

# Outline

---

Speculation



Spectre

Blade Approach

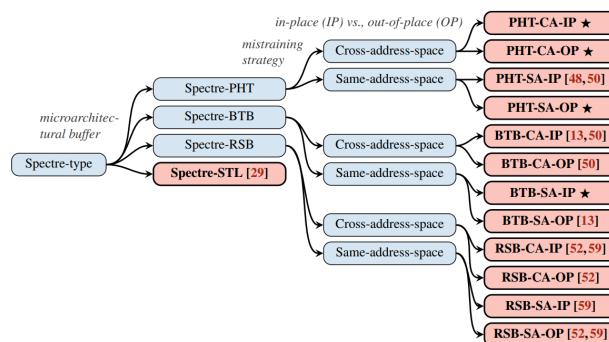
**Spectre** exploits spec ex to leak secrets.

# We are now ready for



**Spectre** is an attack exploiting **speculative execution** to leak secrets

- It is part of a recent wave of **micro-architectural attacks**, i.e., attacks using side-channels induced by the micro-architecture of a processor (including **cache**, timers, virtual memory, ...)
- “**Spectre** [...] transiently bypasses software-defined security policies (e.g., bounds checking, function call/return abstractions, memory stores) to leak secrets out of the program’s intended code/data paths.” [Canella et al., 2019] (← bit outdated, but a good survey on Spectre!)
- Lot of variants:



There are a lot of different variants.

[Canella et al., 2019] Canella, Claudio, et al. "A systematic evaluation of transient execution attacks and defenses." 28th USENIX Security Symposium (USENIX Security 19). 2019. URL: <https://arxiv.org/abs/1811.05441v3>

# What is a gadget?

A **gadget** is a piece of existing code in an (unmodified) existing program binary

- A malicious actor can influence program control flow to cause gadget code to execute
- Gadget code can perform some action of interest to the attacker

That a malicious action can influence to implement a attack

# Spectre v1/v1.1 (special-cases of PHT)

**Spectre**, also called **Spectre-PHT**: the **attacker misstrains** the branch predictor, by poisoning the **Pattern History Table**

How? For example, as follows

◦ The attacker:

- 1. Looks for a piece of “vulnerable code” (**code gadget**) including a condition (next slide)
- 2. Runs that code multiple times with an input such that the condition always hold (so training the branch predictor to take the “true” branch)
- 3. Then, runs the code with a specially-crafted input making the condition false: now the CPU mis-speculates and executes the “true” branch with wrong data (!!?)
- 4. Finally, it looks for information left behind after the CPU discovered it mis-specified and rolled-back the computation (e.g., part of the contents of caches)

↳ To force guess to one branch or the other.

We focus on **Spectre v1** and **v1.1**: due to the fact that modern microprocessors may speculate beyond a **bounds check** condition

# Spectre v1 (Bounds Check Bypass)

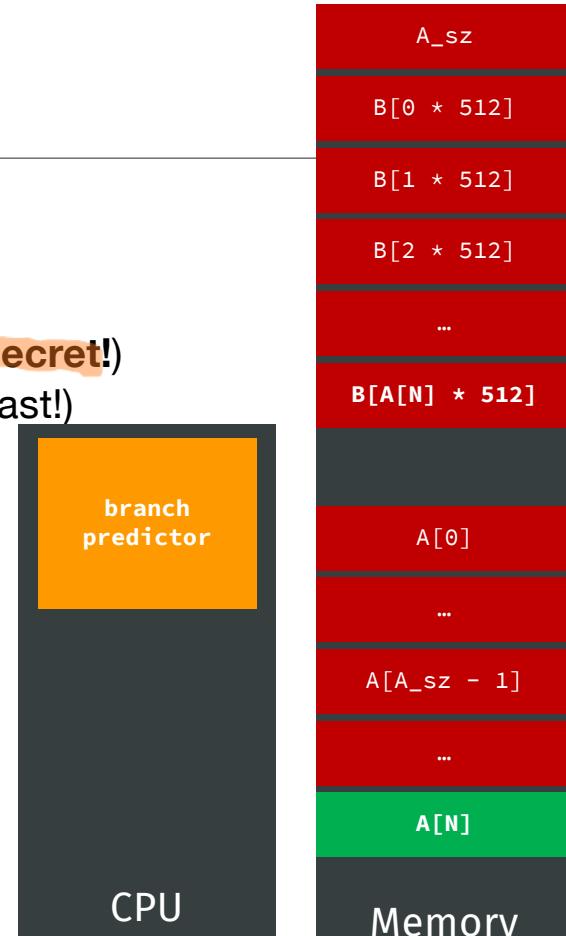
1. Find a code gadget (i.e., a vulnerable piece of code) in victim
2. Mis-train the branch predictor
3. Suppose B and A\_sz not in cache
4. Call the gadget with  $x = N$  with  $N > A_{sz}$  ( $A[N]$  cached, value secret!)
5. CPU mis-speculates on the bound check and accesses  $A[N]$  (fast!)
6. At that point,  $B[A[N] * 512]$  is loaded in the cache
7. CPU discovers that the condition was false: rollback!
8. When control returns to caller (attacker) it just goes over  $B[i * 512]$  for each  $i$  and measures access time:
  - if  $i \neq A[N]$  the access will be slow 🐘
  - if  $i = A[N]$  the access will be fast, the secret is leaked! 🐔

```
if (x < A_sz)
    y = B[A[x] * 512];
```

You have a bound check! ⚡

The attacker can successfully read past the end of the array A

Not in cache      In cache



① And thus is exploited at lower level. We mis-train the predictor that bound check holds.  
Attacker forces program to access  $A[m]$  and  $B[A[m]]$ . We want to discover  $A[m]$ :

You can try to access  $B[n * 512]$  until you have a fast access (already uncached). A attacker can access  $B[n]$ .  
A attacker explores this to read beyond A array.

# Spectre v1.1 (Bounds Check Bypass on stores)

1. Find a **code gadget** in victim
2. Mis-train the branch predictor
3. Let:
  - A, A\_sz, B as above
  - A[N] contains a **secret**
  - y (supplied by the attacker) s.t. A[y] points on the return address on the stack
  - z (supplied by the attacker) to point to an instruction accessing B[A[N] \* 512]
4. If function returns during mis-speculation, the CPU loads B[A[N] \* 512] in cache
5. Despite rollback, when control returns to caller (attacker) it just goes over B[i \* 512] and measures access time:
  - if i ≠ A[N] the access will be slow 🐘
  - if i = A[N] the access will be fast, **the secret is leaked!** 🐚

```
if (y < A_sz)
    A[y] = z;
```

# Spectre: recap

→ Speculation is used here!

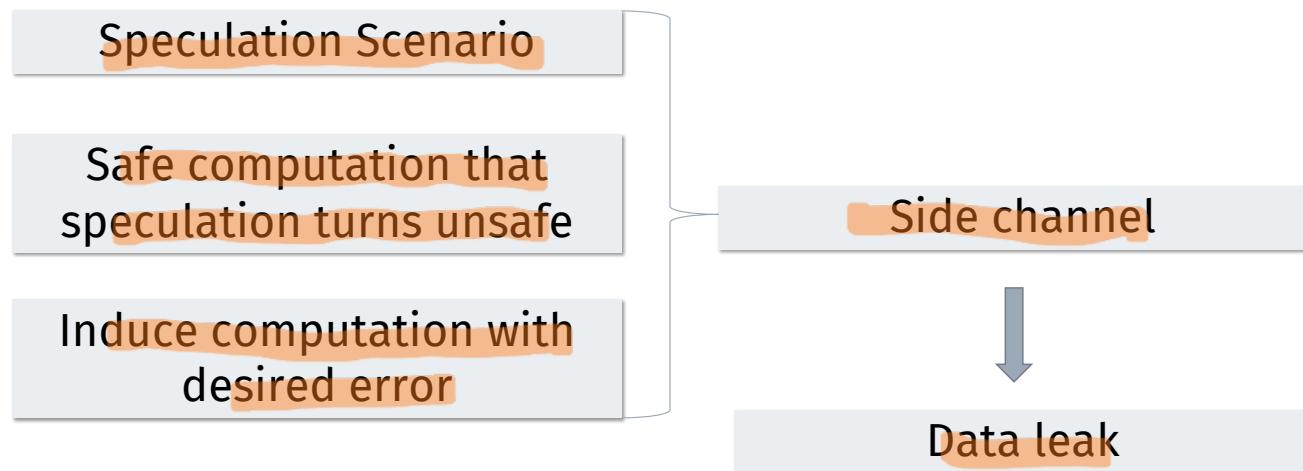
- Branch condition evaluation might require (slow) loads from memory
- Rather than wait for branch resolution, predict outcome of the branch



Then Spectre

- uses side channels to obtain the information from the accessed memory location
- allows an attacker to trick error-free programs into leaking their secrets
- breaks the isolation between different application

# Spectre pattern



Not a bug: everything complies w/Architc. spec.

# Is Spectre a bug? No, it is a symptom

Architecture is ambiguous under this regard.  
Priority given to performance instead of security

- No, everything complies with the architecture specifications
- It is a symptom of excessive architectural ambiguity: typical architectures' guarantees are not enough for security
- Priority given to performance
- In summary, Spectre reveals an architecture-software security gap (SW and HW developers make different assumptions)

# Mitigating Spectre

- Mitigations for Spectre are hard
- Also, the number of different variants does not help
- One simple idea is to stop speculation from accessing to data it should not:

1. Give up speculation, but its costly

## Fences

insertion of a serializing instruction

↳ insert points in which you  
stop speculation, you force CPU to  
serialize instructions

## Speculative Load Hardening (SLH)

conditional masking of the untrusted\_offset

↳ creates artificial data dependencies to  
stop speculation.

# Mitigating Spectre: fences

A **fence** or speculation barrier is a processor instruction that inhibits load speculation

- On x86\_64, we use lfence [load fence] on all branches to inhibit speculation;
- Despite the name, lfence serializes execution completely, and not just loads
- In pseudo-C:
  -

lfence works against load  
speculation and direct  
branch prediction

```
if (x < A_sz) {  
    __asm__("lfence");  
    y = B[A[x] * 512];  
}  
else {  
    __asm__("lfence");  
    /* ... */  
}
```

**Load speculation** is when the CPU tries to predict the outcome of a memory load instruction *before* knowing whether it's actually safe or legal to do so. It does this to run faster, but this can accidentally expose data if the speculation goes wrong.

# Mitigating Spectre: fences

- **The Good:**
  - It works (mostly :)
  - Probably nothing else!
- **The Bad:**
  - Highly inefficient (it stops speculation completely)
  - Requires re-compilation

```
if (x < A_sz) {  
    __asm__ ("lfence");  
    y = B[A[x] * 512];  
}  
else {  
    __asm__ ("lfence");  
    /* ... */  
}
```

- You have to decide where to put them, you can use heuristics but it is not able to cover everything
- Heuristics placement is risky: no guarantees
  - To make this a bit better one could use static analysis to decide where fences should go
    - E.g., Microsoft MSVC does that by detecting known problematic patterns

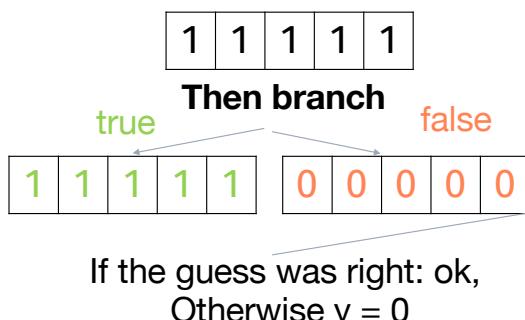
# Mitigating Spectre: SLH

→ Introduce a mask before branches

- Speculative Load Hardening (SLH) is a bit more sophisticated than fences

- Idea:

- It may be more efficient to introduce an **artificial data dependency** between the condition of a jump and the pointer:
- Simplifying, for Spectre v1 (assuming “ ?: ” to be implemented without branching, e.g., using a cmov in x86)



```
uint mask = ALL_ONES; /* all bits set to 1 */ You set all 16 1.  
if (cond) {  
    mask = !cond ? ALL_ZEROES : mask;  
    // ...  
    y = B[(A[x] * 512) & mask];  
}  
else {  
    mask = cond ? ALL_ZEROES : mask;  
    /* ... */  
}
```

Conditional update **cannot** be speculated

before doing anything you set mask that changes bits to 0 if cond is not satisfied.

So we won't execute speculatively

# Mitigating Spectre: Software

- **The Good:**

- It **works** (with some changes also for other Spectre variants)
- **More efficient than fences**
- No need of **static analysis code** (but still could help)

- **The Bad:**

- Still **slow** (mask must be known when accessing to B, also when speculating!)
- Must be **carefully implemented**
- Still requires re-compilation

```
if (cond) {  
    mask = !cond ? ALL_ZEROES : mask;  
    // ...  
    y = B[(A[x] * 512) & mask];  
}  
else {  
    mask = cond ? ALL_ZEROES : mask;  
    /* ... */  
}
```

# Fixing Spectre: other ideas

Other possible fixes:

→ isolate things

- Place secrets in different processes, to make them isolated and non reachable:
  - Just like **Chrome** and **Firefox**
- Reduce the bandwidth of side-channels:
  - e.g., via shadow micro-arch state that can be discarded, less accurate timers, adding noise, ...

Make life of the attacker harder

Is Spectre fixed?

Only if you are willing to re-compile/re-design your system and make it slow!

# Outline

Speculation

Spectre

Blade Approach



# Automatically Eliminating Speculative Leaks from Cryptographic Code with Blade: a short guide

---

[Vassena et al., 2021] M. Vassena, C. Disselkoen, K. von Gleissenthall, S. Cauligi, R. Gökhan Kıcı, R. Jhala, D. Tullsen, and D. Stefan. "Automatically eliminating speculative leaks from cryptographic code with blade." POPL 2021. URL: <https://arxiv.org/abs/2005.00294>

We use formal methods

# Blade approach

Blade



They provide ms w/ static

Transient-Flow Type System

Static detection of Spectre vulnerabilities

Technique through a type sys

With minimal number of fences

Automatic repair of vulnerable programs

↳ Automatic repair alg. to transform programs  
in safe ones.

stops spec. in precise and minimum # of pts  
then introduces spec.

Insight:

- Blade stops speculation only along specific paths
- It cuts the dataflow from expressions that speculatively introduce secrets (**sources**) to those that leak them through the cache (**sinks**), rather than prohibit speculation altogether



How? With a **static type system** that types each expression as

- **transient** (possibly containing speculative secrets) or
- **stable**, and

prohibits speculative leaks by requiring that all **sink** expressions are **stable**

- Blade relies on an abstract primitive, **protect**, to halt speculation at fine granularity, which is implemented using existing architectural mechanisms
- Blade's type system can automatically synthesize a minimal number of protects to provably eliminate speculative leaks

in your words

# The new primitive: protect

---

x := **protect(e)**

primitive as a sort of fine-grained fence

It is a sort of fine-grained fence based on SLH

**Semantics:** evaluate e and assign it to x, only after its value is non speculative, i.e., it is **stable**

# Insertion of protect

---

**Blade** inserts protect to stop speculation in specific points

```
if (i < A_sz)
    x = A[i];
    y = B[x];
...
...
```

# Insertion of protect

---

**Blade** inserts protect to stop speculation in specific points

1. Identify **sources** and **sinks**

```
if (i < A_sz)
    x = A[i];
    y = B[x];
...
```

# Insertion of protect

---

**Blade** inserts protect to stop speculation in specific points

1. Identify **sources** and **sinks**
2. Detect data flow from sources to sinks

```
if (i < A_sz)
    x = A[i];
    y = B[x];
...
```

# Insertion of protect

---

**Blade** inserts protect to stop speculation in specific points

1. Identify **sources** and **sinks**
2. Detect **data flow** from sources to sinks

```
if (i < A_sz)
    x = A[i];
    y = B[i]; green
...
...
```

# Protect to the rescue

Blade inserts protect to stop speculation in specific points

1. Identify sources and sinks
2. Detect data flow from sources to sinks
3. Insert protect to cut these data flows

```
if (i < A_sz)
    x = A[i];
    y = B[x];
...
```



There is a formal fw to do this

Sources are  $x$  and  $A[i]$ , the sink is  $x$ .

```
if (i < A_sz)
    x = protect(A[i]);
    y = B[x];
...
```

Only instructions depending on  $x$  are stopped!



- Type sys statically approximates the runtime propg. of spec. secrets.

# A type system for speculation

Static approximation of the runtime propagation of **speculative secrets**

Expressions are

→ Think of transient like taint

- **transient** possibly containing speculative secrets
- **stable** not containing speculative secrets

Transient = tainted

Propagate taints in expressions

## Static Taint Tracking

Judgment of exp. → Judgment  
 $\Gamma \vdash e : \tau$  ①       $\Gamma \vdash c$  ② for command.  
 ↳ Now has a type

Restricts using transient data in commands

Well-typed expressions are free of speculation leaks

① Exp e has type  $\tau$ .    ② Command is well typed

# Typing rules

→ Crucial in speculation: even though both use stable result vs transient; every time you access memory through array you can leak something.

Array rule

$$\frac{\Gamma \vdash e_1 : S \quad \Gamma \vdash e_2 : S}{\Gamma \vdash e_1 [e_2] : T}$$

Array and index must be stable

Array reads can introduce transient data

No branching on transient data

Condition needs to be stable. You want to avoid misspeculation

$$\frac{\Gamma \vdash e : S \quad \Gamma \vdash c_1 \quad \Gamma \vdash c_2}{\Gamma \vdash \text{if } e \text{ then } c_1 \text{ else } c_2}$$

# Automatic repair algorithm

---

Well-typed programs are free of speculation leaks

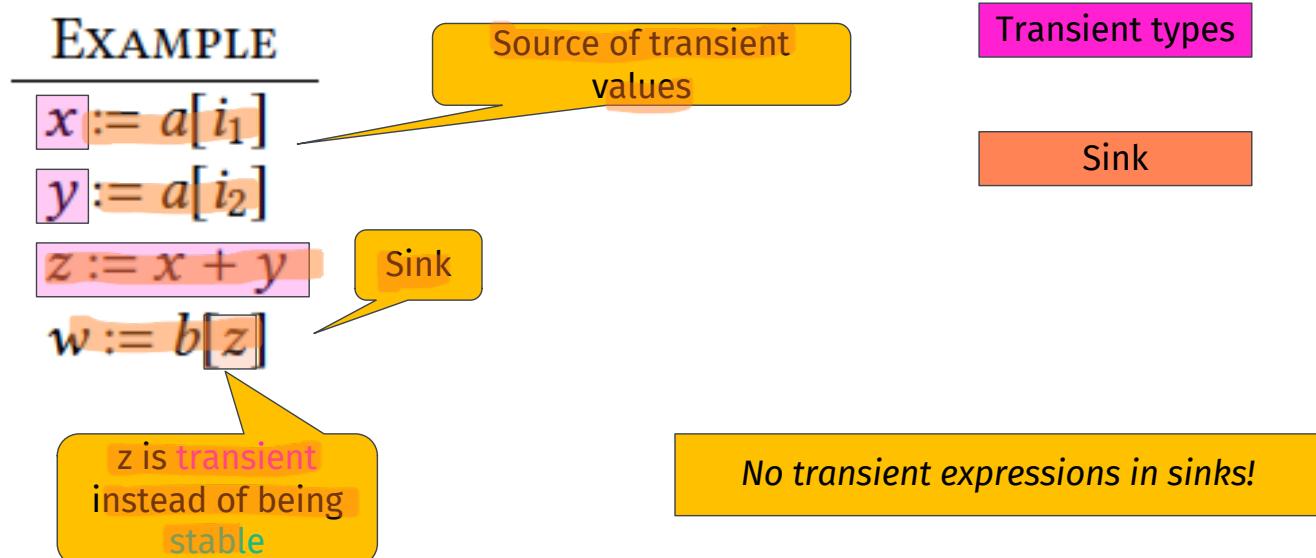
And ill-typed programs?

Ill-typed programs can be repaired by inserting **protect** statements, which **cut** the data-flow between **sources** and **sinks**

# Transient-Flow Type System by Example

1. Build the def-use graph that captures the data dependencies

↳ definition-use variable: show when you define and when you use variables.



# Repair by Example

---

1. Build the **def-use graph** that captures the data dependencies

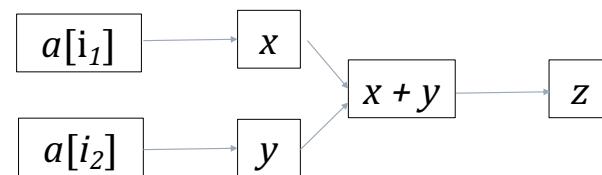
## EXAMPLE

$x := a[i_1]$

$y := a[i_2]$

$z := x + y$

$w := b[z]$



# Repair by Example

1. Build the **def-use graph** that captures the data dependencies
2. Add **source** and **sink** nodes and connect

## EXAMPLE

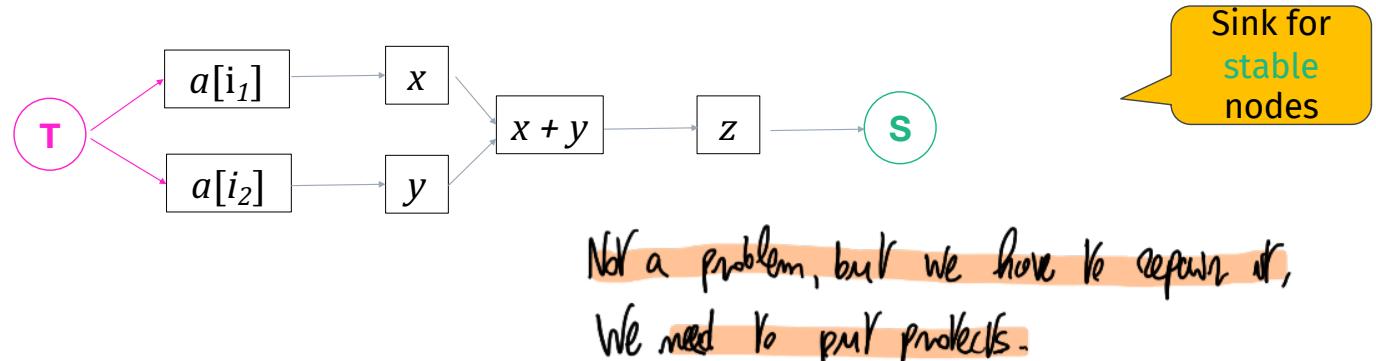
$x := a[i_1]$

$y := a[i_2]$

$z := x + y$

$w := b[z]$

- the **source node** to expressions that read **memory**
- expressions that must be stable to the **sink node**



# Repair by Example

1. Build the **def-use graph** that captures the data dependencies
2. Add **source** and **sink** nodes and connect

## EXAMPLE

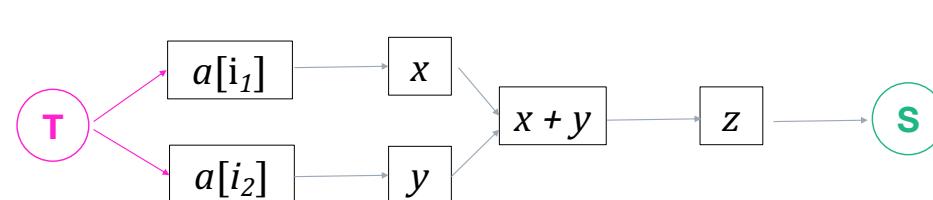
$x := a[i_1]$

$y := a[i_2]$

$z := x + y$

$w := b[z]$

- the **source node** to expressions that read memory
- expressions that must be stable to the **sink node**



Sink for  
stable  
nodes

A path from **source** to **sink** indicates that the program may be **leaky**



# Repair by Example

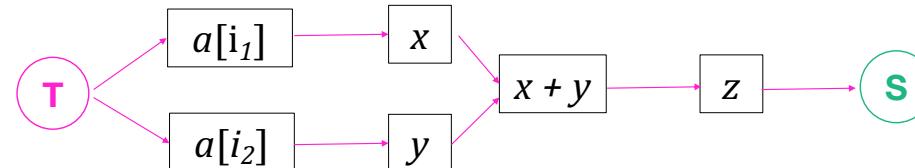
1. Cut the path from **sink** to **source**, i.e., find a set of variables that disconnect **source** and **sink**

## EXAMPLE

$x := a[i_1]$   
 $y := a[i_2]$   
 $z := x + y$   
 $w := b[z]$

Equivalent to the problem of finding a **cut set**  
A **cut set** is set of variables, the removal of which from the graph eliminates all paths from **T** to **S**

↳ find variables whose removal eliminates  
the path you don't want.



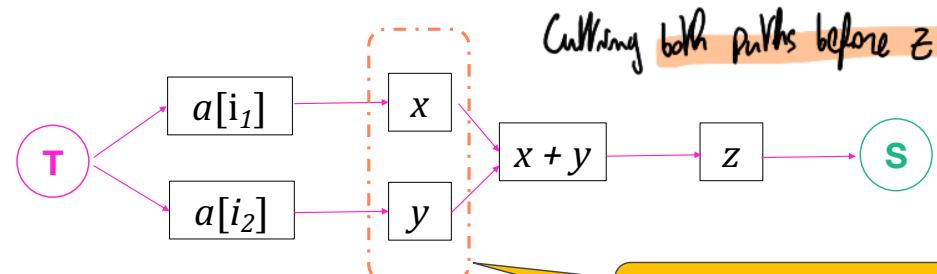
# Repair by Example

1. Cut the path from **sink** to **source**, i.e., find a set of variables that disconnect **source** and **sink**

## EXAMPLE

$x := a[i_1]$   
 $y := a[i_2]$   
 $z := x + y$   
 $w := b[z]$

Equivalent to the problem of finding a **cut set**  
A **cut set** is set of variables, the removal of which from the graph eliminates all paths from **T** to **S**



Insert a protect  
 $x := \text{protect}(a[i_1])$   
 $y := \text{protect}(a[i_2])$

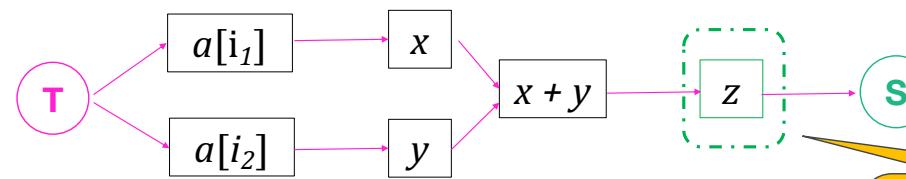
# Repair by Example

1. Cut the path from **sink** to **source**, i.e., find a set of variables that disconnect **source** and **sink**

## EXAMPLE

$x := a[i_1]$   
 $y := a[i_2]$   
 $z := x + y$   
 $w := b[z]$

Equivalent to the problem of finding a **cut set**  
A **cut set** is set of variables, the removal of which from the graph eliminates all paths from **T** to **S**



$z := \text{protect}(x+y)$

You can exploit a  
lot of known solutions; you  
can use

**Minimal cut**  
Classic max-flo min cut  
problem

max-flow min-cut polynomial algorithm! Efficient.

# Challenges

You can use Open Semantics to model low level behavior, but you will need to introduce abstractions for it.

Soundness of the type system

Speculative out-of-order execution

High level reasoning

Low level features

Source-level semantics model of speculation

The problem is semantics of speculation wasn't concerned with security. You have to think about semantics of commands in memory w/ cache

# Modelling Speculation

A relatively simple language is chosen. Why? It

1. Simplifies the reasoning
2. Allows to reason directly at source level

We will use a simple language

Not too many details about the architecture must be known

- A Just In Time (JIT) semantics taking high-level source programs, translating them on-the-fly to low-level processor instructions and executing them
- Two structures:
  - Instruction reorder buffer
  - Command Stack

↑ takes high level source prop. and

Thus is the semantics we will use (one is high level, one low).

# Modeling speculation

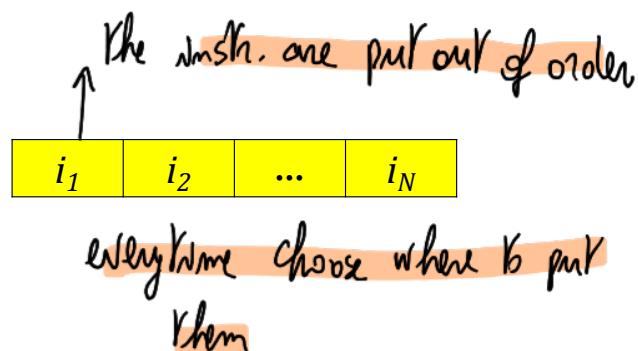
Used structures:

- Instruction reorder buffer (is)
- Command Stack (cs)
- memory ( $\mu$ )
- map from variables to values ( $\rho$ )

Three stages (fetch, exec, retire)

The branch prediction is not modeled:

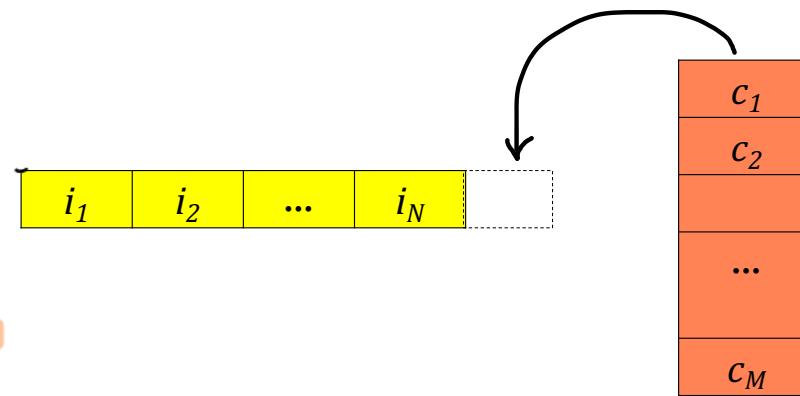
- the attacker can control the processor with a set of directives We assume this. That might impose processor to take branches



# Fetch stage

- Take a command from the top of  $cs$ ,
- flatten it into a sequence of instructions, and
- store them in  $is$

Both  $is$  and  $cs$  are updated in  $is'$  and  $cs'$

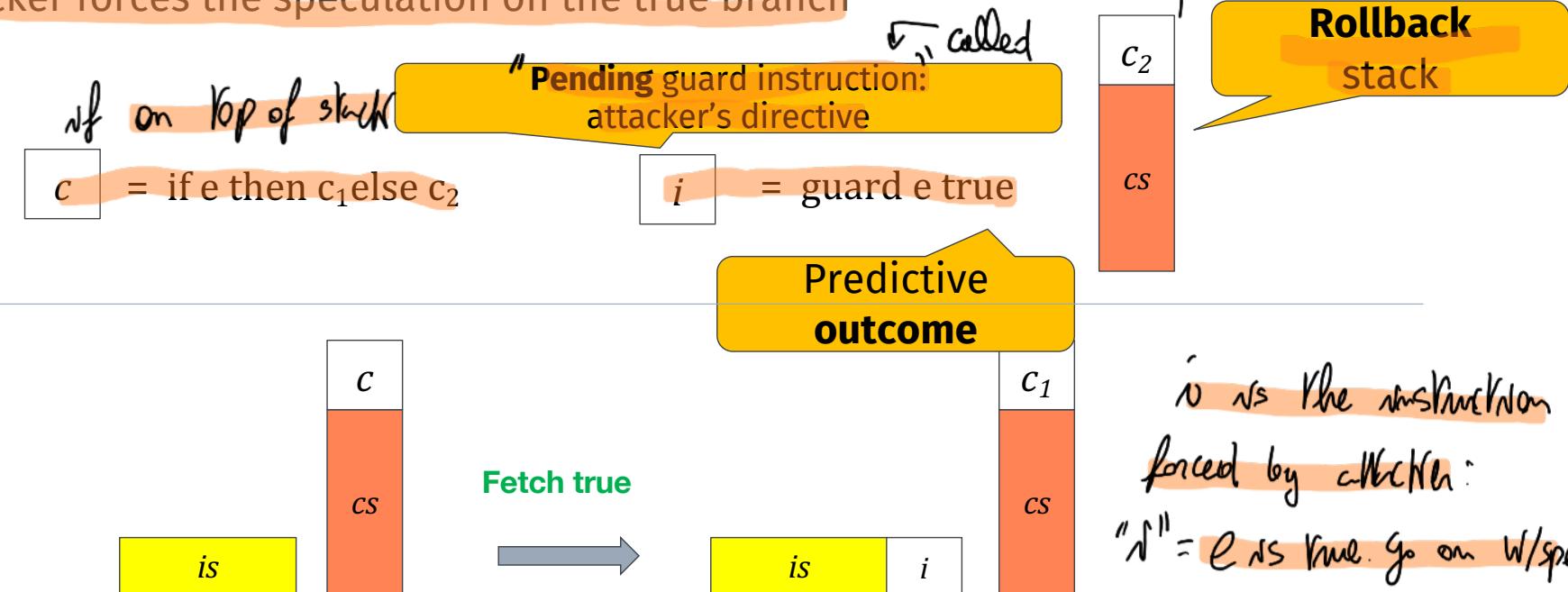


You start from one config and update it into  $(cs', is')$



# Fetch-if-true

The attacker forces the speculation on the true branch



# Execute stage

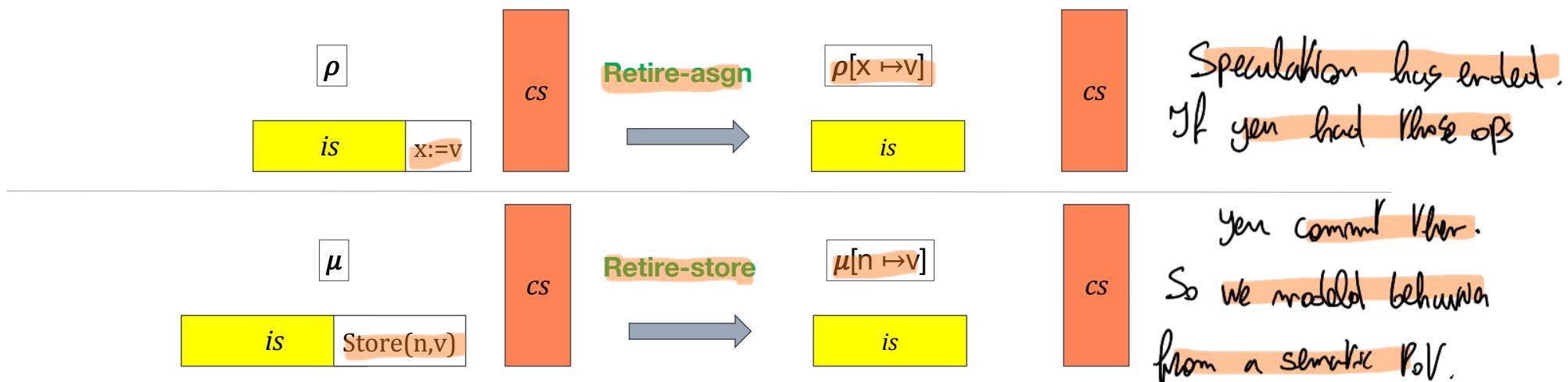
- take the  $n$ -th instruction from  $is$  (provided by the attacker) and
- execute it **speculatively**:  $\rightarrow$  Since you are going ahead, until you commit no changes
  - memory  $\mu$  and map from variables to values  $\rho$  **do not change**

Intuitively **operands of instruction** are evaluated just using **non-transient information**  $\delta$



# Retire stage

- Remove complete instructions from  $is$  and
- commit changes to memory  $\mu$  and variable map  $\rho$  that are modified in  $\mu'$  and  $\rho'$  (done in-order, so giving the illusion of no instruction-level parallelism!)



Using formal method gives guarantees.

# Formal guarantees

- **Consistency:** the JIT semantics agree with the standard sequential (in-order) semantics → agrees with a standard semantics of command, you can commit or rollback too.  
We will see it
- **Correctness:** the repair algorithm synthesizes well-typed programs (provably secure repair)
- **Soundness:** well-typed constant time programs are also speculative constant-time<sup>①</sup>

Constant time programs do not allow attackers to infer info on exec. time. with this approach this is true for spec too.  
↳ with Specie you need to reason over SCT<sup>①</sup>

# The “source” language

---

Values  $v ::= n \mid b \mid a$

Expr.  $e ::= v \mid x \mid e + e \mid e < e$   
|  $e \otimes e \mid e?e:e$   
|  $length(e) \mid base(e)$

Rhs.  $r ::= e \mid *e \mid a[e]$

Cmd.  $c ::= \text{skip} \mid x := r \mid *e = e$   
|  $a[e] := e \mid \text{fail} \mid c;c$   
|  $\text{if } e \text{ then } c \text{ else } c$   
|  $\text{while } e \text{ do } c$   
|  $x := \text{protect}(r)$

- This is **easy**, just the WHILE-language!
- Sure?
  - **Question:** how do you think **protect** works?
    - Hint: SLH prevent “wrong” accesses. Can you “abstract” SLH in a single command?
  - **Answer:** completes the assignment only when the evaluation of  $r$  is **stable**, i.e., no more rollbacks are possible, and  $r$  is the **final, non speculative value**



# The “target” language

Instructions  $i ::= \text{nop} \mid x := e \mid x := \text{load}(e)$   
 $\mid \text{store}(e, e) \mid x := \text{protect}(e)$   
 $\mid \text{guard}(e^b, cs, p) \mid \text{fail}(p)$

↓ This models a Watch

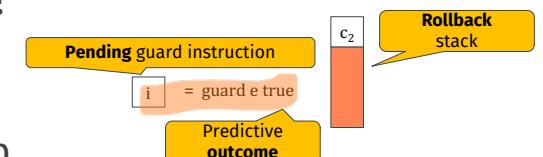
choosing the branch, slide SG.

## move

- Again, we have protect!

- guard** ( $e^b, cs, p$ ):

- expresses a pending jump
- $e$  is the condition,  $b$  is its predicted outcome,  $cs$  is what to do if a rollback is needed (ignore  $p$ , an identifier useful for the security analysis)



The source language can be translated on-the-fly to a target one, using a JIT semantics that models speculation

# The «Compiler»: A JIT semantics for speculation

The semantics formalizes the execution of source programs on a pipelined processor

$$C \xrightarrow{o}{^d} C'$$

Read as: “under the directive  $d$ , the configuration  $C$  moves to  $C'$  with observation  $o$ ” (what the attacker can see), where

~~for general execute a group of instructions.~~

- $d ::= \text{fetch} \mid \text{fetch } b \mid \text{exec } n \mid \text{retire}$ , attacker-supplied directives (also each “class” of directives is a **stage** of the processor’s pipeline)

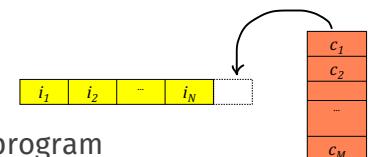
*in this model attacker is choosing what we can do*

- $o ::= \epsilon \mid \text{read } (n, ps) \mid \text{write } (n, ps) \mid \text{fail} \mid \text{rollback } (p)$ , attacker’s observations

- $C, C'$  are of the form  $\langle is, cs, \mu, \rho \rangle$ , where:

*Configurations*

- $is$  is a **reorder buffer**, i.e., a queue of in-flight (target) **instructions** (directives give the order)



- $cs$  represents the **stack** of (source) commands (the current execution path), i.e., the continuation of the program

- $\mu$  is a memory and  $\rho$  is a map from variables to values

# JIT semantics for speculation: the fetch stage

$$\langle is, cs, \mu, \rho \rangle \xrightarrow{o}^{\text{fetch}} \langle is', cs', \mu, \rho \rangle$$

$$\langle is, cs, \mu, \rho \rangle \xrightarrow{o}^{\text{fetch } b} \langle is', cs', \mu, \rho \rangle$$

*either fetch non-blk or force a ground*

- **Idea:** take a command from the top of  $cs$ , flatten it into a sequence of instructions and store them in  $is$

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-ASGN

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

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-ARRAY-LOAD

$$c = x := a[e] \quad e_1 = e < \text{length}(a) \quad e_2 = \text{base}(a) + e \quad c' = \text{if } e_1 \text{ then } x := *e_2 \text{ else fail}$$

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

FETCH-IF-TRUE

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

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

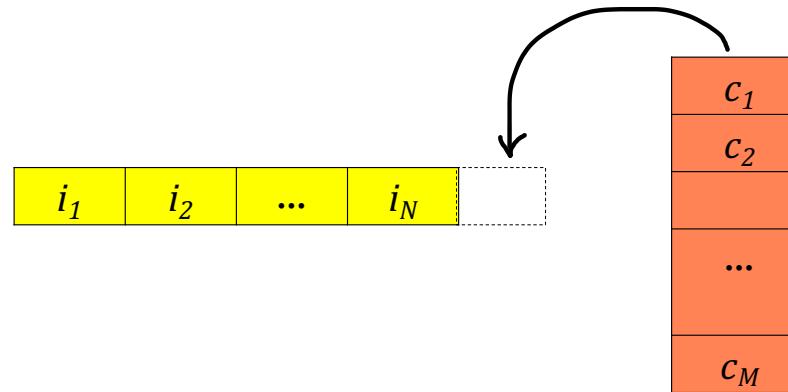
The attacker decides  
which branch

# Fetch stage

$$\langle is, cs, \mu, \rho \rangle \xrightarrow{o}^{\text{fetch}} \langle is', cs', \mu, \rho \rangle$$

- Take a command from the top of  $cs$ ,
- flatten it into a sequence of instructions, and
- store them in  $is$

Both  $is$  and  $cs$  are updated in  $is'$  and  $cs'$

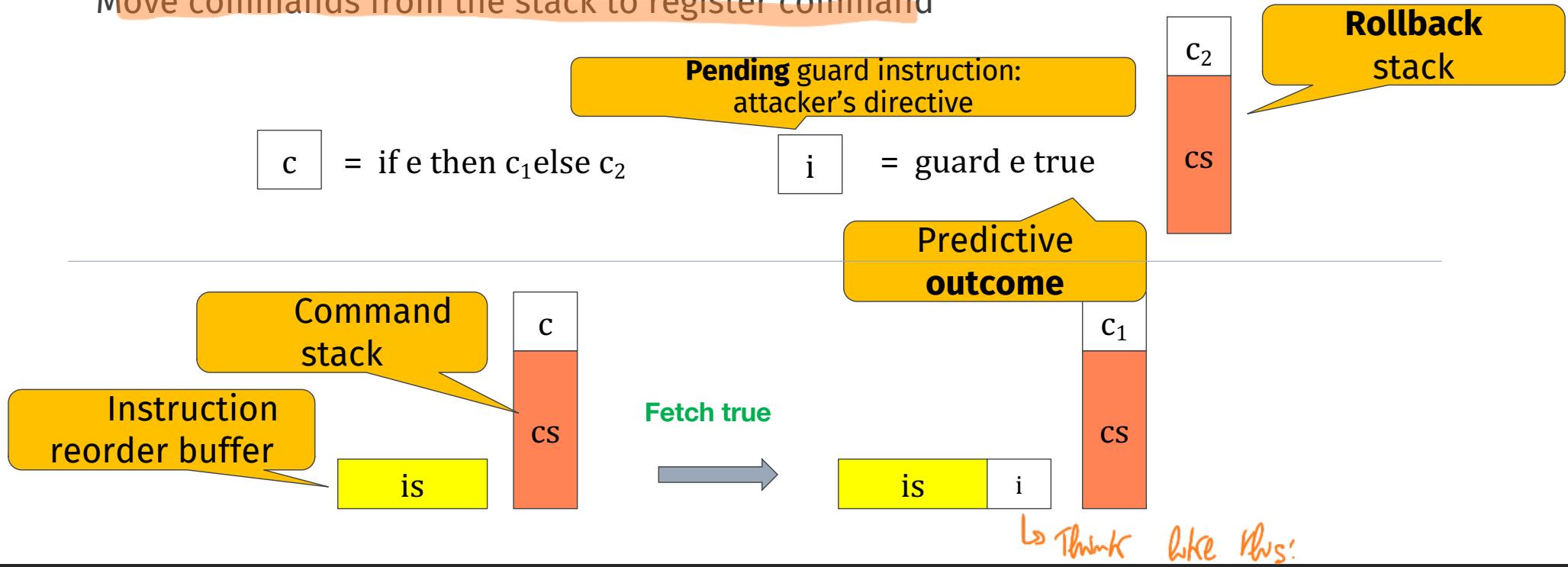


# Fetch-if-true

$$\begin{array}{c}
 \text{FETCH-IF-TRUE} \\
 c = \text{if } e \text{ then } c_1 \text{ else } c_2 \\
 \hline
 \langle is, c : cs, \mu, \rho \rangle \xrightarrow{\text{fetch true}}_e \langle is + [i], c_1 : cs, \mu, \rho \rangle
 \end{array}$$

you keep guard true,  
 you keep the info  
 on  $c_2$ ,  
 longer about NV

Move commands from the stack to register command



When you execute ~ you fetch  $c_1$ ?

# JIT semantics for speculation: the execute stage

$$\langle is, cs, \mu, \rho \rangle \xrightarrow{o}^{\text{exec } n} \langle is', cs', \mu, \rho \rangle$$

You change NS, CS but no  
Touch of memory or env,

- Idea: take the  $n$ -th instruction from  $is$  (provided by the attacker) and execute it **speculatively**
  - Lots of details here, but intuitively **operands of instruction** are evaluated just using **non-transient information** (computed by an ad hoc function  $\phi$  that marks variables from pending assignments)

EXECUTE

$$\frac{|is_1| = n - 1 \quad \rho' = \phi(is_1, \rho) \quad \langle is_1, i, is_2, cs \rangle \xrightarrow{(\mu, \rho', o)} \langle is', cs' \rangle}{\langle is_1 + [i] + is_2, cs, \mu, \rho \rangle \xrightarrow{\text{exec } n} \langle is', cs', \mu, \rho \rangle}$$

$$\begin{array}{c} \text{EXEC-LOAD} \\ i = (x := \text{load}(e)) \quad \text{store}(\_, \_) \notin is_1 \quad n = \llbracket e \rrbracket^\rho \quad ps = (\llbracket is_1 \rrbracket) \quad i' = (x := \mu(n)) \\ \hline \langle is_1, i, is_2, cs \rangle \xrightarrow{(\mu, \rho, \text{read}(n, ps))} \langle is_1 + [i'] + is_2, cs \rangle \end{array}$$

$$\begin{array}{c} \text{EXEC-BRANCH-OK} \\ i = \text{guard}(e^b, cs', p) \quad \llbracket e \rrbracket^\rho = b \\ \hline \langle is_1, i, is_2, cs \rangle \xrightarrow{(\mu, \rho, \epsilon)} \langle is_1 + [\text{nop}] + is_2, cs \rangle \end{array}$$

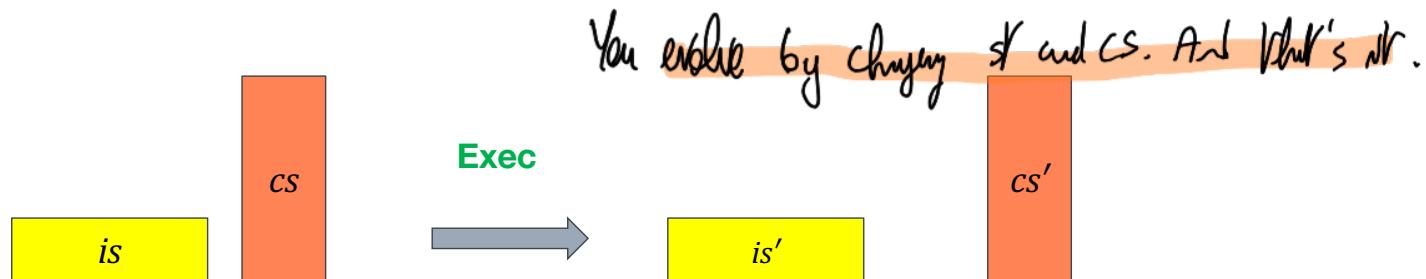
$$\begin{array}{c} \text{EXEC-BRANCH-MISPREDICT} \\ i = \text{guard}(e^b, cs', p) \quad b' = \llbracket e \rrbracket^\rho \quad b' \neq b \\ \hline \langle is_1, i, is_2, cs \rangle \xrightarrow{(\mu, \rho, \text{rollback}(p))} \langle is_1 + [\text{nop}], cs' \rangle \end{array}$$

# Execute stage

$$\langle is, cs, \mu, \rho \rangle \xrightarrow{o}^{\text{fetch}} \langle is', cs', \mu, \rho \rangle$$

- take the  $n$ -th instruction from  $is$  (provided by the attacker) and
- execute it **speculatively**:
  - memory  $\mu$  and map from variables to values  $\rho$  **do not change**

Intuitively **operands of instruction** are evaluated just using **non-transient information**



# JIT semantics for speculation: the retire stage

---

$$\langle is, cs, \mu, \rho \rangle \xrightarrow{o}^{\text{retire}} \langle is', cs', \mu', \rho' \rangle$$

- **Idea:** remove complete instructions from  $is$  and commit changes to memory and variable map
  - Do that in-order, so giving the illusion of no instruction-level parallelism!

RETIRE-NOP

$$\langle \text{nop} : is, cs, \mu, \rho \rangle \xrightarrow{\text{retire}}_{\epsilon} \langle is, cs, \mu, \rho \rangle$$

RETIRE-ASGN

$$\langle x := v : is, cs, \mu, \rho \rangle \xrightarrow{\text{retire}}_{\epsilon} \langle is, cs, \mu, \rho[x \mapsto v] \rangle$$

RETIRE-STORE

$$\langle \text{store}(n, v) : is, cs, \mu, \rho \rangle \xrightarrow{\text{retire}}_{\epsilon} \langle is, cs, \mu[n \mapsto v], \rho \rangle$$

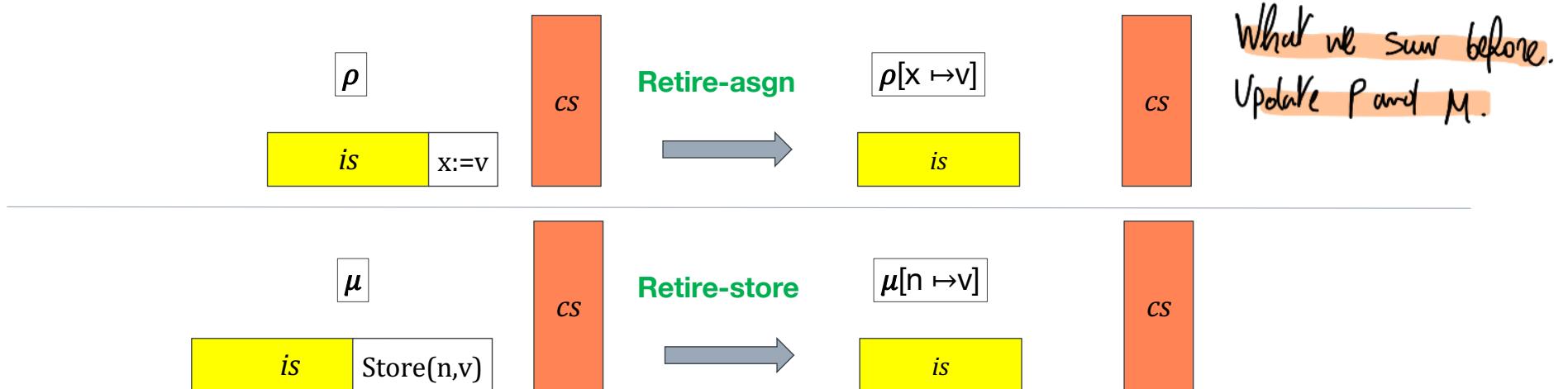
RETIRE-FAIL

$$\langle \text{fail}(p) : is, cs, \mu, \rho \rangle \xrightarrow{\text{retire}}_{\text{fail}(p)} \langle [ ], [ ], \mu, \rho \rangle$$

# Retire stage

$$\langle is, cs, \mu, \rho \rangle \xrightarrow{o}{\text{retire}} \langle is', cs', \mu', \rho' \rangle$$

- Remove complete instructions from  $is$  and
- **commit changes to memory  $\mu$  and variable map  $\rho$**  that are modified in  $\mu'$  and  $\rho'$  (done in-order, so giving the illusion of no instruction-level parallelism!)



You can implement protect in offish ways

# A few words about protect

---

- protect can be implemented both in software and in hardware
- **Software implementation: no support from hardware needed**
  - exactly as SLH for array access (i.e., computes a **mask** depending on the array's **bound check**, and then executes the access depending using the mask)  
↳ **Masks**
- **Hardware implementation:** Can you guess? What needs to be added?
  - **Fetch stage:** translates protect into protect, easy
  - **Exec aux. relation:** reduces the parameter to a value, then if no guard is pending removes the 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

# Protect: software implementation

FETCH-PROTECT-SLH

$$\frac{\begin{array}{c} c = x := \mathbf{protect}(a[e]) \quad e_1 = e < \text{length}(a) \quad e_2 = \text{base}(a) + e \\ c_1 = m := e_1 \quad c_2 = m := m ? \mathbf{1} : \mathbf{0} \quad c_3 = x := *(e_2 \otimes m) \quad c' = c_1; \text{if } m \text{ then } c_2; c_3 \text{ else fail} \end{array}}{\langle is, c : cs, \mu, \rho \rangle \xrightarrow[\epsilon]{\text{fetch}} \langle is, c' : cs, \mu, \rho \rangle}$$

(b) Software implementation of  $\mathbf{protect}(a[e])$ .

```

uint mask = ALL_ONES; /* all bits set to 1 */
if (cond) {
    mask = !cond ? ALL_ZEROES : mask;
    // ...
    y = B[(A[x] * 512) & mask];
}
else {
    mask = cond ? ALL_ZEROES : mask;
    /* ... */
}

```

# Example: the source

---

```
if (i1 < length(a))
{x := a[i1];}
if (i2 < length(a))
{y := a[i2];}
z := x + y
if (z < length(b))
{w := b[z]}
```

every load instruction is the then branch  
of a bound check (we work w/ Bound check bypass)

# Example: the target

Reformulate a program in a way  
that complies with semantics.

```
if (i1 < length(a))  
{x := a[i1];}  
if (i2 < length(a))  
{y := a[i2];}  
z := x + y  
if (z < length(b))  
{w := b[z]}
```

x := load(e)  
guard(eb, cs, p)

Attacker's directives

- 1 guard( $(i_1 < \text{length}(a))$  true, [fail], 1) ↗able
- 2  $x := \text{load}(\text{base}(a) + i_1)$
- 3 guard( $(i_2 < \text{length}(a))$  true, [fail], 2)
- 4  $y := \text{load}(\text{base}(a) + i_2)$
- 5  $z := x + y$  ↗ base addr. + offset
- 6 guard( $(z < \text{length}(b))$  true, [fail], 3)
- 7  $w := \text{load}(\text{base}(b) + z)$

every time you  
have an if you  
have a guard

The attacker instructs the processor to speculatively execute the load instructions (2,4,5,7) and the assignment, but not the guards

# Example of leaking execution

Evolution of the buffer is after each attacker directive (execute 2,4,5,7); the memory  $\mu$  and variable map  $\rho$  are shown above

$C \xrightarrow{o} C'$   
 $C, C'$  are of the form  
 $<is, cs, \mu, \rho>$

Memory Layout		Variable Map
$b[0]$	$b[0]$	$\rho(i_1) = 1$
$a[0]$	$a[0]$	$\rho(i_2) = 2$
$a[1]$	$a[1]$	$\dots$
<b>SECRET</b> $s[0]$	$s[0]$	
	$\dots$	

Reorder Buffer     $is$

```
1 guard(( $i_1 < length(a)$ )true, [ fail ], 1)
2 x := load(base(a) +  $i_1$ ) This is an ok move:  $i_1=1$ ,  $a[1]$  is okay
3 guard(( $i_2 < length(a)$ )true, [ fail ], 2)
4 y := load(base(a) +  $i_2$ )
5 z := x + y
6 guard(( $z < length(b)$ )true, [ fail ], 3)
7 w := load(base(b) + z)
```

Observations:  $o$

# Example: exec 2

Directive **exec 2** executes the first load by computing the memory address 2 and replacing the instruction with the assignment  $x := \mu(2)$  containing the loaded value 0

$C \xrightarrow{o} C'$   
 $C, C'$  are of the form  
 $<is, cs, \mu, \rho>$

Memory Layout

$\mu(0) = 0$	$b[0]$
$\mu(1) = 0$	$a[0]$
$\mu(2) = 0$	$a[1]$
$\mu(3) = 42$	$s[0]$
...	...

Variable Map

$$\begin{aligned}\rho(i_1) &= 1 \\ \rho(i_2) &= 2 \\ &\dots\end{aligned}$$

Reorder Buffer	is	exec 2
1	guard( $(i_1 < \text{length}(a))^\text{true}$ , [fail], 1)	
2	x := load(base(a) + i <sub>1</sub> )	x := $\mu(2)$
3	guard( $(i_2 < \text{length}(a))^\text{true}$ , [fail], 2)	
4	y := load(base(a) + i <sub>2</sub> )	
5	z := x + y	
6	guard( $(z < \text{length}(b))^\text{true}$ , [fail], 3)	
7	w := load(base(b) + z)	

Observations: o      read(2, [1])

# Example: exec 4

Directive **exec 4** transiently reads public 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

Memory Layout		Variable Map	
$\mu(0) = 0$	$b[0]$	$\rho(i_1) = 1$	
$\mu(1) = 0$	$a[0]$	$\rho(i_2) = 2$	
$\mu(2) = 0$	$a[1]$	$\dots$	
$\mu(3) = 42$	$s[0]$		
$\dots$	$\dots$		

Reorder Buffer	$i_S$	exec 2	exec 4	
1	<code>guard(<math>(i_1 &lt; \text{length}(a))^\text{true}</math>, [ fail ], 1)</code>			
2	<code><math>x := \text{load}(\text{base}(a) + i_1)</math></code>	$x := \mu(2)$		
3	<code>guard(<math>(i_2 &lt; \text{length}(a))^\text{true}</math>, [ fail ], 2)</code>			
4	<code><math>y := \text{load}(\text{base}(a) + i_2)</math></code>	<i><math>i_2=2</math>, so base + 2 = 2nd pos of the memory</i>		$y := \mu(3)$
5	<code><math>z := x + y</math></code>			
6	<code>guard(<math>(z &lt; \text{length}(b))^\text{true}</math>, [ fail ], 3)</code>			
7	<code><math>w := \text{load}(\text{base}(b) + z)</math></code>			

Observations: $o$	$\text{read}(2, [1])$	$\text{read}(3, [1, 2])$	
-------------------	-----------------------	--------------------------	--

# Example: exec 5

The processor forwards the values of  $x$  and  $y$  through the transient variable map  $\rho$  [ $x \mapsto \mu(2)$ ,  $y \mapsto \mu(3)$ ] to compute their sum in the fifth instruction, ( $z := 42$ )

Memory Layout		Variable Map		
$\mu(0) = 0$	$b[0]$	$\rho(i_1) = 1$		
$\mu(1) = 0$	$a[0]$	$\rho(i_2) = 2$		
$\mu(2) = 0$	$a[1]$	$\dots$		
$\mu(3) = 42$	$s[0]$			
$\dots$	$\dots$			

Reorder Buffer	is	exec 2	exec 4	exec 5
1	<code>guard(<math>(i_1 &lt; length(a))^{true}</math>, [ fail ], 1)</code>			
2	<code><math>x := load(base(a) + i_1)</math></code>	$x := \mu(2)$		
3	<code>guard(<math>(i_2 &lt; length(a))^{true}</math>, [ fail ], 2)</code>			
4	<code><math>y := load(base(a) + i_2)</math></code>		$y := \mu(3)$	
5	<code><math>z := x + y</math></code>			$z := 42$
6	<code>guard(<math>(z &lt; length(b))^{true}</math>, [ fail ], 3)</code>			
7	<code><math>w := load(base(b) + z)</math></code>			

Observations: $o$	$read(2, [1])$	$read(3, [1, 2])$	$\epsilon$
-------------------	----------------	-------------------	------------

$\hookrightarrow$  active guards

# Example: exec 7

The value of z is then used as an index in the last instruction and leaked to the attacker via observation  $\text{read}(42, [1, 2, 3])$

Memory Layout		Variable Map			
$\mu(0) = 0$	$b[0]$	$\rho(i_1) = 1$			
$\mu(1) = 0$	$a[0]$	$\rho(i_2) = 2$			
$\mu(2) = 0$	$a[1]$	$\dots$			
$\mu(3) = 42$	$s[0]$				
$\dots$	$\dots$				

Reorder Buffer	is	exec 2	exec 4	exec 5	exec 7
1	$\text{guard}((i_1 < \text{length}(a))^{\text{true}}, [\text{fail}], 1)$				
2	$x := \text{load}(\text{base}(a) + i_1)$	$x := \mu(2)$			
3	$\text{guard}((i_2 < \text{length}(a))^{\text{true}}, [\text{fail}], 2)$		$y := \mu(3)$		
4	$y := \text{load}(\text{base}(a) + i_2)$			$z := 42$	
5	$z := x + y$				
6	$\text{guard}((z < \text{length}(b))^{\text{true}}, [\text{fail}], 3)$				
7	$w := \text{load}(\text{base}(b) + z)$				$w := \mu(42)$

Observations: $o$	$\text{read}(2, [1])$	$\text{read}(3, [1, 2])$	$\epsilon$	$\text{read}(42, [1, 2, 3])$
-------------------	-----------------------	--------------------------	------------	------------------------------

Sector 42 was leaked through  
a read!

# Almost there...

---

The language(s) we just saw are nice, but

- Can we be sure that the programs are *speculatively secure*?

Two possibilities

We prove manually that a program  
is not vulnerable to Spectre

We let the machine to prove that a  
program is not vulnerable to  
Spectre

Classify vars in transient and stable.

# Let the machine prove things...

... via type checking!

## Expressions

Types map vars  
 $\Gamma \vdash e : \tau$

Transient-flow type system that statically rejects programs that can potentially leak through transient execution attacks; i.e. that exhibit any source-to-sink data flows

## Lattice

- We use a lattice

$$S \sqsubseteq T$$

Reads as: "expression  $e$  has type  $\tau$  under  $\Gamma$ "

- $\Gamma$  maps variables to types
- $\tau$  is a type, can be either
  - Stable  $S$ , i.e., contains no transient values during executions
  - Transient  $T$ , otherwise
  - $S \sqsubseteq T$  (can-flow-to relation) not viceversa

This order  
should be  
respected

if not, type sys rejects program

## Commands

$$\Gamma, \text{Prot} \vdash c$$

Reads as: "command  $c$  is well-typed under  $\Gamma$  and  $\text{Prot}$ "

- $\Gamma$  as before,  $\text{Prot}$  a set of implicitly protected vars ①
- Idea: if all assignments to variables in  $\text{Prot}$  are protected, then  $c$  does not leak

① You have to prove if sys is well typed when fixed system.

Vars over which you use protect or they can be vars assumed to be protected.

# Let the machine prove things... (cont)

... via a type **inference** (type constraints restrict the types)!

**Expressions:**

$$\Gamma \vdash e : \tau \Rightarrow k$$

In both cases, the algorithm

- i. generates a set of **type constraints**  $k$ , involving concrete types and atoms (i.e., **program variables and type variables**)
- ii. builds a **def-use** graph: nodes are types and atoms, edges their relations; If there is a path from  $\textcolor{red}{T}$  to  $\textcolor{teal}{S}$  the program is leaky (Thm:  $k$  is satisfiable iff no path from  $\textcolor{red}{T}$  to  $\textcolor{teal}{S}$  exists) and reconstructs type information
- iii. cuts all the  $\textcolor{red}{T}$  to  $\textcolor{teal}{S}$  paths so to make  $k$  **satisfiable**, and protects each program variable in the cut

You have a further part in Judgment  
that allows you to also infer the  
type constraints as you are analyzing  
We record which variables are protected

**Commands:**

$$\Gamma, \text{Prot} \vdash c \Rightarrow k \quad \text{programs.}$$

# Let the machine prove things... (cont)

---

This rule disallows assignment from  $\textcolor{red}{T}$  to  $\textcolor{blue}{S}$ , since  $\textcolor{red}{T} \sqsubseteq \textcolor{blue}{S}$  is not true

$$\text{ASGN} \quad \frac{\Gamma \vdash r : \tau \Rightarrow k \quad \tau \sqsubseteq \Gamma(x)}{\Gamma, \text{Prot} \vdash x := r \Rightarrow k \cup (r \sqsubseteq x)}$$

The rule generates the constraint  $r \sqsubseteq x$  disallowing transient to stable assignments

# Let the machine prove things... (cont)

In case of operation with an expression of type  $\textcolor{red}{T}$  and one of type  $\textcolor{blue}{S}$ , since  $\textcolor{red}{T} \sqsubseteq \textcolor{blue}{S}$ , the overall type will be  $\textcolor{red}{T}$

BOP

$$\frac{\Gamma \vdash e_1 : \tau_1 \Rightarrow k_1 \quad \Gamma \vdash e_2 : \tau_2 \Rightarrow k_2 \quad \tau_1 \sqsubseteq \tau \quad \tau_2 \sqsubseteq \tau}{\Gamma \vdash e_1 \oplus e_2 : \tau \Rightarrow k_1 \cup k_2 \cup (e_1 \sqsubseteq e_1 \oplus e_2) \cup (e_2 \sqsubseteq e_1 \oplus e_2)}$$

the unknown type of  $e_1 \oplus e_2$  should be at least as transient as the  
(unknown) type of  $e_1$  and  $e_2$

Why is propagate types of your exp.

# Let the machine prove things... (cont)

Array can be indexed out of bounds during speculation, hence the rule assigns the transient type  $\mathbf{T}$  to array reads.  $\mathbf{E}$  must be of type to avoid leaks

$$\text{ARRAY-READ} \quad \frac{\Gamma \vdash e : \mathbf{S} \Rightarrow k \quad \begin{matrix} \nearrow \text{To prevent lengths } e \text{ must be of type } \mathbf{S}. \\ \Gamma \vdash a[e] : \mathbf{T} \end{matrix}}{\Gamma \vdash a[e] : \mathbf{T} \Rightarrow k \cup (e \sqsubseteq \mathbf{S}) \cup (\mathbf{T} \sqsubseteq a[e])} \quad \begin{matrix} \nwarrow \text{access as transient no matter what.} \\ \text{(a) Typing rules for expressions and arrays.} \end{matrix}$$

The rule generates constraint  $e \sqsubseteq \mathbf{S}$  for the unknown type of the array index, thus forcing it to typed be  $\mathbf{S}$  and forces the type of  $a[e]$  to be  $\mathbf{T}$

# Let the machine prove things... (cont)

---

The indexed expression used in array writes must be  $\mathbf{S}$ , while the stored value can have any type

$$\begin{array}{c} \text{ARRAY-WRITE} \\ \hline \Gamma \vdash e_1 : \mathbf{S} \Rightarrow k_1 \quad \Gamma \vdash e_2 : \tau \Rightarrow k_2 \\ \hline \Gamma, \text{Prot} \vdash a[e_1] := e_2 \Rightarrow k_1 \cup k_2 \cup (e_1 \sqsubseteq \mathbf{S}) \end{array}$$

# Let the machine prove things... (cont)

To prevent leaks, branch condition must be  $\mathbf{S}$

To prevent implicit flow leaks: expression has to be evaluated

IF-THEN-ELSE

$$\frac{\Gamma \vdash e : \mathbf{S} \Rightarrow k \quad \Gamma, \text{Prot} \vdash c_1 \Rightarrow k_1 \quad \Gamma, \text{Prot} \vdash c_2 \Rightarrow k_2}{\Gamma, \text{Prot} \vdash \text{if } e \text{ then } c_1 \text{ else } c_2 \Rightarrow k \cup k_1 \cup k_2 \cup (e \sqsubseteq \mathbf{S})}$$

# Example: constraint generation

---

## EXAMPLE

$x := a[i_1]$

$y := a[i_2]$

$z := x + y$

$w := b[z]$

# Example: constraint generation

## EXAMPLE

$x := a[i_1]$

$y := a[i_2]$

$z := x + y$

$w := b[z]$

## ARRAY-READ

$$\frac{\Gamma \vdash e : \mathbf{S} \Rightarrow k}{\Gamma \vdash a[e] : \mathbf{T} \Rightarrow k \cup (e \sqsubseteq \mathbf{S}) \cup (\mathbf{T} \sqsubseteq a[e])}$$

## ASGN

$$\frac{\Gamma \vdash r : \tau \Rightarrow k \quad \tau \sqsubseteq \Gamma(x)}{\Gamma, \text{Prot} \vdash x := r \Rightarrow k \cup (r \sqsubseteq x)}$$

# Example: constraint generation

## EXAMPLE

$x := a[i_1]$

$y := a[i_2]$

$z := x + y$

$w := b[z]$

- Reading memory creates  $\mathbf{T}$  value
- Types are propagated with assignment:  $\mathbf{T} \sqsubseteq x$ , i.e.  $x$  can contain transient values at runtime
- Constant values like  $a1$  do not produce constraints

CONSTRAINTS
$\mathbf{T} \sqsubseteq a[i_1]$
$a[i_1] \sqsubseteq x$

## ARRAY-READ

$$\frac{\Gamma \vdash e : \mathbf{S} \Rightarrow k}{\Gamma \vdash a[e] : \mathbf{T} \Rightarrow k \cup (e \sqsubseteq \mathbf{S}) \cup (\mathbf{T} \sqsubseteq a[e])}$$

## ASGN

$$\frac{\Gamma \vdash r : \tau \Rightarrow k \quad \tau \sqsubseteq \Gamma(x)}{\Gamma, \text{Prot} \vdash x := r \Rightarrow k \cup (r \sqsubseteq x)}$$

# Example: constraint generation

**EXAMPLE**

---

$$x := a[i_1]$$
$$y := a[i_2] \quad \text{(circled)}$$
$$z := x + y$$
$$w := b[z]$$

CONSTRAINTS
$T \sqsubseteq a[i_1]$
$a[i_1] \sqsubseteq x$
$T \sqsubseteq a[i_2]$
$a[i_2] \sqsubseteq y$

ARRAY-READ

$$\frac{\Gamma \vdash e : S \Rightarrow k}{\Gamma \vdash a[e] : T \Rightarrow k \cup (e \sqsubseteq S) \cup (T \sqsubseteq a[e])}$$

ASGN

$$\frac{\Gamma \vdash r : \tau \Rightarrow k \quad \tau \sqsubseteq \Gamma(x)}{\Gamma, \text{Prot} \vdash x := r \Rightarrow k \cup (r \sqsubseteq x)}$$

# Example: constraint generation

---

**EXAMPLE**

---


$$x := a[i_1]$$

$$y := a[i_2]$$

$$z := x + y \quad \text{(circled)}$$

$$w := b[z]$$

CONSTRAINTS
$T \sqsubseteq a[i_1]$
$a[i_1] \sqsubseteq x$
$T \sqsubseteq a[i_2]$
$a[i_2] \sqsubseteq y$
$x \sqsubseteq x+y$
$y \sqsubseteq x+y$
$x+y \sqsubseteq z$

**BOP**

$$\frac{\Gamma \vdash e_1 : \tau_1 \Rightarrow k_1 \quad \Gamma \vdash e_2 : \tau_2 \Rightarrow k_2 \quad \tau_1 \sqsubseteq \tau \quad \tau_2 \sqsubseteq \tau}{\Gamma \vdash e_1 \oplus e_2 : \tau \Rightarrow k_1 \cup k_2 \cup (e_1 \sqsubseteq e_1 \oplus e_2) \cup (e_2 \sqsubseteq e_1 \oplus e_2)}$$

**ASGN**

$$\frac{\Gamma \vdash r : \tau \Rightarrow k \quad \tau \sqsubseteq \Gamma(x)}{\Gamma, \text{Prot} \vdash x := r \Rightarrow k \cup (r \sqsubseteq x)}$$

x and y used to compute x+y

# Example: constraint generation

---

## EXAMPLE

$x := a[i_1]$   
 $y := a[i_2]$   
 $z := x + y$   
 $w := b[z]$

CONSTRAINTS
$T \sqsubseteq a[i_1]$
$a[i_1] \sqsubseteq x$
$T \sqsubseteq a[i_2]$
$a[i_2] \sqsubseteq y$
$x \sqsubseteq x+y$
$y \sqsubseteq x+y$
$x+y \sqsubseteq z$
$z \sqsubseteq S$
$T \sqsubseteq b[z]$
$b[z] \sqsubseteq w$

## ARRAY-READ

$$\frac{\Gamma \vdash e : S \Rightarrow k}{\Gamma \vdash a[e] : T \Rightarrow k \cup (e \sqsubseteq S) \cup (T \sqsubseteq a[e])}$$

## ASGN

$$\frac{\Gamma \vdash r : \tau \Rightarrow k \quad \tau \sqsubseteq \Gamma(x)}{\Gamma, \text{Prot} \vdash x := r \Rightarrow k \cup (r \sqsubseteq x)}$$

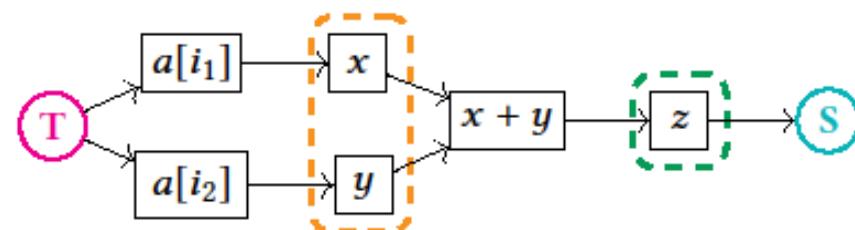
# The dataflow

## EXAMPLE

$x := a[i_1]$   
 $y := a[i_2]$   
 $z := x + y$   
 $w := b[z]$

CONSTRAINTS
$T \sqsubseteq a[i_1]$
$a[i_1] \sqsubseteq x$
$T \sqsubseteq a[i_2]$
$a[i_2] \sqsubseteq y$
$x \sqsubseteq x+y$
$y \sqsubseteq x+y$
$x+y \sqsubseteq z$
$z \sqsubseteq S$
$T \sqsubseteq b[z]$
$b[z] \sqsubseteq w$

## Subset of the Def-Use Graph



- $T$  is the source of transient values of the program
- $S$  is the sink of stable values of the program, i.e., values of a program that must be stable to avoid transient execution attacks

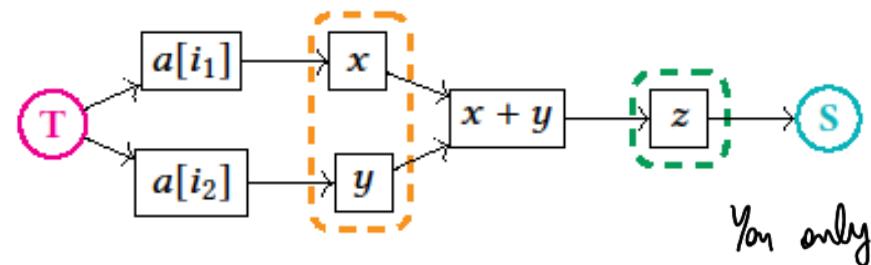
# The dataflow

## EXAMPLE

$x := a[i_1]$   
 $y := a[i_2]$   
 $z := x + y$   
 $w := b[z]$

CONSTRAINTS
$T \sqsubseteq a[i_1]$
$a[i_1] \sqsubseteq x$
$T \sqsubseteq a[i_2]$
$a[i_2] \sqsubseteq y$
$x \sqsubseteq x+y$
$y \sqsubseteq x+y$
$x+y \sqsubseteq z$
$z \sqsubseteq S$
$T \sqsubseteq b[z]$
$b[z] \sqsubseteq w$

## Subset of the Def-Use Graph



you only

- A set of constraints  $k$  is **satisfiable** if and only if there is no  $T-S$  path in  $k$
- Here we have two  $T-S$  paths: transient data flows through data dependencies into (what should be) a stable index expression and thus the program may be leaky

# Cut the dataflow

## EXAMPLE

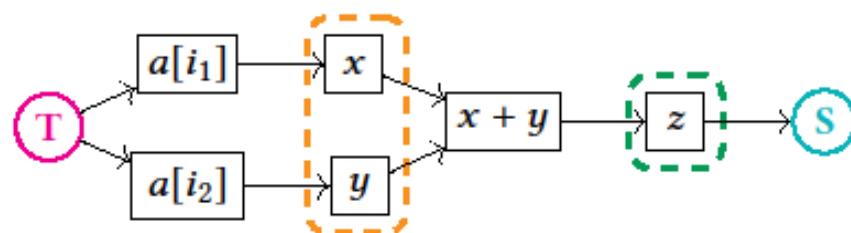
$x := a[i_1]$

$y := a[i_2]$

$z := x + y$

$w := b[z]$

- If a set of constraints is unsatisfiable, as here, we can make it satisfiable by removing some of the nodes or equivalently protecting some of the variables
- A set of variables  $A$  is a cut-set for a set of constraints  $k$ , if  $A$  cuts all  $T-S$  paths in  $k$



CUT-SETS

{x,y}

{z}

{z} is a minimal cut

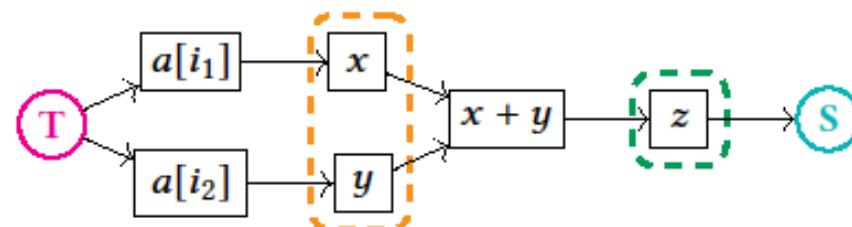
# Program repair

## Program Repair

As a final step, our repair algorithm traverses program  $c$  and inserts a **protect** for each variable in the **cut-set**

### EXAMPLE

```
x := a[i1]  
y := a[i2]  
z := x + y  
w := b[z]
```



### EXAMPLE PATCHED

```
x := protect(a[i1])  
y := protect(a[i2])  
z := protect(x + y)  
w := b[z]
```

The orange patch is sub-optimal because it requires more protect statements than the optimal green patch

# Putting everything together

---

- **Theorem 1 (Consistency):** Programs executed speculatively produce the same result as if they were executed sequentially (no out of order : result is the same : protect when needed)
- **Constant-time (CT) program** if and only if its branches and memory accesses do not depend on secrets (e.g., crypto keys) – I know, it's strange You can define CT programs under speculations,
- **Speculative constant-time (SCT) program** if and only if its observables and final configuration do not depend on secrets (assume the sequence of attacker-provided can be lengthy, directives as fixed before execution)
- **Theorem 2 (Soundness):** If a program  $c$  is CT and well-typed  $\Gamma, \emptyset \vdash c$  then it is also SCT
  - If you want: when compiling  $c$  from the source to the target using the JIT “compiler”, a secure source program (CT & well-typed) is mapped into a secure target program (SCT)

# Summing up

---

- We worked with the “bare-metal”: caches, speculation, side-channels
- We saw how that can be abstracted elegantly
  - A source language translated on-the-fly to a target one, using a JIT semantics
  - The JIT semantics has nice properties: programs deemed secure at the source are such also at the target
  - Also: it is possible to make programs secure by synthesizing protections
  - Nice theoretical framework 😊

# Bibliography

---

- M. Vassena, C. Disselkoen, K. von Gleissenthall, S. Cauligi, R. Gökhan Kıcı, R. Jhala, D. Tullsen, and D. Stefan. *Automatically eliminating speculative leaks from cryptographic code with blade*. POPL 2021. URL:  
<https://arxiv.org/abs/2005.00294>  
<https://www.youtube.com/watch?v=mCzbjGRIVeQ>
- Canella, Claudio, et al. *A systematic evaluation of transient execution attacks and defenses*. 28th USENIX Security Symposium (USENIX Security 19). 2019. URL: <https://arxiv.org/abs/1811.05441v3>

# Take away message

Focus on:

- security vs. performance is required
- hardware and software
- Changes to how we design hardware are required
  - all hardware developers should know which software abstractions are adopted
- Changes to the way we design software are required
  - all software engineers should have some notion of how processors behave and deal with microarchitectural complexity and, in general with all layers
- Also: no more HW and SW people!

## Take away message (cont.)

- Software is eating the world
- Insecurity costs scale faster and faster
- Security is going to eat the software

# End