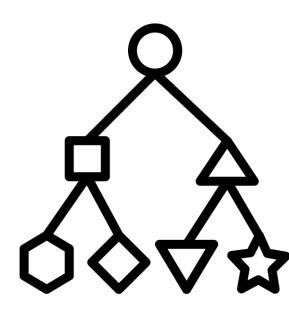# Summary

- Background on races

- Classification on races
- **Unexploitable races**

- New technique turning **unexploitable races to exploitable races**
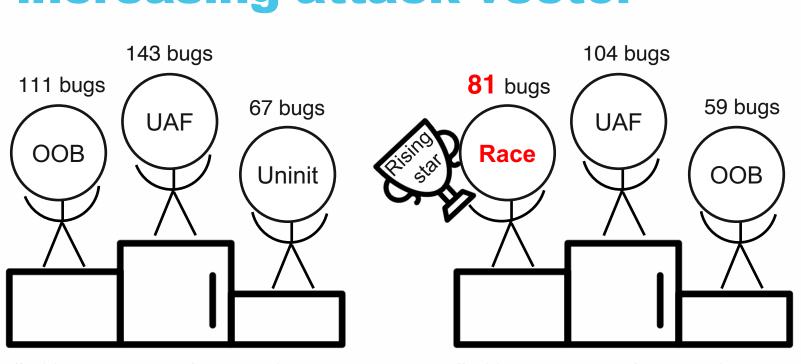
# Race condition is an increasing attack vector

30 bugs

15 bugs

UAF

7 bugs

OOB

Race

# of fixed bugs that Syzkaller found in **2017**

143 bugs

111 bugs

UAF

67 bugs

OOB

Uninit

# of fixed bugs that Syzkaller found in **2018**

104 bugs

**81** bugs

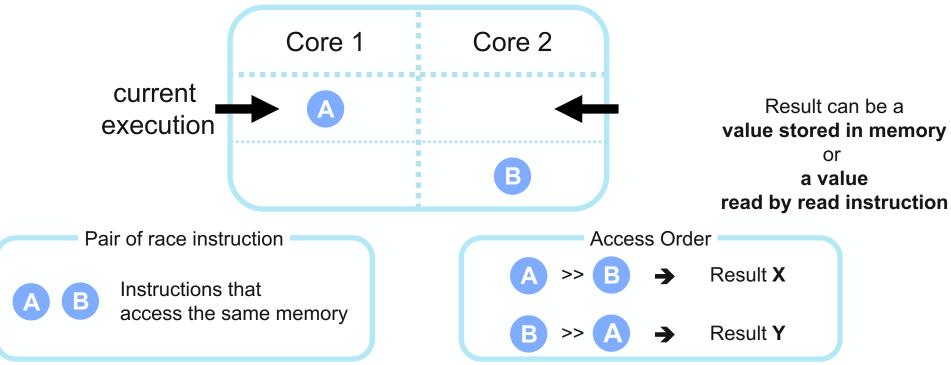Rising star

**Race**

UAF

59 bugs

OOB

# of fixed bugs that Syzkaller found in **2019**

- Razzer, IEEE S&P 2019, found more than **30 race bugs**.

- KCSAN, developed by Google 2019, found more than **300 race bugs**.

# Background : Race condition

Core 1    Core 2

current execution ➡ **A**

**B** ⬅

Result can be a
**value stored in memory**
or
**a value
read by read instruction**

Pair of race instruction

**A** **B**    Instructions that
access the same memory

Access Order

**A** >> **B** ➔ Result **X**

**B** >> **A** ➔ Result **Y**

- **Accessing the same memory** location from two processors

➔ **Execution results are different** depending on the access order.
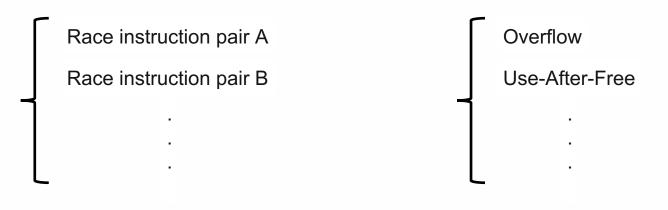
# Background : Race Condition Vulnerability

**Race Condition Vulnerability** = **Race Condition + Memory Corruption**

Race instruction pair A

Race instruction pair B
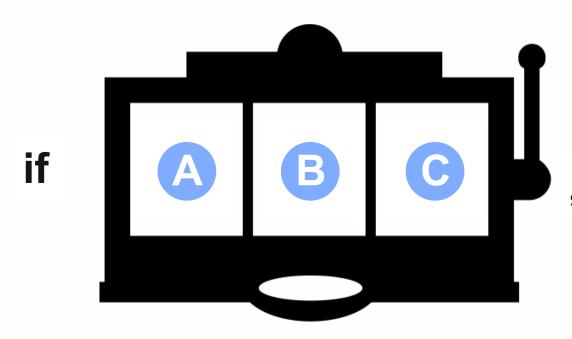
.
.
.

Overflow

Use-After-Free

.
.
.

# Background : to trigger Race Condition Vulnerability

if [A B C] , then **memory corruption** occurs.

Brute forcing :
**Try until success**

# Background : Exploitability of Race Condition Vulnerability

**Exploitable Races?** = **A very specific memory access order** + Availability of Memory Corruption

(e.g., if (A) >> (B) >> (C) , then)

# Classification of Race Condition Vulnerability

Race Condition Vulnerability

**Single Variable**
Race Condition

Race instruction pair 1 for **M1**

Race instruction pair 2 for **M1**

**Single variable**

…

**Multi Variable**
Race Condition

Race instruction pair 1 for **M1**

Race instruction pair 2 for **M2**

**Multi variable**

…

# Single-variable Race Condition



Core 1    Core 2

Pair

A

**Time Window**

B

C

Pair

→ Control Flow Dependency

······▶ Data Flow Dependency

```
raw_sendmsg()
{
    …
A   if ( ! inet->hdrincl ) {
        // initialize rfv variable
        rfv.msg = msg;
        …
    }
C   if ( ! inet->hdrincl ) {
        memcpy(to, rfv->hdr.c, … );
    }
    …
}
```

inet->hdrincl is 1 → A

inet->hdrincl is 0 → C

```
do_ip_setsockopt()
{
    …
B   inet->hdrincl = 0;    ←
    …
}
```

Case study : CVE-2017-17712

if  A  >>  B  >>  C  , then uninitialized buffer use occurs.

# Exploitability of Single-variable Race



- Brute-forcing would somehow trigger the race

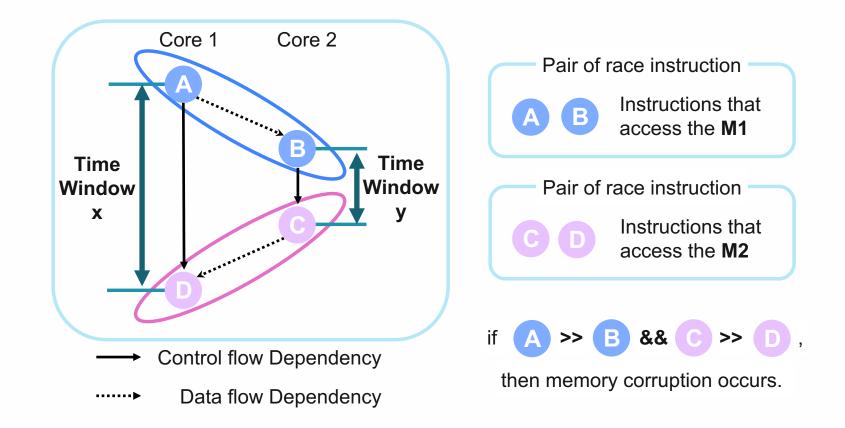  ➔ if B can be executed within the time window

- The smaller the time window is, the lower the probability of successful races.

# Multi-variable Race Condition



Core 1    Core 2

Time Window x

Time Window y

→ Control flow Dependency

┄┄► Data flow Dependency

Pair of race instruction

A B    Instructions that access the **M1**

Pair of race instruction

C D    Instructions that access the **M2**

if A >> B && C >> D ,

then memory corruption occurs.
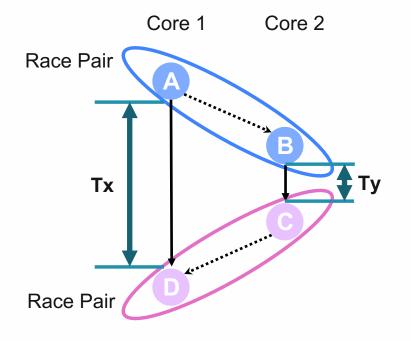
# Multi-variable Race Condition

# Exploitability of Inclusive Multi-variable Race
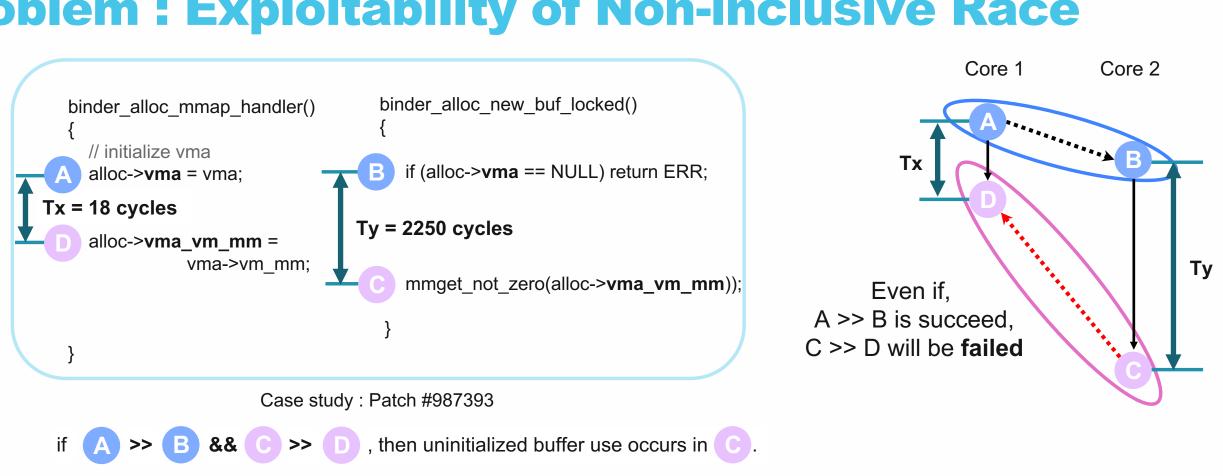


- Brute-force somehow works.

- The more similar the two time windows are, the lower the probability that a race will occur.

# Problem : Exploitability of Non-inclusive Race

```
binder_alloc_mmap_handler()        binder_alloc_new_buf_locked()
{                                   {
    // initialize vma
    alloc->vma = vma;                   if (alloc->vma == NULL) return ERR;   (B)

(A)                                 (B)
    Tx = 18 cycles                      Ty = 2250 cycles

(D) alloc->vma_vm_mm =
            vma->vm_mm;                 mmget_not_zero(alloc->vma_vm_mm));   (C)

                                    }
}
```

Case study : Patch #987393

if (A) >> (B) && (C) >> (D) , then uninitialized buffer use occurs in (C) .

Core 1    Core 2

Tx

Ty

Even if,
A >> B is succeed,
C >> D will be **failed**

- Brute-force never works.

- **impossible to execute** with the order of (A) >> (B) && (C) >> (D) .

# Problem : Exploitability of Non-inclusive Race

|  | Tx | Ty |
|---|---|---|
| CVE-2017-15265 | 35 | 450 |
| CVE-2019-1999 | 150 | 1,800 |
| CVE-2019-2025 | 50 | 600 |
| CVE-2019-6974 | 18 | 1,210 |
| #1035566 | 1,153 | 13,121 |
| #987393 | 18 | 2,250 |
| #759959 | 120 | 730 |
| . |  |  |
| . |  |  |
| . |  |  |

Non-inclusive race vulnerabilities found in linux kernel



Core 1    Core 2

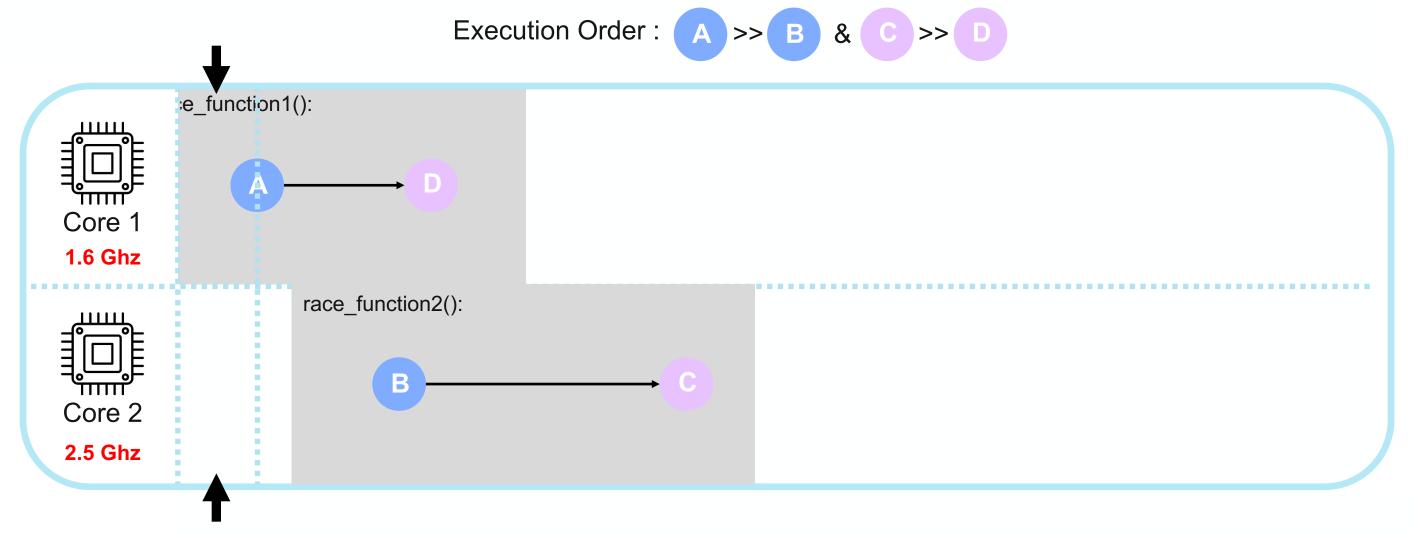Tx

Ty

Even if,
A >> B is succeed,
C >> D will be **failed**

- Brute-force never works.

- **impossible to execute** with the order of A >> B && C >> D .

15

# Previous method : Using Different Core Latency

Execution Order :  A >> B  &  C >> D

e_function1():

A ——→ D

Core 1
**1.6 Ghz**

race_function2():

B ——————————→ C

Core 2
**2.5 Ghz**

- e.g., Qualcomm Snapdragon 845 4x 2.5GHz, 4x 1.6GHz

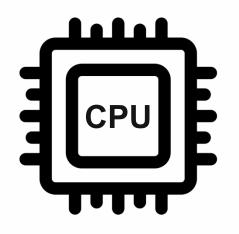# Previous method : Using Different Core Latency

Execution Order : A >> B & C >> D

Core 1
1.6 Ghz

D

Core 2
2.5 Ghz

C

- e.g., Qualcomm Snapdragon 845 4x 2.5GHz, 4x 1.6GHz
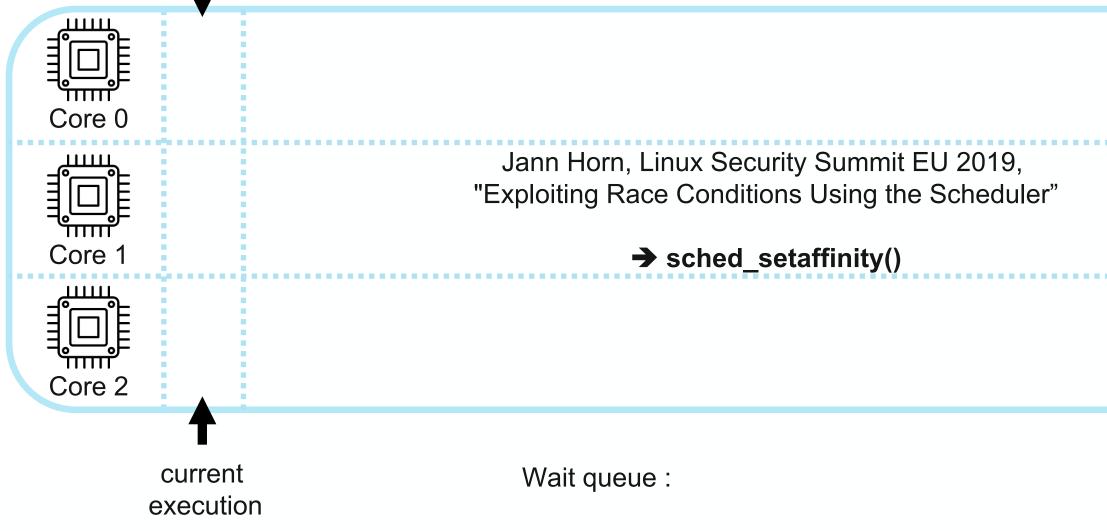
# Limitations of Use Different Core Latency

**CPU**

**CPU dependency**

- **Must use the CPU** that latency between the cores are different.

- Not applicable to vulnerabilities with large time window differences

# Previous Approach : Using scheduler (CONFIG_PREEMPT)
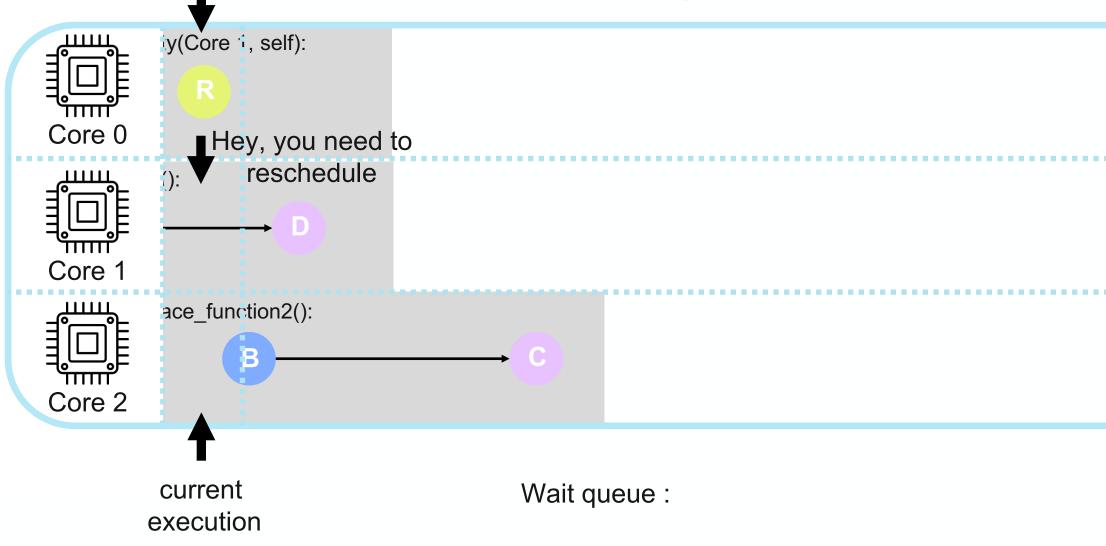
Execution Order :   A  >>  B  &  C  >>  D

Core 0

Core 1

Core 2

Jann Horn, Linux Security Summit EU 2019,
"Exploiting Race Conditions Using the Scheduler"

➔ **sched_setaffinity()**

current
execution

Wait queue :

# Previous Approach : Using scheduler (CONFIG_PREEMPT)

Execution Order : **A**



Core 0

**R**

ty(Core 1, self):

Hey, you need to reschedule

Core 1

():

**D**

Core 2

ace_function2():

**B** → **C**

current execution

Wait queue :

# Previous Approach : Using scheduler (CONFIG_PREEMPT)

Execution Order :  A  >>  B  &  C

Core 0

ty(Core 1, self):

R

Core 1

C

Core 2

current
execution

race_function1():

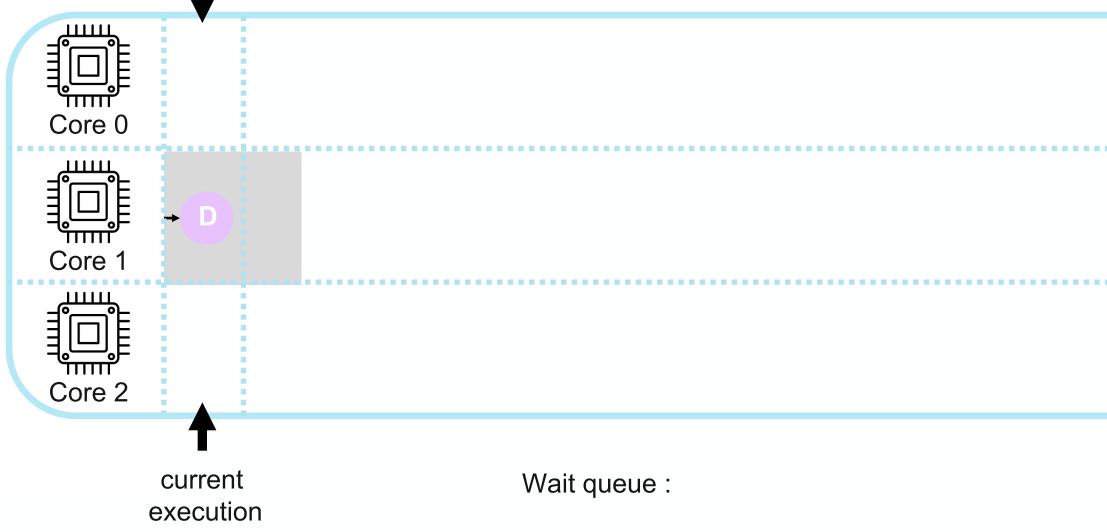Wait queue :  →  D
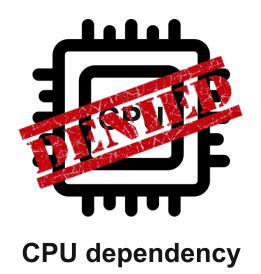
# Limitation of Using scheduler



**Configuration dependency**

- Can be used when COFIG_PREEMPT option is applied.

- Linux uses **CONFIG_PREEMPT_VOLUTARY** option **by default.**

# Each of methods has obvious limitations



**CPU dependency**
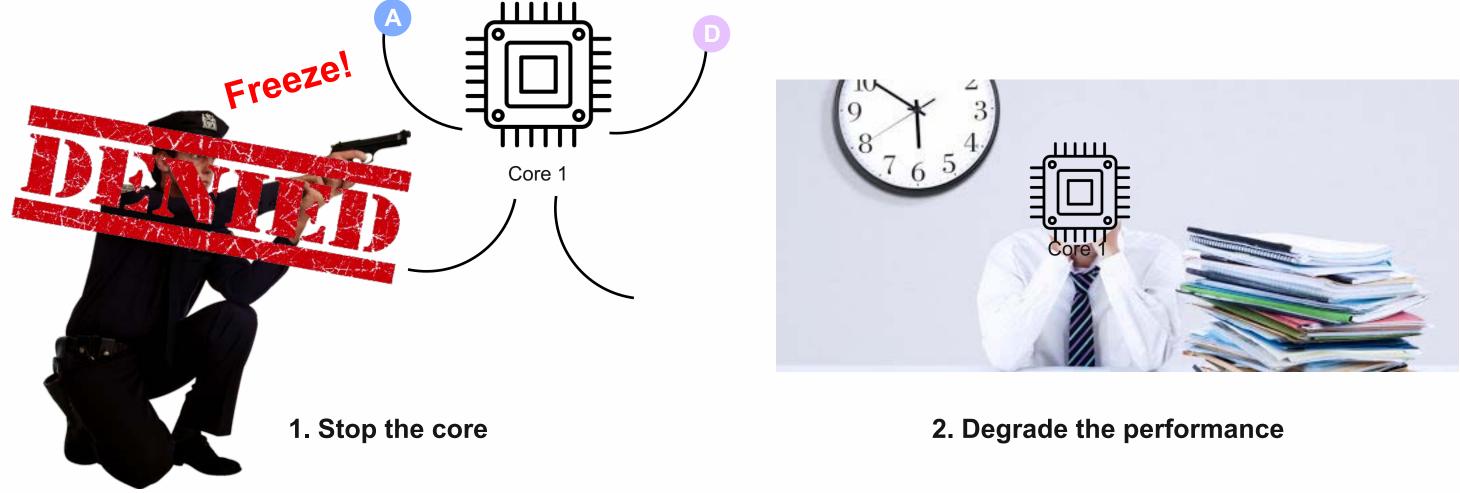


**Configuration dependency**

- All previous methods are hard to be used in general.

- We need **a new method** that extends the time window.

# How to extend the time window?

Freeze!

**DENIED**

A

D

Core 1

Core 1

**1. Stop the core**

**2. Degrade the performance**

# ExpRace

Performance :

Slow                  Fast

A

Interrupt handler!
Interrupt handler!
Interrupt handler!

D

Core 1

interrupt
interrupt
interrupt
interrupt

ExpRace

Attacker

**Bullets**

- Inter-processor interrupt
- Hardware Interrupt
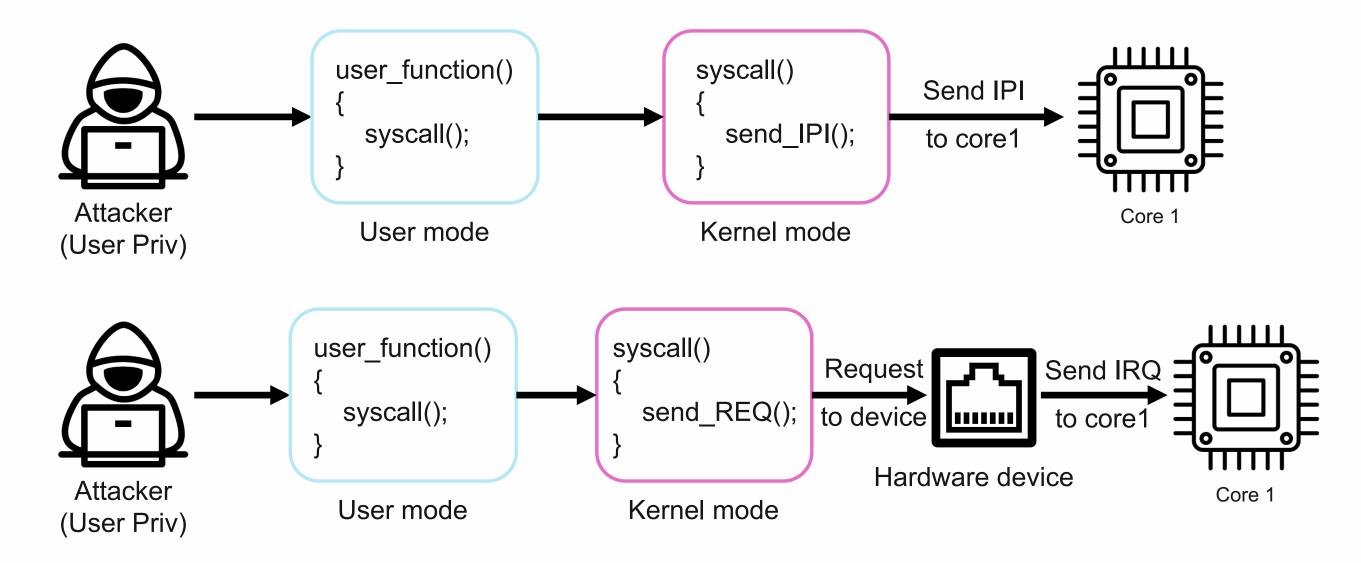
- The key idea of ExpRace is **to keep raising interrupts** to indirectly alter kernel thread's interleaving.

# ExpRace : How to send IPI & IRQ with user priv

# ExpRace : TLB Shootdown



| ~ | ~ | ~ |
|---|---|---|
| 0xABC0 | | 0xABE0 |

**cache**

**munmap(0xABD0)**

**Core 1**

**IPI**
Flush 0xABD0

| ~ | ~ | ~ |
|---|---|---|
| 0xABC0 | | 0xABE0 |

**cache**

**IPI_handler()**

**Core 2**

- Modern OSs implement a TLB shootdown mechanism to ensure that TLB entries are synchronized across different cores.

- Syscalls that either modify the permission of the page (e.g., mprotect()) or unmap (e.g., munmap()) the page use IPI for TLB shootdown.

# ExpRace : IPI Environment setting



If 3 processes have **same mm**

If process A and C have **same mm,** and B have **different mm**

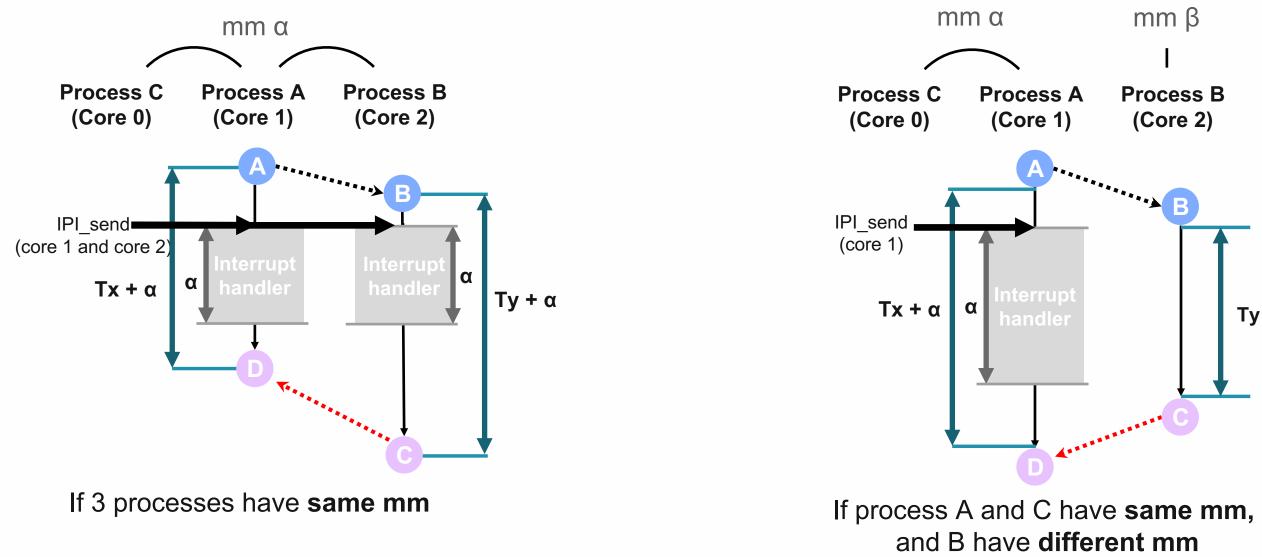# ExpRace : Hardware Interrupt Environment Setting

1. Check **irq's core affinity**.

   (In our environment, ethernet device (IRQ 122) have affinity to core 11)

   ```
   yoochan@compsec:~$ cat /proc/irq/122/smp_affinity_list
   11
   ```

2. Pin the thread to **corresponding core** using sched_setaffinity().

# ExpRace : How many cycles are extended?

Core 1          Core 2

A
B

1,500 ~ 20,000 cycles

IPI
handler

Ty

C

D

**TLB Shootdown**

Core 1          Core 2

A
B

About
15,000 cycles

ISR
handler

Ty

C

D

**Hardware Interrupt**

# ExpRace : Advanced Technique



- IPI and IRQ can be used simultaneously.

- The time window is extended up to 200,000 cycles

# Case Study : CVE-2017-15265

```
snd_seq_create_port()
{

   …
   port = kzalloc();
   list_add_tail(&port->list, &p->list);  ············· Ⓐ


   …                              Tx = 110 cycles


   strlcpy(port->name, info->name,  ············· Ⓓ
                    sizeof(port->name));

}
```

```
snd_seq_delete_port()
{
   list_for_each_entry( … p->list)
   {
      if (p->addr.port == port) {
         found = p;  ····················· Ⓑ
            …
      }
   }
   …                              Ty = 450 cycles
   kfree(found);  ····················· Ⓒ
}
```

**Problems to exploit**

1. Non-inclusive Multi-variable Race

2. No time to reallocate

if Ⓐ **>>** Ⓑ **&&** Ⓒ **>>** Ⓓ , then **Use-After-Free Write occurs**.

# ExpRace can solve two problems at once

```
snd_seq_create_port()
{
    …
    port = kzalloc();
    list_add_tail(&port->list, &p->list);  ········· (A)
```

**Interrupt Handler**

**Tx' = 110 + 15000 cycles**

```
    strlcpy(port->name, info->name,  ········· (D)
            sizeof(port->name));
}
```

```
snd_seq_delete_port()
{
    list_for_each_entry( … p->list)
    {
        if (p->addr.port == port) {
            found = p;  ··························· (B)
            …
        }
    }
    …
    kfree(found);  ························· (C)
}
```

**Ty = 450 cycles**

```
syscall_for_reallocte()
{
    kmalloc();  ···························· ●
}
```

**It takes about 3000 cycles**

if (A) >> (B) && (C) >> (D) , then **Use-After-Free Write occurs**.

# Brief introduction about memory corruption exploit

- Spray struct file pointer using SCM_RIGHT

- Partially overwrite the pointer in reallocated structure for kernel address leak.

- Use iovec structure for arbitrary memory write and read.

1st Use-After-Free Write

Leak : **struct file pointer**

2nd Use-After-Free Write

AAR : **file->f_cred pointer**

3rd Use-After-Free Write

AAW : **f_cred -> uid = 0**

We totally trigger the vulnerability **3 times**

# DEMO

## Conclusion

- Introduced **unexploitable** race types.

- ExpRace can turn **unexploitable** races into **exploitable** races.