# Exploiting modern microarchitectures

Meltdown, Spectre, and other hardware security vulnerabilities in modern processors

Jon Masters - Computer Architect
CRob - Red Hat Product Security
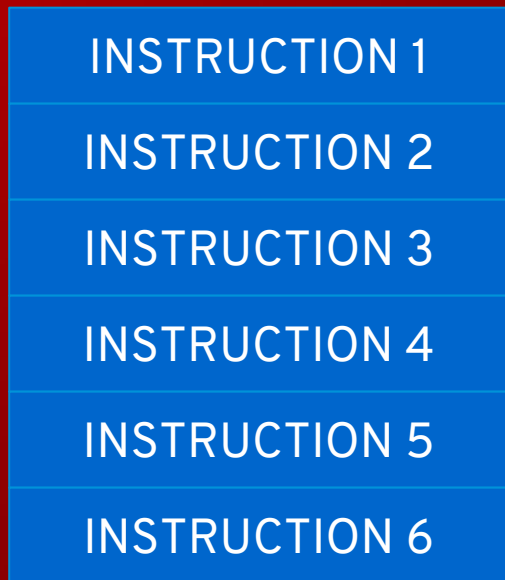10May2018

# Who are we?



JCM



CRob

redhat.

Out-of-Order (OoO) execution converts sequential "in order" machines (CPUs) as conceptualized by programmers into out-of-order "data flow" machines in which computation is performed as soon as dependent data is available*

* Instructions issue and retire "in-order", but are executed using an Out-of-Order "backend"

redhat.

# Sequential vs data flow

## Sequential model

| INSTRUCTION 1 |
|:---:|
| INSTRUCTION 2 |
| INSTRUCTION 3 |
| INSTRUCTION 4 |
| INSTRUCTION 5 |
| INSTRUCTION 6 |

## Data flow model

D1  D2
D3
D4  D5
D6

# Sequential machines

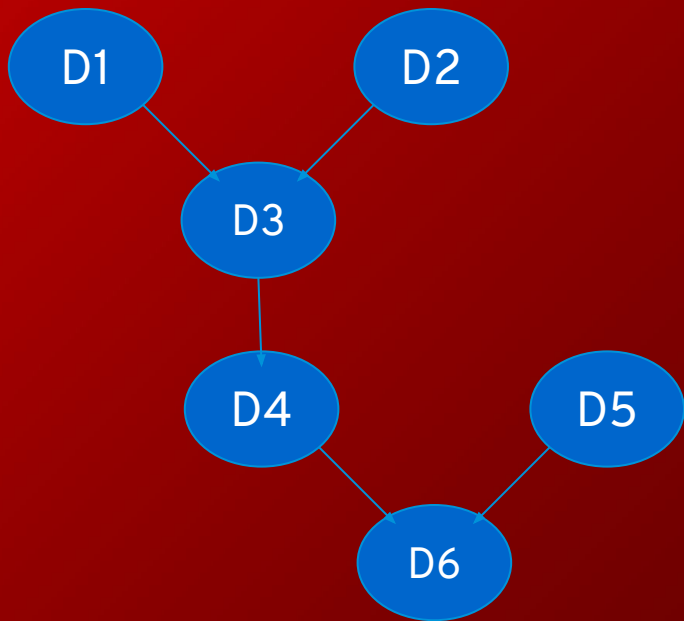| |
|---|
| INSTRUCTION 1 |
| INSTRUCTION 2 |
| INSTRUCTION 3 |
| INSTRUCTION 4 |
| INSTRUCTION 5 |
| INSTRUCTION 6 |

Programmers intuitively think in terms of sequential program flow ("I do this, then I do that")

Expectation is that the machine does everything in the order written down

Debugging involves tracing what each individual instruction did, one after the other, in the order written

redhat.

# Data flow machines



Data flow operates in terms of data. Each operation proceeds when its dependent data operands are ready

Instructions do not execute in the order written by the programmer

Debugging a real world data flow machine would involve understanding all of the internal machine state

redhat.

Robert Tomasulo described how to convert an in-order sequential machine (as understood by programmers) into an Out-of-Order data flow model machine
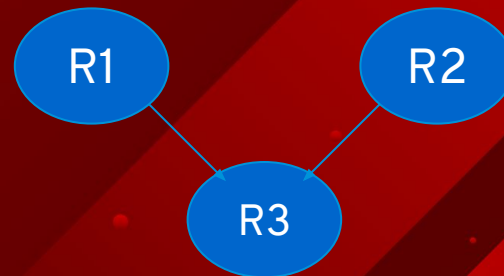
# Tomasulo's Algorithm
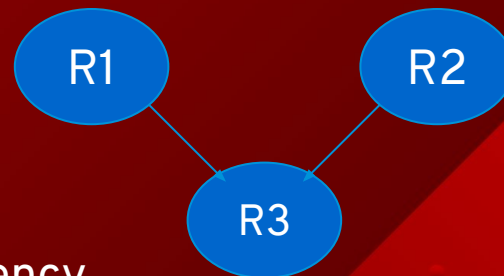
R1 = LOAD A

R2 = LOAD B

R3 = R1 + R2

No data dependency

R1 = 1

R2 = 1

R3 = R1 + R2

# Tomasulo's Algorithm

**Program Order**

| R1 = LOAD A |
| R2 = LOAD B |
| R3 = R1 + R2 |
| R1 = 1 |
| R2 = 1 |
| R3 = R1 + R2 |

**Re-Order Buffer (ROB)**

| Entry | RegRename | Instruction | Deps | Ready? |
|-------|-----------|-------------|------|--------|
| 1 | P1 = R1 | P1 = LOAD A | X | Y |
| 2 | P2 = R2 | P2 = LOAD B | X | Y |
| 3 | P3 = R3 | P3 = R1 + R2 | 1,2 | N |
| 4 | P4 = R1 | P4 = 1 | X | Y |
| 5 | P5 = R2 | P5 = 1 | X | Y |
| 6 | P6 = R3 | P6 = P4 + P5 | 4,5 | N |

redhat.

Tomasulo removed the need for chips to wait around for dependent calculations. Deep reorder buffers allow very high performance gains in contemporary servers, laptops, and even phones

redhat.

...time passed, and computer architects realized it was possible to gain even more performance...

redhat.

"What if we combine Out-of-Order Execution with Branch Prediction?"

# Branches

```
for (i=0;i<10;i++)
{
    // something here
}
// post-loop code here
```

```
        movq $0, %rax
loop:
        incq %rax
        cmpq $10, %rax
        jle loop
```

```
// pack clothes
if (raining) {
    // pack umbrella
}
// more packing
```

```
        movq $raining, %rax
        cmpq $1, %rax
        jne more_packing
        // pack umbrella
more_packing:
        // more packing
```
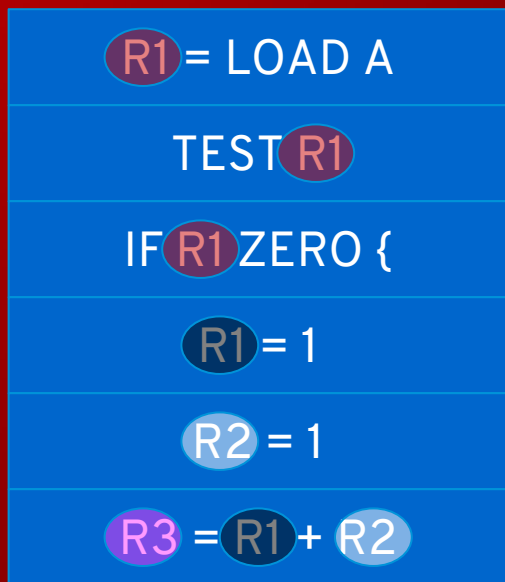
# Branches

```
// pack clothes
if (raining) {
    // pack umbrella
}
// more packing
```

→

```
movq $raining, %rax
cmpq $1, %rax
jne more_packing

more_packing:
...
```

Branch conditions take time to "resolve" (do I take the branch or not?)

# Speculative Execution

R1 = LOAD A

TEST R1

IF R1 ZERO {

R1 = 1

R2 = 1

R3 = R1 + R2

| Entry | RegRename | Instruction | Deps | Ready? | Spec? |
|-------|-----------|-------------|------|--------|-------|
| 1 | P1 = R1 | P1 = LOAD A | X | Y | N |
| 2 | | TEST R1 | 1 | Y | N |
| 3 | | IF R1 ZERO { | 1 | N | N |
| 4 | P2 = R1 | P4 = 1 | X | Y | Y* |
| 5 | P3 = R2 | P5 = 1 | X | Y | Y* |
| 6 | P4 = R3 | P4 = P2 + P3 | 4,5 | Y | Y* |

* Speculatively execute the branch before the condition is known ("resolved")

Branch prediction allows machines to guess which path code is going to take before that is fully known ("resolved")

redhat.

When combined with an Out-of-Order backend, "Speculative Execution" will try to execute code before we know whether we need to do so. If we guess wrong, we throw away our prediction

All is well...unless someone can see what we did speculatively and use it to learn things they should not know...

redhat.

# The Speculation Diner

# Look at what we found!

Flash forward nearly 40 years…. 3 groups of researchers independently find….

Speculation is not an unobservable black box. We can abuse speculation to leak information we shouldn't access using a "side channel"

redhat.

Classic example is "differential power analysis" being used to monitor voltages on the pins of a chip during crypto operations - power varies according to the calculation*

* Your "chip-and-pin" bank card is very carefully designed not to do this

redhat.

Cache side channels exploit the shared nature of memory on a local system to infer what other programs are doing

redhat.

But first...let's talk about memory...recall how modern running programs use "virtual memory"...

redhat.

# Virtual Memory

$ cat /proc/self/maps  →  /bin/cat process

### Virtual Memory

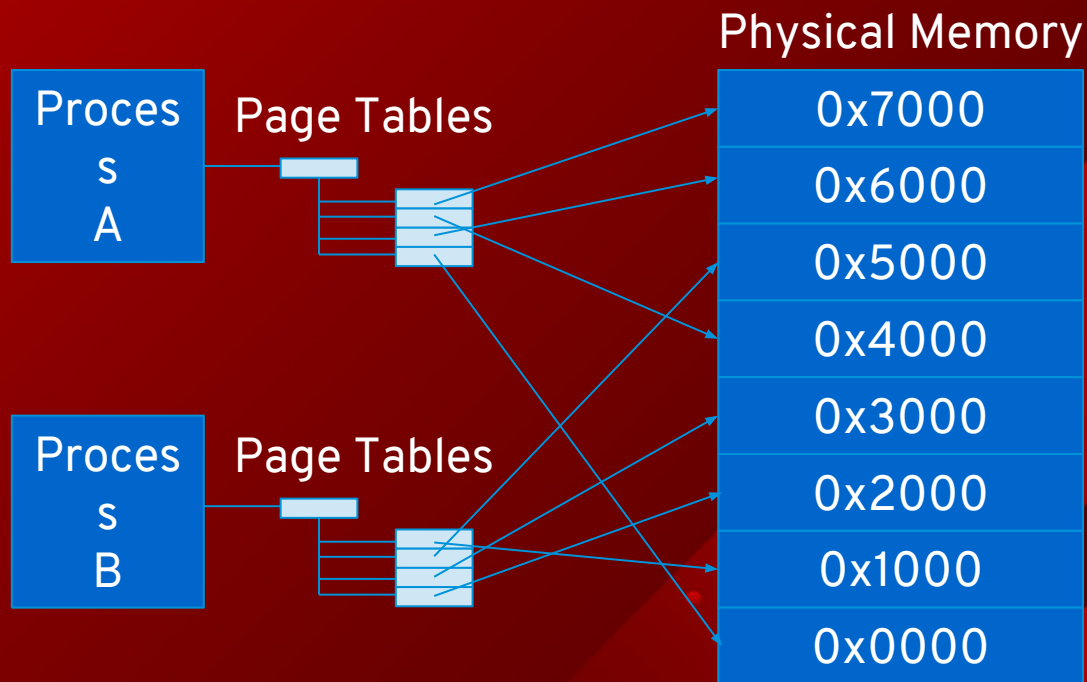| |
|---|
| 0xffff_ffff_81a0_00e0 |
| ... |
| 0xffff_ffff_8100_0000 |
| ... |
| 0x7ffc683f9000* |
| ... |
| 0x7ffc683a6000 |
| ... |
| 0x55d776036000 |

* Special case kernel VDSO (Virtual Dynamic Shared Object)

redhat.

# Virtual Memory

# Virtual Memory

Process A

Page Tables

Physical Memory

| 0x7000 |
| 0x6000 |
| 0x5000 |
| 0x4000 |
| 0x3000 |
| 0x2000 |
| 0x1000 |
| 0x0000 |

## Translation Lookaside Buffer (TLB)

| 0x4000 | → | 0x0000 |
| 0x3000 | → | 0x6000 |
| 0x1000 | → | 0x4000 |
| 0x0000 | → | 0x7000 |

redhat.

Caches are used to store the actual data read by a running program from this virtual memory

# Caches

# Caches

**Virtual Memory**

| |
|---|
| ... |
| 0xf080 |
| 0xf040 |
| 0xf000 |
| ... |
| 0x0080 |
| 0x0040 |
| 0x0000 |

**Cache (L1/L2/etc.)**

| 0xf040 | ksecret |
|---|---|
| 0x0040 | usecret |

**Physical Memory**

| |
|---|
| ... |
| 0x0180 |
| 0x0140 |
| 0x0100 |
| 0x00c0 |
| 0x0080 |
| 0x0040 |
| 0x0000 |

\* For readability priviluged kernel addresses are shortened to begin 0xf instead of 0xfffffffff...

redhat.

Caches are shared resources. The time taken to access data is proportional to whether it is located in the cache

redhat.

# Caches

```
time = rdtsc();
maccess(&data[0x300]);        ⬅
delta3 = rdtsc() - time;

time = rdtsc();
maccess(&data[0x200]);
delta2 = rdtsc() - time;
```
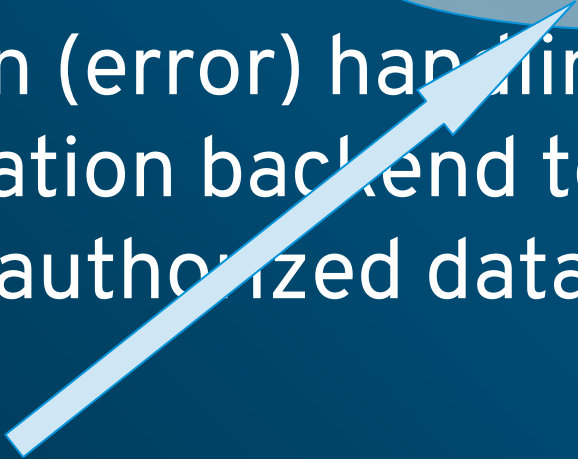
Execution time taken for instruction is proportional to whether it is in cache(s)

"Meltdown" abuses at-retirement exception (error) handling in the speculation backend to read unauthorized data...

"Meltdown" abuses at-retirement exception (error) handling in the speculation backend to read unauthorized data...

fancy way of saying
"at the very end"

...then uses the cache access trick to reconstruct what that unauthorized data was (hence "side channel")

We will get a permission check failure but we won't handle it until too late - after the speculation window completes (we already did the cache access trick)

redhat.

# Meltdown

- A malicious attacker arranges for exploit code similar to the following to speculatively execute:

```
if (spec_cond) {
        unsigned char value = *(unsigned char *)ptr;
        unsigned long index2 = (((value>>bit)&1)*0x100)+0x200;
        maccess(&data[index2]);
}
```

- "data" is a user controller array to which the attacker has access, "ptr" contains privileged data

# Meltdown

char value = *SECRET_KERNEL_PTR;

⬇
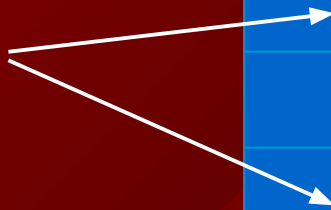
mask out bit I want to read

⬇

calculate offset in "data"
(that I do have access to)

char data[];

| | |
|---|---|
| 0x000 | |
| 0x100 | |
| 0x200 | |
| 0x300 | |

# Meltdown

- Access to "data" element 0x100 pulls the corresponding entry into the cache

char data[];

| 0x000 |      |
|-------|------|
| 0x100 |      |
| 0x200 |      |
| 0x300 | DATA |

Cache

| 0x100 |  |
|-------|--|

redhat.

# Meltdown

- Access to "data" element 0x300 pulls the corresponding entry into the cache

char data[];

| 0x000 | |
|-------|------|
| 0x100 | |
| 0x200 | |
| 0x300 | DATA |

Cache

| 0x300 | |
|-------|---|

# Meltdown: Speculative Execution

| Entry | RegRename | Instruction | Deps | Ready? | Spec? |
|-------|-----------|-------------|------|--------|-------|
| 1 | P1 = R1 | R1 = LOAD SPEC_CONDITION | X | Y | N |
| 2 |  | TEST SPEC_CONDITION | 1 | Y | N |
| 3 |  | IF (SPEC_CONDITION) { | 1 | N | N |
| 4 | P2 = R2 | R2 = LOAD KERNEL_ADDRESS | X | Y | Y* |
| 5 | P3 = R3 | R3 = (((R2&1)*0x100)+0x200) | 2 | Y | Y* |
| 6 | P4 = R4 | LOAD USER_BUF[R3] | 3 | Y | Y* |

**flags for future exception**

redhat

# Meltdown: Speculative Execution

| Entry | RegRename | Instruction | Deps | Ready? | Spec? |
|-------|-----------|-------------|------|--------|-------|
| 1 | P1 = R1 | R1 = LOAD SPEC_CONDITION | X | Y | N |
| 2 | | TEST SPEC_CONDITION | 1 | Y | N |
| 3 | | IF (SPEC_CONDITION) { | 1 | N | N |
| 4 | P2 = R2 | R2 = LOAD KERNEL_ADDRESS | X | Y | Y* |
| 5 | P3 = R3 | R3 = (((R2&1)*0x100)+0x200) | 2 | Y | Y* |
| 6 | P4 = R4 | LOAD USER_BUF[R3] | 3 | Y | Y* |

**should kill speculation here**

redhat.

# Meltdown: Speculative Execution

| Entry | RegRename | Instruction | Deps | Ready? | Spec? |
|-------|-----------|-------------|------|--------|-------|
| 1 | P1 = R1 | R1 = LOAD SPEC_CONDITION | X | Y | N |
| 2 | | TEST SPEC_CONDITION | 1 | Y | N |
| 3 | | IF (SPEC_CONDITION) { | 1 | N | N |
| 4 | P2 = R2 | R2 = LOAD KERNEL_ADDRESS | X | Y | Y* |
| 5 | P3 = R3 | R3 = (((R2&1)*0x100)+0x200) | 2 | Y | Y* |
| 6 | P4 = R4 | LOAD USER_BUF[R3] | 3 | Y | Y* |

**really bad thing(™)**

redhat.

# We mitigated "Meltdown" through "KPTI" (Kernel Page Table Isolation)*

redhat.

"Meltdown" requires that we have valid kernel address translations (TLBs) so ensure we never do by spitting kernel and application page tables

redhat.

Performance impact comes from the "trampoline" code that is needed to switch page tables on every kernel entry and exit back to a program

redhat.

"Spectre" v2 abuses (indirect) branch predictors that don't fully disambiguate between two different process "contexts"...

Spectre-v2 (branch predictor poisoning)

...we train the branch predictor in one context (program) such that it will guess a particular way in another context (the victim/kernel we want to control)

The attack relies upon finding "gadget" code already in the victim program

redhat.

We mitigated "Spectre-v2" through a combination of firmware and OS level branch predictor control interfaces*

* Shipping in a currently supported RHEL kernel near you

redhat.

CPUs have an ability to tweak certain behavior (e.g. "chicken bits") or patch new operations into a small on-chip RAM via "microcode" updates

# We use microcode for two things:

1. To restrict speculation across privilege boundaries (kernel entry)
2. To invalidate the branch predictor entries when switching programs

A later optimization ("retpolines") converts indirect function calls into fake return "trampolines". Avoids using the standard branch predictor hardware*

redhat.

Another variant of Spectre (v1) requires minor code tweaks (barriers) that were added to the kernel source directly

# Why these problems mattered

Nearly every more CPU affected (to differing extents)

A skilled attacker could have the ability to read any memory they desire

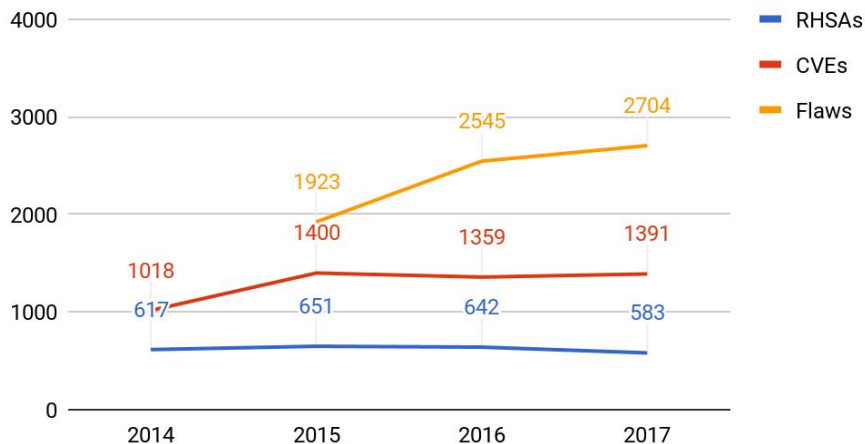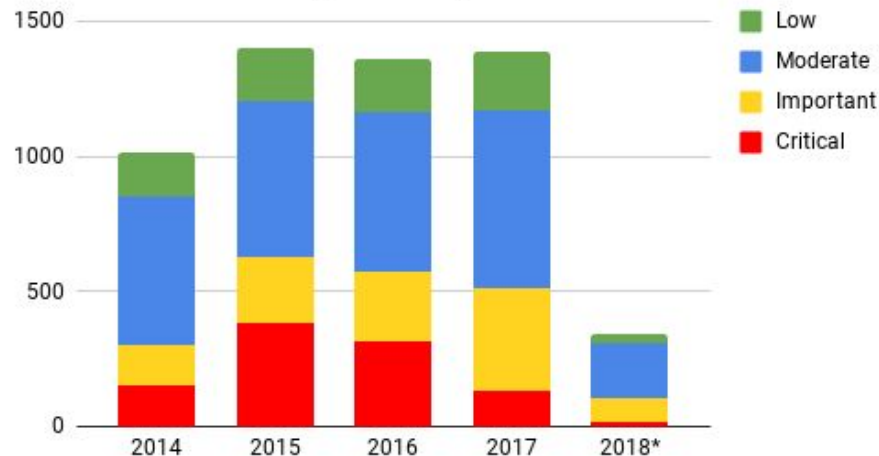Some attacks cross guest/host boundaries

Virtually undetectable

redhat.

# Why Red Hat was included

We have a very particular set of skills, skills we've acquired over our very long career

(the last 25 years and counting)

redhat.

# A snapshot of Red Hat Product Security over the years



**Total CVEs and RHSAs by Year**

RHSAs
CVEs
Flaws

2704
2545
1923
1400
1359
1391
1018
617
651
642
583

2014  2015  2016  2017



**Total CVEs Fixed by Severity**

Low
Moderate
Important
Critical

2014  2015  2016  2017  2018*
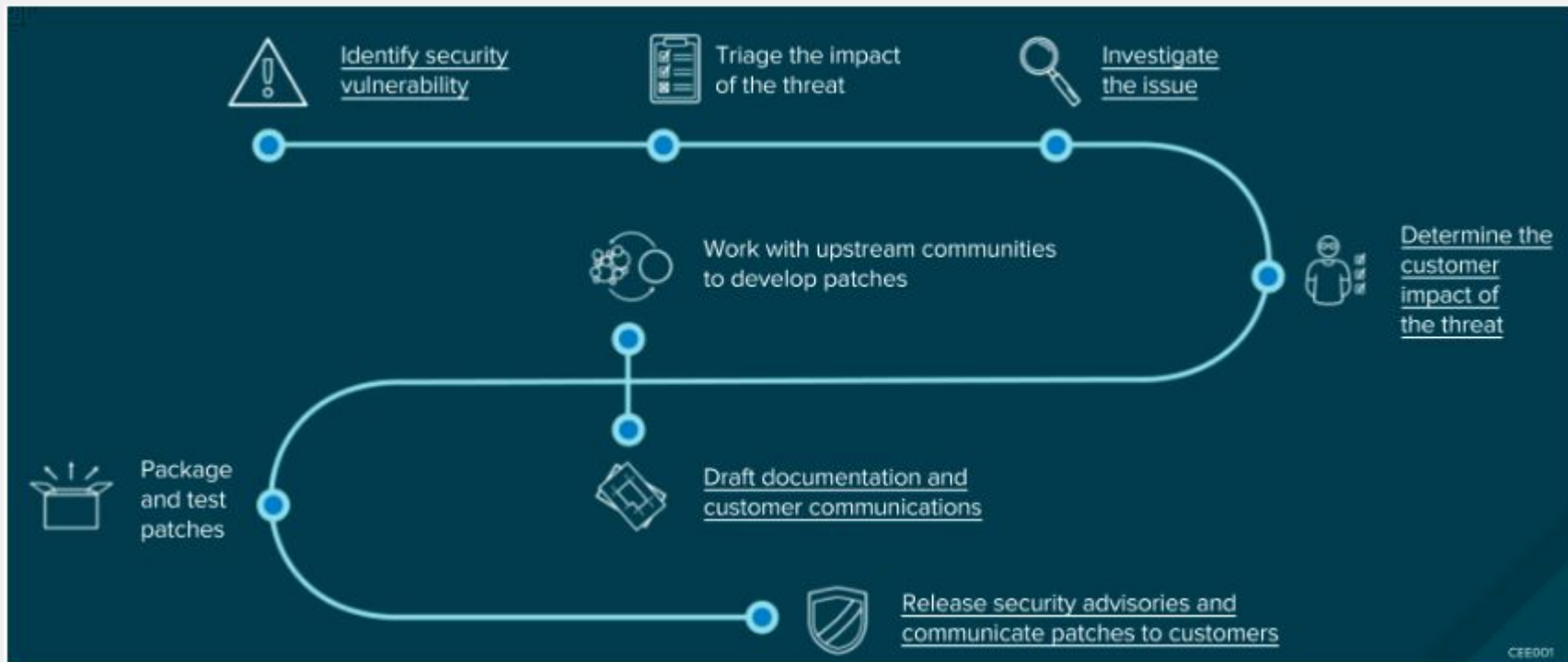
# Coordinated Vulnerability Disclosure

Red Hat is part of a large group of vendor and community security teams

We use a process called Coordinated Vulnerability Disclosure(1) - The goal is to protect customers and the larger global computing community

Red Hat respects the wishes of the issue reporter on how they want the issue to be handled and how long to keep it secret

(1)    https://resources.sei.cmu.edu/asset_files/SpecialReport/2017_003_001_503340.pdf

# Red Hat's Customer Security Awareness Program

Ride, ride my seesaw    Take this place    On this trip    Just for me

# What we delivered to subscribers

15 kernels...3 times (4 if you count retpolines)

59 total patches (plus 15 more for retpolines) PLUS our complete OpenShift Online infrastructure (3 versions), plus all AMIs, ISOs, and base images

Over 60 engineers involved in developing, testing, and packaging the errata

Education video, 1 vulnerability article (translated into 6 languages), 4 kcs articles on specialized topics, 3 forms of customer emails, TAM T3, Performance webinar, and Portal Banners & Alerts

10,000+ engineering hours dedicate to resolving issues

redhat.

CVE-2017-5753
CVE-2017-5715
CVE-2017-5754
3Jan2018

Red Hat Product Security

2017 CVE Response

#redhat  #rhsummit

# To learn more….

Vulnerability article - https://access.redhat.com/security/vulnerabilities/speculativeexecution

3 minute Video - https://youtu.be/syAdX44pokE

Blog - https://www.redhat.com/en/blog/what-are-meltdown-and-spectre-heres-what-you-need-know

CSAw process - https://access.redhat.com/articles/2968471

To sign up for RHSA announcements - https://www.redhat.com/wapps/ugc/protected/notif.html

redhat.