 9, completing the attack. Even though pointer `dst3` is incorrect due to misprediction and the write will therefore be squashed, its side-effects persist, and therefore remain visible to the attacker. The attacker can then extract the target address — and thereby the secret via cache timing measurements [16].

**Preventing the Attack: Memory Fences.** Since the attack exploits the fact that input validation is speculatively skipped, we can prevent it by making sure that the buffer in line 9 is not written until the input has been validated. To mitigate these class of attacks, Intel [19] and AMD [5] recommend inserting a speculation barrier after critical validation check-points. Following this strategy, we would place a *memory fence* on line 6. This fence stops all speculative execution past the fence, *i.e.*, no statements after the fence are executed until all previous statements (including input validation) have been resolved. While the effects of the fence prevent the attack, they are more restrictive than necessary and incur high performance cost [33].

**Preventing the Attack Efficiently.** We propose an alternative way to stop speculation from reaching the write in line 9 through a new primitive called `protect`. Rather than eliminate *all* speculation, `protect` only stops speculation along *a particular data-path*. We use `protect` to patch the program in line 7. Instead of assigning pointer `dst3` directly as in line 5, the expression that computes the address is guarded by a `protect` statement. This ensures that the value assigned

```

111 1 void SHA2_update_last(int *input_len, ...)
112 2 {
113 3     if (! valid(input_len)) { ... }
114 4     int len = *(input_len);
115 5     ...
116 6     int len_safe = protect(*input_len)
117 7     for ( i = 0; i < len_safe + ...)
118 8         dst2[i] = 0;
119 9     ...
120 10 }

```

**Figure 2.** SHA2 code fragment containing a potential speculative execution vulnerability that leaks *implicitly* through a control-flow dependency.

to `dst3_safe` is always guaranteed to use `len`'s final, non-speculative value. Therefore, writing to `dst3_safe` in line 9 prevents any invalid secret-tainted address from speculatively reaching the store, where it could be leaked to the attacker.

The `protect` primitive offers an abstract interface for fine grained control of speculation. There are a number of possible implementations for this interface. For example, `protect` could be implemented in hardware. While unfortunately, today's hardware does not offer an equivalent instruction to `protect`, similar functionalities have been proposed in recent work [2, 32]. Alternatively, `protect` can be implemented in software (a similar proposal has been made in [30]). In general, `protect` can be implemented through a fence instruction. However, better solutions exist for reading arrays. For example, Speculative Load Hardening (SLH), a mitigation deployed in the code generated by Clang [10], stalls individual array reads until the corresponding bounds-check condition gets resolved. We model software implementations of `protect` through a restricted primitive called `safe_read`, which can only be applied to array reads. We then formalize an implementation of `safe_read` via SLH in the supplementary material, and evaluate the number of `protect` and `safe_read` needed to patch HACl\* and their overhead in Section 6.

**2. Leaking Through a Control-Flow Dependency.** Figure 2 shows a code fragment taken from the same function as in Figure 1. The code contains a second potential vulnerability, but in contrast to Figure 1 the vulnerability leaks secrets *implicitly*, through a control-flow dependency.

The function reads from memory a (public) integer `len` (line 4), which determines the number of initialization rounds in the condition of the for-loop (line 7). Like the previous vulnerability, the function is harmless under sequential execution, but leaks under speculation. As before, the attacker manipulates the pointer `input_len` to point to a secret after mistraining the branch predictor to skip validation. But instead of leaking the secret directly through the data cache, they

can leak the value indirectly through a control-flow dependency, e.g., via the instruction cache and non-secret dependent lines of the data cache. In particular, the secret determines how often the initialization loop (line 7) is executed during speculation, and therefore an attacker can make secret dependent observations via instruction- and data-cache timing attacks. Like the previous vulnerability, this vulnerability can be fixed via the `protect` primitive, as shown in lines 6 and 7.

## 2.2 Computing Fixes Via Minimal Fence Inference

BLADE automatically infers the placement of these `protect` statements. We illustrate this process using a simple running example Ex1 shown in Figure 3. The code reads two values from an array ( $x := a[i_1]$  and  $y := a[i_2]$ ), adds them ( $z := x + y$ ), and indexes another array with the result ( $w := b[z]$ ). We assume that all array operations are implicitly bounds-checked and thus no explicit validation code is needed.

Like the examples above, Ex1 contains a speculative execution vulnerability: the array reads may skip their bounds check and so  $x$  and  $y$  can contain transient secrets (i.e., secrets introduced by misspeculation). This secret data then flows to  $z$ , and finally leaks through the data cache by the array read  $b[z]$ .

**Def-Use Graph.** To secure the program, we need to *cut the dataflow* between the array reads which could introduce *transient* secret values into the program, and the index in the array read where they are leaked through the cache. For this, we first build a def-use graph whose nodes and directed edges capture the data dependencies between the expressions and variables of a program. For example, consider the def-use graph of program Ex1 in Figure 4. In the graph, the edge  $x \rightarrow x + y$  indicates that  $x$  is used to compute  $x + y$ .<sup>1</sup> To track how transient values propagate in the def-use graph, we extend the graph with the special circle node **T**, which represents the *source* of *transient* values of the program. Since reading memory creates transient values, we connect the **T** node to all nodes containing expressions that explicitly read memory, e.g.,  $\mathbf{T} \rightarrow a[i_1]$ . Following the data dependencies along the edges of the def-use graph, we can see that node **T** is transitively connected to node  $z$ , which indicates that  $z$  can contain transient data at run-time. To detect insecure uses of transient values, we then extend the graph with the special circle node **S**, which represents the *sink* of *stable* (i.e., non-transient) values of a program. Intuitively, this node draws all the values of a program that *must* be stable to avoid transient execution attacks. Therefore, we connect all expression used as array indices in the program to the **S** node, e.g.,  $z \rightarrow \mathbf{S}$ . The fact that the graph in Figure 3 contains a *path* from **T** to **S** indicates that transient data flows through data dependencies into (what should be) a stable index expression and thus the program is insecure.

**Cutting the Dataflow.** In order to make the program safe, we need to *cut* the data-flow between **T** and **S** by introducing

<sup>1</sup>To avoid ambiguities in the graph, we assume that each variable is assigned at most *once*, i.e., the code is in static single assignment form.

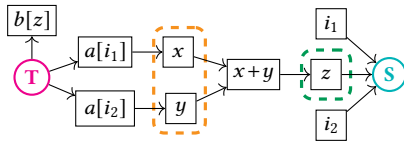


```

221   x := a[i1]           x := protect(a[i1])
222   y := a[i2]           y := protect(a[i2])
223   z := x + y           z := protect(x + y)
224   w := b[z]
225

```

**Figure 3.** Ex1: Running Example. The optimal patch computed by BLADE is shown in **green**. A sub-optimal patch is shown in **orange**.



**Figure 4.** Def-use Graph of Ex1. We omit some irrelevant edges for readability. The Figure contains two choices of cut-sets, shown as dashed lines. The left cut requires removing two nodes and thus, inserting two `protect` statements. The right cut shows a minimal solution, which only requires removing a single node.

as few `protect` statements as necessary. This problem can be equivalently restated as follows: find a *minimal cut-set*, i.e., a minimal set of variables, such that removing the variables from the graph eliminates all paths from **T** from **S**. Each choice of cut-set defines a way to repair the program: simply add a `protect` statement for each variable in the set. Figure 4 contains two choices of cut-sets, shown as dotted lines. The cut-set on the left requires two `protect` statements, for variables  $x$  and  $y$  respectively, corresponding to the **orange** patch in Figure 3. The cut-set on the right is minimal, it requires only a single `protect`, for variable  $z$ , and corresponds to the **green** patch in Figure 3. In general, the a minimal cut-set can be computed as a solution to the Min-Cut/Max-Flow problem, for which efficient polynomial-time algorithms exist [1].

### 2.3 Proving Correctness Via Transient-Flow Types

To formalize and verify the correctness of the patch computed by cutting the def-use graph, we define a transient-flow type system and construct the def-use graph for a given program from the type-constraints generated during type inference.

**Typing Judgement.** The type system statically assigns a transient-flow type to each variable: a variable is typed as *transient* (written as **T**), if it can contain transient data (i.e., potential secrets) at run-time, and as *stable* (written as **S**), otherwise. Given a typing environment  $\Gamma$  which assigns a transient flow type to each variable, and a command  $c$ , the type system defines a judgement  $\Gamma \vdash c$  saying that  $c$  is free of speculative execution bugs. The type system enforces that transient expressions may not be used in positions that may leak their value by affecting memory reads and writes, e.g.,

they may not be used as array indices and in loop conditions. Additionally, it requires that transient expressions may not be assigned to stable variables, except through the use of `protect`. To show that our type system indeed prevents speculative execution attacks, we define a semantics for speculative execution of a while language (Section 3) and prove that well-typed programs do not leak speculatively more than *sequentially*, that is by executing their statements in-order and without speculation (see Section 5).

**Type Inference.** Given an input program, we construct the corresponding def-use graph by collecting the type constraints generated during type inference. Type inference is formalized by a typing-inference judgement  $\Gamma, \text{Prot} \vdash c \Rightarrow k$ , which extends the typing judgement from above with (1) a set of protected variables  $\text{Prot}$  (the *cut-set*), and (2) a set of type-constraints  $k$  (the *def-use graph*). At a high level, type inference has 3 steps: (i) generate a set of constraints under an initial typing environment and protected set that allow any program to type-check, (ii) construct the def-use graph from the constraints and find a cut-set, and (iii) compute the resulting typing environment. To characterize the security of a still *unrepaired* program after type inference, we define a typing judgement  $\Gamma, \text{Prot} \vdash c$ , where unprotected variables are explicitly accounted for in the  $\text{Prot}$  set.<sup>2</sup> Intuitively, the program is secure if we promise to insert a `protect` statement for each variable in  $\text{Prot}$ .

To repair programs, we simply honor the promise of inserting `protect` statements for each for each variable in the protected set of the typing judgement obtained above. Once repaired, the program type checks under an empty protected set and with the same typing environment.

### 2.4 Attacker Model

Before moving to the details of our semantics and transient type system, we discuss the attacker model considered in this work. The attacker runs cryptographic code on a speculative out-of-order processor and, as usual, can choose the values of public inputs and observe public outputs, but may not read secret data (e.g., cryptographic keys) in registers and memory. Additionally, the attacker can influence how programs are speculatively executed through the branch predictor and choose the instructions execution order in the processor pipeline. The effects of these actions are observable through the cache and are otherwise invisible at the ISA level. In particular, while programs run, the attacker can take precise timing measurements through the data- and instruction-cache with a cache-line granularity, which may disclose secret data covertly. These features allow the attacker to mount Spectre-PHT [20, 21] and Spectre-STL [9] attacks and leak data through FLUSH+RELOAD [43] and PRIME+PROBE [34] cache side-channels attacks. We do not consider speculative attacks that rely on the Return Stack Buffer (e.g., Ret2Spec [25])

<sup>2</sup>The judgement  $\Gamma \vdash c$  is just a short-hand for  $\Gamma, \emptyset \vdash c$ .

```

331 Value  $v ::= n \mid b \mid a$ 
332 Expr.  $e ::= v \mid x \mid e_1 + e_2 \mid e_1 \leq e_2$ 
333          $\mid \text{length}(e) \mid \text{base}(e)$ 
334 Rhs.  $r ::= e \mid *e \mid e[e]$ 
335 Cmd.  $c ::= \text{skip} \mid x := r \mid *e = e \mid e_1[e_2] := e_3$ 
336          $\mid \text{if } e \text{ then } c_1 \text{ else } c_2$ 
337          $\mid \text{while } e \text{ do } c \mid \text{fail} \mid c_1; c_2$ 
338          $\mid x := \text{stable\_read}(e_1, e_2)$ 
339          $\mid x := \text{protect}(r)$ 

```

Figure 5. Surface Syntax.

and [22]) or the Branch Target Buffer (Spectre-BTB [21]). We similarly do not consider attacks that do not use the cache to exfiltrate data, e.g., port contention (SMoTherSpectre [7]) and Meltdown attacks [9, 24], since hardware fixes address them.

### 3 A Semantics for Speculation

We now formalize the concepts presented in the overview. We start by giving a formal semantics for a while language with speculative execution. Figure 5 presents the language's surface syntax. Values consist of Booleans  $b$ , pointers  $n$  represented as natural numbers, and arrays  $a$ . Array length and base address are given by functions  $\text{length}(\cdot)$  and  $\text{base}(\cdot)$ . In addition to variable assignments, pointer dereferences, array stores, conditionals and loops, our language features two special commands that help prevent transient execution attacks. Command  $x := \text{protect}(r)$  evaluates  $r$  and assigns its value to  $x$ , only after the value is *stable* (i.e., non-transient). Command  $x := \text{stable\_read}(e_1, e_2)$  is a restricted version of  $\text{protect}(\cdot)$  that only applies to array reads (see Section 3.4) Lastly, **fail** triggers a memory violation error (caused by reading or writing an array out-of-bounds) and aborts the program.

**Processor Instructions.** Our semantics translates the surface syntax into an abstract set of processor instructions shown in Figure 6. Our processor instructions do not contain branching, they represent a *single* predicted path through the control flow. The prediction choices are represented by a sequence of *guard* instructions representing pending branch points. Guard instructions have form  $\text{guard}(e^b, cs, p)$ , which records the branch condition  $e$ , its predicted truth value  $b$  and a unique guard identifier  $p$ , used in our security analysis (Section 5). Each guard attests the fact that the current execution is valid only if the branch condition gets resolved as predicted. In order to enable a roll-back in case of a missprediction, guards additionally record the set of commands  $cs$  along the alternative branch.

**Directives and Observations.** Instructions do not have to be executed in sequence, they can be executed in any order, enabling out-of-order execution. We use a simple three stage processor pipeline: the execution of each instruction is split into **fetch**, **exec**, and **retire**. We do not fix the order in which

```

386 Instr.  $i ::= \text{nop} \mid x := e \mid x := \text{load}(e)$ 
387          $\mid \text{store}(e_1, e_2) \mid x := \text{protect}(e)$ 
388          $\mid \text{guard}(e^b, cs, p) \mid \text{fail}$ 
389 Dir.  $d ::= \text{fetch} \mid \text{fetch } b \mid \text{exec } n$ 
390          $\mid \text{retire}$ 
391 Obs.  $o ::= \epsilon \mid \text{load}(n, ps) \mid \text{store}(n, ps)$ 
392          $\mid \text{fail} \mid \text{rollback}(p)$ 
393 Prediction  $b \in \{\text{true}, \text{false}\}$ 
394 Guard Id.  $p \in \mathbb{N}$ 
395 Reorder Buffer  $is ::= i:is \mid []$ 
396 Cmd Stack  $cs ::= c:cs \mid []$ 
397 Memory Store  $\mu \in \mathbb{N} \rightarrow \text{Value}$ 
398 Variables Map  $\rho \in \text{Var} \rightarrow \text{Value}$ 
399 Configuration  $C ::= \langle is, cs, \mu, \rho \rangle$ 

```

Figure 6. Processor Syntax.

instructions, and their individual stages are executed, nor do we supply a model of the branch predictor to decide which control flow path to follow. Instead, we let the attacker supply those decisions through a set of *directives* [11] shown in Fig. 6. For example, directive **fetch true** fetches the **true** branch of a conditional and **exec  $n$**  executes the  $n$ th instruction in the reorder buffer. Executing an instruction generates an *observation* (Fig. 6) which records attacker observable behavior. Observations include *speculative* memory reads and writes (i.e.,  $\text{load}(n, ps)$  and  $\text{store}(n, ps)$  issued while guards  $ps$  are pending), rollbacks (i.e.,  $\text{rollback}(p)$  due to misspeculation of guard  $p$ ), and memory violations (**fail**). Most instructions generate the empty observation  $\epsilon$ .

**Configurations and Reduction Relation.** We formally specify our semantics as a reduction relation between processor configurations. A configuration  $\langle is, cs, \mu, \rho \rangle$  consists of a queue of in-flight instructions  $is$  called the reorder buffer, a stack of commands  $cs$ , a memory  $\mu$ , and map from variables to values  $\rho$ . A reduction step  $C \xrightarrow{d}_o C'$  denotes that, under directive  $d$ , configuration  $C$  is transformed into  $C'$  and generates observation  $o$ . To execute a program  $c$  with initial memory  $\mu$  and variable map  $\rho$ , the processor initializes the configuration with an empty reorder buffer and inserts the program into the command stack, i.e.,  $\langle [], [c], \mu, \rho \rangle$ . Then, the execution proceeds until both the reorder buffer and the stack in the configuration are empty, i.e., we reach a configuration of the form  $\langle [], [], \mu', \rho' \rangle$ , for some final memory store  $\mu'$  and variable map  $\rho'$ .

We now discuss the semantics rules of each execution stage and then those for our security primitives.

#### 3.1 Fetch Stage

The fetch stage flattens the input command into a sequence of instructions which it stores in the reorder buffer. Figure 7 presents selected rules; the remaining rules are in Appendix A. Rule [FETCH-SEQ] pops command  $c_1; c_2$  from the commands

$$\begin{array}{l}
\text{FETCH-SEQ} \\
\langle is, (c_1; c_2) : cs, \mu, \rho \rangle \xrightarrow{\text{fetch}}_{\epsilon} \langle is, c_1 : c_2 : cs, \mu, \rho \rangle \\
\text{FETCH-ASGN} \\
\langle is, x := e : cs, \mu, \rho \rangle \xrightarrow{\text{fetch}}_{\epsilon} \langle is ++ [x := e], cs, \mu, \rho \rangle \\
\text{FETCH-PTR-LOAD} \\
\langle is, x := *e : cs, \mu, \rho \rangle \xrightarrow{\text{fetch}}_{\epsilon} \langle is ++ [x := \text{load}(e)], cs, \mu, \rho \rangle \\
\text{FETCH-ARRAY-LOAD} \\
\frac{c = x := e_1[e_2] \quad e = e_2 < \text{length}(e_1) \quad \text{fresh}(p) \quad e' = \text{base}(e_1) + e_2 \quad c' = \text{if } e \text{ then } x := *e' \text{ else fail}}{\langle is, c : cs, \mu, \rho \rangle \xrightarrow{\text{fetch}}_{\epsilon} \langle is, c' : cs, \mu, \rho \rangle} \\
\text{FETCH-IF-TRUE} \\
\frac{c = \text{if } e \text{ then } c_1 \text{ else } c_2 \quad \text{fresh}(p) \quad i = \text{guard}(e^{\text{true}}, c_2 : cs, p)}{\langle is, c : cs, \mu, \rho \rangle \xrightarrow{\text{fetch true}}_{\epsilon} \langle is ++ [i], c_1 : cs, \mu, \rho \rangle}
\end{array}$$

Figure 7. Fetch stage (selected rules).

stack and pushes the two sub-commands for further processing. [FETCH-ASGN] pops an assignment from the commands stack and appends the corresponding processor instruction ( $x := e$ ) at the end of the reorder buffer.<sup>3</sup> Rule [FETCH-PTR-LOAD] is similar and simply translates pointer dereferences to the corresponding load instruction. Arrays provide a memory-safe interface to read and write memory: the processor injects bounds-checks when fetching commands that read and write arrays. For example, rule [FETCH-LOAD-TRUE] expands command  $x := e_1[e_2]$  into the corresponding pointer dereference, but guards the command with a bounds-check condition. First, the rule generates the condition  $e = e_2 < \text{length}(e_1)$  and calculates the address of the indexed element  $e' = \text{base}(e_1) + e_2$ . Then, it replaces the array read on the stack with command **if  $e$  then  $x := *e'$  else fail** to abort the program and prevent the buffer overrun if the bounds check fails. Later, we show that speculative out-of-order execution can simply ignore the bounds check guard and cause the processor to transiently read memory at an invalid address. Rule [FETCH-IF-TRUE] fetches a conditional branch from the stack and, following the prediction provided in directive **fetch true**, speculates that the condition  $e$  will evaluate to **true**. Thus, the processor inserts the corresponding instruction **guard**( $e^{\text{true}}, c_2 : cs, p$ ) with a fresh guard identifier  $p$  in the reorder buffer and pushes the then-branch  $c_1$  onto the stack  $cs$ . Importantly, the guard instruction stores the else-branch together with a copy of

<sup>3</sup>Notation  $[i_1, \dots, i_n]$  represents a list of  $n$  elements,  $is_1 ++ is_2$  denotes list concatenation, and  $|is|$  computes the length of the list  $is$ .

$$\begin{array}{l}
\text{EXECUTE} \\
\frac{\rho' = \phi(is_1, \rho) \quad |is_1| = n-1 \quad \langle is_1, i, is_2, cs \rangle \xrightarrow{(\mu, \rho', o)} \langle is', cs' \rangle}{\langle is_1 ++ [i] ++ is_2, cs, \mu, \rho \rangle \xrightarrow{\text{exec } n}_o \langle is', cs', \mu, \rho \rangle} \\
\text{EXEC-ASGN} \\
\frac{i = (x := e) \quad v = \llbracket e \rrbracket^{\rho} \quad i' = (x := v)}{\langle is_1, i, is_2, cs \rangle \xrightarrow{(\mu, \rho, \epsilon)} \langle is_1 ++ [i'] ++ is_2, cs \rangle} \\
\text{EXEC-BRANCH-OK} \\
\frac{i = \text{guard}(e^b, cs', p) \quad \llbracket e \rrbracket^{\rho} = b}{\langle is_1, i, is_2, cs \rangle \xrightarrow{(\mu, \rho, \epsilon)} \langle is_1 ++ [\text{nop}] ++ is_2, cs \rangle} \\
\text{EXEC-BRANCH-MISPREDICT} \\
\frac{i = \text{guard}(e^b, cs', p) \quad \llbracket e \rrbracket^{\rho} \neq b}{\langle is_1, i, is_2, cs \rangle \xrightarrow{(\mu, \rho, \text{rollback}(p))} \langle is_1, cs' \rangle} \\
\text{EXEC-LOAD} \\
\frac{i = (x := \text{load}(e)) \quad \text{store}(-, -) \notin is_1 \quad n = \llbracket e \rrbracket^{\rho} \quad ps = (|is_1|) \quad i' = (x := \mu(n))}{\langle is_1, i, is_2, cs \rangle \xrightarrow{(\mu, \rho, \text{read}(n, ps))} \langle is_1 ++ [i'] ++ is_2, cs \rangle}
\end{array}$$

Figure 8. Execute stage (selected rules).

the current commands stack (*i.e.*,  $c_2 : cs$ ) as a rollback stack to restart the execution in case of misprediction.

### 3.2 Execute Stage

In the execute stage, the processor evaluates the operands of instructions in the reorder buffer and rolls back the program state whenever it detects a misprediction.

**Transient Variable Map.** To evaluate operands in the presence of out-of-order execution, we need to take into account how previous, possibly unresolved assignments in the reorder buffer affect the variable map. In particular, we need to ensure that an instruction cannot execute if it depends on a preceding assignment whose value is still unknown. To update variable map  $\rho$  with the pending assignments in reorder buffer  $is$ , we define a function  $\phi(is, \rho)$ , called the *transient variable map*. The function walks through the reorder buffer, registers each resolved assignment instruction ( $x := v$ ) in the variable map (through function update  $\rho[x \mapsto v]$ ) and marks variables from pending assignments (*i.e.*,  $x := e$ ,  $x := \text{load}(e)$ , and  $x := \text{protect}(r)$ ) as *undefined* ( $\rho[x \mapsto \perp]$ ), making their respective values unavailable to following instructions.

**Execute Rule and Auxiliary Relation.** Step rules for the reduction relation are shown in Figure 8. Rule [EXECUTE] executes the  $n$ -th instruction in the reorder buffer, following the directive **exec  $n$** . For this, the rule splits the reorder buffer into

prefix  $is_1$ ,  $n$ -th instruction  $i$  and suffix  $is_2$ . Next, it computes the transient variable map  $\phi(is_1, \rho)$  and executes a transition step under the new map using an auxiliary relation  $\rightsquigarrow$ . Notice that [EXECUTE] does not update the store or the variable map (the transient map is simply discarded). These changes are performed later in the retire stage.

The rules for the auxiliary relation are shown in Fig. 8. The relation transforms a tuple  $\langle is_1, i, is_2, cs \rangle$  consisting of prefix, suffix and current instruction  $i$  into a tuple  $\langle is', cs' \rangle$  specifying the reorder buffer and command stack obtained by executing  $i$ . For example, rule [EXEC-ASGN] evaluates the right-hand side of the assignment  $x := e$  where  $\llbracket e \rrbracket^\rho$  denotes the value of  $e$  under  $\rho$ . The premise  $v = \llbracket e \rrbracket^\rho$  ensures that the expression is defined *i.e.*, it does not evaluate to  $\perp$ . Then, the rule substitutes the computed value into the assignment ( $x := v$ ), and reinserts the instruction back into its original position in the reorder buffer.

**Guards and Rollback.** Rules [EXEC-BRANCH-OK] and [EXEC-BRANCH-MISPREDICT] resolve guard instructions. In rule [EXEC-BRANCH-OK], the predicted and computed value of the guard expression match, and the processor only has to replace the guard with a **nop**. In contrast, in rule [EXEC-BRANCH-MISPREDICT] the predicted and computed value differ ( $\llbracket e \rrbracket^\rho \neq b$ ). This causes the processor to revert the program state and issue a rollback observation. For the rollback, the processor discards the instructions *past* the guard (*i.e.*,  $is_2$ ) and substitutes the current commands stack  $cs$  with the rollback stack  $cs'$  which causes execution to revert to the alternative branch.

**Loads.** Rule [EXEC-LOAD] executes a memory load. The rule computes the address ( $n = \llbracket e \rrbracket^\rho$ ), retrieves the value at that address from memory ( $\mu(n)$ ) and rewrites the load into an assignment ( $x := \mu(n)$ ). Inserting the assignment into the reorder buffer allows transiently forwarding the loaded value to later instructions. The premise  $\text{store}(\_, \_) \notin is_1$  prevents the processor from reading stale data from memory: if the load aliases with a preceding (but pending) store, ignoring the store would produce a stale read. To record that the load is issued *speculatively*, the observation  $\text{read}(n, ps)$  stores list  $ps$  containing the identifiers of the guards still pending in the reorder buffer. Function  $\langle is \rangle$  simply extracts the identifiers of the guard instructions in the buffer  $is$ .

### 3.3 Retire Stage

The retire stage removes completed instructions from the reorder buffer and propagates their changes to variable map and memory store. While instructions are executed out-of-order, they are retired in-order to preserve the illusion of sequential execution to the user. Figure 9 presents the rules for the retire stage. Rule [RETIRE-NOP] removes **nop**. Rules [RETIRE-ASGN] and [RETIRE-STORE] remove the resolved assignment  $x := v$  and instruction  $\text{store}(n, v)$  from the reorder buffer and update the variable map ( $\rho[x \mapsto v]$ ) and the memory store ( $\mu[n \mapsto v]$ )

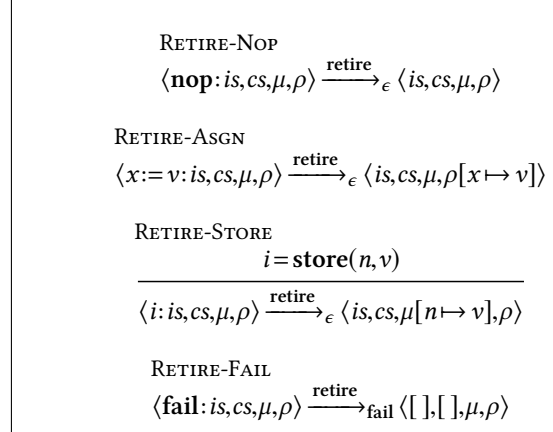


Figure 9. Retire stage.

respectively. Rule [RETIRE-FAIL] aborts the program by emptying reorder buffer and command stack and generates a **fail** observation, simulating a processor raising an exception (e.g., a page fault).

We demonstrate how the attacker can leak a secret from program Ex1 (Fig. 3) in our model. First, the attacker instructs the processor to fetch all the instructions, supplying prediction **true** for all bounds-check conditions. Figure 10 shows the resulting buffer and how it evolves after each attacker directive, which instruct the processor to speculatively execute the load instructions and the assignment (but not the guard instructions). Memory  $\mu$  and variable map  $\rho$  are shown on the right. Directive **exec 4** transiently reads array  $a$  past its bound, at index 2, reading into the memory ( $\mu(3) = 42$ ) of secret array  $s[0]$  and generates the corresponding observation. Finally, the processor forwards the values of  $x$  and  $y$  to compute their sum in the fifth instruction, ( $z := 42$ ), which is then used as an index in the last instruction and leaked to the attacker via observation  $\text{read}(42, [1, 2, 3])$ .

### 3.4 Security Primitives

Next, we turn to the rules describing our security primitives.

**Protect.** Instruction  $x := \text{protect}(r)$  assigns the value of  $r$ , only after all previous **guard** instructions have been executed, *i.e.*, when the value has become stable and no more rollbacks are possible. Figure 11 formalizes this intuition. Rule [FETCH-PROTECT-EXPR] fetches protect commands involving simple expressions ( $x := \text{protect}(e)$ ) and inserts the corresponding protect instruction in the reorder buffer. Rule [FETCH-PROTECT-ARRAY] piggy-backs on the previous rule by splitting a protect of an array read ( $x := \text{protect}(e_1[e_2])$ ) into a separate assignment of the array value ( $x := e_1[e_2]$ ) and protect of the variable ( $x := \text{protect}(x)$ ). Rules [EXEC-PROTECT<sub>1</sub>] and [EXEC-PROTECT<sub>2</sub>] extend auxiliary relation  $\rightsquigarrow$ . Rule [EXEC-PROTECT<sub>1</sub>] evaluates the expression ( $v = \llbracket e \rrbracket^\rho$ )



Reorder Buffer		exec 2	exec 4	exec 5	exec 7	Memory Layout
1	$\text{guard}((i_1 < \text{length}(a))^{\text{true}}, [\text{fail}], 1)$					$\mu(0)=0$   $b[0]$
2	$x := \text{load}(\text{base}(a) + i_1)$	$x := \mu(2)$				$\mu(1)=0$   $a[0]$
3	$\text{guard}((i_2 < \text{length}(a))^{\text{true}}, [\text{fail}], 2)$					$\mu(2)=0$   $a[1]$
4	$y := \text{load}(\text{base}(a) + i_2)$		$y := \mu(3)$			$\mu(3)=42$   $s[0]$
5	$z := x + y$			$z := 42$		...
6	$\text{guard}((z < \text{length}(b))^{\text{true}}, [\text{fail}], 3)$					...
7	$w := \text{load}(\text{base}(b) + z)$				$w := \mu(42)$	<b>Variable Map</b>
Observations:		$\text{read}(2, [1])$	$\text{read}(3, [1, 2])$	$\epsilon$	$\text{read}(42, [1, 2, 3])$	$\rho(i_1) = 1$
						$\rho(i_1) = 2$

Figure 10. Leaking execution of running example Ex1.

FETCH-PROTECT-ARRAY	
$c = (x := \text{protect}(e_1[e_2]))$	
$c_1 = (x := e_1[e_2])$	$c_2 = (x := \text{protect}(x))$
$\langle is, c, cs, \mu, \rho \rangle \xrightarrow{\text{fetch}}_{\epsilon} \langle is, c_1 : c_2 : cs, \mu, \rho \rangle$	
FETCH-PROTECT-EXPR	
$c = (x := \text{protect}(e)) \quad i = (x := \text{protect}(e))$	
$\langle is, c, cs, \mu, \rho \rangle \xrightarrow{\text{fetch}}_{\epsilon} \langle is ++ [i], cs, \mu, \rho \rangle$	
EXEC-PROTECT <sub>1</sub>	
$i = (x := \text{protect}(e)) \quad v = \llbracket e \rrbracket^{\rho} \quad i' = (x := \text{protect}(v))$	
$\langle is_1, i, is_2, cs \rangle \xrightarrow{(\mu, \rho, \epsilon)} \langle is_1 ++ [i'] ++ is_2, cs \rangle$	
EXEC-PROTECT <sub>2</sub>	
$i = (x := \text{protect}(v)) \quad \text{guard}(\_, \_, \_) \notin is_1 \quad i' = (x := v)$	
$\langle is_1, i, is_2, cs \rangle \xrightarrow{(\mu, \rho, \epsilon)} \langle is_1 ++ [i'] ++ is_2, cs \rangle$	

Figure 11. Semantics of  $\text{protect}(\cdot)$  (selected rules).

and reinserts the instruction in the reorder buffer as if it were a normal assignment. However, the processor leaves the value wrapped inside the  $\text{protect}$  instruction in the reorder buffer, i.e.,  $x := \text{protect}(v)$ , to prevent forwarding the value to the later instructions via the transient variable map. When no guards are pending in the reorder buffer ( $\text{guard}(\_, \_, \_) \notin is_1$ ), rule [EXEC-PROTECT<sub>2</sub>] transforms the instruction into a normal assignment, so that the processor can propagate and commit its value.

**Example.** Consider again Ex1 and the execution shown in Figure 10. In the repaired program,  $x + y$  is wrapped in a  $\text{protect}$  statement. As a result, directive  $\text{exec 5}$  produces value  $z := \text{protect}(42)$ , instead of  $z := 42$  which prevents instruction 7 from executing (as its target address is undefined), until all guards are resolved. This in turn prevents the leaking of the transient value.

**Stable Read.** Unfortunately, current processors do not provide the means to implement  $\text{protect}$  in its full generality. Our

semantics therefore contains a primitive  $\text{stable\_read}(e_1, e_2)$  that implements a restricted version of  $\text{protect}(e_1[e_2])$  for array reads. While  $\text{protect}(\cdot)$  prevents forwarding loaded values until all pending branches are resolved,  $\text{stable\_read}(\cdot)$  stalls memory loads until individual bounds-check conditions have been resolved.  $\text{stable\_read}(\cdot)$  can be implemented using today's hardware, for example through speculative Load Hardening (SLH) [10], the spectre mitigation proposed by and deployed in the Clang compiler. We provide formal semantics in Appendix B.

**Example.** Consider again Ex1. Instead of using  $\text{protect}(\cdot)$ , we can repair the example by inserting  $\text{stable\_read}$ . Instead of a single  $\text{protect}(\cdot)$  for expression  $x + y$ , we however need to insert two  $\text{stable\_read}$  for  $a[i_1]$  and  $a[i_2]$ , respectively.

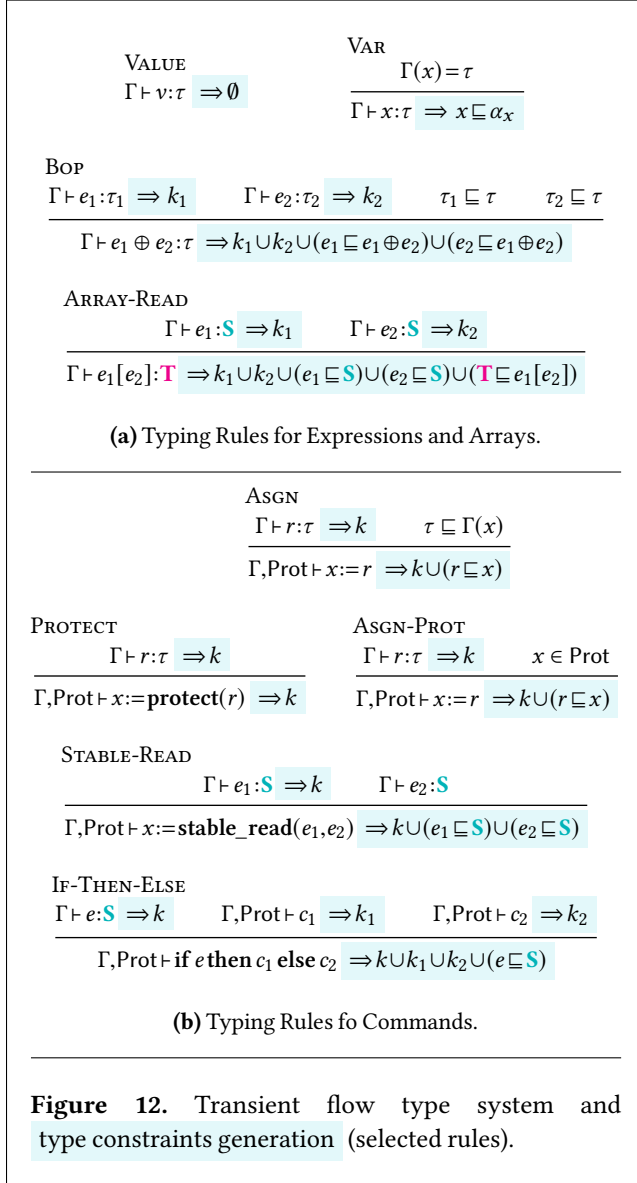
## 4 Type System and Inference

In Section 4.1, we present a transient-flow type system which statically rejects programs that can potentially leak through transient execution attacks. Given an unannotated program, we apply constraint-based type inference [3, 27] to generate its use-def graph and reconstruct type information (Section 4.2). Then, reusing off-the-shelf Max-Flow/Min-Cut algorithms, we analyze the graph and locate potential speculative vulnerabilities in the form of a variable min-cut set. Finally, using a simple program repair algorithm we patch the program by inserting a minimum number of  $\text{protect}$  so that it does not leak speculatively anymore (Figure 13).

### 4.1 Type System

Our transient-flow type system prevents programs from leaking transient values via cache timing channels. To this end, the type system assigns a *transient-flow type* to expressions and tracks how transient values propagate within programs, rejecting programs in which transient values reach commands which may leak them. An expression can either be typed as *stable* (**S**) indicating that it cannot contain transient values during execution, or as *transient* (**T**) indicating that it can. These types form a 2-point lattice [23], which allows stable expressions to be typed as transient, but not vice versa, i.e., we define a can-flow-to relation  $\sqsubseteq$  such that  $\mathbf{S} \sqsubseteq \mathbf{T}$ , but  $\mathbf{T} \not\sqsubseteq \mathbf{S}$ .





## 4.2 Type Inference

We now present our type inference algorithm.

**Constraints.** We start by collecting a set of constraints  $k$  via typing judgement  $\Gamma, \text{Prot} \vdash s \Rightarrow k$ . For this, we define a dummy environment  $\Gamma^*$  and protected set  $\text{Prot}^*$ , such that  $\Gamma^*, \text{Prot}^* \vdash c \Rightarrow k$  holds for any command  $c$ , (i.e., we let  $\Gamma^* = \lambda x. \mathbf{S}$  and include *all* variables in the cut-set) and use it to extract the set of constraints  $k$ . The syntax for constraints is shown in Figure 21. The constraints relate *atoms* which represent the unknown type of variables, i.e.,  $\alpha_x$  for  $x$ , and expression, i.e.,  $r$ . Constraints record can-flow-to relationships between the atoms and lattice values  $\mathbf{T}$  and  $\mathbf{S}$ . They are accumulated via operator  $\cup$ , where we identify  $k_1 \cup \dots \cup k_n$  with the set  $\{k_1, \dots, k_n\}$ .

**Solutions and Satisfiability.** We define the solution to a set of constraints as a function  $\sigma$  from atoms to flow types, i.e.,  $\sigma \in \text{ATOMS} \mapsto \{\mathbf{T}, \mathbf{S}\}$ , and extend solutions to map  $\mathbf{T}$  and  $\mathbf{S}$  to themselves. For a set of constraints  $k$  and a solution function  $\sigma$ , we write  $\sigma \vdash k$  to say that the constraints  $k$  are satisfied under solution  $\sigma$ . A solution  $\sigma$  satisfies  $k$ , if all can-flow-to constraints hold, when the atoms are replaced by their values under  $\sigma$ . We say that a set of constraints  $k$  is satisfiable, if there is a solution  $\sigma$  such that  $\sigma \vdash k$ .

**Def-Use Graph & Paths.** The constraints generated by our type system give rise to the def-use graph of the type-checked program. For a set of constraints  $k$ , we call a sequence of atoms  $a_1 \dots a_n$  a *path* in  $k$ , if  $a_i \sqsubseteq a_{i+1} \in k$  for  $i \in \{1, \dots, n-1\}$  and say that  $a_1$  is the path's entry and  $a_n$  its exit. A  $\mathbf{T-S}$  path is a path with entry  $\mathbf{T}$  and exit  $\mathbf{S}$ . A set of constraints  $k$  is satisfiable if and only if there is no  $\mathbf{T-S}$  path in  $k$ , as such a path would correspond to a derivation of false. If  $k$  is satisfiable, we can compute a solution  $\sigma(k)$  by letting  $\sigma(k)(a) = \mathbf{T}$ , if there is a path with entry  $\mathbf{T}$  and exit  $a$ , and  $\mathbf{S}$  otherwise.

**Cuts.** If a set of constraints is unsatisfiable, we can make it satisfiable by removing some of the nodes in its graph or equivalently protecting some of the variables. A set of atoms  $A$  *cuts* a path  $a_1 \dots a_n$ , if some  $a \in A$  occurs along the path, i.e., there exists  $a \in A$  and  $i \in \{1, \dots, n\}$  such that  $a_i = a$ . We call  $A$  a cut-set for a set of constraints  $k$ , if  $A$  cuts all  $\mathbf{T-S}$  paths in  $k$ . A cut-set  $A$  is minimal for  $k$ , if all other cut-sets  $A'$  contain as many or more atoms than  $A$ , i.e.,  $\#A \leq \#A'$ .

**Extracting Types From Cuts.** From a set of variables  $A$  such that  $A$  is a cut-set of constraints  $k$ , we can extract a typing environment  $\Gamma(k, A)$  as follows: for an atom  $\alpha_x$ , we define  $\Gamma(k, A)(x) = \mathbf{T}$ , if there is a path with entry  $\mathbf{T}$  and exit  $\alpha_x$  in  $k$  that is not cut by  $A$ , and let  $\Gamma(k, A)(x) = \mathbf{S}$  otherwise.

**Proposition 1 (Type Inference).** *If  $\Gamma^*, \text{Prot}^* \vdash c \Rightarrow k$  and  $A$  is a set of variables that cut  $k$ , then  $\Gamma(k, A), A \vdash s$ .*

**Remark.** To infer a repair using `stable_read` instead of `protect`, we can restrict our cut-set to only include variables that are assigned from an array read.

Atom	$a$	$::=$	$\alpha_x \mid r$
Constraint	$k$	$::=$	$a \sqsubseteq \mathbf{S} \mid \mathbf{T} \sqsubseteq a \mid a \sqsubseteq a \mid k \cup k \mid \emptyset$
Solution	$\sigma$	$\in$	$\text{ATOMS} \mapsto \{\mathbf{S}, \mathbf{T}\}$

Figure 13. Constraint Syntax.

**Example.** Consider again Ex1 in Figure 3. The graph defined by the constraints  $k$ , given by  $\Gamma^*, \text{Prot}^* \vdash \text{Ex1} \Rightarrow k$  is shown in Figure 4, where we have omitted  $\alpha$ -nodes. The constraints are not satisfiable, since there are  $\mathbf{T-S}$  paths. Both  $\{x, y\}$  and  $\{z\}$  are cut-sets, since they cut each  $\mathbf{T-S}$  path, however, the set  $\{z\}$  contains only one element and is therefore minimal. The typing environment  $\Gamma(k, \{x, y\})$  extracted from the sub-optimal cut  $\{x, y\}$  types all variables as  $\mathbf{S}$ , while the typing extracted from the optimal cut, i.e.,  $\Gamma(k, \{z\})$  types  $x$  and  $y$  as  $\mathbf{T}$  and  $z, i_1$  and  $i_2$  as  $\mathbf{S}$ . By Proposition 2 both  $\Gamma(k, \{x, y\}), \{x, y\} \vdash \text{Ex1}$  and  $\Gamma(k, \{z\}), \{z\} \vdash \text{Ex1}$  hold.

## 4.3 Program Repair

As a final step, our repair algorithm `repair(c, Prot)` traverses program  $c$  and inserts a `protect(·)` statement for each variable in the cut-set  $\text{Prot}$ . Since we assume that programs are in static single assignment form, there is a single assignment  $x := r$  for each variable  $x \in \text{Prot}$ , and our repair algorithm simply replaces it with  $x := \text{protect}(r)$ .

## 5 Consistency and Security

We now present two formal results about our speculative semantics and the security of the type system. Our full definitions and proofs can be found in Appendix D.

**Consistency.** We write  $C \Downarrow_O^D C'$  for the complete speculative execution of configuration  $C$  to final configuration  $C'$ , which generates a trace of observations  $O$  under list of directives  $D$ . Similarly, we write  $\langle \mu, \rho \rangle \Downarrow_O^c \langle \mu', \rho' \rangle$  for the sequential execution of program  $c$  with initial memory  $\mu$  and variable map  $\rho$  resulting in final memory  $\mu'$  and variable map  $\rho'$ . To relate speculative and sequential observations, we define a projection function, written  $O\downarrow$ , which removes prediction identifiers, rollbacks, and mispredicted loads and stores.

**Theorem 5.1 (Consistency).** *For all programs  $c$ , initial memory stores  $\mu$ , variable maps  $\rho$ , and directives  $D$ , such that  $\langle \mu, \rho \rangle \Downarrow_O^c \langle \mu', \rho' \rangle$  and  $\langle [], [c], \mu, \rho \rangle \Downarrow_{O'}^D \langle [], [], \mu'', \rho'' \rangle$ , then  $\mu' = \mu'', \rho' = \rho''$ , and  $O \cong O'\downarrow$ .*

The theorem ensures equivalence of the final memory stores, variable maps, and observation traces from the sequential and the speculative execution. Notice that trace equivalence is up to *permutation*, i.e.,  $O \cong O'\downarrow$ , because the processor can execute load and store instructions out-of-order.

**Speculative Non-Interference.** Speculative non-interference is parametric in the security policy that specifies which variables and part of the memory are controlled by the attacker [17]. In the following, we write  $L$  for the set of public

variables and memory locations that are *observable* by the attacker. Two variable maps are *indistinguishable* to the attacker, written  $\rho_1 \approx_L \rho_2$ , if and only if  $\rho_1(x) = \rho_2(x)$  for all  $x \in L$ . Similarly, memory stores are related pointwise, i.e.,  $\mu_1 \approx_L \mu_2$  iff  $\mu_1(n) = \mu_2(n)$  for all  $n \in L$ .

**Definition 1** (Speculative Non-Interference). *A program  $c$  satisfies speculative non-interference if and only if for all directives  $D$ , memory stores and variable maps such that  $\mu_1 \approx_L \mu_2$  and  $\rho_1 \approx_L \rho_2$ , let  $C_i = \langle [c], [\mu_i, \rho_i] \rangle$  for  $i \in \{1, 2\}$ , such that  $C_1 \Downarrow_{O_1}^D C'_1, C_2 \Downarrow_{O_2}^D C'_2$ , if  $O_1 \Downarrow = O_2 \Downarrow$ , then  $O_1 = O_2$ .*

In the definition above, programs *leak* by producing different observations starting from memories and variables indistinguishable to the attacker. Speculative non-interference requires showing absence of leaks for the speculative traces ( $O_1 = O_2$ ) assuming that the program does not already leak sequentially ( $O_1 \Downarrow = O_2 \Downarrow$ ). Notice that here we consider syntactic equivalence for the traces because both executions follow the same list of directives. We now present our soundness theorem: well-typed programs satisfy speculative non-interference.

**Theorem 5.2** (Soundness). *For all programs  $c$ , if  $\Gamma \vdash c$  then  $c$  satisfies speculative non-interference.*

We conclude with a corollary that combines all the components of our protection chain (type inference, type checking and automatic repair via our security primitives) and shows that repaired programs satisfy speculative non-interference.

**Corollary 5.3.** *For all programs  $c$ , there exists a set of constraints  $k$  such that  $\Gamma^*, \text{Prot}^* \vdash c \Rightarrow k$ . Let  $A$  be a set of variables that cut  $k$ , then the repaired program  $\text{repair}(c, A)$  satisfies speculative non-interference.*

## 6 Implementation and Evaluation

We now describe our implementation and evaluate BLADE on an implementation of the Signal secure messaging protocol and various cryptographic algorithms. Our evaluation shows that BLADE can secure existing software systems against speculative execution attacks automatically. Moreover, BLADE introduces two orders of magnitude less fences than a baseline algorithm implemented in Clang. As a result, the repairs computed by BLADE incur only a minimal performance overhead.

### 6.1 Implementation

We implemented BLADE in 3500 lines of Haskell code. BLADE takes as input a WebAssembly program, computes a repaired program that is safe under speculative execution and verifies its correctness via type-checking. Internally, BLADE proceeds in three stages. First, BLADE converts the WebAssembly program into an intermediate representation similar to the While language in Figure 5. This simplifies further processing as WebAssembly is a stack-based language, i.e., arguments are not represented directly, but instead kept on an argument stack.

Second, BLADE builds the use-def graph (§4.1) of the input program, infers a minimal cut-set (§4.2), and computes the repair (§4.3). Finally, in the last stage, BLADE extracts a typing-environment from the use-def graph and type-checks the repaired program (§4). This independent checking step provides extra confidence that the repaired program indeed does not leak more speculatively, than it does sequentially (§5). Source code will be made available under an open source license.

### 6.2 Evaluation

We evaluate BLADE by answering three questions: **(Q1)** Can we apply BLADE to secure existing software? **(Q2)** How many `protect` statements does BLADE have to insert in order to secure those systems? and **(Q3)** How do the inserted fences affect performance?

**(Q1) Applicability.** To evaluate BLADE's applicability, we run it on crypto code, which is already carefully written to eschew cache-timing side channels. Our benchmarks are taken from two main sources: first, a verified implementation [29] of the Signal messaging protocol [15], and second, verified implementations of several crypto primitives taken from [38]. In particular, our benchmarks consist of

- ▶ The messaging algorithm implemented in module Signal Core and common cryptographic constructions implemented in module Signal Crypto and used in Signal.
- ▶ The HAcl\* SHA2 hash, AES block cypher, Curve25519 elliptic curve function, and ED25519 digital signature used in Signal.
- ▶ The SALSA20 stream cypher, SHA2 hash, and TEA block cypher from [38].

The original implementations of our benchmarks are *provably* free from cache and timing side-channel. However, those proofs considered only a sequential execution model and therefore do not account for the speculative execution vulnerabilities addressed in this work.

**Results.** Table 1 shows the code size in WebAssembly text format, and the runtime of BLADE on each benchmark. The runtime includes translation, repair and type-checking. The results are encouraging: the execution time scales proportionally with the code size and the analysis completes fairly quickly, even for large benchmarks (>60k WASM LOC): the runtime is less than 10s for all of our benchmarks.

**(Q2) Number of Fences.** Next, we evaluate how many fences the analysis has to insert to make the programs secure. The results are shown in Table 1. Column **B** contains our baseline, which replaces all non-constant array reads, i.e., reads whose address depends on a variable, with statement `stable_read` (Section 3.4), implementing a SLH-like mitigation that masks the address with the array bounds-check condition. This is the proposed mitigation in the Clang compiler [10]. Column **P** shows the number of `protect` inserted by BLADE. All benchmarks are modified by the baseline, except for TEA, which is a simple, toy encryption algorithm



Name	B	P	S	P/B	LOC	Time
CRYPTO [29]	92	1	2	1.1	3386	181.0 ms
CORE [29]	47	1	2	2.1	6595	347.8 ms
SHA2 [29]	156	18	34	11.5	7310	286.7 ms
AES[29]	48	0	0	0	6284	28.95 ms
CURVE [29]	2214	0	0	0	59921	5.571 s
ED25519 [29]	2403	6	10	0.2	60308	8.797 s
SALSA 20 [38]	7	0	0	0	529	20.20 ms
SHA 256 [38]	23	0	0	0	334	11.23 ms
TEA [38]	0	0	0	-	112	3.036 ms
<b>Total</b>	<b>4990</b>	<b>26</b>	<b>48</b>	<b>0.5</b>	<b>144779</b>	<b>-</b>

**Table 1. (B)** contains our baseline, *i.e.*, the number of `stable_read`, if every non-constant read is protected; **(P)** contains the number of `protect` statements insert by BLADE; **(S)** contains the number of `stable_read` inserted, if `stable_read` is used to implement `protect`; **(P/B)** contains the ratio of `protect` statments to the baseline fences in %; **(LOC)** contains the number of lines of WASM code in text format; **(Time)** shows the mean timing for fence inference, repair, and typechecking over 15 trials; Experiments were run on a 12” Macbook with 8GB RAM.

(that should not be used in practice). In particular, for five of the nine programs, BLADE does not need to insert any fences. Column **P/B** shows the ratio of `protect` statements to baseline read masks in percent. For most benchmarks, our analysis has to insert under 3% of fences compared to the baseline. For the SHA2 implementation of HACL\* this rises to 11.5%. Across all benchmarks, the number of fences is two orders of magnitude lower than the baseline. Since `protect` statements are an idealized primitive that are not available in todays hardware, we show the number of `stable-read` primitives that are needed to implement the `protect` in column **S**. The table shows that using `stable reads` requires inserting more fences by a factor of 1.8x, which underlines the benefits of a hardware implementation of `protect`.

**(Q3) Performance Impact of Fences.** To evaluate the performance impact of our repair, we compared how a *naive* placement of fences—applying speculative load hardening to every load of a non-constant address—compares against our approach. We picked the SHA2-512 hash function for this test, and used inputs of size 4KB. Naive fence placement introduced 44 fences while ours introduced only 5. Our measurements showed that while the naive repair algorithm caused 13.9% overhead, the overhead of our minimal fence replacement algorithm was only 0.42%. We used a sample size of 500, and found the relative margin of error of our measurements were less than 0.07%.

## 7 Related Work

**Transient Execution Attacks.** Since Spectre [21] and Melt-down [24] were announced, many transient execution attacks

exploiting different microarchitectural components and side-channels have been discovered and new ones come to light at a steady pace. These attacks leak data across arbitrary security boundaries, including SGX enclaves [14, 35], hypervisors and virtual machines [40], and even remotely over a network [31]. We refer to [9] for a comprehensive systematization.

**Detection and Repair.** Wu and Wang [41] detect cache side channels via abstract interpretation by augmenting the control-flow to accommodate for speculation. Spectector [17] and Pitchfork [11] use symbolic execution on x86 binaries to detect speculative vulnerabilities. Cheang et al. [13] and Bloem et al. [8] apply bounded model checking to detect potential speculative vulnerabilities respectively via 4-ways self-composition and taint-tracking. Almost all these efforts [8, 11, 13, 17, 41] consider only *in-order* execution (except Pitchfork [11]) for a *fixed* speculation bound, and focus on vulnerability detection but do not propose techniques to *repair* vulnerable programs. In contrast, our type system enforces speculative non-interference even when program instructions are executed *out-of-order* with *unbounded* speculation and automatically synthesizes repairs. Given a set of untrusted input source, oo7 Wang et al. [37] statically analyzes a binary to detect vulnerable patterns and inserts fences in turn. Our tool, BLADE, not only repairs vulnerable programs without user annotation, but ensures that program patches contain a minimum number of fences. Furthermore, BLADE formally guarantees that repaired programs are free from speculation-based attacks.

**Speculative Execution Semantics.** There have been several recent proposals for speculative execution semantics [11, 13, 17, 26]. Of those, [11] is closest to ours, and inspired our semantics (e.g., we share the 3-stages pipeline, attacker-supplied directives and the instruction reorder buffer). However their semantics targets an assembly language with direct jumps, while we reason about speculative execution of imperative programs with structured control-flow.

**Hardware Mitigations and Secure Design.** Both AMD AMD [5] and Intel Intel [19] recommend inserting serializing, fence instructions after bounds checks to protect against Spectre v1 attacks and some compilers followed suit [18, 28]. Unfortunately, these defenses causes significant performance degradation [9]. Taram et al. [32] propose context-sensitive fencing, a hardware-based mitigation that dynamically inserts fences in the instruction stream when dangerous conditions arise. Several secure hardware designs have been studied to remove speculative attacks from future processors. InvisiSpec Yan et al. [42] is a new micro-architecture design that features a special *speculative buffer* to prevent speculative loads from polluting the cache. STT [2] tracks speculative taints inside the processor micro-architecture and prevent speculative values from reaching instructions that could serve as covert channels. We think our approach could be applied to guide such hardware mitigations by pinpointing the program parts that need to be protected.



## References

- [1] *Flows in Networks*. Princeton University Press, 1962.
- [2] Speculative taint tracking (stt): A comprehensive protection for speculatively accessed data. In *MICRO*, 2019.
- [3] Alex Aiken. Constraint-based program analysis. In Radhia Cousot and David A. Schmidt, editors, *Static Analysis*, pages 1–1, Berlin, Heidelberg, 1996. Springer Berlin Heidelberg. ISBN 978-3-540-70674-8.
- [4] José Bacerlar Almeida, Manuel Barbosa, Gilles Barthe, François Dupres-soir, and Michael Emmi. Verifying constant-time implementations. In *Usenix Security*, 2016.
- [5] AMD. Software techniques for managing speculation on AMD processors. <https://developer.amd.com/wp-content/resources/Managing-Speculation-on-AMD-Processors.pdf>, 2018.
- [6] GILLES BARTHE, SANDRINE BLAZY, BENJAMIN GRÉGOIRE, RÉMI HUTIN, VINCENT LAPORTE, DAVID PICHARDIE, and ALIX TRIEU. Formal verification of a constant-time preserving c compiler. In *POPL*, 2020.
- [7] Atri Bhattacharyya, Alexandra Sandulescu, Matthias Neugschwand-ner, Alessandro Sorniotti, Babak Falsafi, Mathias Payer, and Anil Kurmus. Smotherspectre: Exploiting speculative execution through port contention. In *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security, CCS '19*, pages 785–800, New York, NY, USA, 2019. ACM. ISBN 978-1-4503-6747-9. doi: 10.1145/3319535.3363194. URL <http://doi.acm.org/10.1145/3319535.3363194>.
- [8] Roderick Bloem, Swen Jacobs, and Yakir Vizel. Efficient information-flow verification under speculative execution. In Yu-Fang Chen, Chih-Hong Cheng, and Javier Esparza, editors, *Automated Technology for Verification and Analysis*, pages 499–514, Cham, 2019. Springer International Publishing. ISBN 978-3-030-31784-3.
- [9] Claudio Canella, Jo Van Bulck, Michael Schwarz, Moritz Lipp, Benjamin Von Berg, Philipp Ortner, Frank Piessens, Dmitry Evtvyushkin, and Daniel Gruss. A systematic evaluation of transient execution attacks and defenses. In *Proceedings of the 28th USENIX Conference on Security Symposium, SEC'19*, pages 249–266, Berkeley, CA, USA, 2019. USENIX Association. ISBN 978-1-939133-06-9. URL <http://dl.acm.org/citation.cfm?id=3361338.3361356>.
- [10] Chandler Carruth. Speculative load hardening. <https://llvm.org/docs/SpeculativeLoadHardening.html>, 2019.
- [11] Sunjay Cauligi, Craig Disselkoen, Klaus von Gleissenthall, Deian Stefan, Tamara Rezk, and Gilles Barthe. Towards constant-time foundations for the new spectre era. *CoRR*, abs/1910.01755, 2019. URL <http://arxiv.org/abs/1910.01755>.
- [12] Sunjay Cauligi, Gary Soeller, Brian Johannesmeyer, Fraser Brown, Riad S. Wahby, John Renner, Benjamin Gregoire, Gilles Barthe, Ranjit Jhala, and Deian Stefan. FaCT: A dsl for timing-sensitive computation. In *Programming Language Design and Implementation (PLDI)*. ACM SIGPLAN, June 2019.
- [13] Kevin Cheang, Cameron Rasmussen, Sanjit A. Seshia, and Pramod Subramanyan. A formal approach to secure speculation. In *Proceedings of the Computer Security Foundations Symposium (CSF)*, 2019.
- [14] Guoxing Chen, Sanchuan Chen, Yuan Xiao, Yinqian Zhang, Zhiqiang Lin, and Ten-Hwang Lai. Sgxpectre attacks: Leaking enclave secrets via speculative execution. *CoRR*, abs/1802.09085, 2018. URL <http://arxiv.org/abs/1802.09085>.
- [15] Katriel Cohn-Gordon, Cas Cremers, Benjamin Dowling, Luke Garratt, and Douglas Stebila. A formal security analysis of the signal messaging protocol. In *EuroS&P*, 2017.
- [16] Qian Ge, Yuval Yarom, David Cock, and Gernot Heiser. A survey of microarchitectural timing attacks and countermeasures on contemporary hardware. In *Journal of Cryptographic Engineering*, 2018.
- [17] Marco Guarnieri, Boris Koepf, José Francisco Morales, Jan Reineke, and Andrés Sánchez. Spectector: Principled detection of speculative information flows. In *Proc. IEEE Symp. on Security and Privacy, SSP '20*, 2020.
- [18] Intel. Using intel compilers to mitigate speculative execution side-channel issues. <https://software.intel.com/en-us/articles/using-intel-compilers-to-mitigate-speculative-execution-side-channel-issues>. URL <http://arxiv.org/abs/1807.03757>.
- [19] Intel. Intel analysis of speculative execution side channels. <https://newsroom.intel.com/wp-content/uploads/sites/11/2018/01/Intel-Analysis-of-Speculative-Execution-Side-Channels.pdf>, 2018.
- [20] Vladimir Kiriansky and Carl Waldspurger. Speculative buffer overflows: Attacks and defenses. *CoRR*, abs/1807.03757, 2018. URL <http://arxiv.org/abs/1807.03757>.
- [21] Paul Kocher, Jann Horn, Anders Fogh, Daniel Genkin, Daniel Gruss, Werner Haas, Mike Hamburg, Moritz Lipp, Stefan Mangard, Thomas Prescher, Michael Schwarz, and Yuval Yarom. Spectre attacks: Exploiting speculative execution. In *40th IEEE Symposium on Security and Privacy (S&P'19)*, 2019.
- [22] Esmaeil Mohammadian Koruyeh, Khaled N. Khasawneh, Chengyu Song, and Nael Abu-Ghazaleh. Spectre returns! speculation attacks using the return stack buffer. In *Proceedings of the 12th USENIX Conference on Offensive Technologies, WOOT'18*, pages 3–3, Berkeley, CA, USA, 2018. USENIX Association. URL <http://dl.acm.org/citation.cfm?id=3307423.3307426>.
- [23] J. Landauer. A lattice of information. In *CSFW*, 1993.
- [24] Moritz Lipp, Michael Schwarz, Daniel Gruss, Thomas Prescher, Werner Haas, Anders Fogh, Jann Horn, Stefan Mangard, Paul Kocher, Daniel Genkin, Yuval Yarom, and Mike Hamburg. Meltdown: Reading kernel memory from user space. In *27th USENIX Security Symposium (USENIX Security 18)*, 2018.
- [25] Giorgi Maisuradze and Christian Rossow. Ret2spec: Speculative execution using return stack buffers. In *Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security, CCS '18*, pages 2109–2122, New York, NY, USA, 2018. ACM. ISBN 978-1-4503-5693-0. doi: 10.1145/3243734.3243761. URL <http://doi.acm.org/10.1145/3243734.3243761>.
- [26] Ross McIlroy, Jaroslav Sevcik, Tobias Tebbi, Ben L. Titzer, and Toon Verwaest. Spectre is here to stay: An analysis of side-channels and speculative execution. *CoRR*, abs/1902.05178, 2019. URL <http://arxiv.org/abs/1902.05178>.
- [27] Hanne Riis Nielson and Flemming Nielson. Flow logics for constraint based analysis. In Kai Koskimies, editor, *Compiler Construction*, pages 109–127, Berlin, Heidelberg, 1998. Springer Berlin Heidelberg. ISBN 978-3-540-69724-4.
- [28] Andrew Pardoe. Spectre mitigations in msvc. <https://devblogs.microsoft.com/cppblog/spectre-mitigations-in-msvc/>, 2018.
- [29] Jonathan Protzenko, Benjamin Beurdouche, Denis Merigoux, and Karthikeyan Bhargavan. Formally verified cryptographic web applications in webassembly. In *Security and Privacy*, 2019.
- [30] G. Romer and C. Carruth. C++ proposal, 2019. URL <http://www.open-std.org/jtc1/sc22/wg21/docs/papers/2018/p0928r0.pdf>.
- [31] Michael Schwarz, Martin Schwarzl, Moritz Lipp, and Daniel Gruss. Net-spectre: Read arbitrary memory over network. *CoRR*, abs/1807.10535, 2018. URL <http://arxiv.org/abs/1807.10535>.
- [32] Mohammadkazem Taram, Ashish Venkat, and Dean Tullsen. Context-sensitive fencing: Securing speculative execution via microcode customization. In *ASPLOS'19*.
- [33] Vadim Tkachenko. 20-30% performance hit from the spectre bug fix on ubuntu. <https://www.percona.com/blog/2018/01/23/20-30-performance-hit-spectre-bug-fix-ubuntu/>, Jan 2018.
- [34] Eran Tromer, Dag Arne Osvik, and Adi Shamir. Efficient cache attacks on aes, and countermeasures. *J. Cryptol.*, 23(1):37–71, January 2010. ISSN 0933-2790. doi: 10.1007/s00145-009-9049-y. URL <http://dx.doi.org/10.1007/s00145-009-9049-y>.
- [35] Jo Van Bulck, Marina Minkin, Ofir Weisse, Daniel Genkin, Baris Kasikci, Frank Piessens, Mark Silberstein, Thomas F. Wenisch, Yuval Yarom,

1321	and Raoul Strackx. Foreshadow: Extracting the keys to the Intel SGX kingdom with transient out-of-order execution. In <i>Proceedings of the 27th USENIX Security Symposium</i> . USENIX Association, August 2018.	1376
1322	See also technical report Foreshadow-NG [40].	1377
1323	[36] D. Volpano, G. Smith, and C. Irvine. A Sound Type System for Secure Flow Analysis. <i>J. Computer Security</i> , 4(3):167–187, 1996.	1378
1324	[37] Guanhua Wang, Sudipta Chattopadhyay, Ivan Gotovchits, Tulika Mitra, and Abhik Roychoudhury. oo7: Low-overhead defense against spectre attacks via binary analysis. <i>CoRR</i> , abs/1807.05843, 2018. URL <a href="http://arxiv.org/abs/1807.05843">http://arxiv.org/abs/1807.05843</a> .	1379
1325	[38] Conrad Watt, John Renner, Natalie Popescu, Sunjay Cauligi, and Deian Stefan. Ct-wasm: Type-driven secure cryptography for the web ecosystem. In <i>POPL</i> , 2019.	1380
1326	[39] Conrad Watt, John Renner, Natalie Popescu, Sunjay Cauligi, and Deian Stefan. Ct-wasm: Type-driven secure cryptography for the web ecosystem. <i>Proc. ACM Program. Lang.</i> , 3(POPL):77:1–77:29, January 2019. ISSN 2475-1421. doi: 10.1145/3290390. URL <a href="http://doi.acm.org/10.1145/3290390">http://doi.acm.org/10.1145/3290390</a> .	1381
1327	[40] Ofir Weisse, Jo Van Bulck, Marina Minkin, Daniel Genkin, Baris Kasikci, Frank Piessens, Mark Silberstein, Raoul Strackx, Thomas F. Wenisch, and Yuval Yarom. Foreshadow-NG: Breaking the virtual memory abstraction with transient out-of-order execution. <i>Technical report</i> , 2018. See also USENIX Security paper Foreshadow [35].	1382
1328	[41] Meng Wu and Chao Wang. Abstract interpretation under speculative execution. In <i>Proceedings of the 40th ACM SIGPLAN Conference on Programming Language Design and Implementation</i> , PLDI 2019, pages 802–815, New York, NY, USA, 2019. ACM. ISBN 978-1-4503-6712-7. doi: 10.1145/3314221.3314647. URL <a href="http://doi.acm.org/10.1145/3314221.3314647">http://doi.acm.org/10.1145/3314221.3314647</a> .	1383
1329	[42] Mengjia Yan, Jiho Choi, Dimitrios Skarlatos, Adam Morrison, Christopher W. Fletcher, and Josep Torrellas. Invisispec: Making speculative execution invisible in the cache hierarchy. In <i>Proceedings of the 51st Annual IEEE/ACM International Symposium on Microarchitecture</i> , MICRO-51, pages 428–441, Piscataway, NJ, USA, 2018. IEEE Press. ISBN 978-1-5386-6240-3. doi: 10.1109/MICRO.2018.00042. URL <a href="https://doi.org/10.1109/MICRO.2018.00042">https://doi.org/10.1109/MICRO.2018.00042</a> .	1384
1330	[43] Yuval Yarom and Katrina Falkner. Flush+reload: A high resolution, low noise, l3 cache side-channel attack. In <i>23rd USENIX Security Symposium (USENIX Security 14)</i> , pages 719–732, San Diego, CA, August 2014. USENIX Association. ISBN 978-1-931971-15-7. URL <a href="https://www.usenix.org/conference/usenixsecurity14/technical-sessions/presentation/yarom">https://www.usenix.org/conference/usenixsecurity14/technical-sessions/presentation/yarom</a> .	1385
1331		1386
1332		1387
1333		1388
1334		1389
1335		1390
1336		1391
1337		1392
1338		1393
1339		1394
1340		1395
1341		1396
1342		1397
1343		1398
1344		1399
1345		1400
1346		1401
1347		1402
1348		1403
1349		1404
1350		1405
1351		1406
1352		1407
1353		1408
1354		1409
1355		1410
1356		1411
1357		1412
1358		1413
1359		1414
1360		1415
1361		1416
1362		1417
1363		1418
1364		1419
1365		1420
1366		1421
1367		1422
1368		1423
1369		1424
1370		1425
1371		1426
1372		1427
1373		1428
1374		1429
1375		1430

## A Full Calculus

FETCH-SKIP

$$\langle is, \text{skip} : cs, \mu, \rho \rangle \xrightarrow{\text{fetch}}_{\epsilon} \langle is++[\text{nop}], cs, \mu, \rho \rangle$$

FETCH-ASGN

$$\langle is, x := e : cs, \mu, \rho \rangle \xrightarrow{\text{fetch}}_{\epsilon} \langle is++[x := e], cs, \mu, \rho \rangle$$

FETCH-SEQ

$$\langle is, c_1; c_2 : cs, \mu, \rho \rangle \xrightarrow{\text{fetch}}_{\epsilon} \langle is, c_1 : c_2 : cs, \mu, \rho \rangle$$

FETCH-PTR-LOAD

$$\langle is, x := *e : cs, \mu, \rho \rangle \xrightarrow{\text{fetch}}_{\epsilon} \langle is++[x := \text{load}(e)], cs, \mu, \rho \rangle$$

FETCH-PTR-STORE

$$\langle is, *e_1 := e_2 : cs, \mu, \rho \rangle \xrightarrow{\text{fetch}}_{\epsilon} \langle is++[\text{store}(e_1, e_2)], cs, \mu, \rho \rangle$$

FETCH-FAIL

$$\langle is, \text{fail} : cs, \mu, \rho \rangle \xrightarrow{\text{fetch}}_{\epsilon} \langle is++[\text{fail}], cs, \mu, \rho \rangle$$

FETCH-ARRAY-LOAD

$$\begin{array}{l} c = x := e_1[e_2] \quad e = e_2 < \text{length}(e_1) \quad \text{fresh}(p) \\ e' = \text{base}(e_1) + e_2 \quad c' = \text{if } e \text{ then } x := *e' \text{ else fail} \end{array}$$

$$\langle is, c : cs, \mu, \rho \rangle \xrightarrow{\text{fetch}}_{\epsilon} \langle is, c' : cs, \mu, \rho \rangle$$

FETCH-ARRAY-STORE

$$\begin{array}{l} c = e_1[e_2] := e_3 \quad e = e_2 < \text{length}(e_1) \quad \text{fresh}(p) \\ e' = \text{base}(e_1) + e_2 \quad c' = \text{if } e \text{ then } *e' := e \text{ else fail} \end{array}$$

$$\langle is, c : cs, \mu, \rho \rangle \xrightarrow{\text{fetch}}_{\epsilon} \langle is, c' : cs, \mu, \rho \rangle$$

FETCH-IF-TRUE

$$\begin{array}{l} c = \text{if } e \text{ then } c_1 \text{ else } c_2 \\ \text{fresh}(p) \quad i = \text{guard}(e^{\text{true}}, c_2 : cs, p) \end{array}$$

$$\langle is, c : cs, \mu, \rho \rangle \xrightarrow{\text{fetch true}}_{\epsilon} \langle is++[i], c_1 : cs, \mu, \rho \rangle$$

FETCH-IF-FALSE

$$\begin{array}{l} c = \text{if } e \text{ then } c_1 \text{ else } c_2 \\ \text{fresh}(p) \quad i = \text{guard}(e^{\text{false}}, c_1 : cs, p) \end{array}$$

$$\langle is, c : cs, \mu, \rho \rangle \xrightarrow{\text{fetch false}}_{\epsilon} \langle is++[i], c_2 : cs, \mu, \rho \rangle$$

FETCH-WHILE

$$c_1 = c; \text{while } e \text{ c} \quad c_2 = \text{if } e \text{ then } c_1 \text{ else skip}$$

$$\langle is, \text{while } e \text{ c} : cs, \mu, \rho \rangle \xrightarrow{\text{fetch}}_{\epsilon} \langle is, c_2 : cs, \mu, \rho \rangle$$

Figure 14. Fetch stage.

1541  
1542  
1543  
1544  
1545  
1546  
1547  
1548  
1549  
1550  
1551  
1552  
1553  
1554  
1555  
1556  
1557  
1558  
1559  
1560  
1561  
1562  
1563  
1564  
1565  
1566  
1567  
1568  
1569  
1570  
1571  
1572  
1573  
1574  
1575  
1576  
1577  
1578  
1579  
1580  
1581  
1582  
1583  
1584  
1585  
1586  
1587  
1588  
1589  
1590  
1591  
1592  
1593  
1594  
1595

$$\begin{array}{c}
\text{EXECUTE} \\
\frac{|is_1| = n - 1 \quad \rho' = \phi(is_1, \rho) \quad \langle is_1, i, is_2, cs \rangle \xrightarrow{(\mu, \rho', 0)} \langle is', cs' \rangle}{\langle is_1 ++ [i] ++ is_2, cs, \mu, \rho \rangle \xrightarrow{\text{exec } n} \langle is', cs', \mu, \rho \rangle} \\
\text{EXEC-ASGN} \\
\frac{i = (x := e) \quad v = \llbracket e \rrbracket^\rho \quad i' = (x := v)}{\langle is_1, i, is_2, cs \rangle \xrightarrow{(\mu, \rho, \epsilon)} \langle is_1 ++ [i'] ++ is_2, cs \rangle} \\
\text{EXEC-BRANCH-OK} \\
\frac{i = \text{guard}(e^b, cs', p) \quad \llbracket e \rrbracket^\rho = b}{\langle is_1, i, is_2, cs \rangle \xrightarrow{(\mu, \rho, \epsilon)} \langle is_1 ++ [\text{nop}] ++ is_2, cs \rangle} \\
\text{EXEC-BRANCH-MISPREDICT} \\
\frac{i = \text{guard}(e^b, cs', p) \quad \llbracket e \rrbracket^\rho \neq b}{\langle is_1, i, is_2, cs \rangle \xrightarrow{(\mu, \rho, \text{rollback}(p))} \langle is_1, cs' \rangle} \\
\text{EXEC-LOAD} \\
\frac{i = x := \text{load}(e) \quad \text{store}(\_, \_) \notin is_1 \quad n = \llbracket e \rrbracket^\rho \quad ps = \langle is_1 \rangle \quad i' = (x := \mu(n))}{\langle is_1, i, is_2, cs \rangle \xrightarrow{(\mu, \rho, \text{read}(n, ps))} \langle is_1 ++ [i'] ++ is_2, cs \rangle} \\
\text{EXEC-STORE-ADDR} \\
\frac{i = \text{store}(e_1, e_2) \quad n = \llbracket e_1 \rrbracket^\rho \quad i' = \text{store}(n, e_2)}{\langle is_1, i, is_2, cs \rangle \xrightarrow{(\mu, \rho, \epsilon)} \langle is_1 ++ [i'] ++ is_2, cs \rangle} \\
\text{EXEC-STORE-VALUE} \\
\frac{i = \text{store}(n, e) \quad v = \llbracket e \rrbracket^\rho \quad ps = \langle is_1 \rangle \quad i' = \text{store}(n, v)}{\langle is_1, i, is_2, cs \rangle \xrightarrow{(\mu, \rho, \text{write}(n, ps))} \langle is_1 ++ [i'] ++ is_2, cs \rangle}
\end{array}$$

Figure 15. Execute stage.

1596  
1597  
1598  
1599  
1600  
1601  
1602  
1603  
1604  
1605  
1606  
1607  
1608  
1609  
1610  
1611  
1612  
1613  
1614  
1615  
1616  
1617  
1618  
1619  
1620  
1621  
1622  
1623  
1624  
1625  
1626  
1627  
1628  
1629  
1630  
1631  
1632  
1633  
1634  
1635  
1636  
1637  
1638  
1639  
1640  
1641  
1642  
1643  
1644  
1645  
1646  
1647  
1648  
1649  
1650

$$\begin{array}{c}
\text{RETIRE-NOP} \\
\langle \text{nop} : is, cs, \mu, \rho \rangle \xrightarrow{\text{retire}}_\epsilon \langle is, cs, \mu, \rho \rangle \\
\text{RETIRE-ASGN} \\
\langle x := v : is, cs, \mu, \rho \rangle \xrightarrow{\text{retire}}_\epsilon \langle is, cs, \mu, \rho[x \mapsto v] \rangle \\
\text{RETIRE-STORE} \\
\frac{i = \text{store}(n, v)}{\langle i : is, cs, \mu, \rho \rangle \xrightarrow{\text{retire}}_\epsilon \langle is, cs, \mu[n \mapsto v], \rho \rangle} \\
\text{RETIRE-FAIL} \\
\langle \text{fail} : is, cs, \mu, \rho \rangle \xrightarrow{\text{retire}}_{\text{fail}} \langle [], [], \mu, \rho \rangle
\end{array}$$

Figure 16. Retire stage.

1621  
1622  
1623  
1624  
1625  
1626  
1627  
1628  
1629  
1630  
1631  
1632  
1633  
1634  
1635  
1636  
1637  
1638  
1639  
1640  
1641  
1642  
1643  
1644  
1645  
1646  
1647  
1648  
1649  
1650

$$\begin{array}{c}
\phi(\rho, []) = \rho \\
\phi(\rho, (x := v) : is) = \phi(\rho[x \mapsto v], is) \\
\phi(\rho, (x := e) : is) = \phi(\rho[x \mapsto \perp], is) \\
\phi(\rho, (x := \text{load}(e)) : is) = \phi(\rho[x \mapsto \perp], is) \\
\phi(\rho, (x := \text{protect}(e)) : is) = \phi(\rho[x \mapsto \perp], is) \\
\phi(\rho, i : is) = \phi(\rho, is)
\end{array}$$

(a) Transient Variable Map.

1621  
1622  
1623  
1624  
1625  
1626  
1627  
1628  
1629  
1630  
1631  
1632  
1633  
1634  
1635  
1636  
1637  
1638  
1639  
1640  
1641  
1642  
1643  
1644  
1645  
1646  
1647  
1648  
1649  
1650

$$\begin{array}{c}
\llbracket v \rrbracket^\rho = v \\
\llbracket x \rrbracket^\rho = \rho(x) \\
\llbracket \text{length}(e) \rrbracket^\rho = \text{length}(\llbracket e \rrbracket^\rho) \\
\llbracket \text{base}(e) \rrbracket^\rho = \text{base}(\llbracket e \rrbracket^\rho) \\
\llbracket e_1 + e_2 \rrbracket^\rho = \llbracket e_1 \rrbracket^\rho + \llbracket e_2 \rrbracket^\rho \\
\llbracket e_1 \leq e_2 \rrbracket^\rho = \llbracket e_1 \rrbracket^\rho \leq \llbracket e_2 \rrbracket^\rho
\end{array}$$

(b) Evaluation Function.

1621  
1622  
1623  
1624  
1625  
1626  
1627  
1628  
1629  
1630  
1631  
1632  
1633  
1634  
1635  
1636  
1637  
1638  
1639  
1640  
1641  
1642  
1643  
1644  
1645  
1646  
1647  
1648  
1649  
1650

$$\begin{array}{c}
\langle [] \rangle = [] \\
\langle \text{guard}(e^b, cs, p) : is \rangle = p : \langle is \rangle \\
\langle i : is \rangle = \langle is \rangle
\end{array}$$

(c) Pending Guard Identifiers.

Figure 17. Helper functions.



$$\begin{array}{c}
\text{FETCH-PROTECT-PTR} \\
\frac{c = x := \mathbf{protect}(*e)}{c_1 = x := *e \quad c_2 = x := \mathbf{protect}(x)} \\
\langle is, c : cs, \mu, \rho \rangle \xrightarrow{\text{fetch}}_{\epsilon} \langle is, c_1 : c_2 : cs, \mu, \rho \rangle \\
\\
\text{FETCH-PROTECT-ARRAY} \\
\frac{c = x := \mathbf{protect}(e_1[e_2])}{c_1 = x := e_1[e_2] \quad c_2 = x := \mathbf{protect}(x)} \\
\langle is, c : cs, \mu, \rho \rangle \xrightarrow{\text{fetch}}_{\epsilon} \langle is, c_1 : c_2 : cs, \mu, \rho \rangle \\
\\
\text{FETCH-PROTECT-EXPR} \\
\frac{c = x := \mathbf{protect}(e) \quad i = x := \mathbf{protect}(e)}{\langle is, c : cs, \mu, \rho \rangle \xrightarrow{\text{fetch}}_{\epsilon} \langle is ++ [i], cs, \mu, \rho \rangle} \\
\\
\text{EXEC-PROTECT}_1 \\
\frac{i = x := \mathbf{protect}(e) \quad v = \llbracket e \rrbracket^{\rho} \quad i' = x := \mathbf{protect}(v)}{\langle is_1, i, is_2, cs \rangle \xrightarrow{(\mu, \rho, \epsilon)} \langle is_1 ++ [i'] ++ is_2, cs \rangle} \\
\\
\text{EXEC-PROTECT}_2 \\
\frac{i = x := \mathbf{protect}(v) \quad \mathbf{guard}(\rightarrow, \rightarrow) \notin is_1 \quad i' = (x := v)}{\langle is_1, i, is_2, cs \rangle \xrightarrow{(\mu, \rho, \epsilon)} \langle is_1 ++ [i'] ++ is_2, cs \rangle}
\end{array}$$

**Figure 18.** Semantics of  $\mathbf{protect}(\cdot)$ .

## B Semantics of Stable Read

Current processors do not provide a protect primitive instruction nor the means to implement it on top of existing instructions, in its full generality. However, for array reads, it is possible to replicate the effects of protect by exploiting the same data-dependencies tracking capabilities at the core of the processor pipeline. Indeed, Speculative Load Hardening (SLH), a mitigation technique deployed in the code generated by the CLANG compiler, relies on data-dependencies to secure memory loads automatically [10]. Using our formal model, we give rigorous semantics to SLH and show that it can stop transient execution attacks.

At a high level, SLH injects *artificial* data-dependencies between the conditions used in branch instructions and the addresses loaded in the following instructions to transform control-flow dependencies into data-flow dependencies. Intuitively, these data-dependencies validate control-flow decisions at runtime by stalling speculative loads until the processor resolves the conditions. Using branch conditions, SLH *masks* the address of loads instructions in such a way that the processor zeroes out the address if the condition is mispredicted, preventing misloads.

To formalize this mechanism, we extend our processor model as follows. We introduce a new processor instruction  $x := e ? e_1 : e_2$ , which corresponds to the conditional move instruction CMOV on x86 processors. This instruction simply assigns the value of  $e_1$  (resp.  $e_2$ ) to variable  $x$ , if the condition  $e$  evaluates to true (resp. false). Importantly, this instruction is not subject to speculation: the processor must first evaluate the condition before it can resolve the assignment. We also extend expressions with the standard bitwise AND operator ( $\&$ ) and write  $\bar{0}$  and  $\bar{1}$  for bit words consisting of all 0 and 1. As usual bitmask  $\bar{0}$  and  $\bar{1}$  are respectively the zero and identity element for  $\&$ , i.e.,  $\llbracket e \& \bar{0} \rrbracket^\rho = \bar{0}$  and  $\llbracket e \& \bar{1} \rrbracket^\rho = \llbracket e \rrbracket^\rho$ .

Figure 19 presents the semantics rules for CMOV and for the stable read command implemented using SLH. Rule [EXEC-CMOV] evaluates the condition ( $b = \llbracket e \rrbracket^\rho$ ) of the conditional assignment  $x := e ? e_{\text{true}} : e_{\text{false}}$  and assigns the corresponding expressions ( $x := e_b$ ). Rule [FETCH-STABLE-READ-SLH] fetches command  $x := \text{stable\_read}(e_1, e_2)$ , computes the bounds check condition, the address of the indexed element, and push on the stack the following command.

$$\begin{array}{l} r := e_1 \leq \text{length}(e_2) \\ \text{if } r \text{ then} \\ \quad r := r ? \bar{1} : \bar{0}; \\ \quad x := *((\text{base}(e_1) + e_2) \& r); \\ \text{else} \\ \quad \text{fail} \end{array}$$

The code above is similar to the code generated by a regular array read, but additionally stores the result of the bounds-check condition in reserved variable  $r$ . In the then-branch, the condition is then converted into a suitable bitmask using using

$$\begin{array}{l} \text{EXEC-CMOV} \\ i = x := e ? e_{\text{true}} : e_{\text{false}} \quad b = \llbracket e \rrbracket^\rho \quad i' = x := e_b \\ \hline \langle is_1, i, is_2, cs \rangle \xrightarrow{(\mu, \rho, \epsilon)} \langle is_1 + + [i'] + + is_2, cs \rangle \\ \\ \text{FETCH-STABLE-READ-SLH} \\ c = x := \text{stable\_read}(e_1, e_2) \quad e = e_2 < \text{length}(e_1) \\ e' = \text{base}(e_1) + e_2 \quad c_1 = r := e \quad c_2 = r := r ? \bar{1} : \bar{0} \\ c_3 = x := *(e' \& r) \quad c' = c_1; \text{if } r \text{ then } c_2; c_3 \text{ else fail} \\ \hline \langle is, c : cs, \mu, \rho \rangle \xrightarrow{\text{fetch}}_\epsilon \langle is, c' : cs, \mu, \rho \rangle \end{array}$$

Figure 19. Semantics of  $x := \text{stable\_read}(e_1, e_2)$ .

the non-speculative CMOV instruction i.e.,  $r := r ? \bar{1} : \bar{0}$ , which then masks the address loaded, i.e.,  $*((\text{base}(e_1) + e_2) \& r)$ . As a result, the value of the address remains undefined until the processor evaluates the bounds check condition. When the condition resolves, if the index is inbound  $r = \bar{1}$  and the program reads the correct address  $\llbracket e \& \bar{1} \rrbracket^\rho = \llbracket e \rrbracket^\rho$ . If the index is out-of-bounds, instead,  $r = \bar{0}$  and the load can only read speculatively from a constant address ( $x := \mu(0)$ ), thus closing the leak.<sup>4</sup>

**Revisited Example.** Consider again running example Ex1 in Figure 3, where instead of standard array reads, we employ the `stable_read`( $\cdot$ ) primitive from above. After fetching the program, the addresses of the loads are masked with the respective array bounds-check conditions. Assuming the same memory layout and content as in Figure 10 (except for the fact that arrays are shifted by one position since  $\mu(0) = 0$  is reserved), the processor resolves the first bounds check and reads the array within its bounds, i.e.,  $x := \mu(3) = 0$ . The second load attempts to read the array out of bounds ( $y := a[2]$ ), and our countermeasure prevents the buffer overrun by redirecting the load to the dummy value stored at address 0. First, the processor resolves the bounds check, i.e.,  $r := \bar{0}$ , and forwards it to the load  $y := \text{load}((\text{base}(a) + i_2) \& r)$ . Then, the condition zeroes out the address and the processor assigns the dummy value to variable  $y$ , i.e.,  $y := \mu(0)$ . As a result, we always read array  $b$  at index  $z = 0$  and close the leak.

<sup>4</sup>We assume that the first memory cell is reserved to the processor, which initializes it with dummy data, e.g.,  $\mu(0) = 0$ .