

# SOSCON

## Micro-architectural attack and defense on Linux kernel

### Spectre and Meltdown

Samsung Research | Security team | Jinbum Park | [jinb-park.github.io](https://jinb-park.github.io)

Samsung Research | Security team | Joonwon Kang

16-Oct-2019



# Index

---

Overview	<b>01</b>
Attack 1 : Out-of-bounds attack	<b>02</b>
Defense 1 : Bounds check	<b>03</b>
Attack 2 : Bounds check bypass (Spectre v1)	<b>04</b>
Defense 2 : Index sanitization	<b>05</b>
Attack 3 : Indirect branch poisoning (Spectre v2)	<b>06</b>
Defense 3 : Retpoline	<b>07</b>
Attack 4 : Meltdown	<b>08</b>
Defense 4 : Page Table Isolation	<b>09</b>
Attack 5 : L1 Terminal Fault	<b>10</b>
Defense 5 : PTE Inversion	<b>11</b>

**SOSCON 2019**

SAMSUNG OPEN SOURCE CONFERENCE 2019



# Overview

---



SOSCON 2019

SAMSUNG OPEN SOURCE CONFERENCE 2019

**Meltdown?**

**Spectre?**

**Intel hardware vulnerabilities?**

**Micro-architecture attacks?**



## Software

```
int *ptr = ...;

void func(void)
{
    int a = *ptr;
}
```

## Architecture (Intel / ARM ...)

```
mov %ebx, (%eax)

// eax == ptr
```

## Micro-architecture (i7-1, i7-2, i5-3, A12, ...)

```
pa = convert_va_to_pa(%eax);
read pa;
```



## Software

```
int *ptr = NULL;  
  
void func(void)  
{  
    int a = *ptr;  
}
```

## Architecture (Intel / ARM ...)

```
mov %ebx, (%eax)  
  
// eax == ptr (NULL)
```

## Micro-architecture (i7-1, i7-2, i5-3, A12, ...)

```
pa = convert_va_to_pa(%eax);  
if (is_valid_address(pa))  
    read pa;  
else  
    error;
```

Segmentation fault (core dumped)



## Software

```
int *ptr = &obj;  
  
void func(void)  
{  
    int a = *ptr;  
}
```

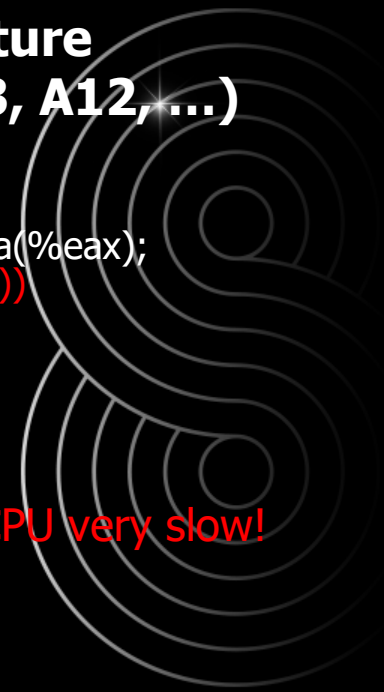
## Architecture (Intel / ARM ...)

```
mov %ebx, (%eax)  
  
// eax == ptr (NULL)
```

## Micro-architecture (i7-1, i7-2, i5-3, A12, ...)

```
pa = convert_va_to_pa(%eax);  
if (is_valid_address(pa))  
    read pa;  
else  
    error;
```

Checking on every memory access makes CPU very slow!  
How can we make it better?



## Software

```
int *ptr = &obj;

void func(void)
{
    int a = *ptr;
}
```

## Architecture (Intel / ARM ...)

```
mov %ebx, (%eax)

// eax == ptr (NULL)
```

## Micro-architecture

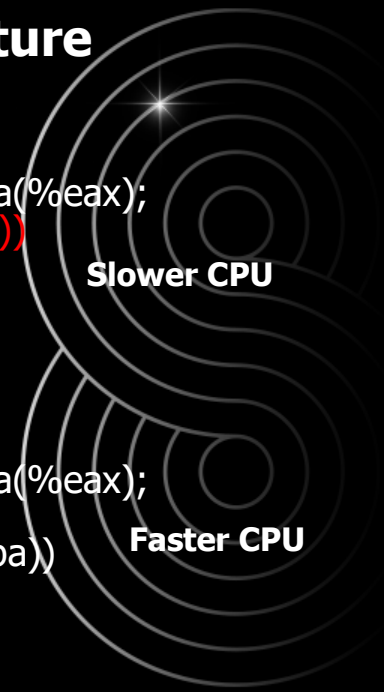
```
pa = convert_va_to_pa(%eax);
if (is_valid_address(pa))
    read pa;
else
    error;
```

Slower CPU

```
pa = convert_va_to_pa(%eax);
read pa;
if (is_invalid_address(pa))
    error;
```

Faster CPU

Speculate "pa" is not NULL  
Read "pa" first!  
It called "Speculative execution"





## Software

```
int *ptr = &obj;

void func(void)
{
    int a = *ptr;
}
```

## Architecture (Intel / ARM ...)

```
mov %ebx, (%eax)

// eax == ptr (NULL)
```

## Micro-architecture

```
pa = convert_va_to_pa(%eax);
if (is_valid_address(pa))
    read pa;
else
    error;
```

Slower CPU

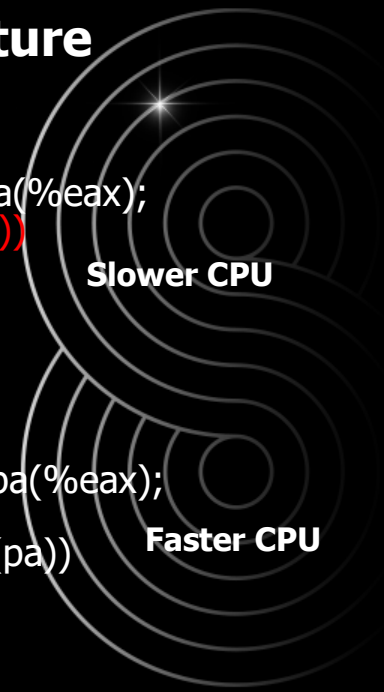
```
pa = convert_va_to_pa(%eax);
read pa;
if (is_invalid_address(pa))
    error;
```

Faster CPU

**Security implication in it**

**:: Even though it occurs a fault,**

**We can read an invalid memory before the fault!**



## Software

```
int *ptr = &obj;

void func(void)
{
    int a = *ptr;
}
```

## Architecture (Intel / ARM ...)

```
mov %ebx, (%eax)

// eax == ptr (NULL)
```

## Micro-architecture

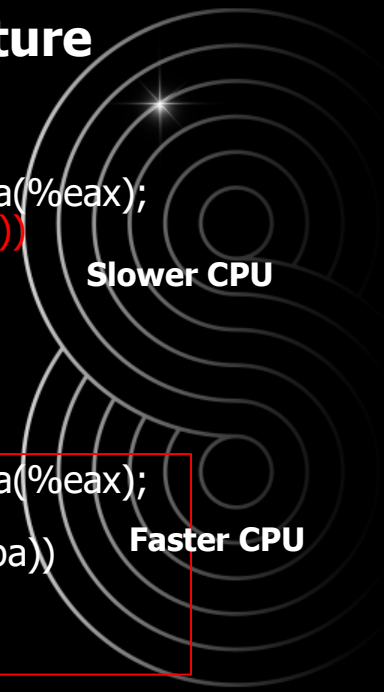
```
pa = convert_va_to_pa(%eax);
if (is_valid_address(pa))
    read pa;
else
    error;
```

Slower CPU

```
pa = convert_va_to_pa(%eax);
read pa;
if (is_invalid_address(pa))
    error;
```

Faster CPU

If you can exploit this implementation,  
we call this "Hardware Vulnerability"  
Call the attack on this "Micro-architecture attack"



## Intel Desktop CPUs Affected By Meltdown + Spectre

### Intel Desktop CPUs Vulnerable To Meltdown + Spectre

**Affected Variants :** These Intel CPUs are affected by **all three variants** of the **speculative execution CPU bug**. They are vulnerable to the Meltdown and both Spectre exploits.

#### Intel Coffee Lake-S (2017)

- Intel Core i7-8700K
- Intel Core i7-8700
- Intel Core i5-8600K
- Intel Core i5-8400
- Intel Core i3-8350K
- Intel Core i3-8100

#### Intel Gemini Lake (2017)

- Intel Pentium Silver J5005
- Intel Celeron J4105
- Intel Celeron J4005

#### Intel Denverton (2017)

- Intel Celeron C3958
- Intel Celeron C3955

#### Contents

1. [The Complete List Of CPUs Affected By Meltdown + Spectre](#)
2. [AMD Workstation, Desktop & Mobile CPUs Vulnerable To Spectre](#)
3. [AMD Mobile CPUs Vulnerable To Spectre](#)
4. [Apple, ARM & Intel CPUs Affected By Meltdown & Spectre](#)
5. [Intel Server / Workstation CPUs Vulnerable To Meltdown + Spectre](#)
6. [Intel Desktop CPUs Affected By Meltdown + Spectre](#)
7. [Intel Mobile CPUs Affected By Meltdown + Spectre](#)
8. [VIA Desktop + Mobile CPUs Vulnerable To Meltdown + Spectre](#)



## 4 Hardware Vulnerabilities

- Spectre Variant1, Variant2, Meltdown, L1TF

How to exploit them for user to read kernel-memory?

How does Linux kernel defend against them?



# Spectre Attack & Defense

---



SOSCON 2019

SAMSUNG OPEN SOURCE CONFERENCE 2019

## 2. Attack 1: Out-of-bounds attack

SOSCON 2019

### Intuition

```
#include <stdio.h>
#include <stdlib.h>

char *secret = "This is victim's secret.";
char ary1[10];

int main(int argc, char **argv)
{
    int idx = atoi(argv[1]);
    char var = ary1[idx];

    return 0;
}
```

**Q. What if execute it with `idx = -100`?**

**A. `var == 'T'`**

address	value
100	'T'
101	'h'
...	...
200	ary1[0]

**Can access out of ary with the attacker chosen index without crash!**



## 2. Attack 1: Out-of-bounds attack (Cont.)

SOSCON 2019

### Attack steps

```
#include <stdio.h>
#include <stdlib.h>

extern char ary2[];

char *secret = "This is victim's secret.";
char ary1[10];

int main(int argc, char **argv)
{
    int idx = atoi(argv[1]);

    char var = ary2[ary1[idx]];

    return 0;
}
```

Shared memory between processes.

Goal: What is in victim's ary1[-100] ?



## 2. Attack 1: Out-of-bounds attack (Cont.)

SOSCON 2019

### Attacker

1. Flush `ary2[0...255]`.
2. Set `idx = -100` and Execute the victim.

5. Access `ary2[0...255]`.

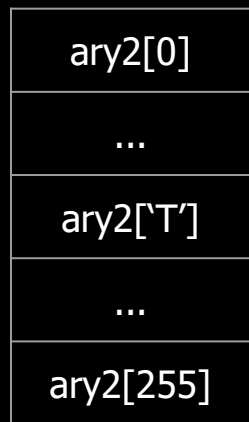
6. Find `idx` which takes **the shortest time**.  
**Found! `ary1[-100] == 'T'`.**

### Cache



4. **Cached!**

### Shared memory



### Victim

3. **var = ary2[ary1[-100]]**

address	variable
<b>100</b>	<b>'T'</b>
101	'h'
...	...
200	ary1[0]

Called "Flush+Reload" analysis.



## 2. Attack 1: Out-of-bounds attack (Cont.)

SOSCON 2019

How can we defend against this attack?

```
#include <stdio.h>
#include <stdlib.h>

extern char ary2[];

char *secret = "This is victim's secret.";
char ary1[10];

int main(int argc, char **argv)
{
    int idx = atoi(argv[1]);

    char var = ary2[ary1[idx]];

    return 0;
}
```

← **Let's check bounds here!**



### 3. Defense 1: Bounds check

Check bounds before the gadget.

```
#include <stdio.h>
#include <stdlib.h>

extern char ary2[];

char *secret = "This is victim's secret.";
char ary1[10];
int size = 10;

int main(int argc, char **argv)
{
    int idx = atoi(argv[1]);
    char var;

    if (idx < size)
        var = ary2[ary1[idx]];

    return 0;
}
```

**Out-of-bounds attack is blocked!**  
**Really..?**



# 4. Attack 2: Bounds check bypass (Spectre v1)

Can we bypass bounds check?

Software

```
if (idx < size)
    var = ary2[ary1[idx]];
```

Nope... →

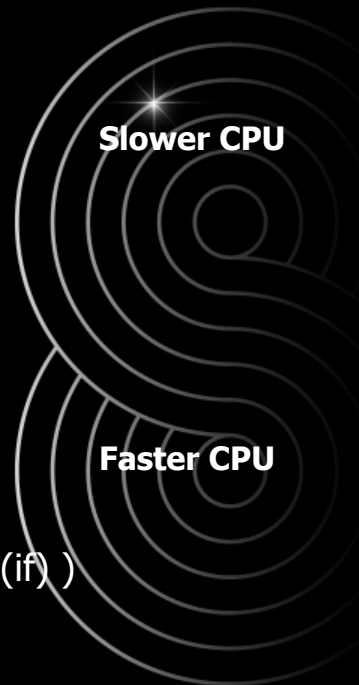
Micro-architecture

```
read idx
read size
res = idx < size
if ( res )
    read ary2[ary1[idx]]
```

Yes! but we have to make  
**likely\_taken(if) == true.** →

```
if ( likely_taken(if) )
    read ary2[ary1[idx]]
read idx
read size
res = idx < size
if ( res == likely_taken(if) )
    commit
else
    revert
...

```



## 4. Attack 2: Bounds check bypass (Cont.)

SOSCON 2019

### Attack steps


#### Software

```
if (idx < size)
    var = ary2[ary1[idx]];
```

**Goal: What is in victim's ary1[-100] ?**

#### Micro-architecture

```
if ( likely_taken(if) )
    read ary2[ary1[idx]]
read idx
read size
res = idx < size
if ( res == likely_taken(if) )
    commit
else
    revert
...
```



## 4. Attack 2: Bounds check bypass (Cont.)

SOSCON 2019

### Attack steps

#### Software


```
if (idx < size)
    var = ary2[ary1[idx]];
```

1. Execute victim with legitimate idx multiple times.

→ **Conditional branch predictor will learn that**  
**likely\_taken(if) == true.**

#### Micro-architecture

```
if ( likely_taken(if) )
    read ary2[ary1[idx]]
read idx
read size
res = idx < size
if ( res == likely_taken(if) )
    commit
else
    revert
...
```



## 4. Attack 2: Bounds check bypass (Cont.)

SOSCON 2019

### Attack steps

#### Software


```
if (idx < size)
    var = ary2[ary1[idx]];
```

2. Execute victim with `idx = -100`.

→ `ary2[ary1[-100]]` will be read and cached!

#### Micro-architecture

```
if ( likely_taken(if) )
    read ary2[ary1[idx]]
read idx
read size
res = idx < size
if ( res == likely_taken(if) )
    commit
else
    revert
...
```



## 4. Attack 2: Bounds check bypass (Cont.)

SOSCON 2019

### Attack steps

#### Software

```
if (idx < size)
    var = ary2[ary1[idx]];
```


3. Do Flush+Reload analysis on ary2.

→ **ary2['T'] will take the shortest time.**

**Thus, ary1[-100] == 'T' !!**

#### Micro-architecture

```
if ( likely_taken(if) )
    read ary2[ary1[idx]]
read idx
read size
res = idx < size
if ( res == likely_taken(if) )
    commit
else
    revert
...
```



How can we defend against this attack?


### Software

```
if (idx < size)
    var = ary2[ary1[idx]];
```

← **Let's sanitize idx here!**

### Micro-architecture

```
if ( likely_taken(if) )
    read ary2[ary1[idx]]
read idx
read size
res = idx < size
if ( res == likely_taken(if) )
    commit
else
    revert
...
```





### How can we defend against this attack?

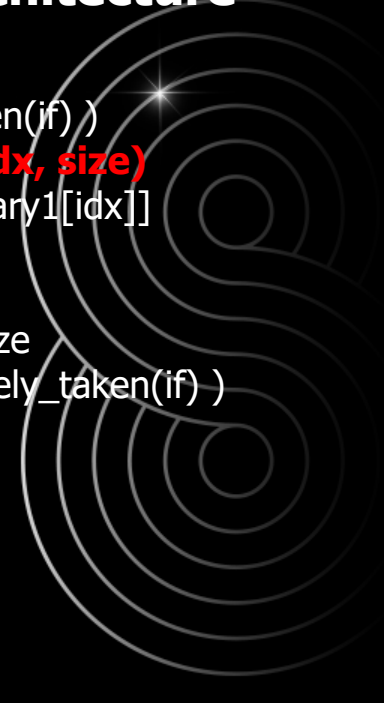
#### Software

```
if (idx < size)
    sanitize(idx, size)
    var = ary2[ary1[idx]];
```

**Clamp idx in [0, size).**

#### Micro-architecture

```
if ( likely_taken(if) )
    sanitize(idx, size)
    read ary2[ary1[idx]]
read idx
read size
res = idx < size
if ( res == likely_taken(if) )
    commit
else
    revert
...
```



# 5. Defense 2: Index sanitization

## Example: idx sanitization in Linux kernel.

```
static inline struct file *__fcheck_files(struct files_struct *files, unsigned int fd)
{
    struct fdtable *fdt = rcu_dereference_raw(files->fdt);

    if (fd < fdt->max_fds) {
        fd = array_index_nospec(fd, fdt->max_fds);
        return rcu_dereference_raw(fdt->fd[fd]);
    }
    return NULL;
}
```

```
static struct perf_event *ptrace_hbp_get_event(unsigned int note_type,
                                              struct task_struct *tsk,
                                              unsigned long idx)
{
    struct perf_event *bp = ERR_PTR(-EINVAL);

    switch (note_type) {
    case NT_ARM_HW_BREAK:
        if (idx >= ARM_MAX_BRP)
            goto out;
        idx = array_index_nospec(idx, ARM_MAX_BRP);
        bp = tsk->thread.debug.hbp_break[idx];
        break;
    case NT_ARM_HW_WATCH:
        if (idx >= ARM_MAX_WRP)
            goto out;
        idx = array_index_nospec(idx, ARM_MAX_WRP);
        bp = tsk->thread.debug.hbp_watch[idx];
        break;
    }
```

```
/*
 * array_index_nospec - sanitize an array index after a bounds check
 */
#define array_index_nospec(index, size) \
({ \
    typeof(index) _i = (index); \
    typeof(size) _s = (size); \
    unsigned long _mask = array_index_mask_nospec(_i, _s); \
    BUILD_BUG_ON(sizeof(_i) > sizeof(long)); \
    BUILD_BUG_ON(sizeof(_s) > sizeof(long)); \
    (typeof(_i)) (_i & _mask); \
})
```

linux/include/linux/nospec.h

### How can we defense against this attack?

#### Software

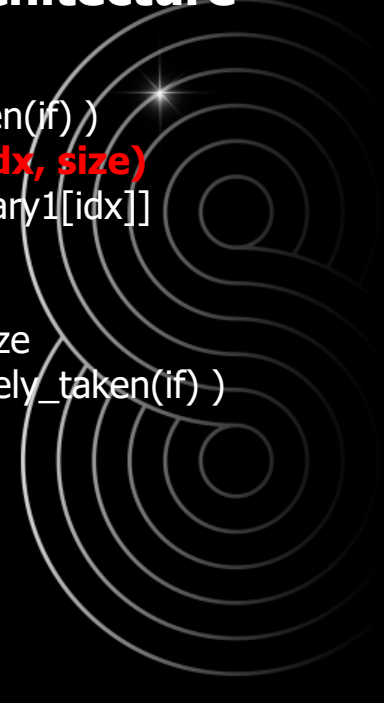
```
if (idx < size)
    sanitize(idx, size)
    var = ary2[ary1[idx]];
```

**Out-of-bounds attack is blocked!**

**You sure...?**

#### Micro-architecture

```
if ( likely_taken(if) )
    sanitize(idx, size)
    read ary2[ary1[idx]]
read idx
read size
res = idx < size
if ( res == likely_taken(if) )
    commit
else
    revert
...
```



## 6. Attack 3: Indirect branch poisoning (Spectre v2) SOSCON 2019

---

Can we execute the gadget only?

### Software

```
if (idx < size)
    sanitize(idx, size)
    var = ary2[ary1[idx]];
```

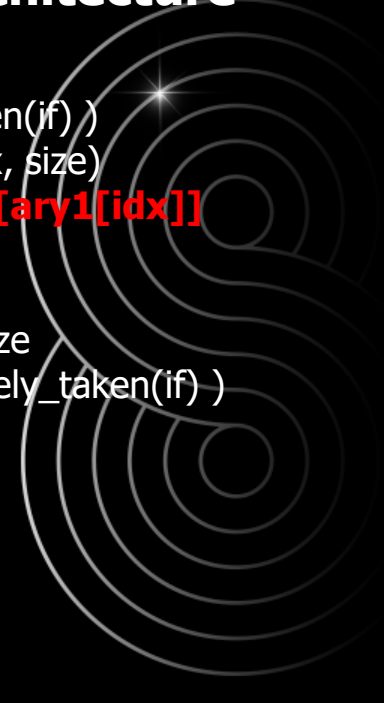
Yes, if we can **jump to the gadget** like ROP.

**But.. How can we make jump?**

**Maybe through function pointer?**

### Micro-architecture

```
if ( likely_taken(if) )
    sanitize(idx, size)
    read ary2[ary1[idx]]
read idx
read size
res = idx < size
if ( res == likely_taken(if) )
    commit
else
    revert
...
```



# 6. Attack 3: Indirect branch poisoning (Cont.)

Can we jump to the gadget?

## Software

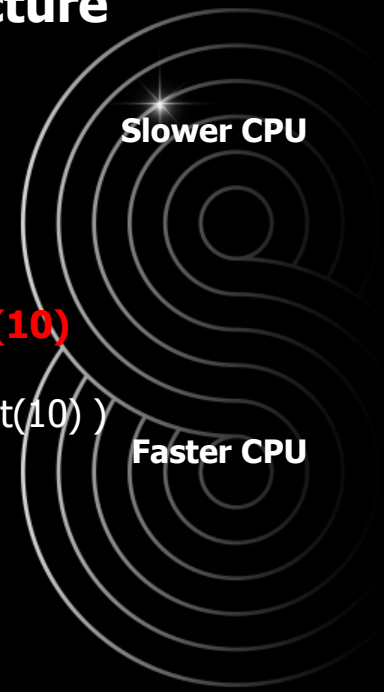
```
void (*fp)(void);  
  
void func(void)  
{  
    fp();    ← VA: 10  
}
```

Nope. →

## Micro-architecture

```
read fp  
jump fp
```

```
jump likely_target(10)  
read fp  
if ( fp == likely_target(10) )  
    commit  
else  
    revert  
jump fp
```



Yes! but we have to make  
**likely\_target(10) == gadget addr.**

## 6. Attack 3: Indirect branch poisoning (Cont.)

SOSCON 2019

### Attack steps

#### Software

```
void (*fp)(void);

void func(void)
{
    fp();          ← VA: 10
}
```

**Goal: What is in victim's ary1[-100] ?**

#### Micro-architecture

```
jump likely_target(10)
read fp
if ( fp == likely_target(10) )
    commit
else
    revert
jump fp
```

**Faster CPU**



## 6. Attack 3: Indirect branch poisoning (Cont.)

SOSCON 2019

### Attack steps

#### Victim

```
void (*fp)(void);

void func(void)
{
    fp();    ← VA: 10
}
```

#### Attacker

```
void (*fp)(void);

void func(void)
{
    jump gadget; ← VA: 10
}
```

#### Micro-architecture

```
jump likely_target(10)
read fp
if ( fp == likely_target(10) )
    commit
else
    revert
    jump fp
```

**Faster CPU**



1. Allocate the same memory as victim's and  
Modify fp() to **jump gadget**.

## 6. Attack 3: Indirect branch poisoning (Cont.)

SOSCON 2019

### Attack steps

#### Victim

```
void (*fp)(void);  
  
void func(void)  
{  
    fp();    ← VA: 10  
}
```

#### Attacker

```
void (*fp)(void);  
  
void func(void)  
{  
    jump gadget; ← VA: 10  
}
```

#### Micro-architecture

→ **jump likely\_target(10)**  
read fp  
if ( fp == likely\_target(10) )  
 commit  
else  
 revert  
 jump fp

Faster CPU



2. Execute the attacker's program multiple times.

→ Indirect branch predictor will learn that  
**likely\_target(10) == gadget addr.**



## 6. Attack 3: Indirect branch poisoning (Cont.)

SOSCON 2019

### Attack steps

#### Victim

```
void (*fp)(void);  
  
void func(void)  
{  
    fp();    ← VA: 10  
}
```

**ary2[ary1[-100]] is executed!**

#### Micro-architecture

```
jump likely_target(10)  
read fp  
if ( fp == likely_target(10) )  
    commit  
else  
    revert  
    jump fp
```

**Faster CPU**



3. Execute the victim program with **idx = -100**.

**→ ary2[ary1[-100]] will be read and cached!**

## 6. Attack 3: Indirect branch poisoning (Cont.)

SOSCON 2019

### Attack steps

#### Victim

```
void (*fp)(void);

void func(void)
{
    fp();      ← VA: 10
}
```

4. Do Flush+Reload analysis on ary2.

→ **ary2['T'] will take the shortest time.**

**Thus, ary1[-100] == 'T' !!**

#### Micro-architecture

```
jump likely_target(10)
read fp
if ( fp == likely_target(10) )
    commit
else
    revert
    jump fp
```

**Faster CPU**



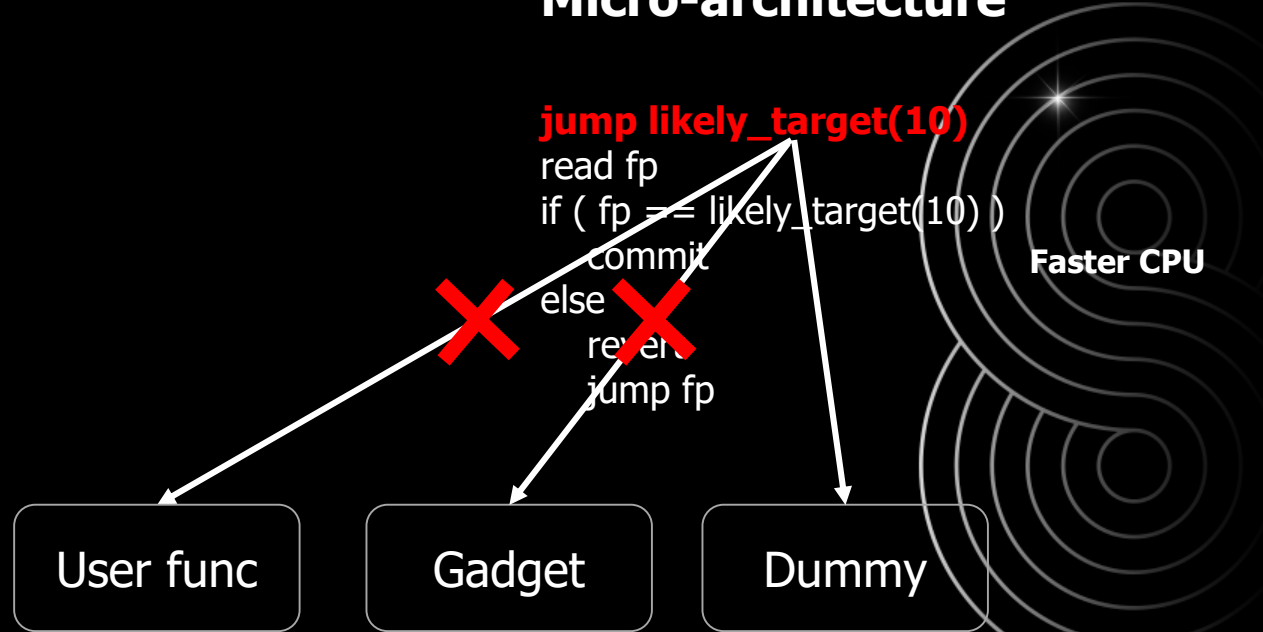
# 6. Attack 3: Indirect branch poisoning (Cont.)

How can we defend against this attack?

## Victim

```
void (*fp)(void);  
  
void func(void)  
{  
    fp();    ← VA: 10  
}
```

## Micro-architecture



How about making it like this?

## 7. Defense 3: Retpoline

SOSCON 2019

### Indirect call replacement with retpoline.

```
jmp *%rax    →  
  
capture_ret_spec:  
    call load_label  
    pause ; lfence  
    jmp capture_ret_spec  
  
load_label:  
    mov %rax, (%rsp)  
    ret
```



## 7. Defense 3: Retpoline (Cont.)

SOSCON 2019

**Goal: Capture speculative execution for indirect branch into a dummy loop.**

Details are hard to understand for beginner.

To get more details, please refer to

- Appendix-D: Retpoline details
- [intel's guidance](#)



# 7. Defense 3: Retpoline (Cont.)

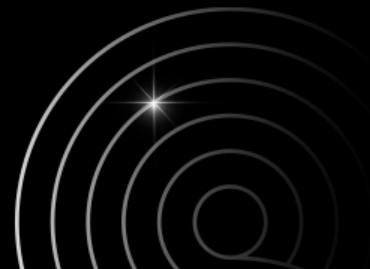
## Example: Retpoline in Linux kernel.

```
config RETPOLINE
    bool "Avoid speculative indirect branches in kernel"
    default y
    select STACK_VALIDATION if HAVE_STACK_VALIDATION
    help
        Compile kernel with the retpoline compiler options to guard against
        kernel-to-user data leaks by avoiding speculative indirect
        branches. Requires a compiler with -mindirect-branch=thunk-extern
        support for full protection. The kernel may run slower.
```

- linux/arch/x86/Kconfig

```
/*
 * JMP_NOSPEC and CALL_NOSPEC macros can be used instead of a simple
 * indirect jmp/call which may be susceptible to the Spectre variant 2
 * attack.
 */
.macro JMP_NOSPEC reg:req
#ifdef CONFIG_RETPOLINE
    ANNOTATE_NOSPEC_ALTERNATIVE
    ALTERNATIVE_2 __stringify(ANNOTATE_RETPOLINE_SAFE; jmp *\reg), \
    __stringify(RETPOLINE_JMP \reg), X86_FEATURE_RETPOLINE, \
    __stringify(lfence; ANNOTATE_RETPOLINE_SAFE; jmp *\reg), X86_FEATURE_RETPOLINE_AMD
#else
    jmp *\reg
#endif
.endm
```

- linux/arch/x86/include/asm/nospec-branch.h



**For actual v1, v2 attacks, there are more prerequisites.**

- Additional cache miss for successful speculative execution.
- Control over specific registers.
- Run programs on the same core to share BTB.
- etc..

**To get more details, please refer to**

- [Spectre Attacks: Exploiting Speculative Execution](#)
- [Google project zero blog](#)



# Attack4: Meltdown

---



SOSCON 2019

SAMSUNG OPEN SOURCE CONFERENCE 2019



# Two phases of Meltdown

---

SOSCON 2019

Goal : Reading kernel memory from user-space

1. Speculatively read kernel memory into register.
2. Stick the register to the leak gadget.



## Software (User)

```
int *ptr = Kernel_Addr;  
void func(void)  
{  
    int a = *ptr;  
}
```

Segmentation  
fault!

What we expect

Actually happens

1. Segmentation Fault
2. Die or Run signal handler

1. int a = \*ptr; (a == kernel value)
2. Segmentation Fault
3. a = 0; (clear kernel value)
4. Die or Run signal handler

**We can read kernel-memory!!**



## Micro-architecture

### Software

```
int *ptr = Kernel_Addr;  
  
void func(void)  
{  
    int a = *ptr;  
}
```

Vulnerable-to-Meltdown

```
pa = convert_va_to_pa(ptr);  
if (is_kernel_address(pa))  
    error;  
else  
    read pa;
```

```
pa = convert_va_to_pa(ptr);  
read pa;  
if (is_kernel_address(pa))  
    error;
```

Let's turn our focus to this!



# Speculatively read kernel memory (Zoom in more)

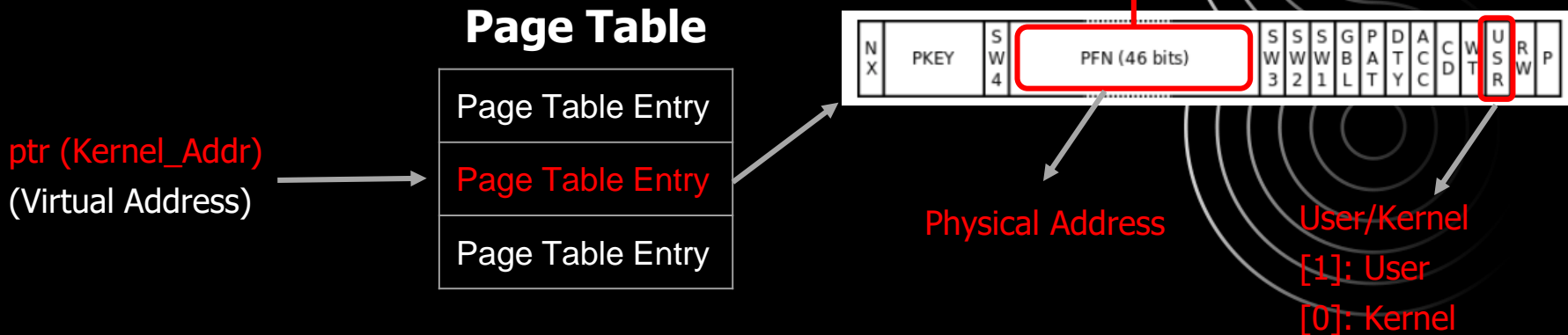
SOSCON 2019

## Software

```
int *ptr = Kernel_Addr;  
  
void func(void)  
{  
    int a = *ptr;  
}
```

## Micro-architecture

- (1) pa = convert\_va\_to\_pa(ptr);
- (2) read pa;
- (3) if (is\_kernel\_address(pa))  
 error;



## Leak gadget + FLUSH-RELOAD + Meltdown = Kernel memory leak!

```
int *ptr = Kernel_Addr;
```

```
void func(void)
```

```
{  
    int a = *ptr;    // (1) a = kernel-value  
    var = ary1[a];  // (2) [leak gadget] kernel-value dependent access → Cache loading  
    // (3) invalid access! fault!  
}
```

```
void fault_handler(...)
```

```
{  
    for (i=0; i<SIZE; i++) {  
        if (time_to_access(ary[i]) < CONSTANT)  
            // (4) i is kernel-value!  
    }  
}
```



# Defense4: Page Table Isolation

---



SOSCON 2019

SAMSUNG OPEN SOURCE CONFERENCE 2019

## Micro-architecture

### Software

```
int *ptr = Kernel_Addr;  
  
void func(void)  
{  
    int a = *ptr;  
}
```

**Fix**

```
pa = convert_va_to_pa(ptr);  
if (is_kernel_address(pa))  
    read pa;  
else  
    error;
```

```
pa = convert_va_to_pa(ptr);  
read pa;  
if (is_kernel_address(pa))  
    error;
```

- (1) It's infeasible to update a lot of deployed CPUs..
- (2) then.. what is the root cause in the view of software?
- (3) Figure out solution to fix Meltdown with software fix! (Linux kernel)



## Software

```
int *ptr = Kernel_Addr;  
  
void func(void)  
{  
    int a = *ptr;  
}
```

## Micro-architecture

```
(1) pa = convert_va_to_pa(ptr);  
(2) read pa;  
(3) if (is_kernel_address(pa))  
    error;
```



**Why does it allow User to access Kernel??**

## Page Table

`ptr (Kernel_Addr)`  
(Virtual Address)



`pa (physical address)`





## Software

```
int *ptr = Kernel_Addr;  
  
void func(void)  
{  
    int a = *ptr;  
}
```

## Micro-architecture

```
(1) pa = convert_va_to_pa(ptr); ⇒ Fail!  
  
(2) read pa;  
  
(3) if (is_kernel_address(pa))  
    error;
```



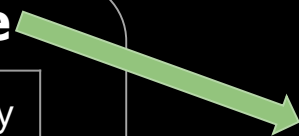
**If no Page Table Entry for Kernel?**

### Page Table

`ptr (Kernel_Addr)`  
(Virtual Address)



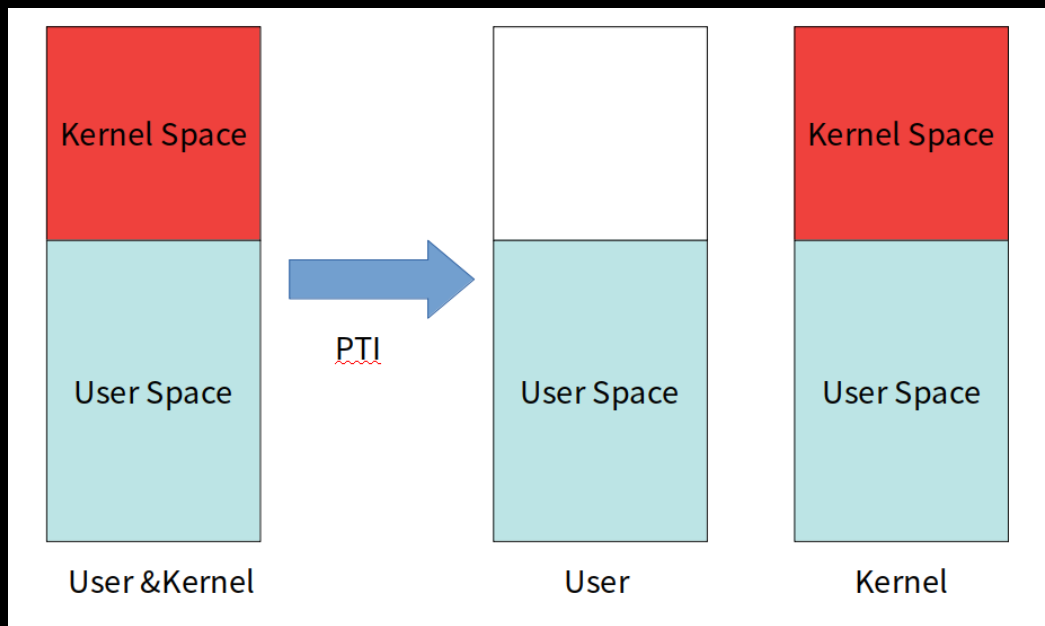
Page Table Entry
-
Page Table Entry



**Managed by Software (Linux)!**



## Unmapping Kernel Space in User Mode



- Too difficult details for a beginner to understand...
- Links to get more details
  - <https://gruss.cc/files/kaiser.pdf>
  - <https://jinb-park.blogspot.com/2019/06/deep-dive-into-page-table-isolation.html>



# Attack5: L1TF (L1 Terminal Fault)



SOSCON 2019

SAMSUNG OPEN SOURCE CONFERENCE 2019

## Software

```
int *ptr = Not_Exist_Addr;  
  
void func(void)  
{  
    int a = *ptr;  
}
```

Vulnerable-to-L1TF

## Micro-architecture

```
pa = convert_va_to_pa(ptr);  
if (!is_present(pa))  
    error;
```

```
read pa;
```

```
pa = convert_va_to_pa(ptr);  
read pa;
```

```
if (!is_present(pa))  
    error;
```



## Software

```
int *ptr = Not_Exist_Addr;  
  
void func(void)  
{  
    int a = *ptr;  
}
```

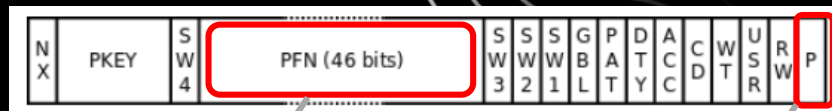
## Micro-architecture

- (1) pa = convert\_va\_to\_pa(ptr);
- (2) read pa;
- (3) if (!is\_present(pa))  
 error;

From address Zero?  
We can't make sure that what's in it.  
Secret? Dummy value?

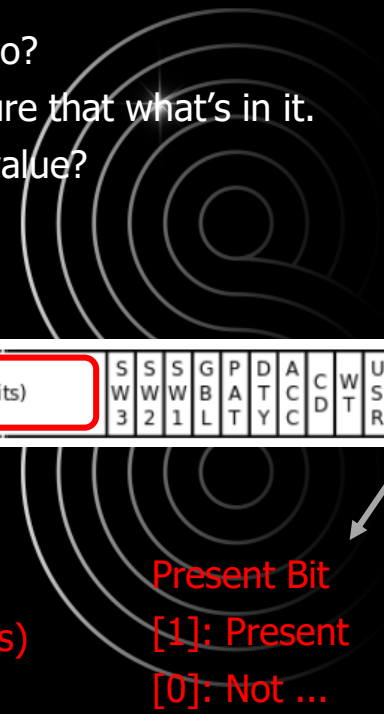
## Page Table

ptr (Not\_Exist\_Addr)  
(Virtual Address)



Normally Zero  
(Physical Address)

Present Bit  
[1]: Present  
[0]: Not ...



## Software

```
int *ptr = Not_Exist_Addr;  
  
void func(void)  
{  
    int a = *ptr;  
}
```

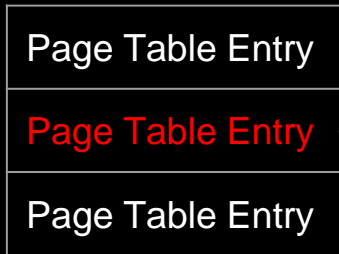
## Micro-architecture

- (1) pa = convert\_va\_to\_pa(ptr);
- (2) read pa;
- (3) if (!is\_present(pa))  
 error;

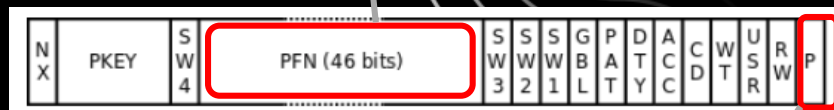
**Kernel Memory**

**Attacker can control!**

## Page Table



ptr (Not\_Exist\_Addr)  
(Virtual Address)



Present Bit  
[1]: Present  
[0]: Not ...

## Software

```
int *ptr = Not_Exist_Addr;  
  
void func(void)  
{  
    int a = *ptr;  
}
```

Virtual Address

## Page Table

Page Table Entry

Page Table Entry  
(Kernel Address)

Page Table Entry

Physical  
Address

: If attackers control the physical address of page table, they can read arbitrary kernel memory!

: But we are user, not kernel. It's infeasible for user to control page table.

: That's why L1TF is not as popular as Meltdown.





## Software

```
int *ptr = Not_Exist_Addr;  
  
void func(void)  
{  
    int a = *ptr;  
}
```

Virtual Address

## Page Table

Page Table Entry
Page Table Entry (Arbitrary Address)
Page Table Entry

Out-Of-RAM

User-1 Memory

User-2 Memory

Kernel Memory

Lucky Case!

Lucky Case!

: There might still be possibility relying on luck! How can we eliminate such possibility?



## Leak gadget + FLUSH-RELOAD + L1TF = Kernel memory leak!

```
int *ptr = Not_Exist_Addr; // Physical address ==> Kernel address
```

```
void func(void)
```

```
{  
    int a = *ptr; // (1) a = kernel-value  
    var = ary1[a]; // (2) [leak gadget] kernel-value dependent access → Cache loading  
    // (3) invalid access! fault!  
}
```

```
void fault_handler(...)
```

```
{  
    for (i=0; i<SIZE; i++) {  
        if (time_to_access(ary[i]) < CONSTANT)  
            // (4) i is kernel-value!  
    }  
}
```



- L1TF has some prerequisites,
  - L1 cache must contain entry for target physical address.
- Links to get more details
  - [https://www.usenix.org/system/files/conference/usenixsecurity18/sec18-van\\_bulck.pdf](https://www.usenix.org/system/files/conference/usenixsecurity18/sec18-van_bulck.pdf)
  - <https://software.intel.com/security-software-guidance/insights/deep-dive-intel-analysis-l1-terminal-fault>



# Defense5: PTE Inversion

---

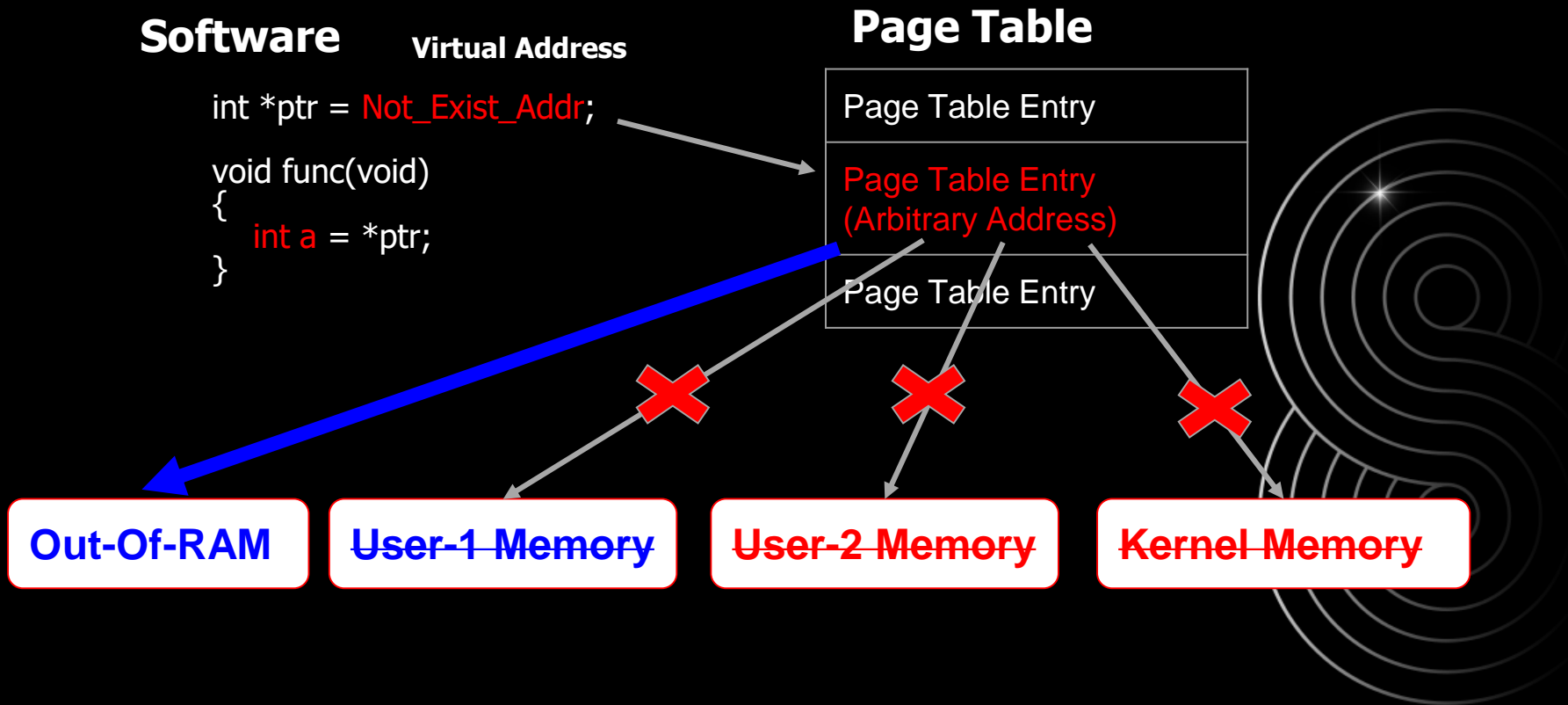


SOSCON 2019

SAMSUNG OPEN SOURCE CONFERENCE 2019

# Forcing PTE to point to Out-Of-RAM

SOSCON 2019



**Out-Of-RAM**

## Page Table

Page Table Entry  
(Physical Address)

: Assume the size of RAM is 16GB. → **0x00 ~ 0x400000000**

: Setting one to the top bit. → **0x00 ~ 0x8000000400000000**  
→ **Over 1024\*1024 TBs (Out-Of-All-Kinds-of-RAM)**

: Performing an Inversion!

→ Original: **0x00000000000001000**

→ Inversion: **0x11111111111110111**



## Page Table

Page Table Entry  
(Physical Address)

For all accesses on PTE,  
pte\_pfn running!

```
// PTE inversion

static inline bool __pte_needs_invert(u64 val)
{
    // NB. PTE is exist, but invalid.
    return val && !(val & _PAGE_PRESENT);
}

static inline u64 protnone_mask(u64 val)
{
    // NB. If PTE inversion needed,
    // return the mask for PTE to point to invalid memory.
    return __pte_needs_invert(val) ? ~0ull : 0;
}

static inline unsigned long pte_pfn(pte_t pte)
{
    phys_addr_t pfn = pte_val(pte);
    // NB. Masking PFN (physical address)
    // after masking, it's pointing to invalid memory.
    pfn ^= protnone_mask(pfn);
    return (pfn & PTE_PFN_MASK) >> PAGE_SHIFT;
}
```

THANK YOU

SOSCON 2019

SAMSUNG OPEN SOURCE CONFERENCE 2019





# Appendix-A

---



SOSCON 2019

SAMSUNG OPEN SOURCE CONFERENCE 2019

# All attacks so far

SOSCON 2019

Vulnerability	N	CVE	Public Vulnerability Name (codename)	Affected CPU architectures and mitigations						
				Intel <sup>[5]</sup>				AMD <sup>[6]</sup>		
				Ice Lake <sup>[7]</sup>	Cascade Lake	Whiskey Lake, Coffee Lake (9th gen) <sup>[8]</sup>	Coffee Lake (8th gen)*	Zen 1 / Zen 1+	Zen 2 <sup>[9]</sup>	
Spectre	1	<a href="#">2017-5753</a>	Bounds Check Bypass	OS/VMM				Firmware + OS/VMM	Hardware + OS/VMM	
Spectre	2	<a href="#">2017-5715</a>	Branch Target Injection	Hardware + OS		Firmware + OS	Firmware + OS			
SpectreRSB <sup>[10]</sup> /ret2spec <sup>[11]</sup>	2	<a href="#">2018-15572</a>	Return Mispredict							
Meltdown	3	<a href="#">2017-5754</a>	Rogue Data Cache Load	Not affected			Firmware	Not affected		
Spectre-NG	3a	<a href="#">2018-3640</a>	Rogue System Register Read	Not affected <sup>[12]</sup>		Firmware				
Spectre-NG	4	<a href="#">2018-3639</a>	Speculative Store Bypass	Hardware + OS/VMM <sup>[12]</sup>		Firmware + OS		OS/VMM	Hardware + OS/VMM	
Foreshadow	5	<a href="#">2018-3615</a>	L1 Terminal Fault (L1TF)	Not affected			Firmware	Not affected		
Spectre-NG		<a href="#">2018-3665</a>	Lazy FP State Restore			OS/VMM <sup>[13]</sup>		Not affected		
Spectre-NG	1.1	<a href="#">2018-3693</a>	Bounds Check Bypass Store			OS/VMM <sup>[14]</sup>				
Spectre-NG	1.2		Read-only Protection Bypass (RPB)							
Foreshadow-OS		<a href="#">2018-3620</a>	L1 Terminal Fault (L1TF)							
Foreshadow-VMM		<a href="#">2018-3646</a>	L1 Terminal Fault (L1TF)							
ZombieLoad		<a href="#">2018-12130</a>	Microarchitectural Fill Buffer Data Sampling (MFBDS)	Not affected			Firmware + OS	Not affected		
RIDL		<a href="#">2018-12127</a> <a href="#">2019-11091</a>	Microarchitectural Data Sampling Uncacheable Memory MLPDS/MDSUM							
Fallout		<a href="#">2018-12126</a>	Microarchitectural Store Buffer Data Sampling (MSBDS)							
Spectre SWAPGS	1	<a href="#">2019-1125</a>	None yet <sup>[15][16][17][18]</sup>	Same as Spectre 1						

: From [https://en.wikipedia.org/wiki/Transient\\_execution\\_CPU\\_vulnerabilities](https://en.wikipedia.org/wiki/Transient_execution_CPU_vulnerabilities)

Discussed here

Worth to look at

- The most useful link : [A dynamic Tree](#)
- [Intel Sandybridge Microarchitecture](#)



# Appendix-B

---



SOSCON 2019

SAMSUNG OPEN SOURCE CONFERENCE 2019

Used to predict indirect branch target.

Stores recent branch history (source & destination address pairs).

Typically, stores **partial** src addr and dst addr.

Shared per cpu core.

Ref: [Jump Over ASLR: Attacking Branch Predictors to Bypass ASLR](#)



# Branch Target Buffer (Cont.)

SOSCON 2019

## BTB

src	dst
0x80100	0x56780A00
0x80200	0x56780B00
0x80300	0x56780C00



Only low 20 bits of src addr are stored!  
(It depends on the CPU design.)

## Code

```
0x56780100: jmp [0x56781010]
...
0x56780200: jmp [0x56781020]
...
0x56780300: jmp [0x56781030]
```

## Data

```
0x56781010: 0x56780A00
0x56781020: 0x56780B00
0x56781030: 0x56780C00
```



# Appendix-C

---



SOSCON 2019

SAMSUNG OPEN SOURCE CONFERENCE 2019

# Spectre v2 detailed steps

SOSCON 2019

## Victim

### Memory

...	
0x56780100	jmp [0x56781010]
...	
0x56780200	jmp [0x56781020]
...	
0x56780300	jmp [0x56781030]
...	
0x56780F00	<b>var = ary2[ary1[idx]]</b>
...	

### BTB

src	dst
0x80100	0x56780A00
0x80200	0x56780B00
0x80300	<b>0x56780C00</b>

Attacker will change this address to **0x56780F00**.

Gadget address!

### Data

addr	value
0x56781010	0x56780A00
0x56781020	0x56780B00
<b>0x56781030</b>	0x56780C00

Attacker will flush this memory for speculative execution.





# Spectre v2 detailed steps (Cont.)

SOSCON 2019

## Attacker: 1. Mistrain BTB.

### Memory

	...
0x56780100	ret
	...
0x56780200	ret
	...
0x56780300	ret
	...
0x56780F00	<b>ret</b>
	...

### BTB

src	dst
0x80100	0x56780A00
0x80200	0x56780B00
0x80300	0x56780C00

- ← 1) Allocate memory same as victim's.  
2) Fill up with ret instructions.

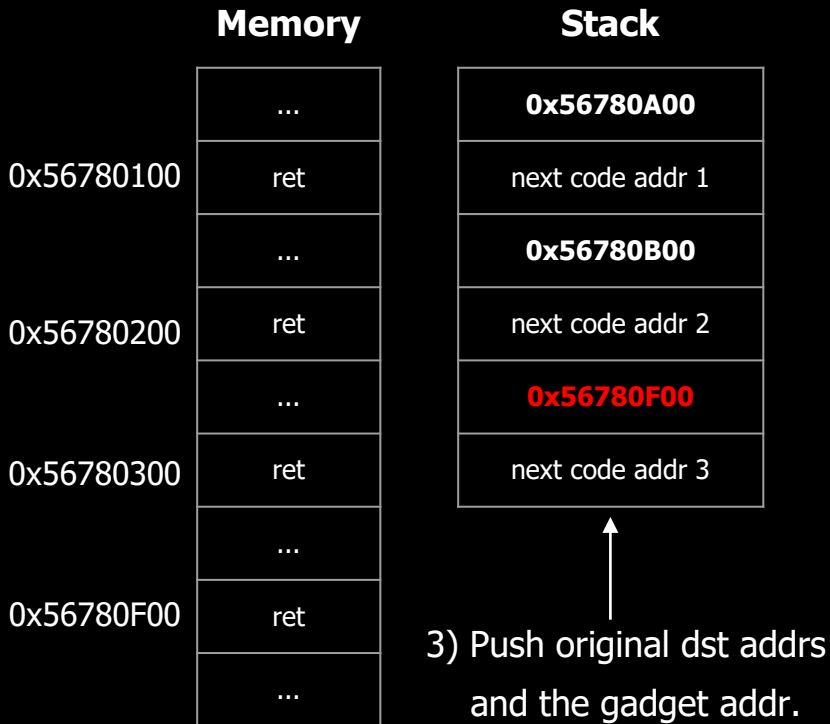


# Spectre v2 detailed steps (Cont.)

SOSCON 2019

## Attacker: 1. Mistrain BTB.

Before: 0x56780300 → 0x56780C00  
Now: 0x56780300 → **0x56780F00**



**BTB**

src	dst
0x80100	0x56780A00
0x80200	0x56780B00
0x80300	<b>0x56780F00</b>

**Changed!!**

4) Jump to original src addrs.

```
jmp 0x56780100  
jmp 0x56780200  
jmp 0x56780300
```

5) Iterate 3) and 4).

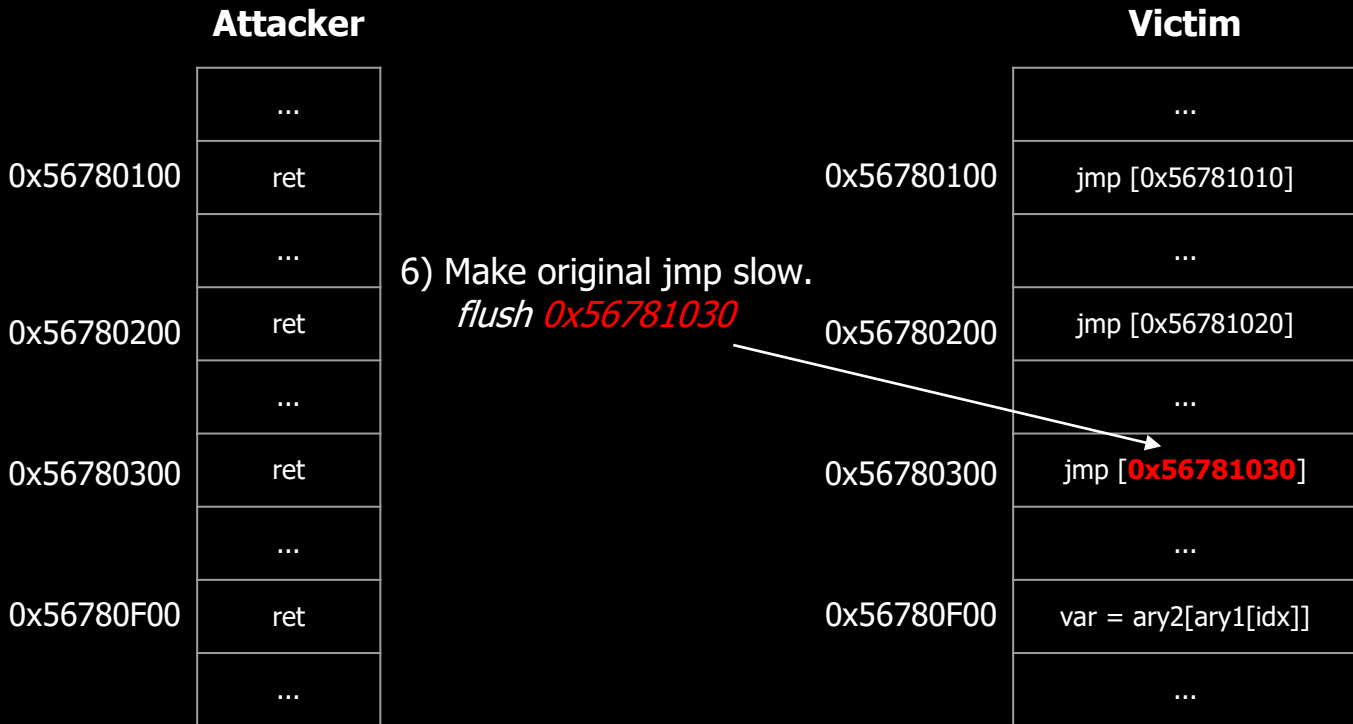


# Spectre v2 detailed steps (Cont.)

SOSCON 2019

## Attacker: 1. Mistrain BTB.

Before: 0x56780300 → 0x56780C00  
Now: 0x56780300 → **0x56780F00**



# Spectre v2 detailed steps (Cont.)

SOSCON 2019

Attacker: 2. Execute victim with malicious idx.

Victim

	...
0x56780100	jmp [0x56781010]
	...
0x56780200	jmp [0x56781020]
	...
0x56780300	jmp [0x56781030]
	...
0x56780F00	<b>var = ary2[ary1[idx]]</b>
	...

Flushed!

BTB

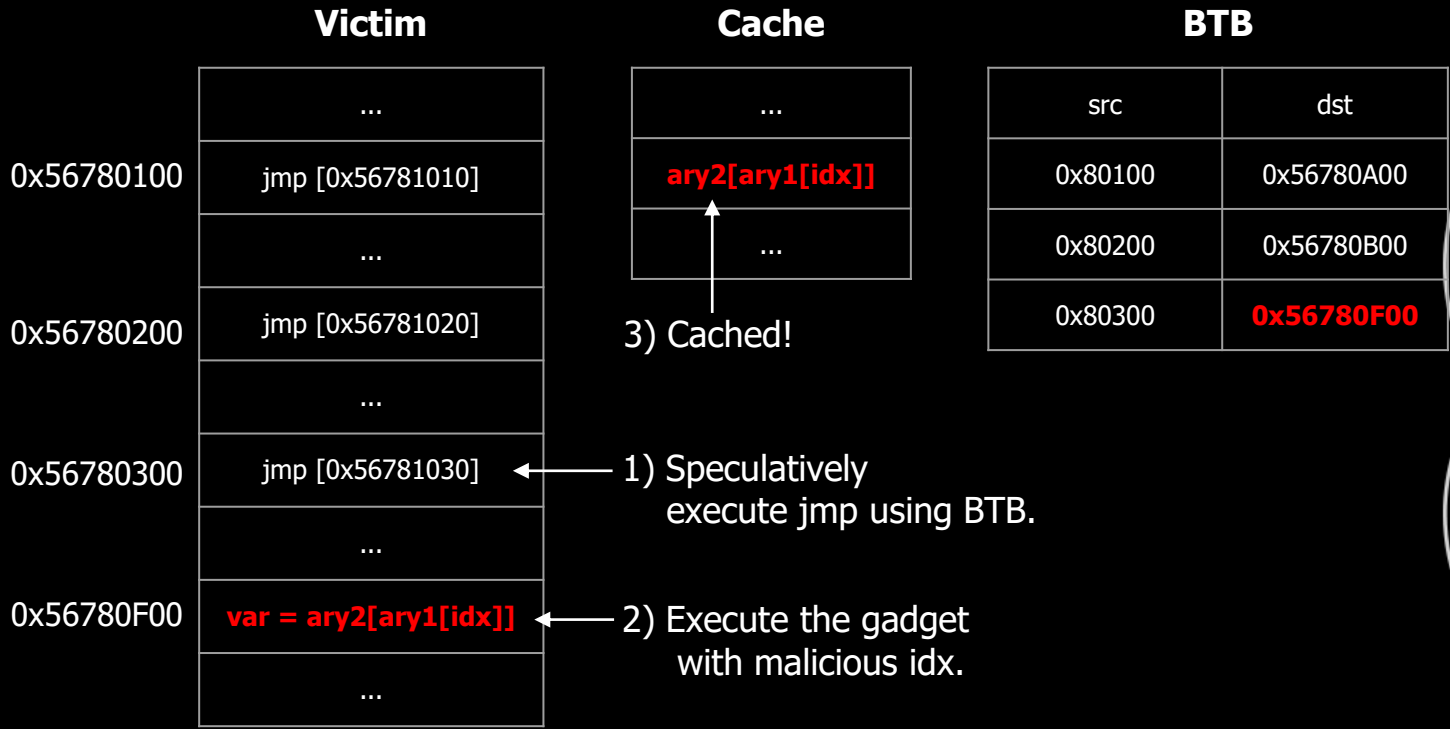
src	dst
0x80100	0x56780A00
0x80200	0x56780B00
0x80300	<b>0x56780F00</b>

Mistrained to the gadget.

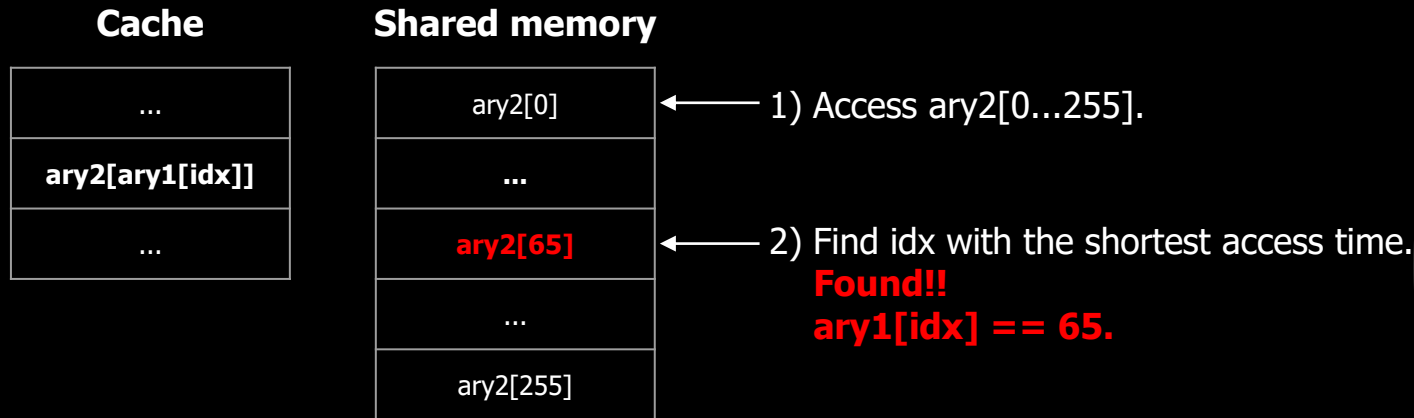


# Spectre v2 detailed steps (Cont.)

Attacker: 2. Execute victim with malicious idx.



## Attacker: 3. Retrieve value via Flush + Reload analysis.



# Appendix-D

---



SOSCON 2019

SAMSUNG OPEN SOURCE CONFERENCE 2019

**Goal: Avoid predictions which use the BTB.**

## Steps

1. Change indirect jmp to ret.

When predicting ret's destination, **RSB(Return Stack Buffer)** is used ahead of BTB.

2. Capture speculative execution into a fake infinite loop.





## Indirect call replacement with retpoline.

```
jmp *%rax    →  
  
capture_ret_spec:  
    pause ; lfence  
    jmp capture_ret_spec  
  
load_label:  
    mov %rax, (%rsp)  
    ret
```



## Original path

```
call load_label <--
capture_ret_spec:
    pause ; lfence
    jmp capture_ret_spec
load_label:
    mov %rax, (%rsp)
    ret
```

rax register

Original jmp dst

Stack

capture\_ret\_spec



## Original path

```
call load_label  
capture_ret_spec:  
    pause ; lfence  
    jmp capture_ret_spec  
load_label:  
    mov %rax, (%rsp) <--  
    ret
```

rax register

Original jmp dst

Stack

Original jmp dst



## Original path

```
capture_ret_spec:
    call load_label

load_label:
    mov %rax, (%rsp)
    ret
```

rax register

Original jmp dst

Stack

Original-jmp-dst

**<-- Jump to the original dst.**



## Speculative path

```
call load_label <--
capture_ret_spec:
    pause ; lfence
    jmp capture_ret_spec
load_label:
    mov %rax, (%rsp)
    ret
```

rax register

Original jmp dst

**RSB (Return Stack Buffer)**

**capture\_ret\_spec**



## Speculative path

call load\_label

capture\_ret\_spec:

pause ; lfence

jmp capture\_ret\_spec

load\_label:

mov %rax, (%rsp) <--

ret

rax register

Original jmp dst

RSB (Return Stack Buffer)

capture\_ret\_spec



## Speculative path

```
call load_label  
capture_ret_spec:  
    pause ; lfence  
    jmp capture_ret_spec  
load_label:  
    mov %rax, (%rsp)  
    ret
```

rax register

Original jmp dst

RSB (Return Stack Buffer)

capture\_ret\_spec

<-- Jump to capture\_ret\_spec.



## Speculative path

```
call load_label  
capture_ret_spec:  
    pause ; lfence  
    jmp capture_ret_spec  
load_label:  
    mov %rax, (%rsp)  
    ret
```

rax register

Original jmp dst

RSB (Return Stack Buffer)



**Captured in the infinite loop!**

