

# バイナリ書き換えを用いた システムコールフック機構

第 159 回 システムソフトウェアとオペレーティング・システム研究会  
2023 年 5 月 16 日

安形 憲一<sup>1</sup>、田崎 創<sup>1</sup>、オブラン ピエールルイ<sup>1</sup>、石黒 健太<sup>2</sup>

<sup>1</sup> IJ 技術研究所

<sup>2</sup> 法政大学

# About this Work

- zpoline is a system call hook mechanism for x86-64 CPUs
- We made this for transparently applying user-space OS subsystems to existing applications
- The source code of zpoline has been publicly available at <https://github.com/yasukata/zpoline> since October 2021

USENIX ATC 2023 で発表予定の内容です

<https://www.usenix.org/conference/atc23/presentation/yasukata>

# System Call

- System calls are the primary interface for user-space programs to communicate with OS kernels

# System Call

- System calls are the primary interface for user-space programs to communicate with OS kernels

User-space program

# System Call

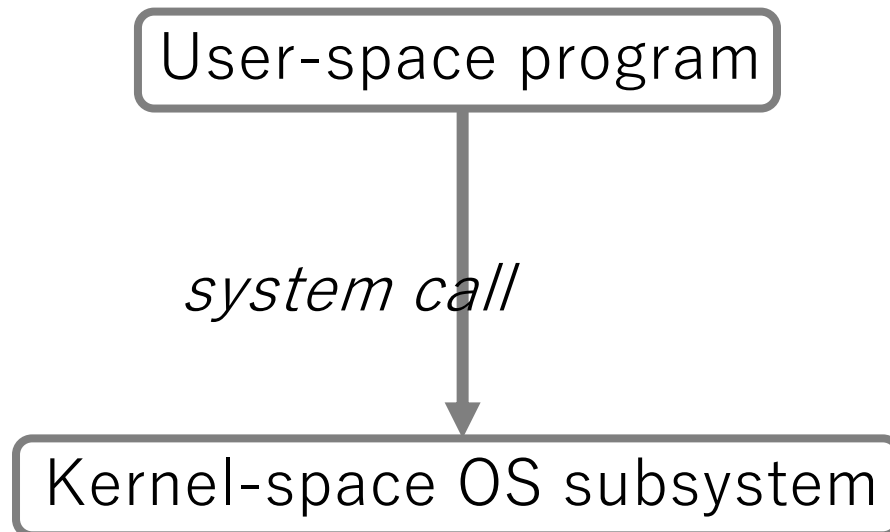
- System calls are the primary interface for user-space programs to communicate with OS kernels

User-space program

Kernel-space OS subsystem

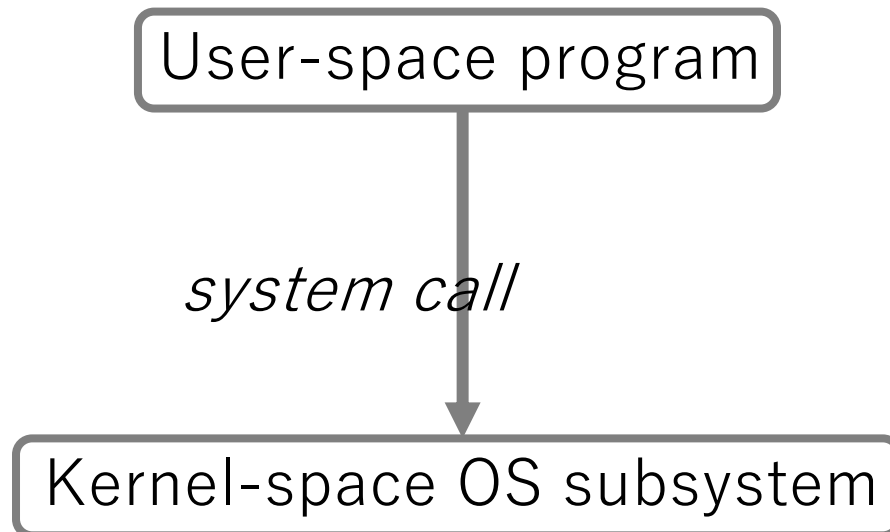
# System Call

- System calls are the primary interface for user-space programs to communicate with OS kernels



# System Call Hook

- System calls are the primary interface for user-space programs to communicate with OS kernels
- A system call hook mechanism intercepts a system call



# System Call Hook

- System calls are the primary interface for user-space programs to communicate with OS kernels
- A system call hook mechanism intercepts a system call

User-space program

intercept

*system call hook*

Kernel-space OS subsystem



# System Call Hook

- System calls are the primary interface for user-space programs to communicate with OS kernels
- A system call hook mechanism intercepts a system call, and redirects the execution to a user-defined hook function

User-space program

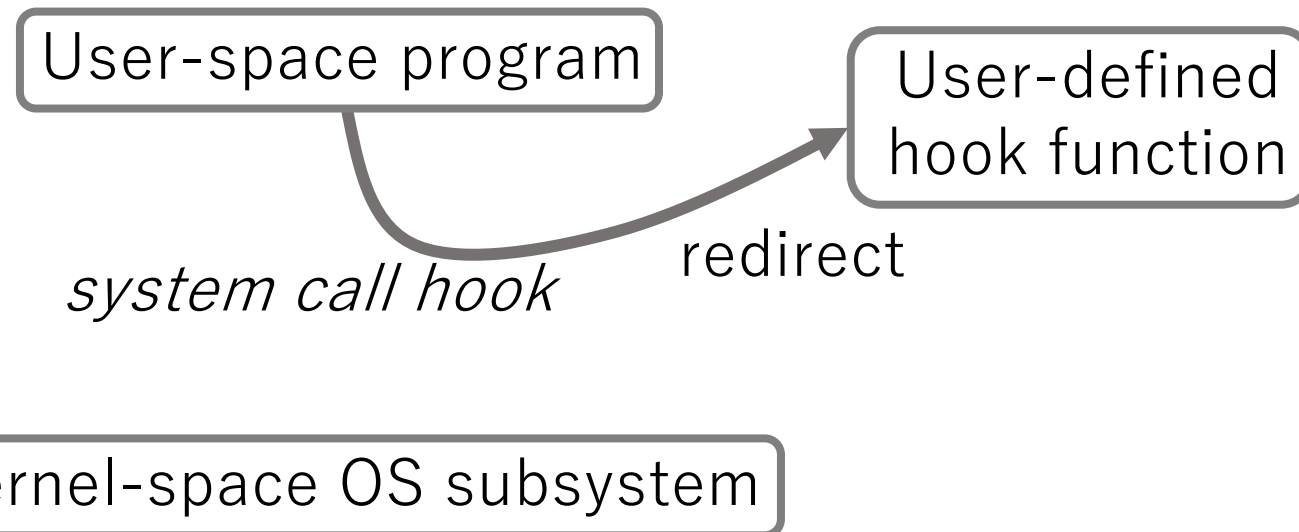
intercept

*system call hook*

Kernel-space OS subsystem

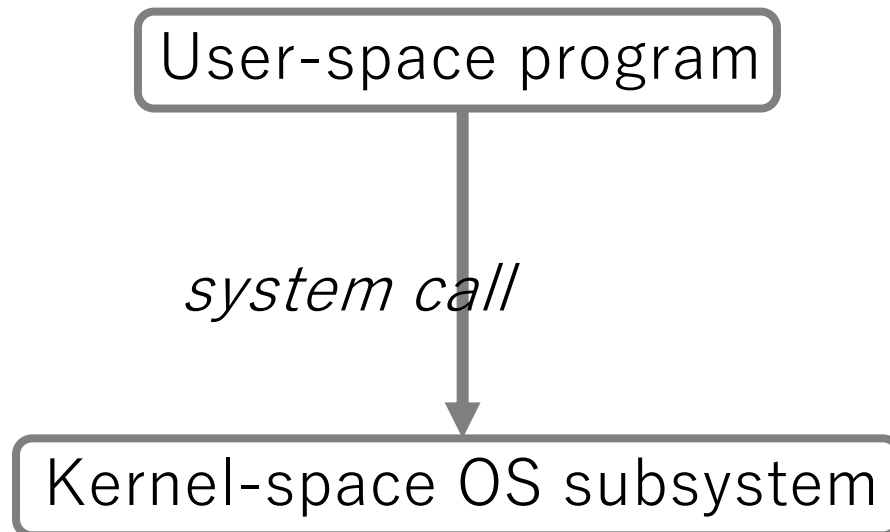
# System Call Hook

- System calls are the primary interface for user-space programs to communicate with OS kernels
- A system call hook mechanism intercepts a system call, and redirects the execution to a user-defined hook function



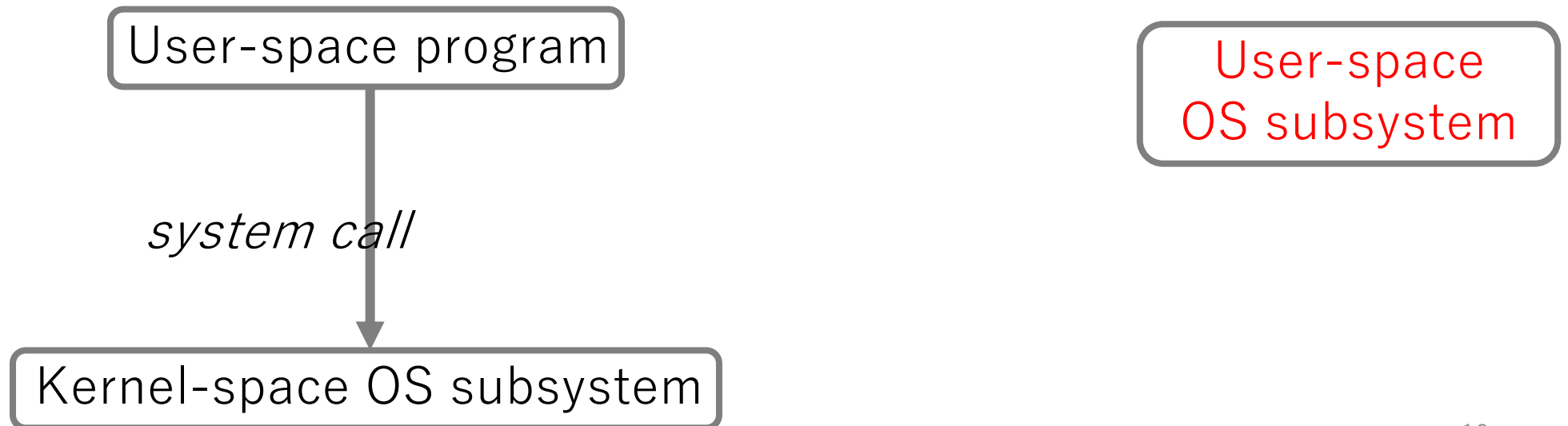
# Motivating Use Case

- System call hook mechanisms allow us to transparently apply user-space OS subsystems to existing applications



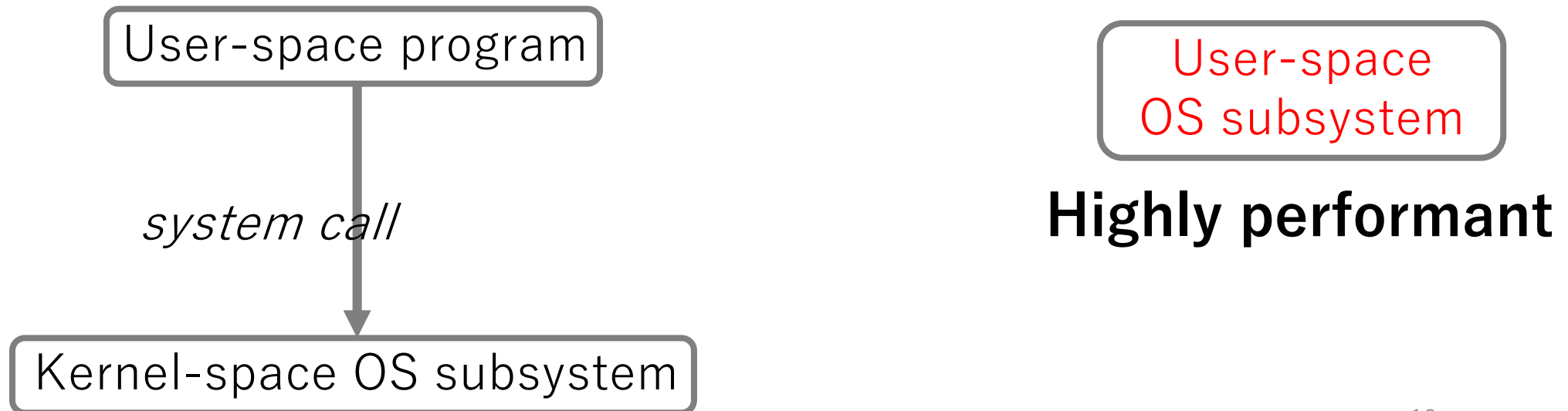
# Motivating Use Case

- System call hook mechanisms allow us to transparently apply user-space OS subsystems to existing applications



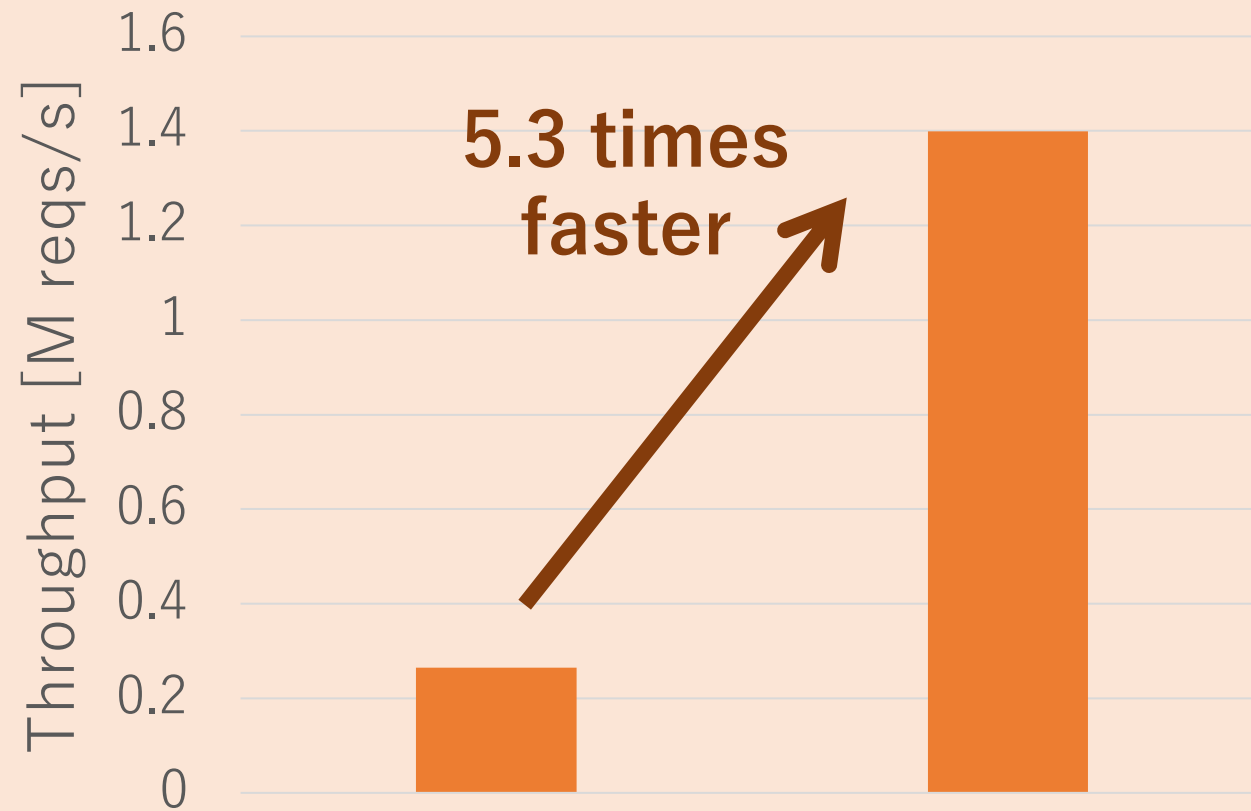
# Motivating Use Case

- System call hook mechanisms allow us to transparently apply user-space OS subsystems to existing applications



# Motivating Use Case

TCP ping-pong performance



Linux TCP stack lwIP on DPDK = user-space network stack

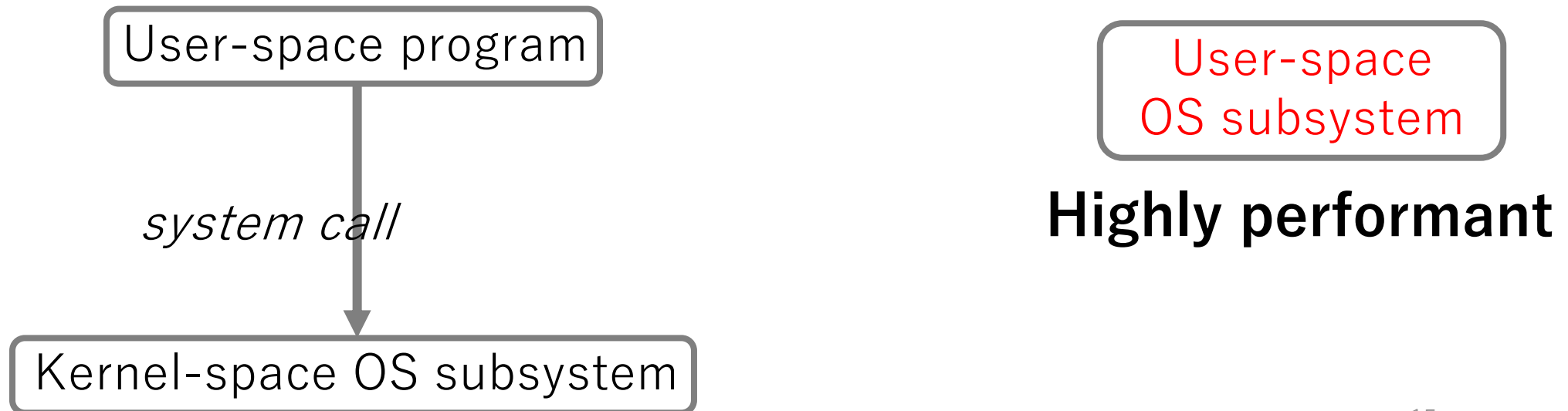
to transparently apply applications

User-space  
OS subsystem

**Highly performant**

# Motivating Use Case

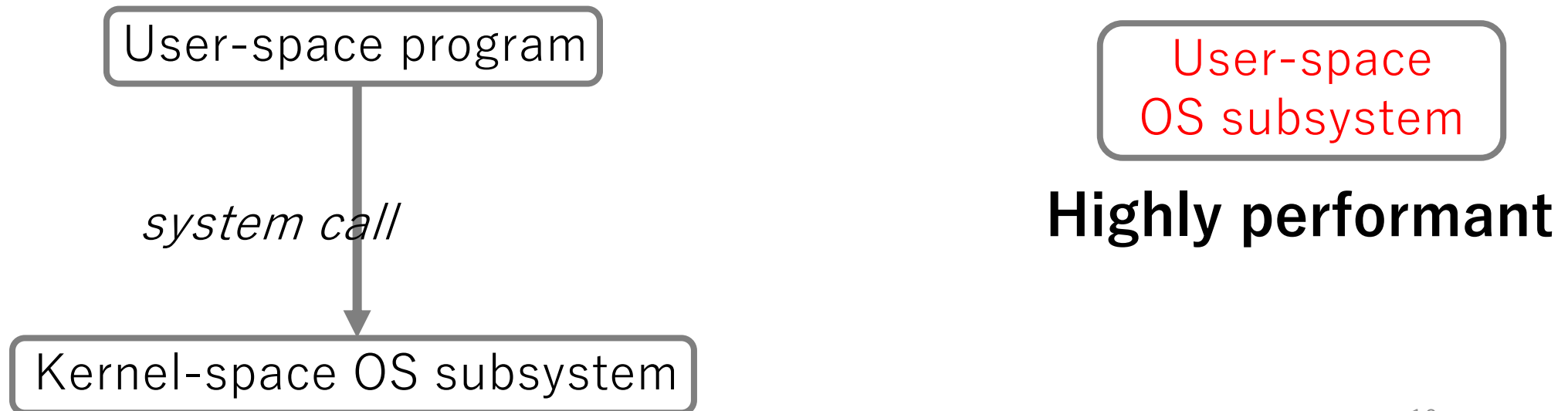
- System call hook mechanisms allow us to transparently apply user-space OS subsystems to existing applications



# Motivating Use Case

- System call hook mechanisms allow us to transparently apply user-space OS subsystems to existing applications

**Normally, adaptation requires changes of a user-space program to apply a specific API of a user-space OS subsystem**





# Motivating Use Case

- System call hook mechanisms allow us to transparently apply user-space OS subsystems to existing applications

**Normally, adaptation requires changes of a user-space program to apply a specific API of a user-space OS subsystem**



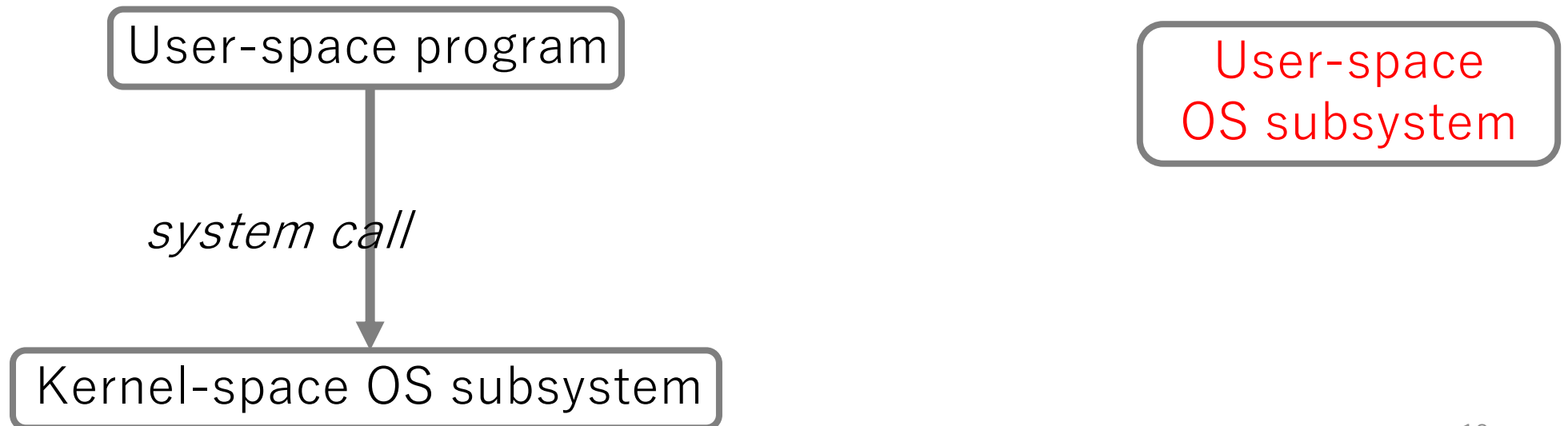
**Highly performant**

Kernel-space OS subsystem

# Motivating Use Case

- System call hook mechanisms allow us to transparently apply user-space OS subsystems to existing applications

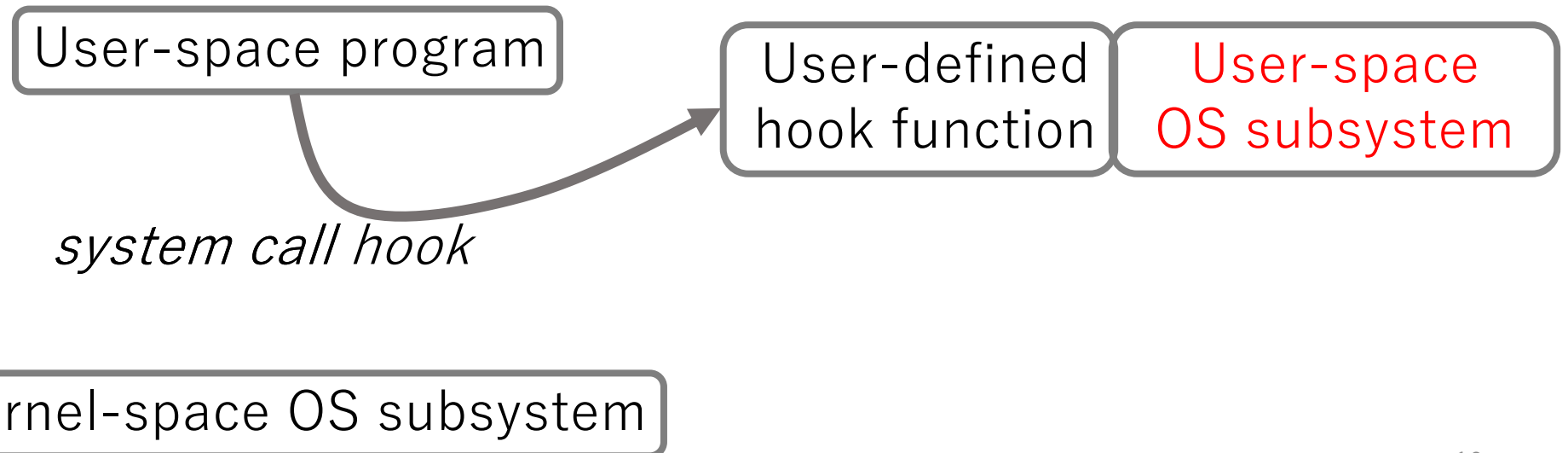
**If we use a system call hook mechanism, ...**



# Motivating Use Case

- System call hook mechanisms allow us to transparently apply user-space OS subsystems to existing applications

**If we use a system call hook mechanism, ...**

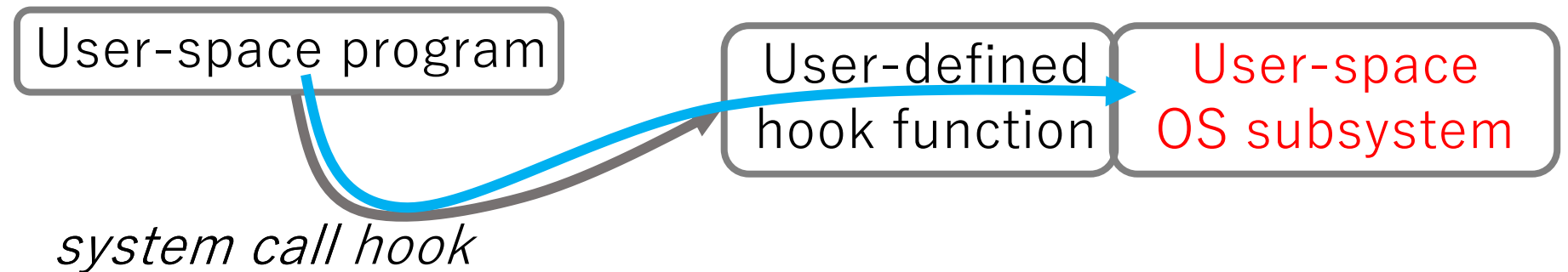


# Motivating Use Case

- System call hook mechanisms allow us to transparently apply user-space OS subsystems to existing applications

**If we use a system call hook mechanism, ...**

**no modification of the user-space program is necessary**



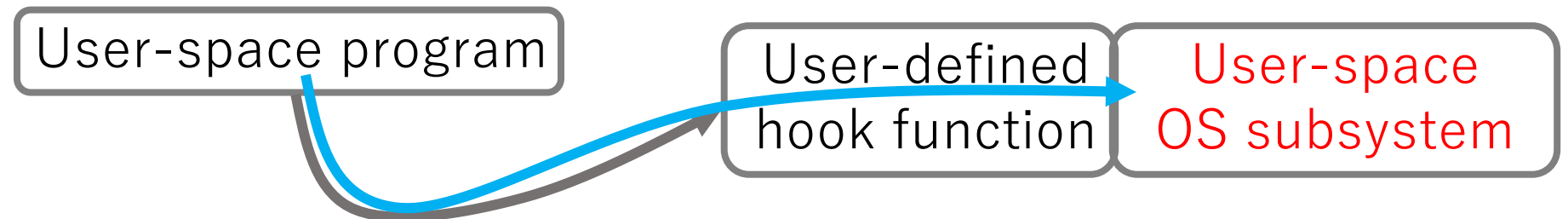
Kernel-space OS subsystem

# Problem

- System call hook mechanisms allow us to transparently apply user-space OS subsystems to existing applications

**If we use a system call hook mechanism, ...**

**no modification of the user-space program is necessary**



*system call hook* ← **There is no perfect hook mechanism**

Kernel-space OS subsystem

# Problem

## Existing Mechanisms

- ptrace
- Syscall User Dispatch (SUD)
- int3 signaling technique
- LD\_PRELOAD trick
- Binary rewriting techniques
- ...

These mechanisms allow us to transparently apply hooks to existing applications

**Hook mechanism, ...**

**user-space program is necessary**

*system call hook*

User-defined  
hook function

User-space  
OS subsystem

Kernel-space OS subsystem

**There is no perfect hook mechanism**

# Problem

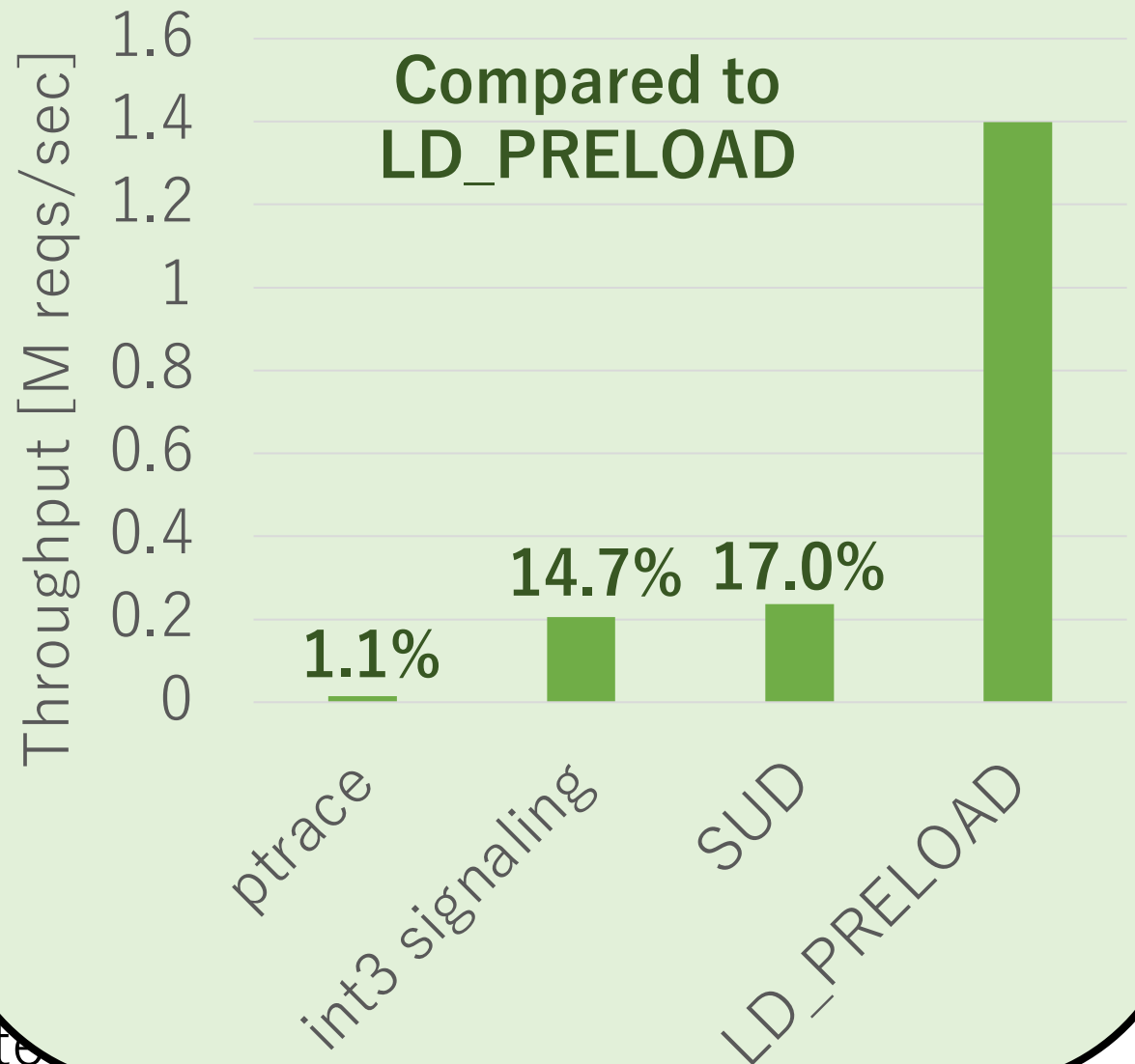
## Existing Mechanisms

- ptrace
- Syscall User Dispatch (SUD)
- int3 signaling technique
- LD\_PRELOAD trick
- Binary rewriting techniques
- ...

*system call hook*

Kernel-space OS subsystem

lwIP on DPDK : TCP ping-pong



# Problem

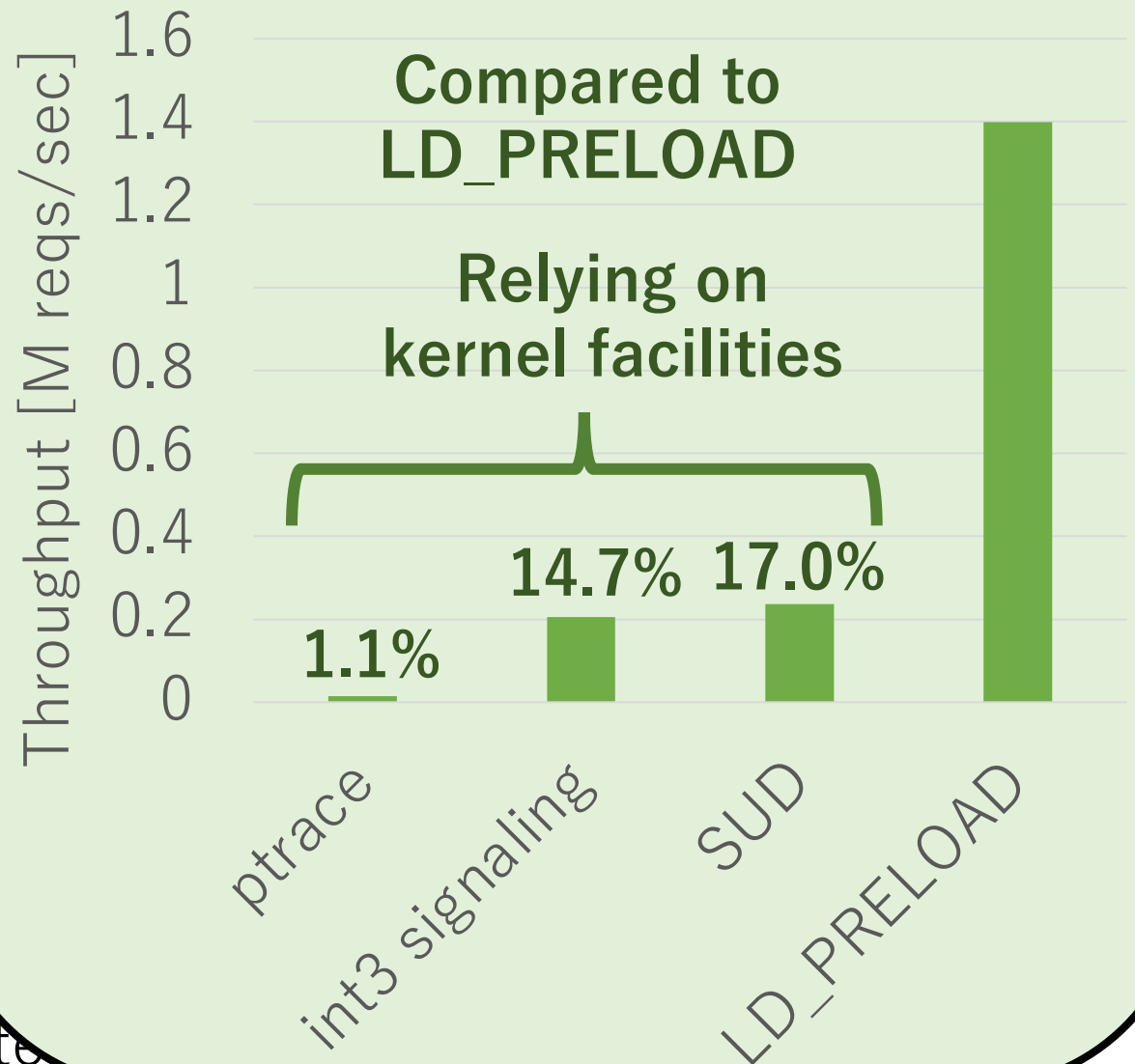
## Existing Mechanisms

- ptrace
- Syscall User Dispatch (SUD)
- int3 signaling technique
- LD\_PRELOAD trick
- Binary rewriting techniques
- ...

*system call hook*

Kernel-space OS subsystem

lwIP on DPDK : TCP ping-pong





Prob

```
app_function(...)  
{
```

```
  ...  
  write(...)
```



libc  
write()  
library call

Ex

- ptra
- Sysc
- int3
- LD\_
- Bina
- ...

WORK : TCP ping-pong

Compared to  
PRELOAD

Relying on  
kernel facilities

14.7% 17.0%

SUD  
LD\_PRELOAD

Problem

Example

- ptrace
- Syscall
- int3
- LD\_PRELOAD
- Binary
- ...

```
app_function(...)  
{  
    ...  
    write(...) →  
    ...  
}
```

libc  
write()  
library call

User-defined  
write()  
library call

**LD\_PRELOAD**

Work : TCP ping-pong

Compared to  
**PRELOAD**

Relying on  
kernel facilities

14.7% 17.0%

SUD  
LD\_PRELOAD

Problem

Example

- ptrace
- Syscall
- int3
- LD\_PRELOAD
- Binary
- ...

```
app_function(...)  
{  
    ...  
    write(...)  
    ...  
}
```

function call  
replacement

libc  
write()  
library call

User-defined  
write()  
library call

**LD\_PRELOAD**

Work : TCP ping-pong

Compared to  
**PRELOAD**

relying on  
kernel facilities

14.7% 17.0%

SUD  
LD\_PRELOAD

Prob

```
app_function(...)  
{
```

...

```
    special_write(...)
```

...

```
}
```

libc  
write()  
library call

User-defined  
write()  
library call

**LD\_PRELOAD**

WORK : TCP ping-pong

Compared to  
**PRELOAD**

relying on  
kernel facilities

14.7% 17.0%

SUD  
LD\_PRELOAD

- Ex
- ptrace
  - Syscall
  - int3
  - LD\_PRELOAD
  - Binary instrumentation
  - ...

Problem

Example

- ptrace
- Syscall
- int3
- LD\_PRELOAD
- Binary
- ...

```
app_function(...)  
{
```

...

```
    special_write(...)
```

...

```
}
```

```
special_write(...)  
{
```

```
    asm volatile (  
        "trigger  
        write syscall  
    ");
```

```
}
```

libc  
write()  
library call

User-defined  
write()  
library call

**LD\_PRELOAD**

Work : TCP ping-pong

Compared to  
**LD\_PRELOAD**

Relying on  
kernel facilities

14.7% 17.0%

SUD  
LD\_PRELOAD

Problem

Example

- ptrace
- Syscall
- int3
- LD\_PRELOAD
- Binary
- ...

```
app_function(...)  
{
```

...

```
    special_write(...)  
    ...
```

```
}
```



```
special_write(...)  
{
```

```
    asm volatile (  
        "trigger  
        write syscall  
    ");
```

```
}
```

libc  
write()  
library call

User-defined  
write()  
library call

**LD\_PRELOAD**

Work : TCP ping-pong

Compared to  
LD\_PRELOAD

Relaying on  
kernel facilities

14.7% 17.0%

SUD  
LD\_PRELOAD

Problem

Example

- ptrace
- Syscall
- int3
- LD\_PRELOAD
- Binary
- ...

```
app_function(...)  
{
```

...

```
    special_write(...)
```

...

```
}
```



```
special_write(...)  
{
```

```
    asm volatile (  
        trigger
```

```
        write syscall  
    )
```

```
}
```

libc  
write()  
library call

User-defined  
write()  
library call

**LD\_PRELOAD**

Work : TCP ping-pong

Compared to  
LD\_PRELOAD

Relying on  
kernel facilities

14.7% 17.0%

SUD  
LD\_PRELOAD

Problem

Example

- ptrace
- Syscall
- int3
- LD\_PRELOAD
- Binary
- ...

```
app_function(...)  
{
```

...

```
    special_write(...)
```

...

```
}
```

libc  
write()  
library call

```
special_write(...)  
{
```

```
    asm volatile (
```

```
        trigger  
        write syscall
```

```
    )
```

User-defined  
write()  
library call

```
} Hook is not applied LD_PRELOAD
```

Work : TCP ping-pong

Compared to  
LD\_PRELOAD

Relying on  
kernel facilities

14.7% 17.0%

SUD  
LD\_PRELOAD



Problem

Example

- ptrace
- Syscall
- int3
- LD\_PRELOAD
- Binary
- ...

```
app_function(...)  
{
```

```
...  
    special_write(...)  
...  
}
```

libc  
write()  
library call

```
special_write(...)  
{
```

because names  
are different

```
    asm volatile (  
        trigger  
        write syscall  
    )  
}
```

User-defined  
write()  
library call

**Hook is not applied** LD\_PRELOAD

Work : TCP ping-pong

Compared to  
LD\_PRELOAD

Relying on  
kernel facilities

14.7% 17.0%

SUD  
LD\_PRELOAD

Prob

We see this case in glibc often

```
app_function(...)
```

```
{
```

```
...
```

```
special_write(...)
```

```
...
```

```
}
```

libc  
write()  
library call

```
special_write(...)
```

```
{
```

```
asm volatile (
```

```
trigger  
write syscall
```

```
)
```

because names  
are different

User-defined  
write()  
library call

```
} Hook is not applied LD_PRELOAD
```

WORK : TCP ping-pong

Compared to  
LD\_PRELOAD

relying on  
kernel facilities

14.7% 17.0%

SUD  
LD\_PRELOAD

- Ex
- ptrace
  - Syscall
  - int3
  - LD\_PRELOAD
  - Binary
  - ...

# Problem

## Existing Mechanisms

- ptrace
- Syscall User Dispatch (SUD)
- int3 signaling technique
- LD\_PRELOAD trick
- Binary rewriting techniques
- ...

These mechanisms allow us to transparently apply hooks to existing applications

**hook mechanism, ...**

**user-space program is necessary**

*system call hook*

User-defined  
hook function

User-space  
OS subsystem

Kernel-space OS subsystem

**There is no perfect hook mechanism**

# Problem

## Existing Mechanisms

- System call intercept (SUD)
- int3 signaling technique
- LD\_PRELOAD trick
- Binary rewriting techniques
- ...

**High performance penalty**

These mechanisms allow us to transparently apply hooks to existing applications

**Hook mechanism, ...**

**user-space program is necessary**

*system call hook*

User-defined  
hook function

User-space  
OS subsystem

Kernel-space OS subsystem

**There is no perfect hook mechanism**

# Problem

## Existing Mechanisms

- System call intercept (SUD)
- int3 signaling technique
- FPU load trick
- Binary rewriting techniques
- ...

**High performance penalty**

**Sometimes fail to hook**

... allow us to transparently apply to existing applications

**hook mechanism, ...**

**user-space program is necessary**

*system call hook*

User-defined  
hook function

User-space  
OS subsystem

Kernel-space OS subsystem

**There is no perfect hook mechanism**

Problem ➡ **Applicability of user-space OS subsystems has been limited regardless of their benefits**

### Existing Mechanisms

- System call intercept (SUD)
- int3 signaling technique
- FPU load trick
- Binary rewriting techniques
- ...

**High performance penalty**

**Sometimes fail to hook**

These mechanisms allow us to transparently apply to existing applications

**Hook mechanism, ...**

**user-space program is necessary**

*system call hook*

User-defined  
hook function

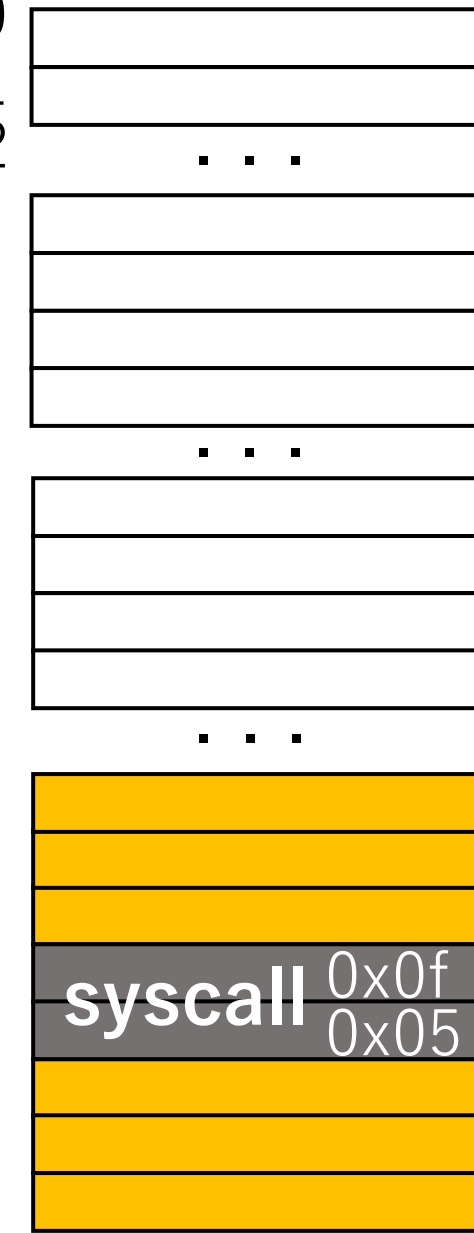
User-space  
OS subsystem

Kernel-space OS subsystem

**There is no perfect hook mechanism**

# Contribution

- zpoline employs binary rewriting and offers 6 advantages
  1. Low hook overhead
  2. Exhaustive hooking
  3. No breakage of user-space program logic
  4. No kernel change and no additional kernel module are necessary
  5. No source code of a user-space program is needed
  6. It can be used for emulating system calls
- None of existing mechanisms achieve them simultaneously

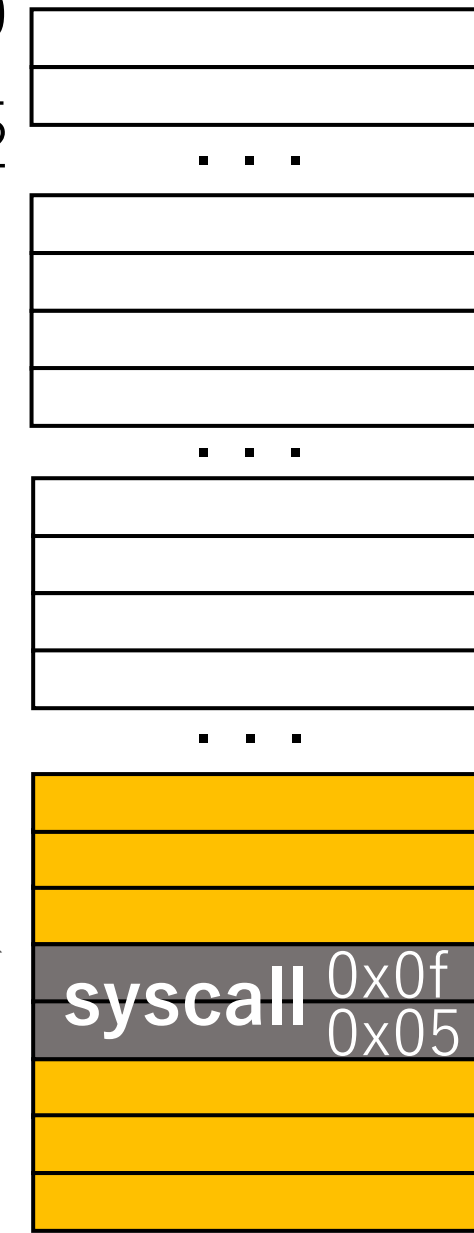


# Binary Rewriting Approach

- On x86-64 CPUs, syscall and sysenter instructions trigger a system call
  - syscall: 0x0f 0x05, sysenter: 0x0f 0x34

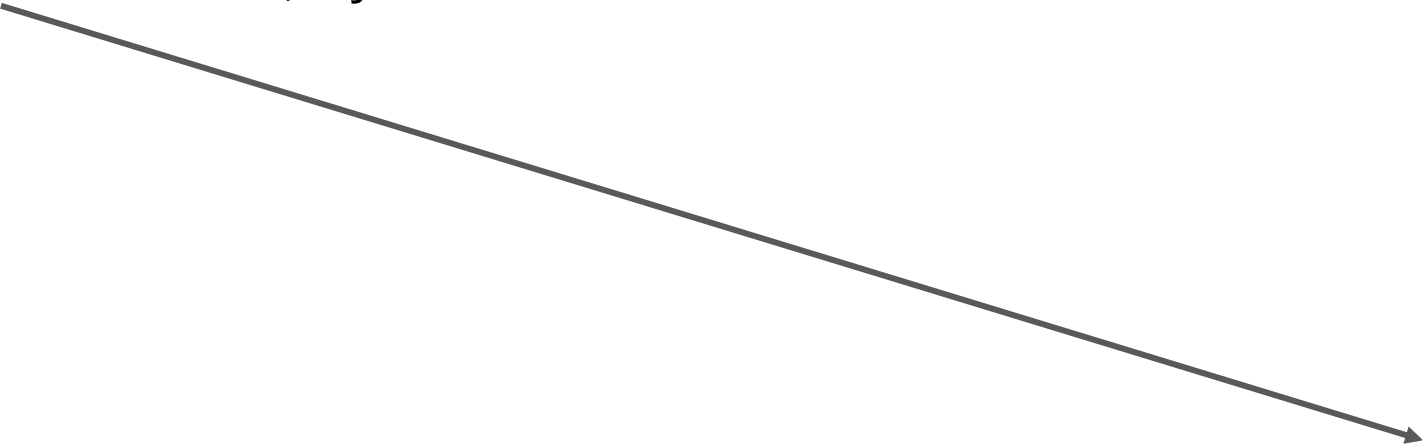
0x0000  
0x0001  
0x0002





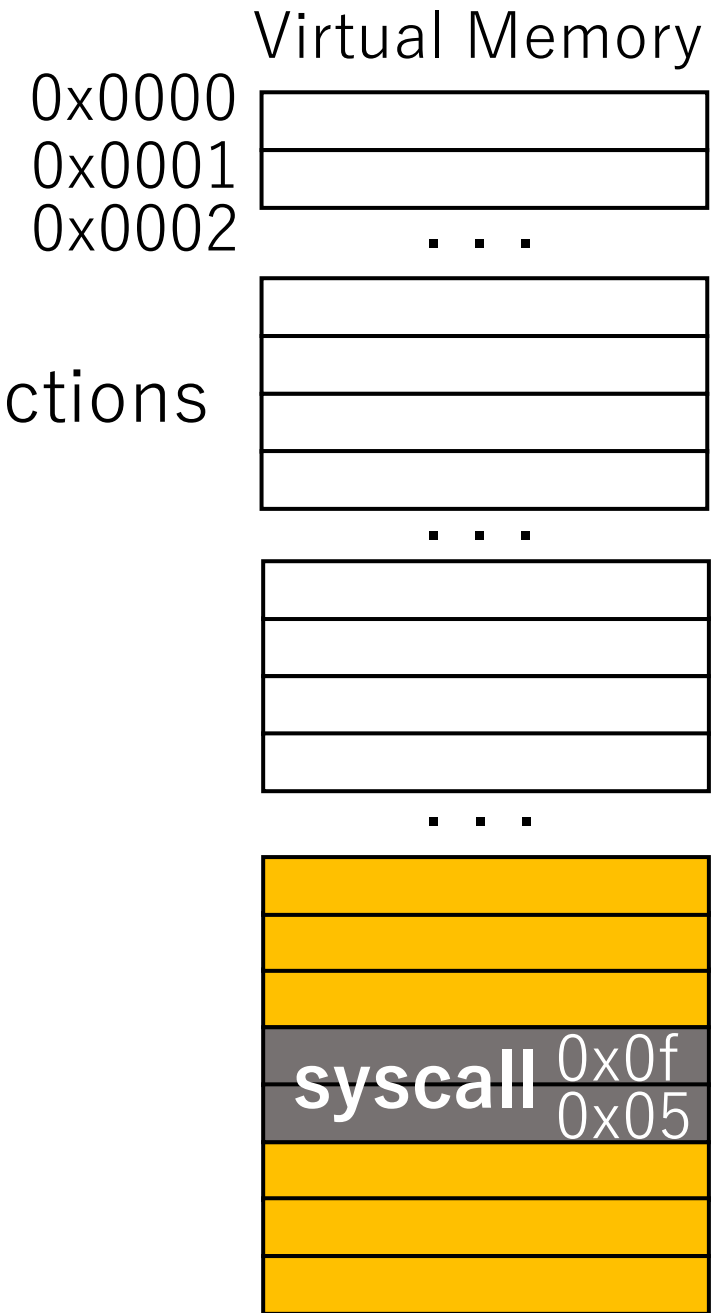
# Binary Rewriting Approach

- On x86-64 CPUs, syscall and sysenter instructions trigger a system call
  - syscall: 0x0f 0x05, sysenter: 0x0f 0x34



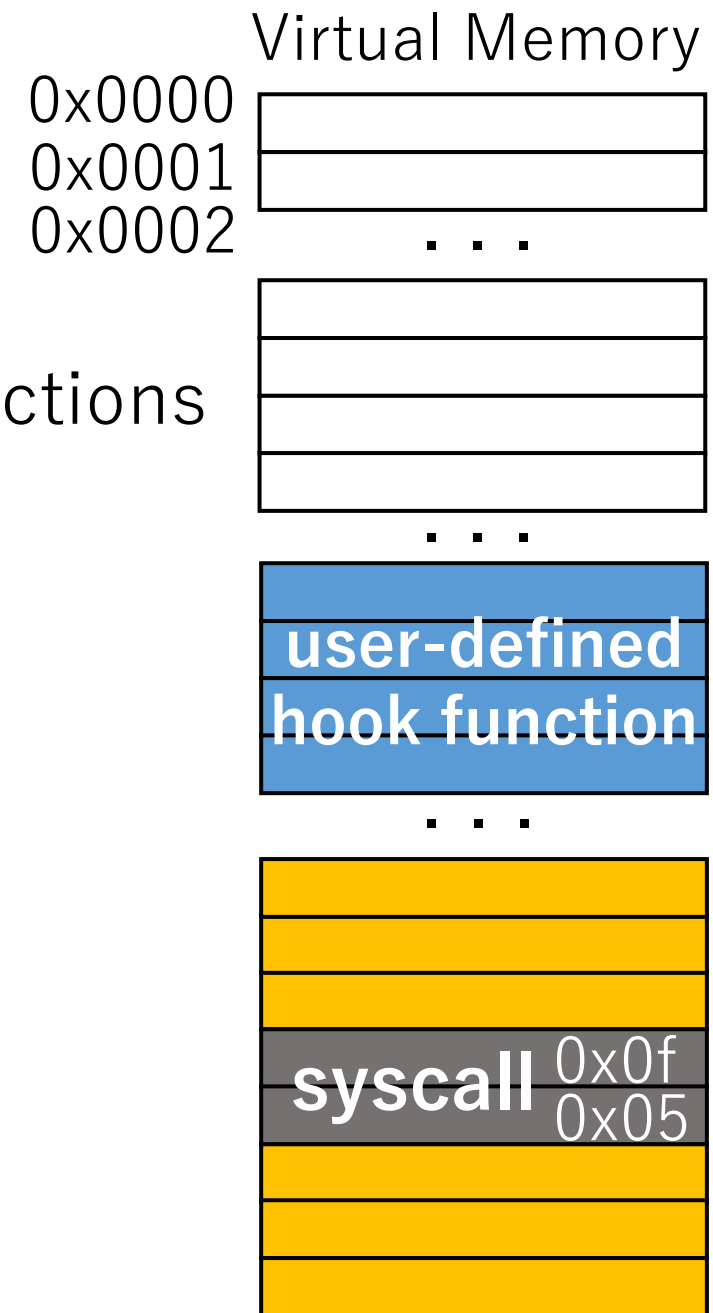
# Binary Rewriting Approach

- On x86-64 CPUs, syscall and sysenter instructions trigger a system call
  - syscall: 0x0f 0x05, sysenter: 0x0f 0x34
- What we wish to achieve



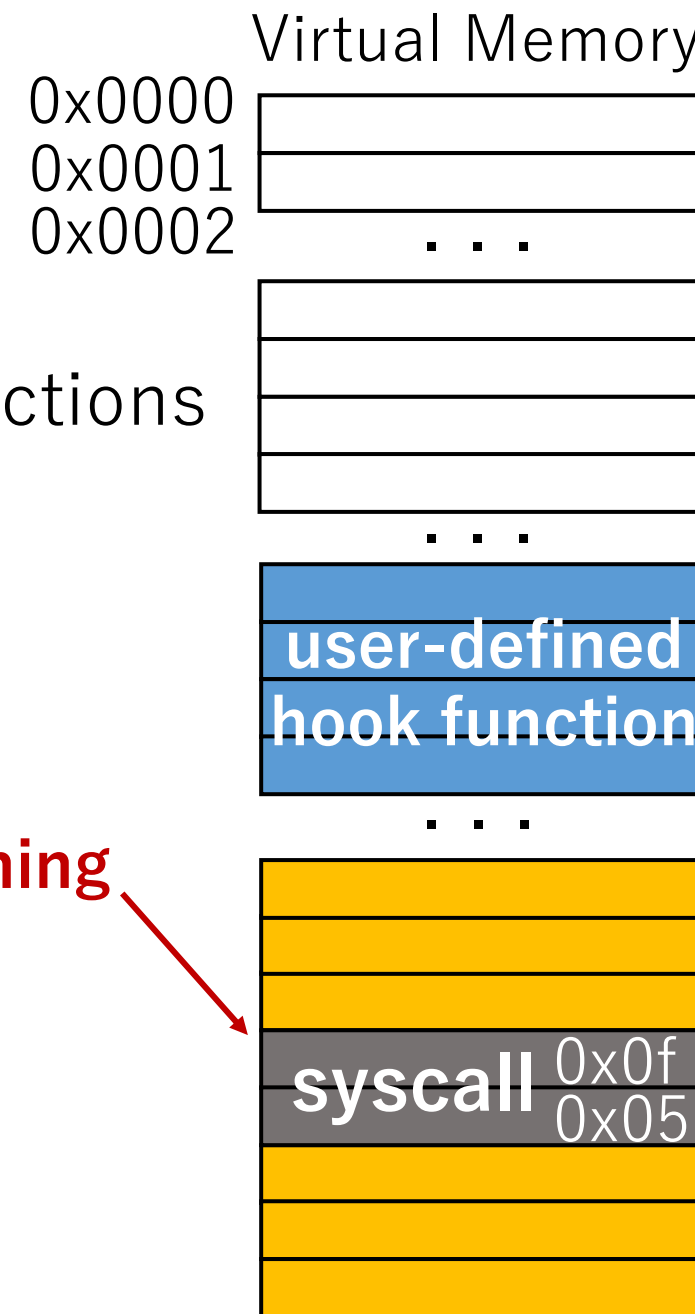
# Binary Rewriting Approach

- On x86-64 CPUs, syscall and sysenter instructions trigger a system call
  - syscall: 0x0f 0x05, sysenter: 0x0f 0x34
- What we wish to achieve



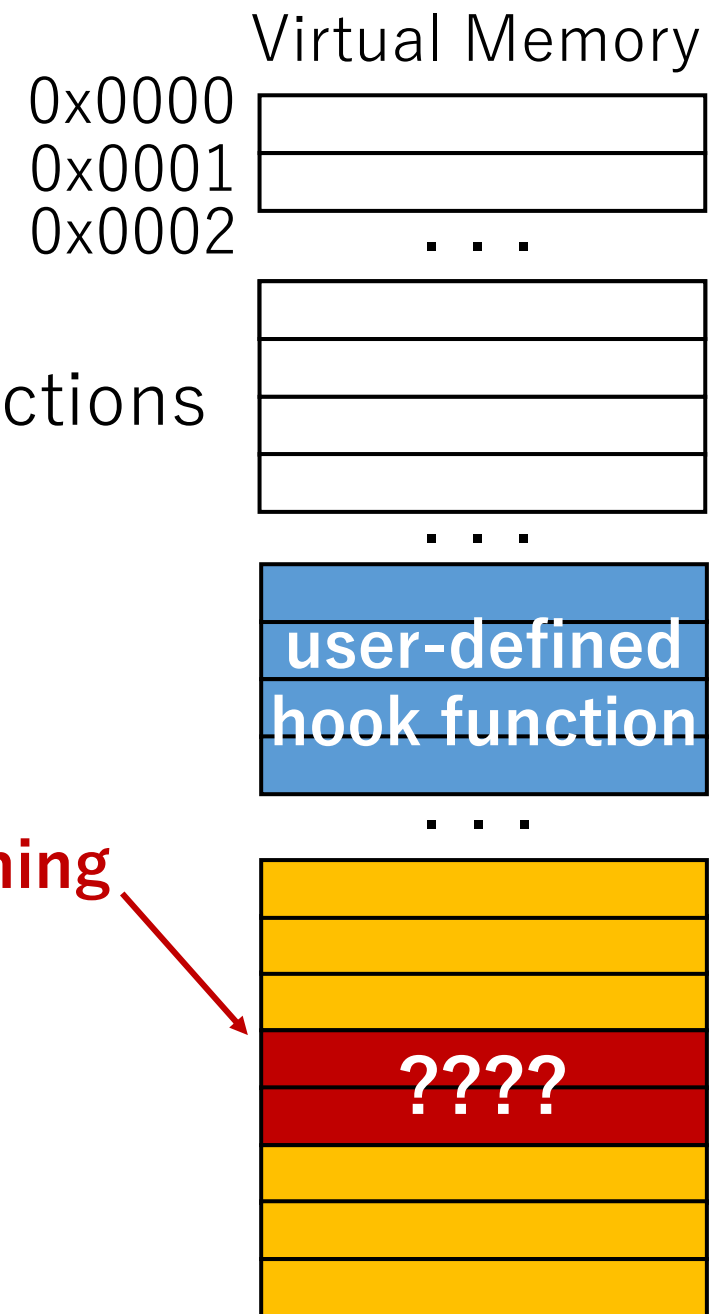
# Binary Rewriting Approach

- On x86-64 CPUs, syscall and sysenter instructions trigger a system call
  - syscall: 0x0f 0x05, sysenter: 0x0f 0x34
- What we wish to achieve
  - replace syscall/sysenter instruction with **something**



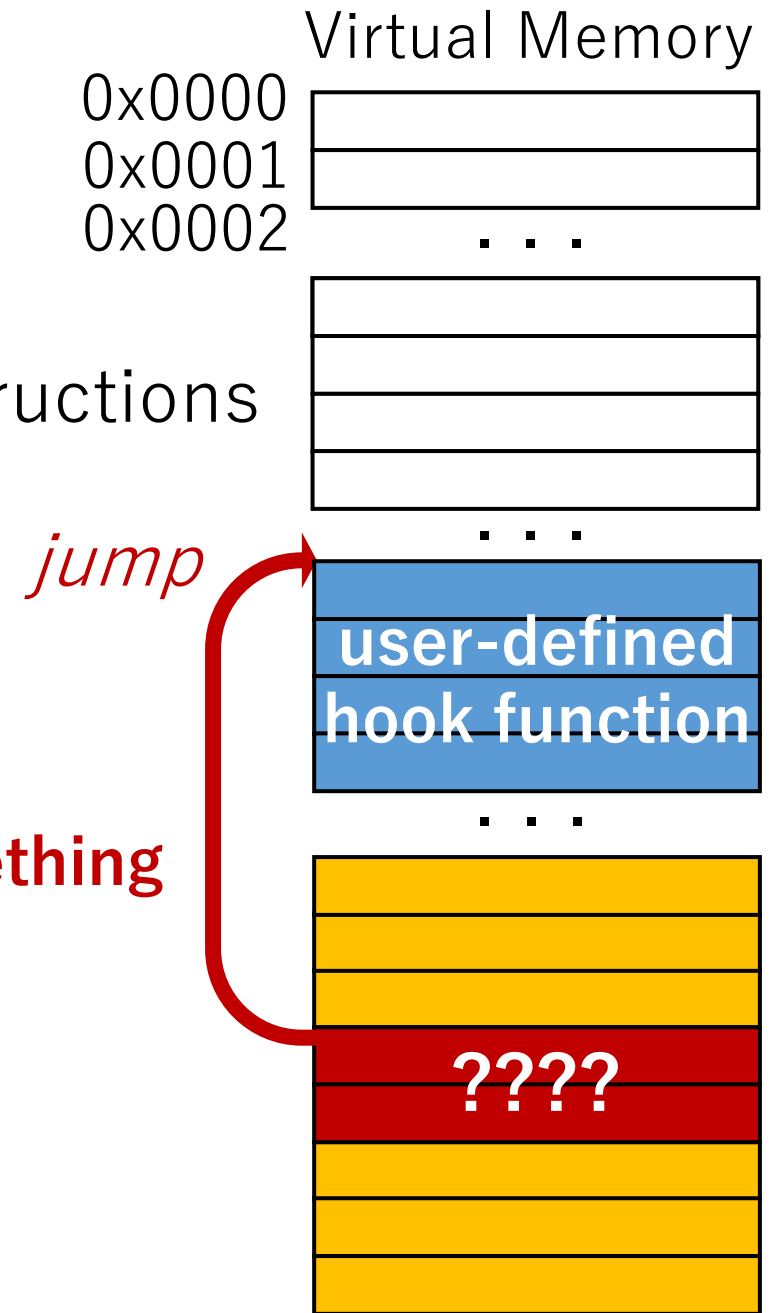
# Binary Rewriting Approach

- On x86-64 CPUs, syscall and sysenter instructions trigger a system call
  - syscall: 0x0f 0x05, sysenter: 0x0f 0x34
- What we wish to achieve
  - replace syscall/sysenter instruction with **something**



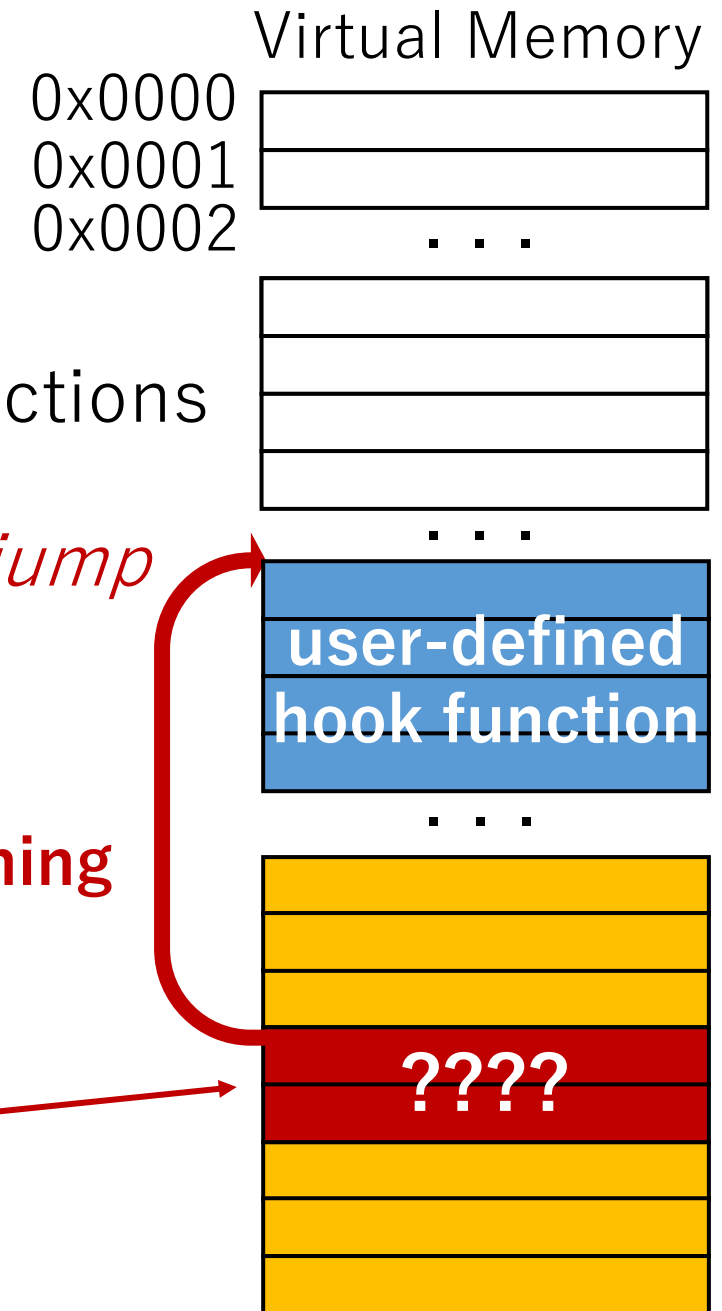
# Binary Rewriting Approach

- On x86-64 CPUs, syscall and sysenter instructions trigger a system call
  - syscall: 0x0f 0x05, sysenter: 0x0f 0x34
- What we wish to achieve
  - replace syscall/sysenter instruction with **something**
  - to **jump** to a user-defined hook function



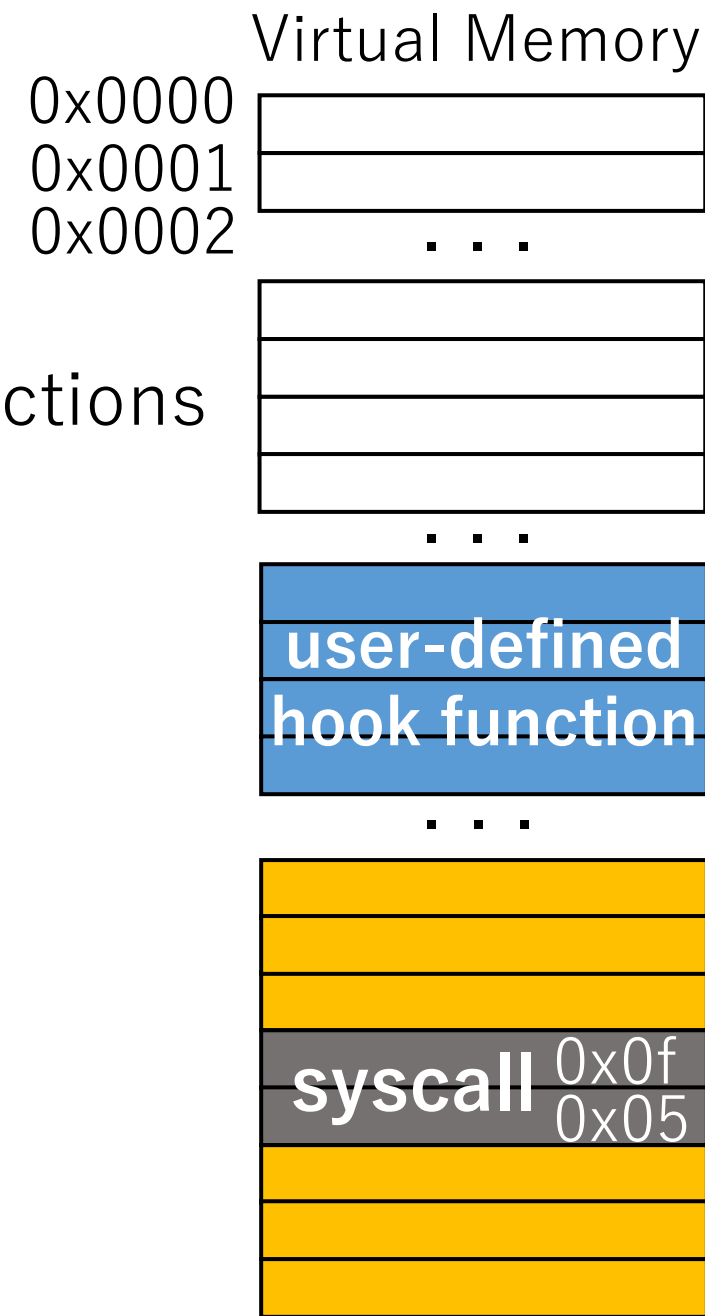
# Binary Rewriting Approach

- On x86-64 CPUs, syscall and sysenter instructions trigger a system call
  - syscall: 0x0f 0x05, sysenter: 0x0f 0x34
- What we wish to achieve
  - replace syscall/sysenter instruction with **something**
  - to **jump** to a user-defined hook function
- Question: what should we put here?



# Challenge

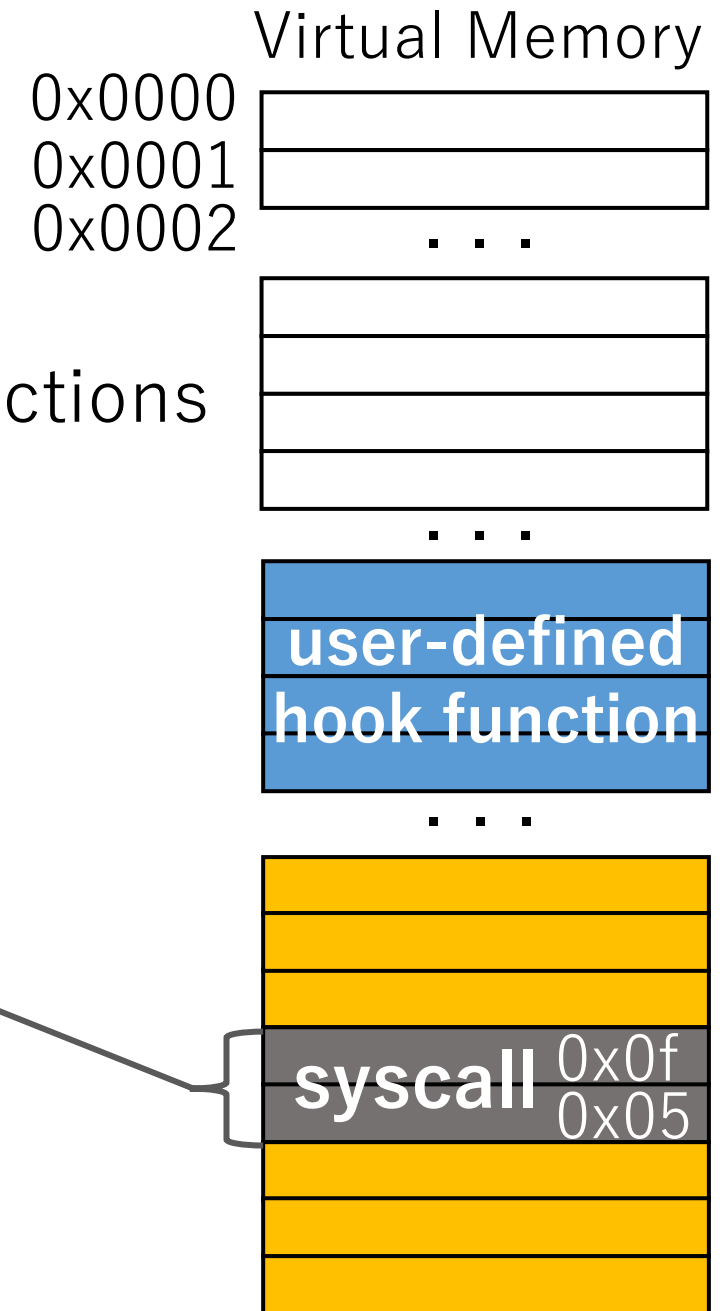
- On x86-64 CPUs, syscall and sysenter instructions trigger a system call
  - syscall: 0x0f 0x05, sysenter: 0x0f 0x34





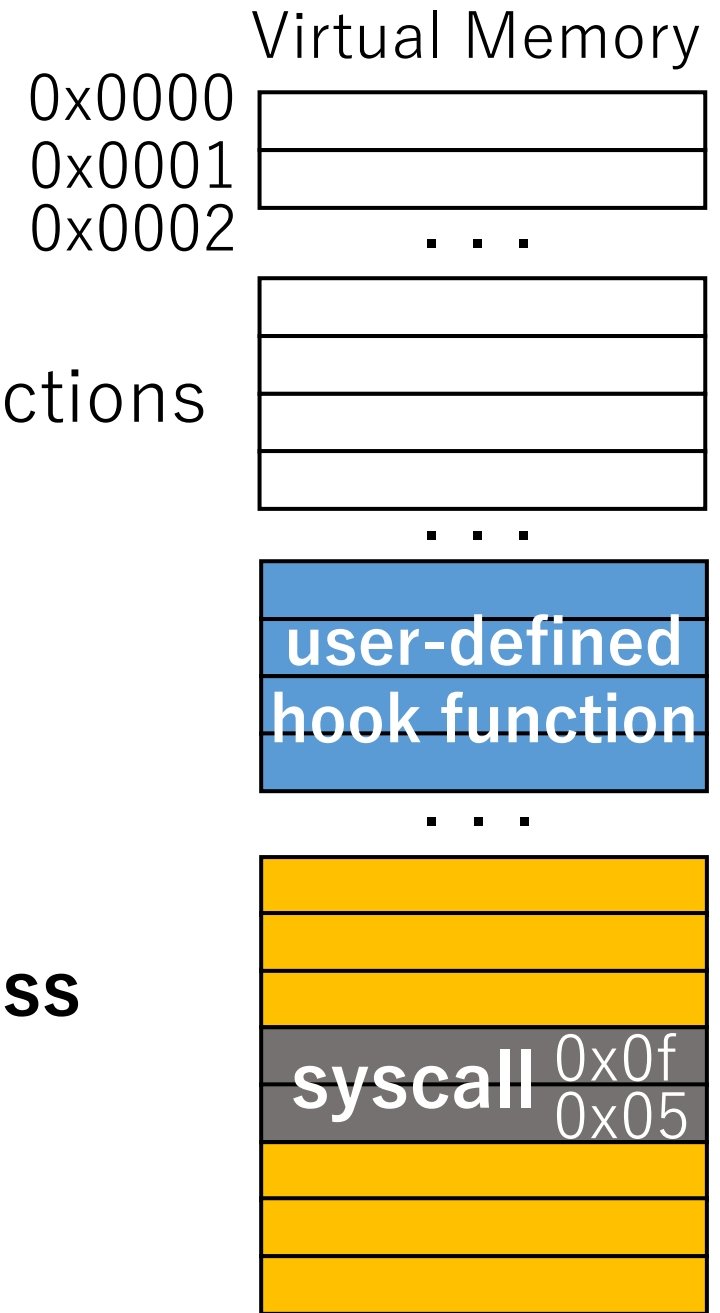
# Challenge

- On x86-64 CPUs, syscall and sysenter instructions trigger a system call
  - syscall: 0x0f 0x05, sysenter: 0x0f 0x34
- syscall and sysenter are **2-byte** instructions



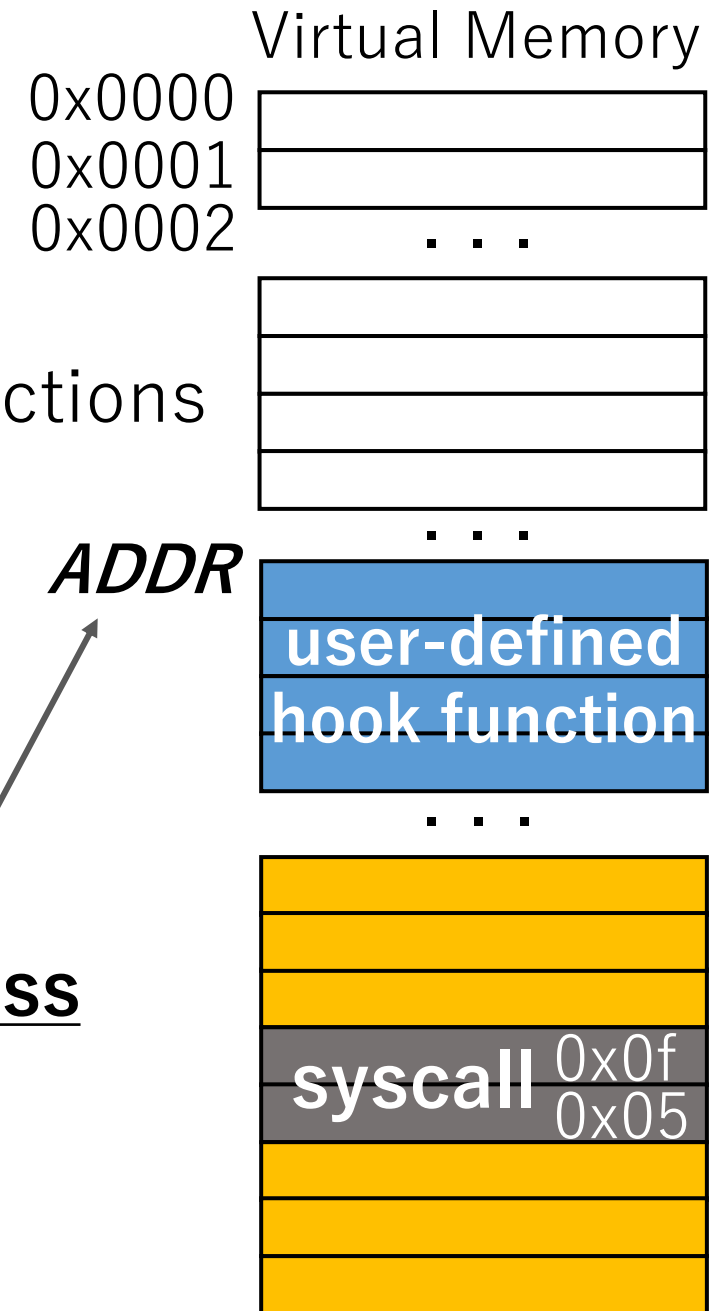
# Challenge

- On x86-64 CPUs, syscall and sysenter instructions trigger a system call
  - syscall: 0x0f 0x05, sysenter: 0x0f 0x34
- syscall and sysenter are **2-byte** instructions
- **Specification for a jump destination address needs more than 2 bytes**



# Challenge

- On x86-64 CPUs, syscall and sysenter instructions trigger a system call
  - syscall: 0x0f 0x05, sysenter: 0x0f 0x34
- syscall and sysenter are **2-byte** instructions
- **Specification for a jump destination address needs more than 2 bytes**



# Challenge

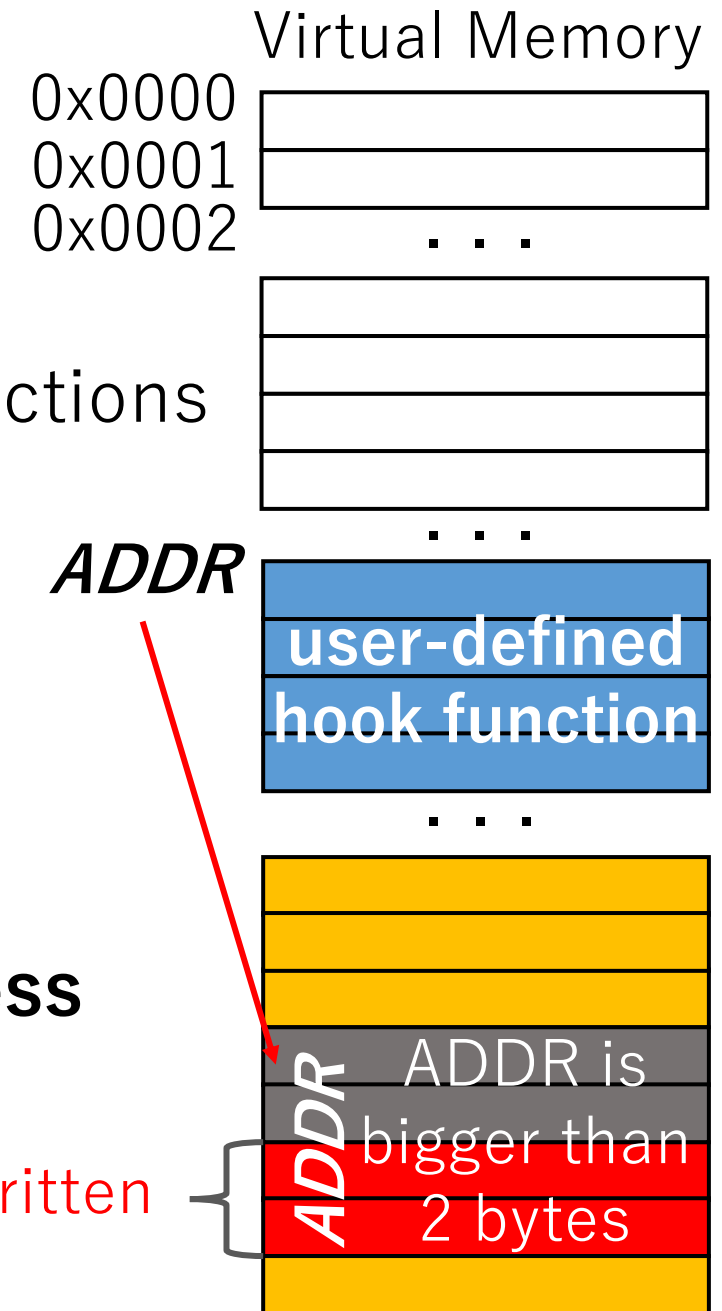
- On x86-64 CPUs, syscall and sysenter instructions trigger a system call

- syscall: 0x0f 0x05, sysenter: 0x0f 0x34

- syscall and sysenter are **2-byte** instructions

- Specification for a jump destination address needs more than 2 bytes**

If we put ADDR, subsequent instructions are overwritten



# Challenge

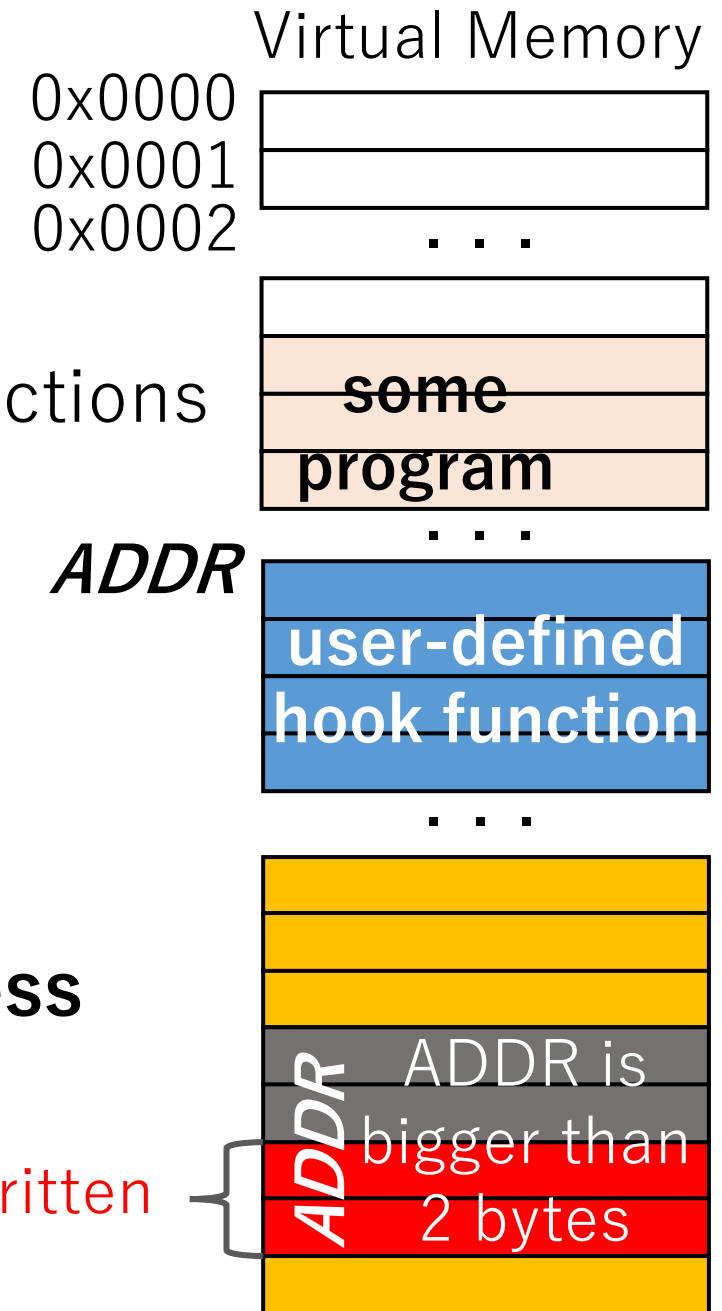
- On x86-64 CPUs, syscall and sysenter instructions trigger a system call

- syscall: 0x0f 0x05, sysenter: 0x0f 0x34

- syscall and sysenter are **2-byte** instructions

- Specification for a jump destination address needs more than 2 bytes**

If we put ADDR, subsequent instructions are overwritten



# Challenge

- On x86-64 CPUs, syscall and sysenter instructions trigger a system call

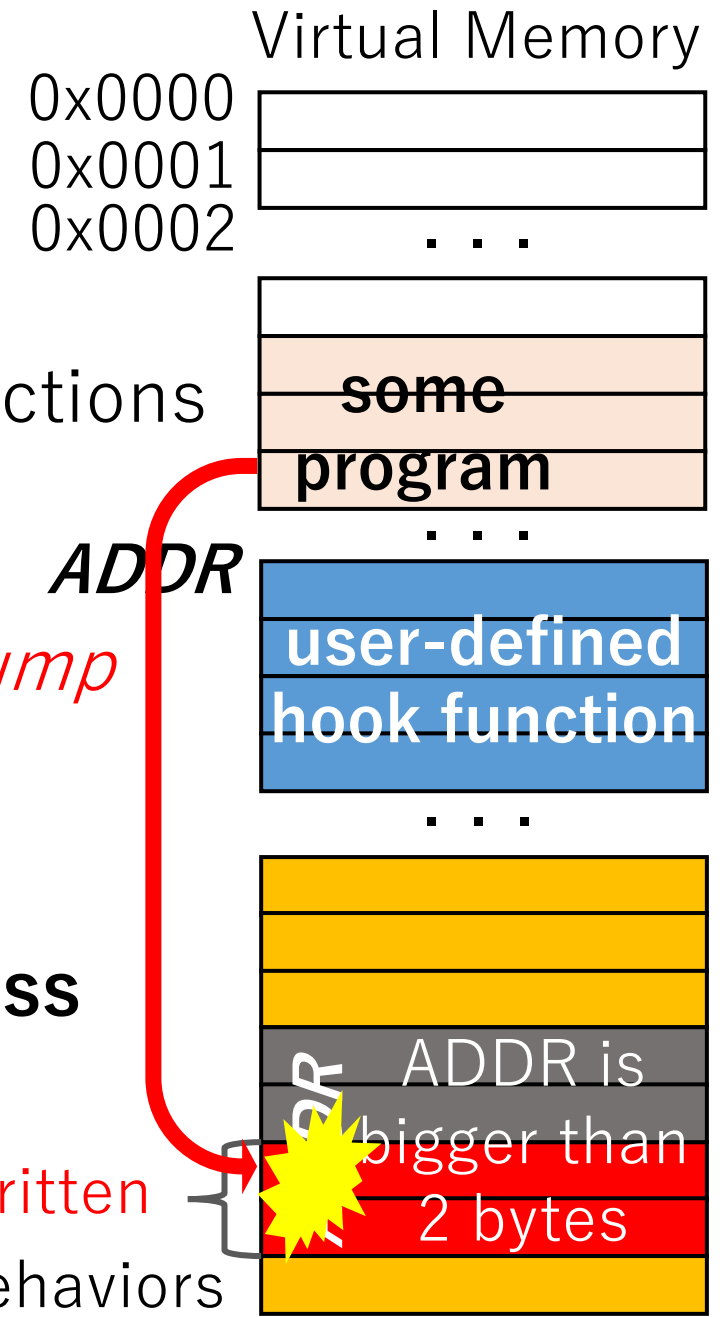
- syscall: 0x0f 0x05, sysenter: 0x0f 0x34

- syscall and sysenter are **2-byte** instructions

- Specification for a jump destination address needs more than 2 bytes**

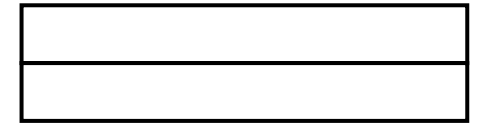
If we put ADDR, subsequent instructions are overwritten

jump to the overwritten part leads to unexpected behaviors



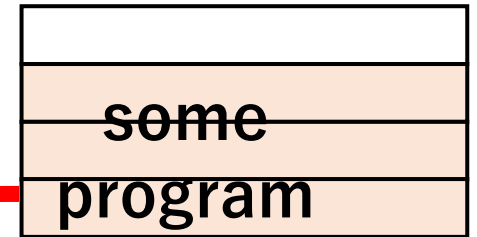
Goal ➡ **jump to a function using only 2 bytes**  
**originally occupied by syscall/sysenter**

0x0000  
 0x0001  
 0x0002  
 . . .



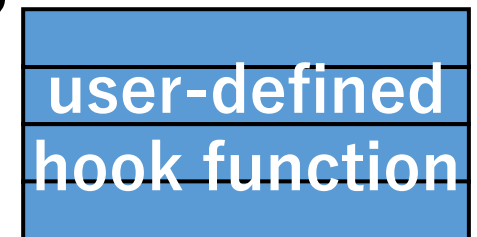
- On x86-64 CPUs, syscall and sysenter instructions trigger a system call

- syscall: 0x0f 0x05, sysenter: 0x0f 0x34



**ADDR**

*jump*

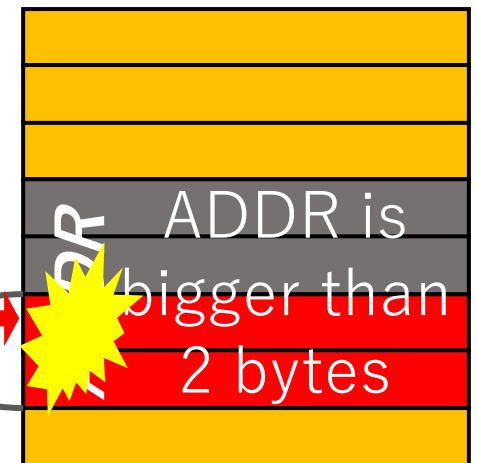


- syscall and sysenter are **2-byte** instructions

- Specification for a jump destination address needs more than 2 bytes**

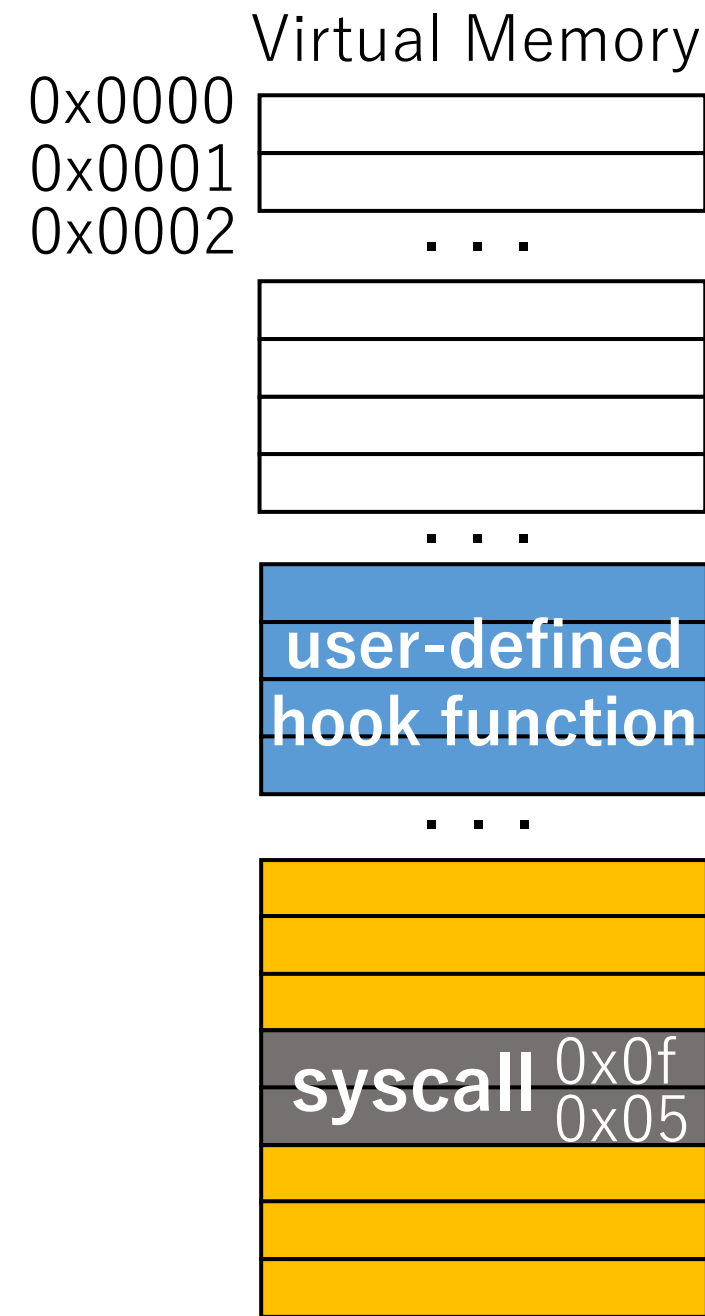
If we put ADDR, subsequent instructions are overwritten

jump to the overwritten part leads to unexpected behaviors



# Calling Convention

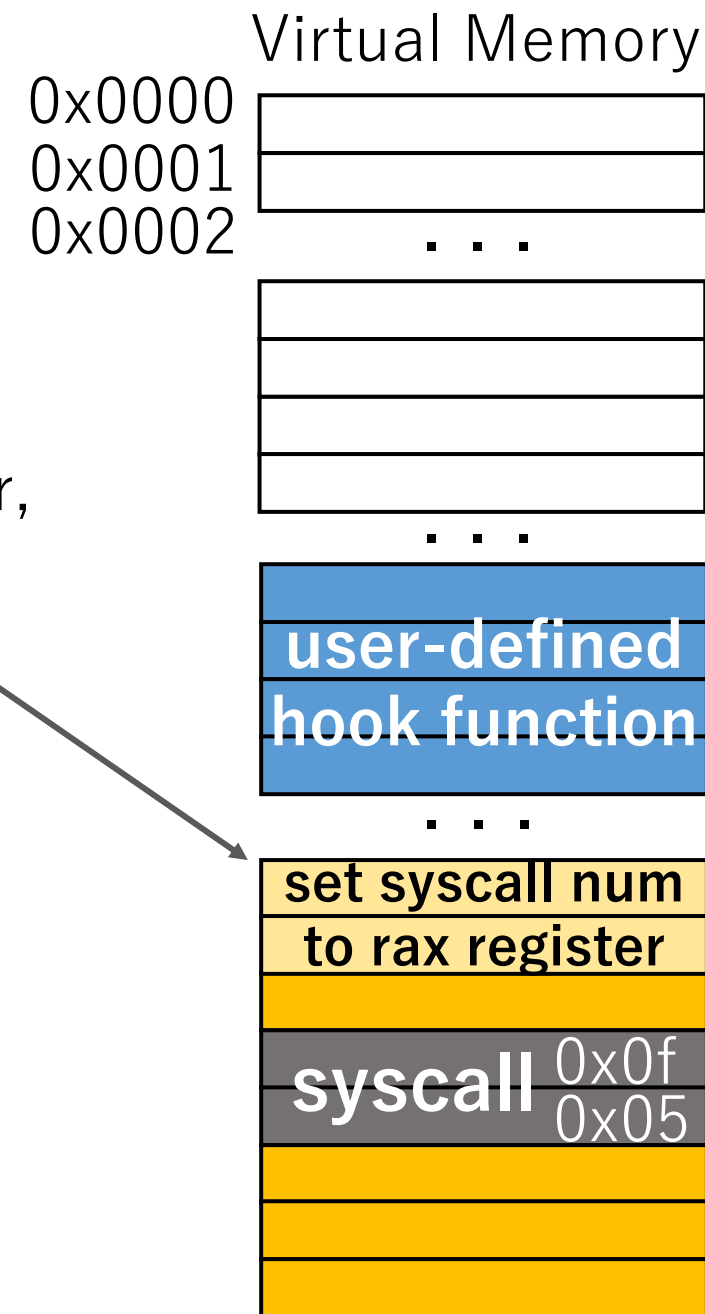
- How to invoke a system call





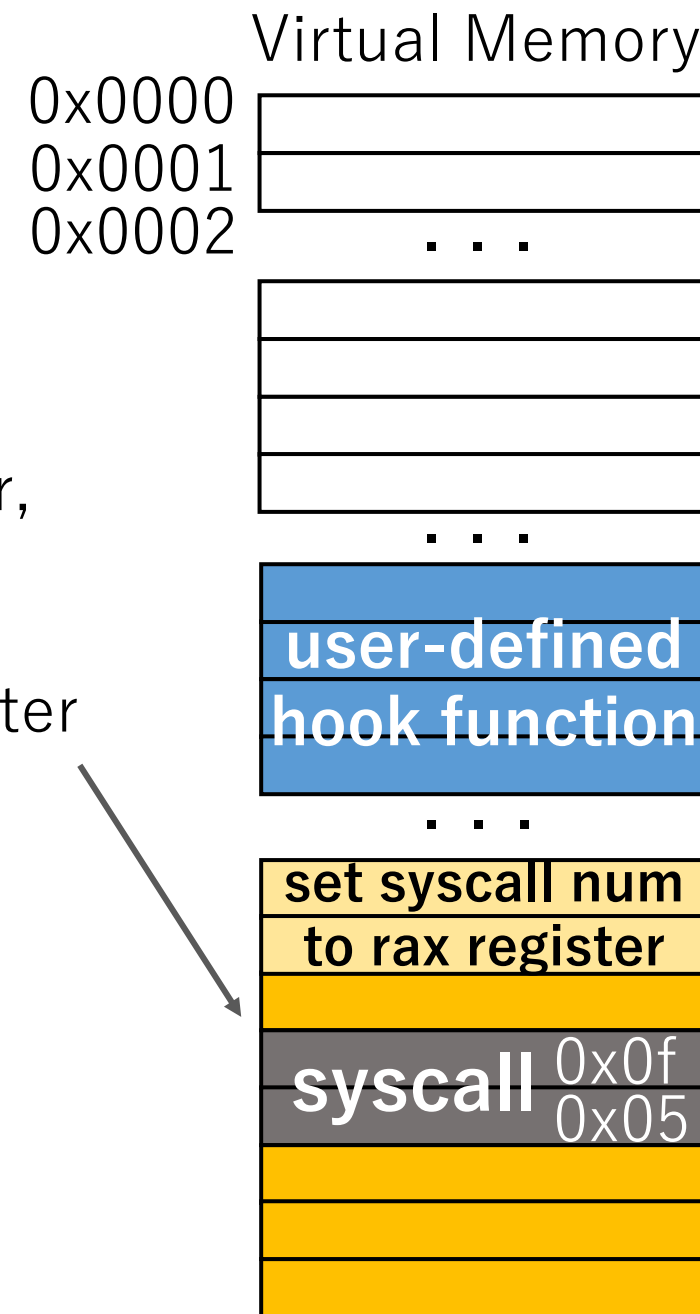
# Calling Convention

- How to invoke a system call
  - A user-space program sets a system call number, predefined by the kernel, to the **rax register**
    - e.g., 0: read(), 1: write(), 2: open(), ...



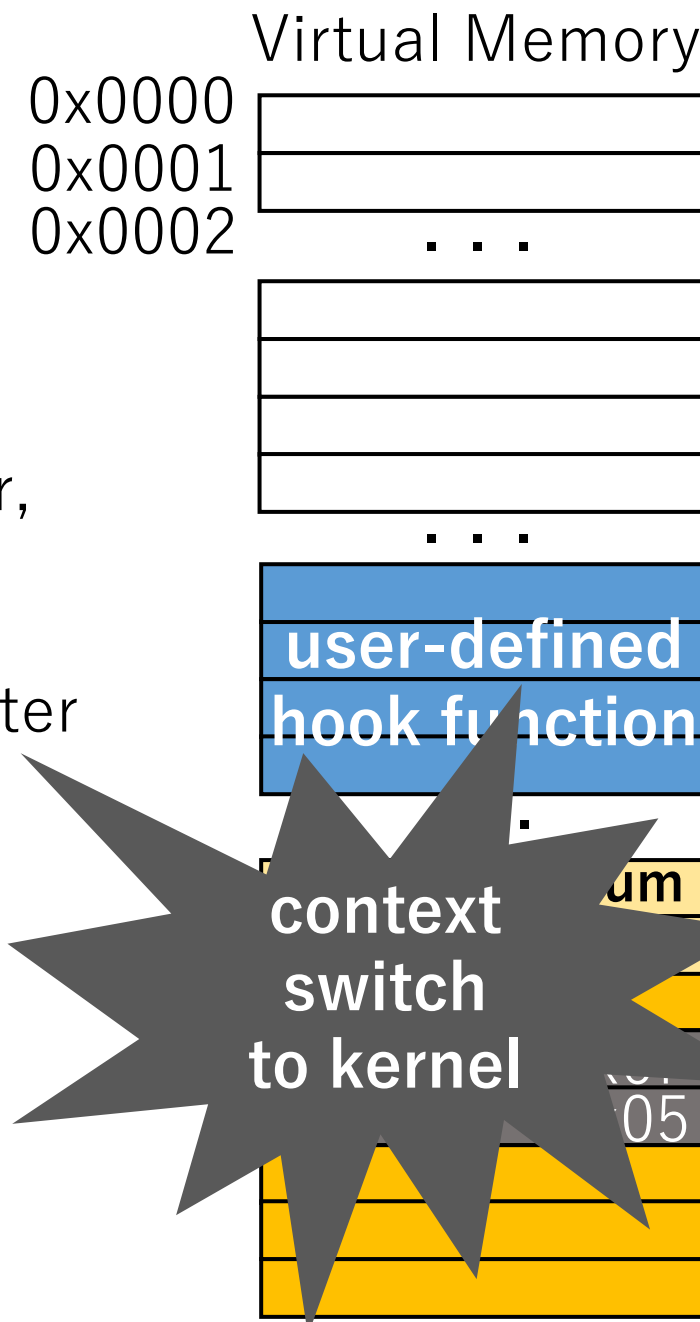
# Calling Convention

- How to invoke a system call
  - A user-space program sets a system call number, predefined by the kernel, to the **rax register**
    - e.g., 0: read(), 1: write(), 2: open(), ...
  - The user-space program executes syscall/sysenter



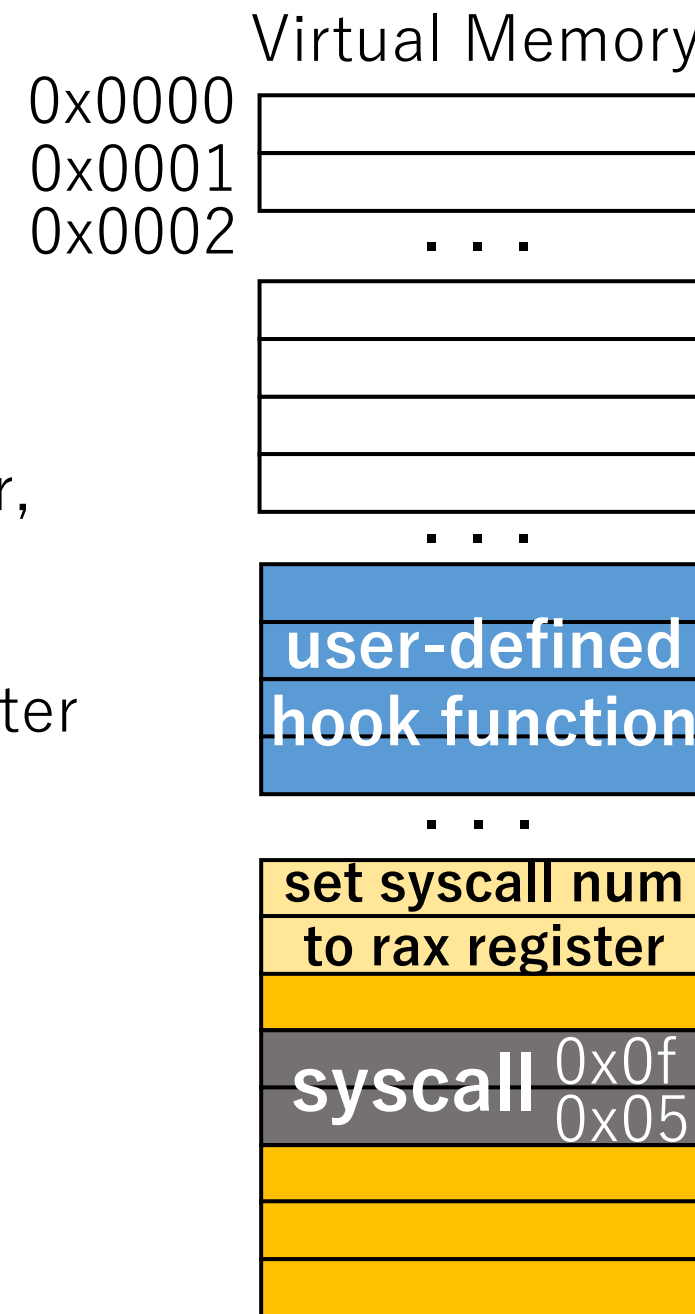
# Calling Convention

- How to invoke a system call
  - A user-space program sets a system call number, predefined by the kernel, to the **rax register**
    - e.g., 0: read(), 1: write(), 2: open(), ...
  - The user-space program executes syscall/sysenter
    - the context is switched to the kernel ----



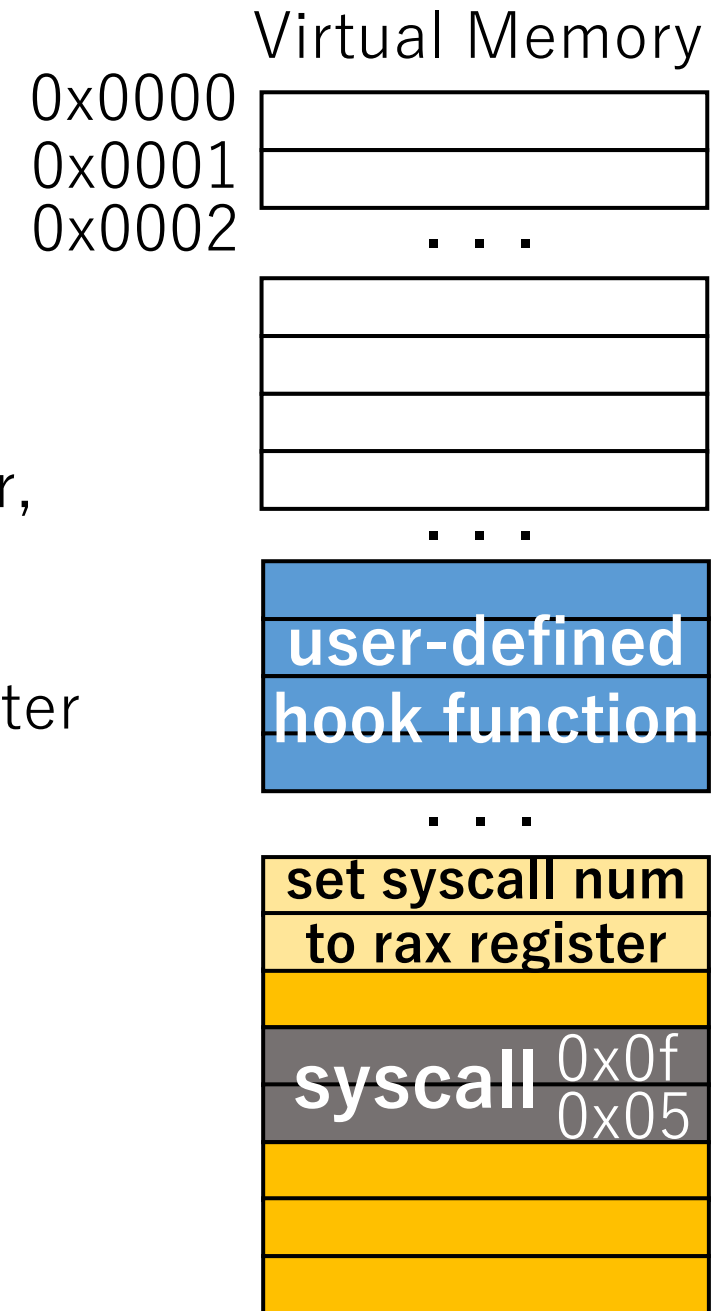
# Calling Convention

- How to invoke a system call
  - A user-space program sets a system call number, predefined by the kernel, to the **rax register**
    - e.g., 0: read(), 1: write(), 2: open(), ...
  - The user-space program executes syscall/sysenter
    - the context is switched to the kernel ----
  - Kernel executes a system call specified through the system call number set to the **rax register**



# Calling Convention

- How to invoke a system call
  - A user-space program sets a system call number, predefined by the kernel, to the **rax register**
    - e.g., 0: read(), 1: write(), 2: open(), ...
  - The user-space program executes syscall/sysenter  
---- the context is switched to the kernel ----
  - Kernel executes a system call specified through the system call number set to the **rax register**
    - if the rax register has 0, kernel executes read()
    - if the rax register has 1, kernel executes write()
    - if the rax register has 2, kernel executes open()

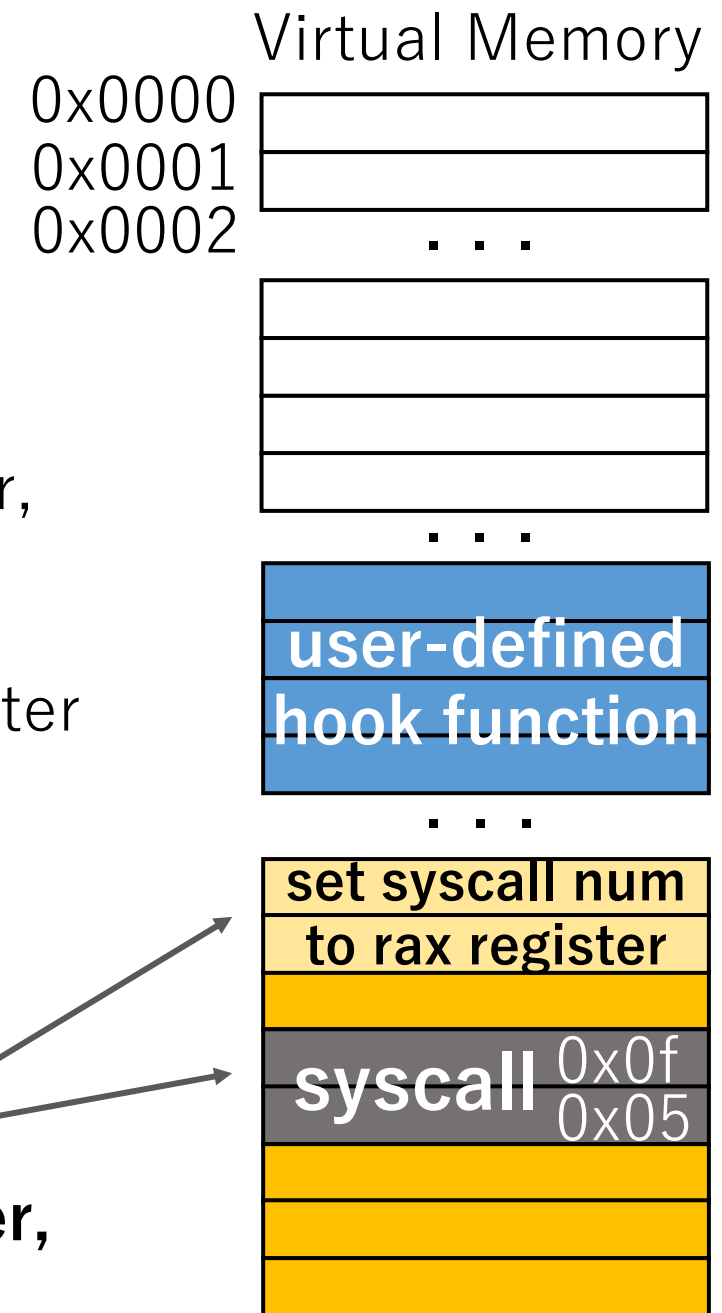


# Calling Convention

- How to invoke a system call
  - A user-space program sets a system call number, predefined by the kernel, to the **rax register**
    - e.g., 0: read(), 1: write(), 2: open(), ...
  - The user-space program executes syscall/sysenter
    - the context is switched to the kernel ----
  - Kernel executes a system call specified through the system call number set to the **rax register**

## Point: Calling Convention

**When syscall/sysenter is executed,  
the rax register always has a system call number,**

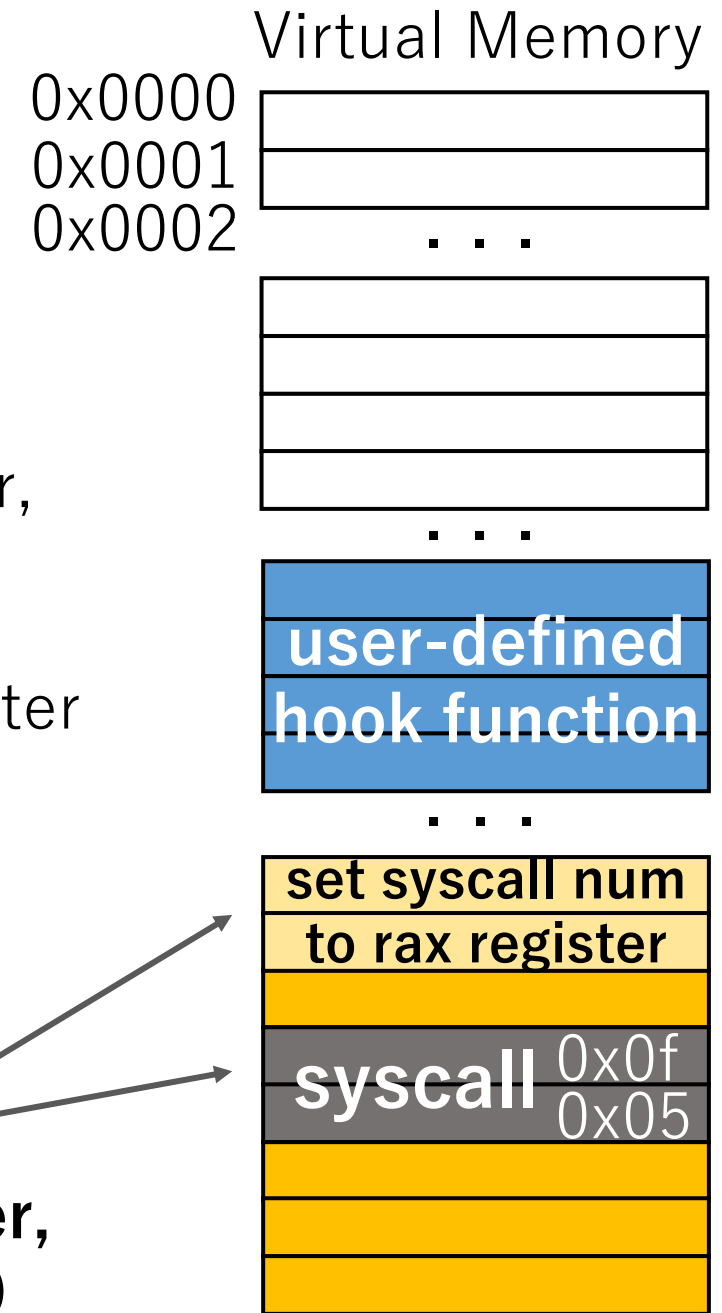


# Calling Convention

- How to invoke a system call
  - A user-space program sets a system call number, predefined by the kernel, to the **rax register**
    - e.g., 0: read(), 1: write(), 2: open(), ...
  - The user-space program executes syscall/sysenter
    - the context is switched to the kernel ----
  - Kernel executes a system call specified through the system call number set to the **rax register**

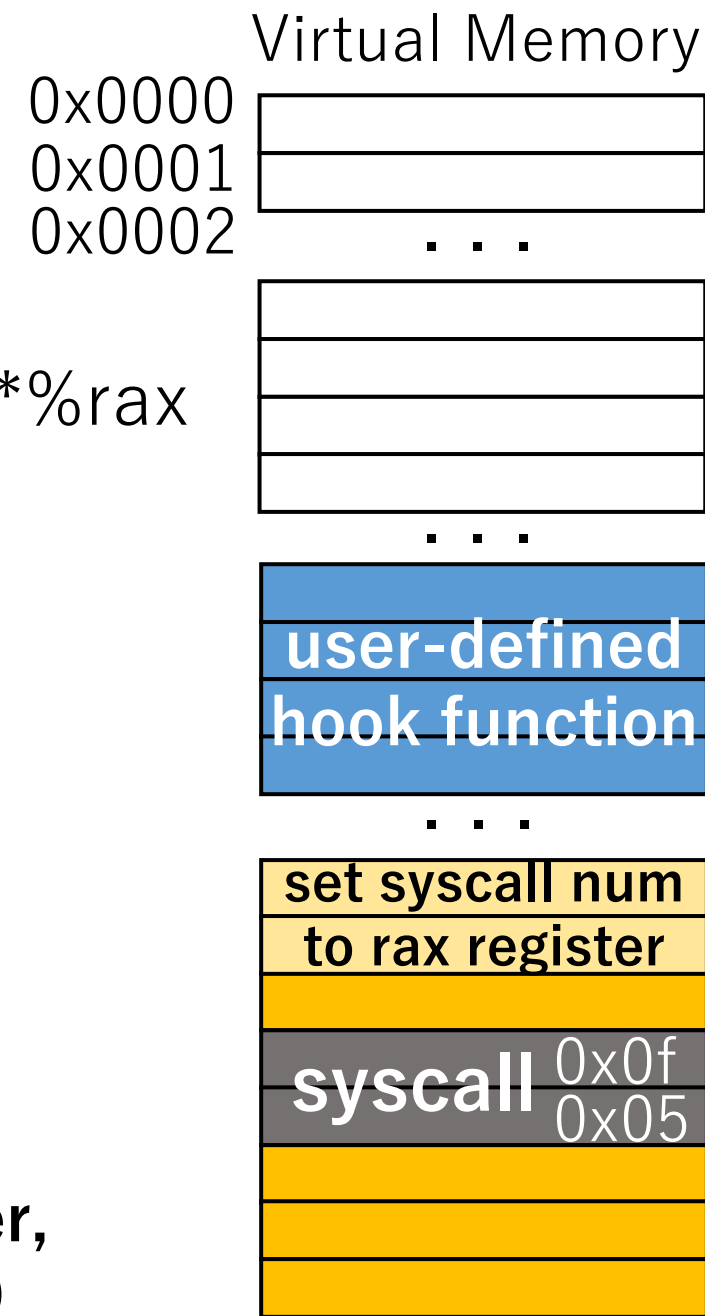
## Point: Calling Convention

**When syscall/sysenter is executed, the rax register always has a system call number, which is 0 ~ around 500 (defined in the kernel)**



zpline

- zpoline replaces syscall/sysenter with callq \*%rax



## Point: Calling Convention

**When syscall/sysenter is executed,  
the rax register always has a system call number,  
which is 0 ~ around 500 (defined in the kernel)**

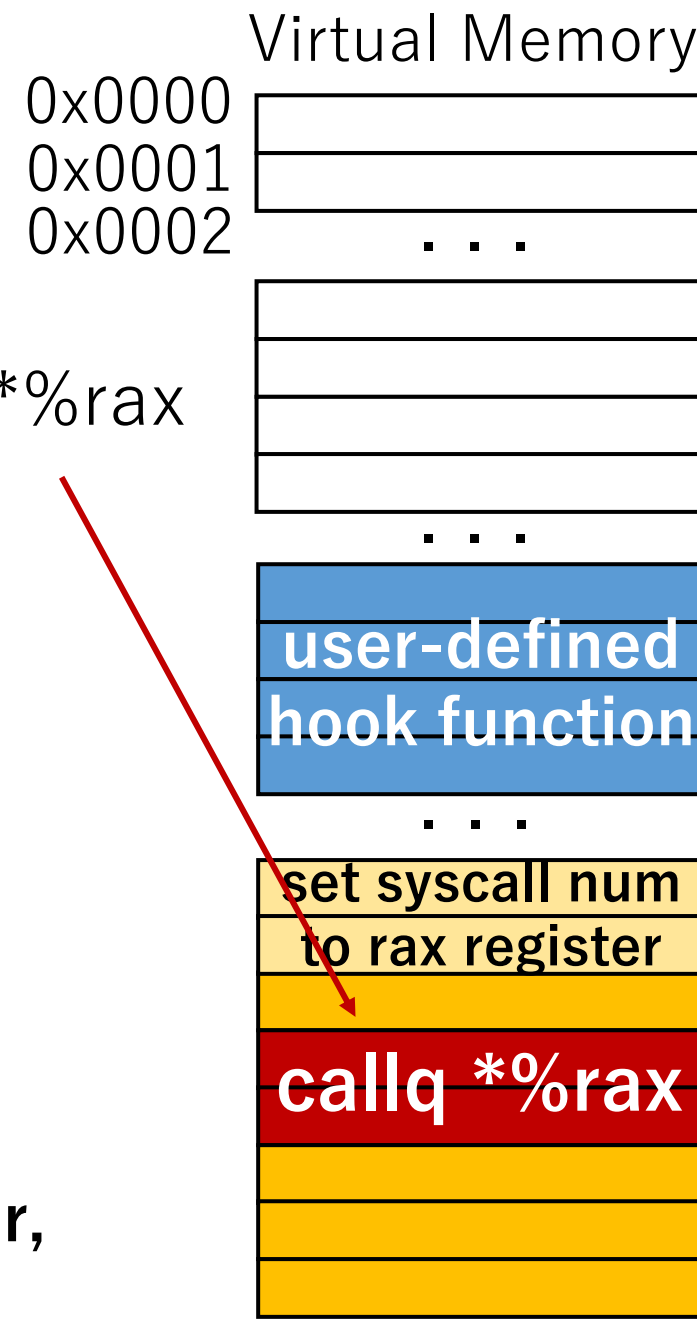


# zpoline

- zpoline replaces syscall/sysenter with callq \*%rax
  - callq \*%rax is a **2-byte** instruction (0xff 0xd0)

## Point: Calling Convention

When syscall/sysenter is executed,  
the rax register always has a system call number,  
which is 0 ~ around 500 (defined in the kernel)

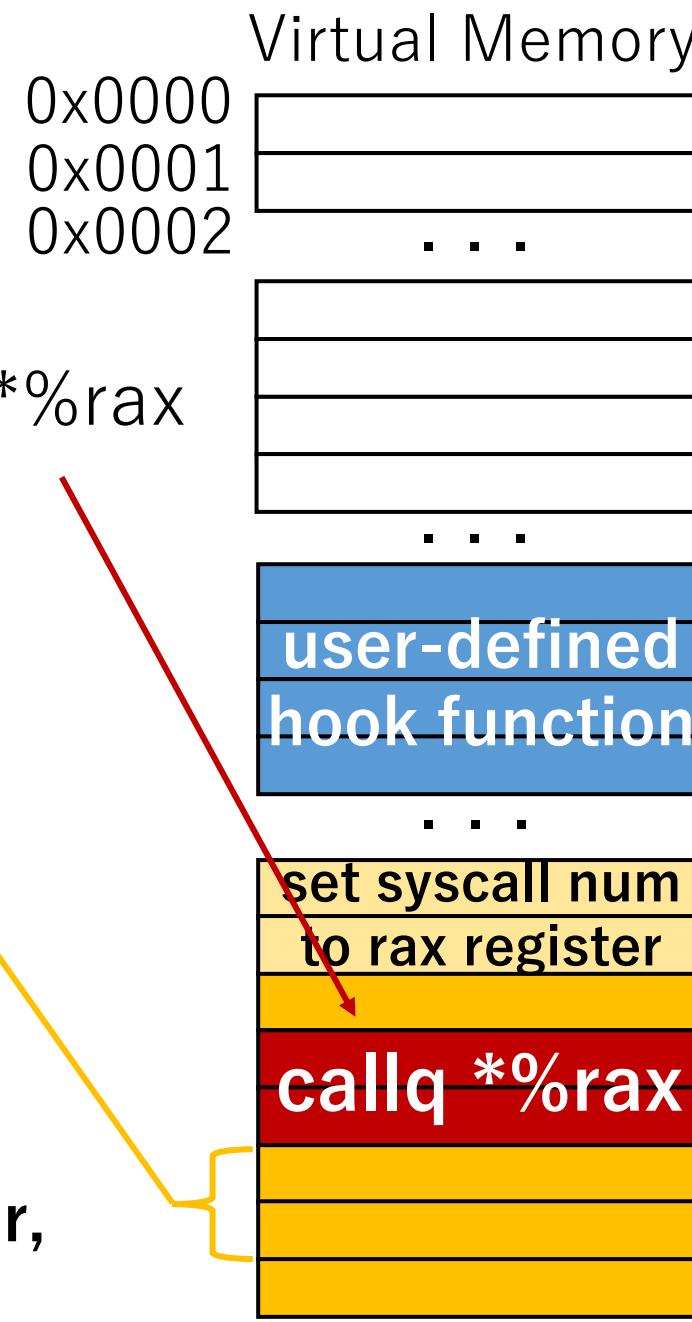


# zpoline

- zpoline replaces syscall/sysenter with callq \*%rax
  - callq \*%rax is a **2-byte** instruction (0xff 0xd0)
  - Neighbour instructions are not overwritten

## Point: Calling Convention

When syscall/sysenter is executed,  
the rax register always has a system call number,  
which is 0 ~ around 500 (defined in the kernel)

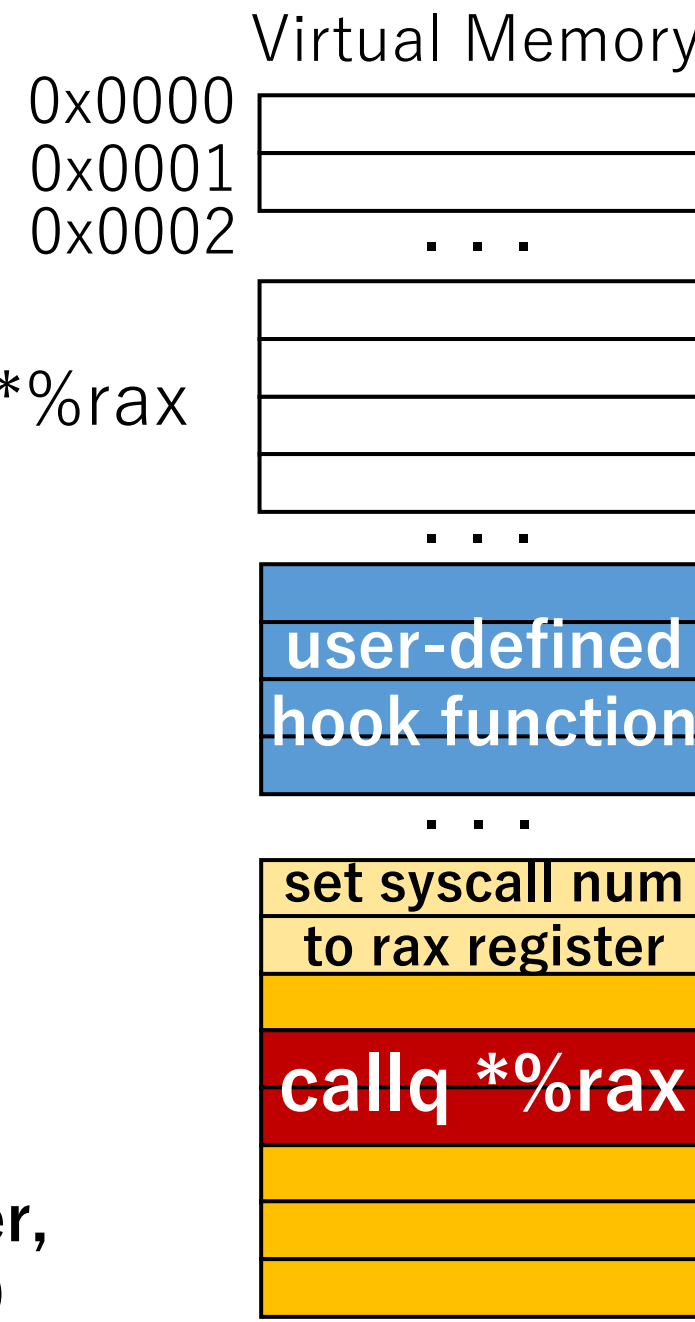


# zpoline

- zpoline replaces syscall/sysenter with callq \*%rax
  - callq \*%rax is a **2-byte** instruction (0xff 0xd0)
    - Neighbour instructions are not overwritten
  - callq \*%rax is an instruction to jump to the address stored in the rax register

## Point: Calling Convention

**When syscall/sysenter is executed,  
the rax register always has a system call number,  
which is 0 ~ around 500 (defined in the kernel)**

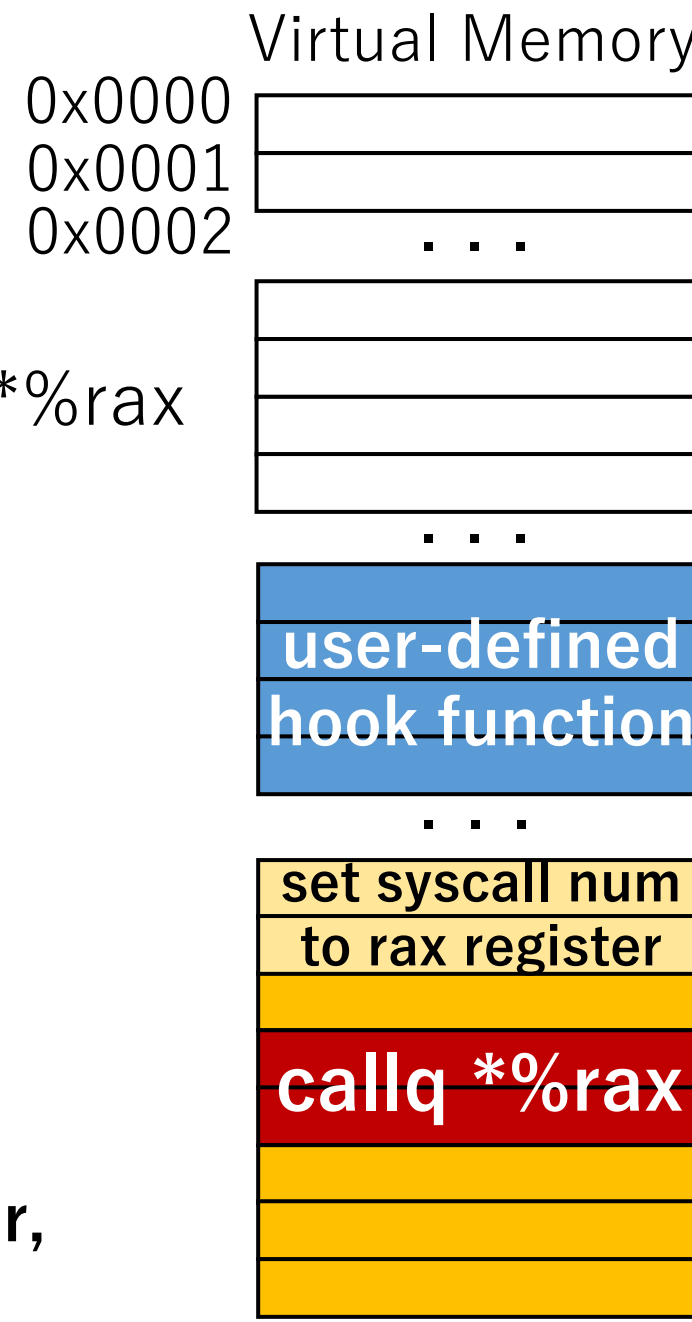


# zpoline

- zpoline replaces syscall/sysenter with callq \*%rax
  - callq \*%rax is a **2-byte** instruction (0xff 0xd0)
    - Neighbour instructions are not overwritten
  - callq \*%rax is an instruction to jump to the address stored in the rax register

## Point: Calling Convention

When syscall/sysenter is executed,  
the rax register always has a system call number,  
which is 0 ~ around 500 (defined in the kernel)



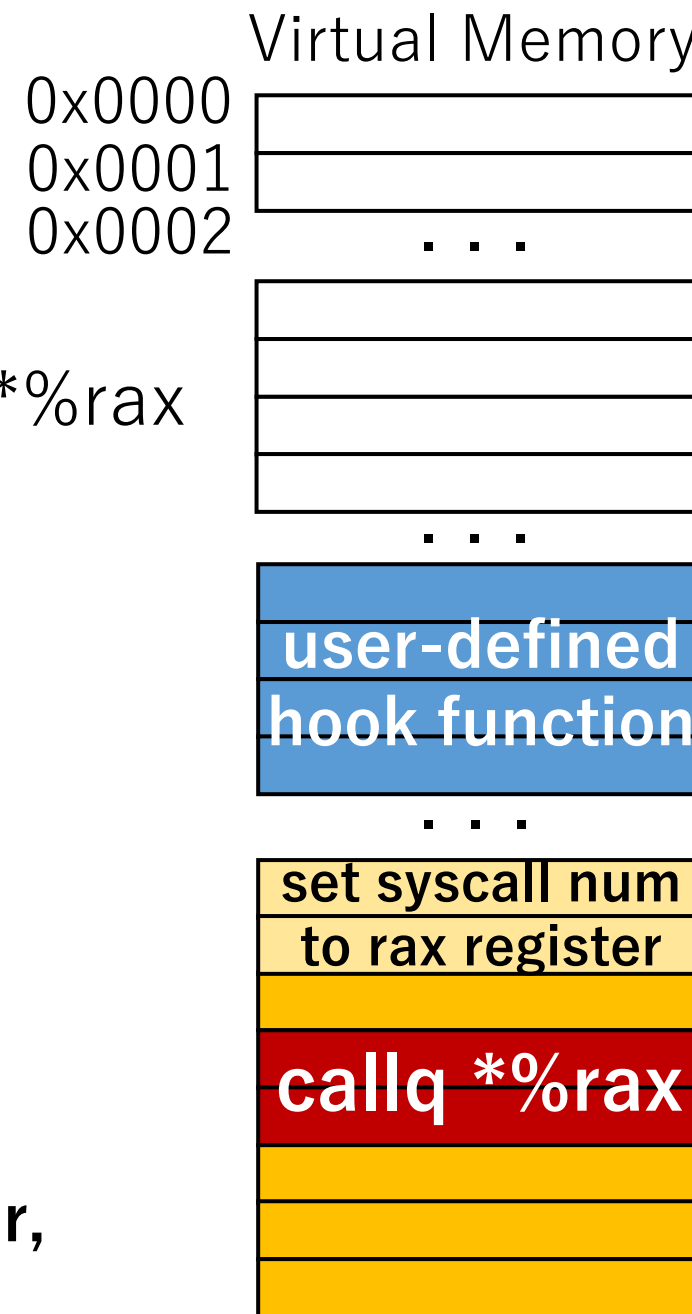
# zpoline

- zpoline replaces syscall/sysenter with callq \*%rax
  - callq \*%rax is a **2-byte** instruction (0xff 0xd0)
    - Neighbour instructions are not overwritten
  - callq \*%rax is an instruction to jump to the address stored in the rax register

**After the binary rewriting**

**Point: Calling Convention**

**When syscall/sysenter is executed,  
the rax register always has a system call number,  
which is 0 ~ around 500 (defined in the kernel)**



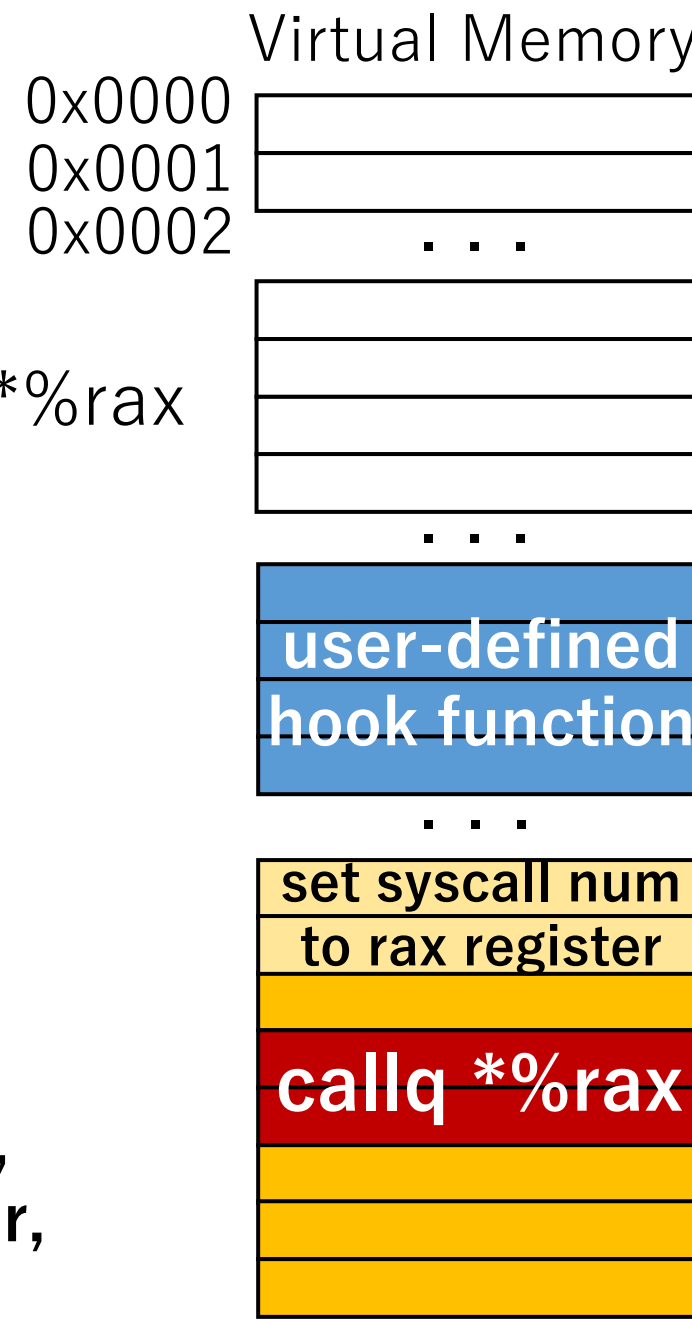
# zpoline

- zpoline replaces syscall/sysenter with callq \*%rax
  - callq \*%rax is a **2-byte** instruction (0xff 0xd0)
    - Neighbour instructions are not overwritten
  - callq \*%rax is an instruction to jump to the address stored in the rax register

**After the binary rewriting**

**Point: Calling Convention**

**When ~~syscall/sysenter~~ **callq \*%rax** is executed, the rax register always has a system call number, which is 0 ~ around 500 (defined in the kernel)**



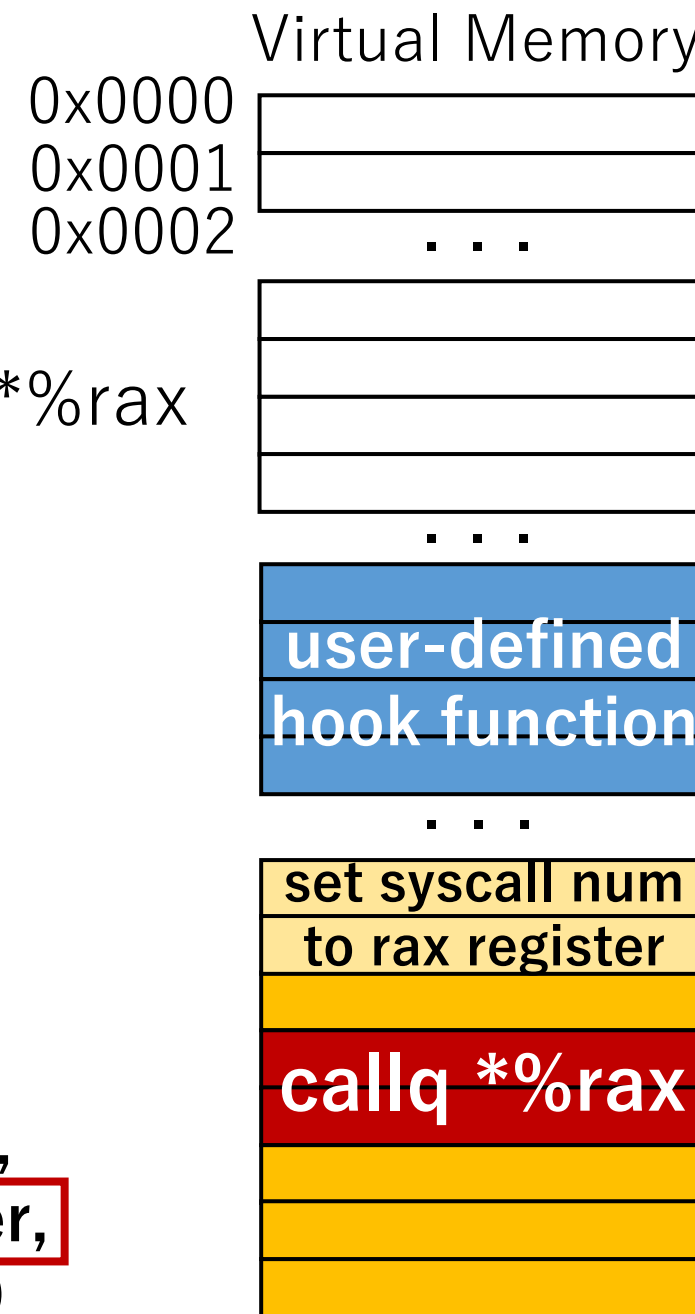
# zpoline

- zpoline replaces syscall/sysenter with callq \*%rax
  - callq \*%rax is a **2-byte** instruction (0xff 0xd0)
    - Neighbour instructions are not overwritten
  - callq \*%rax is an instruction to jump to the address stored in the rax register

**After the binary rewriting**

Point: Calling Convention

When ~~syscall/sysenter~~ **callq \*%rax** is executed,  
the rax register always has a system call number,  
which is 0 ~ around 500 (defined in the kernel)



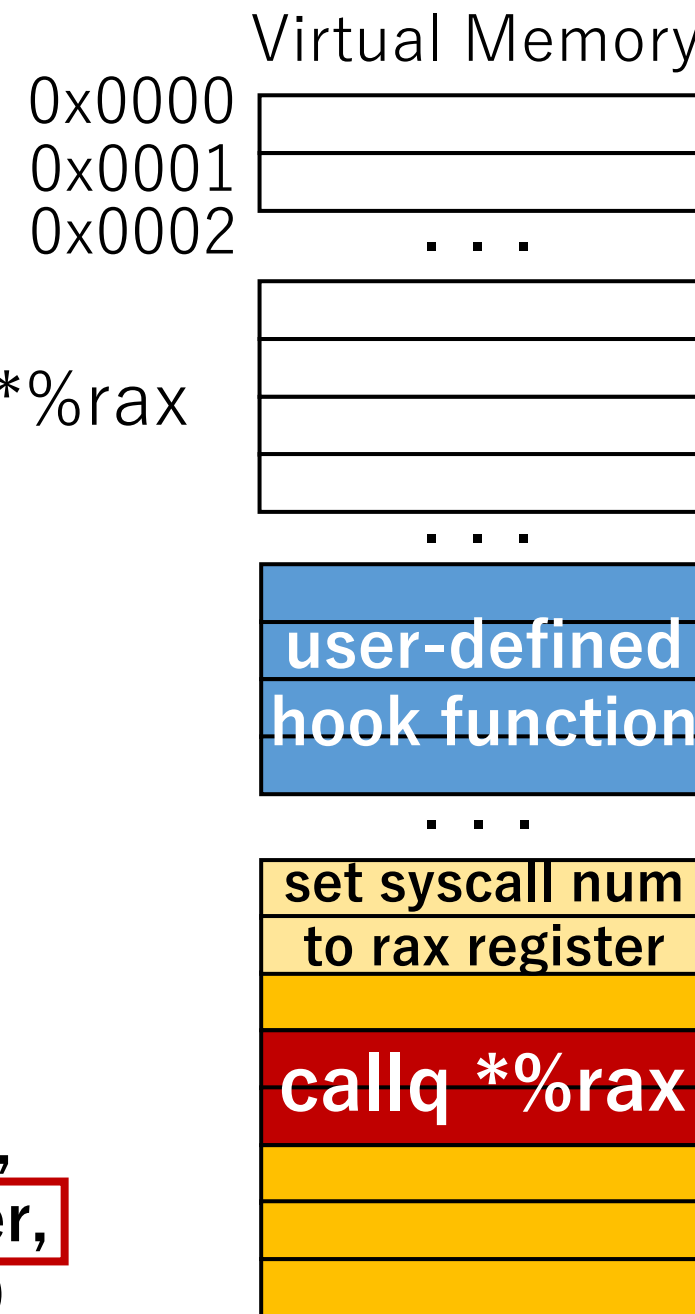
# zpoline

- zpoline replaces syscall/sysenter with callq \*%rax
  - callq \*%rax is a **2-byte** instruction (0xff 0xd0)
    - Neighbour instructions are not overwritten
  - callq \*%rax is an instruction to jump to the address stored in the rax register

**After the binary rewriting**

**Point: Calling Convention**

**When ~~syscall/sysenter~~ callq \*%rax is executed,**  
**the rax register always has a system call number,**  
**which is 0 ~ around 500 (defined in the kernel)**





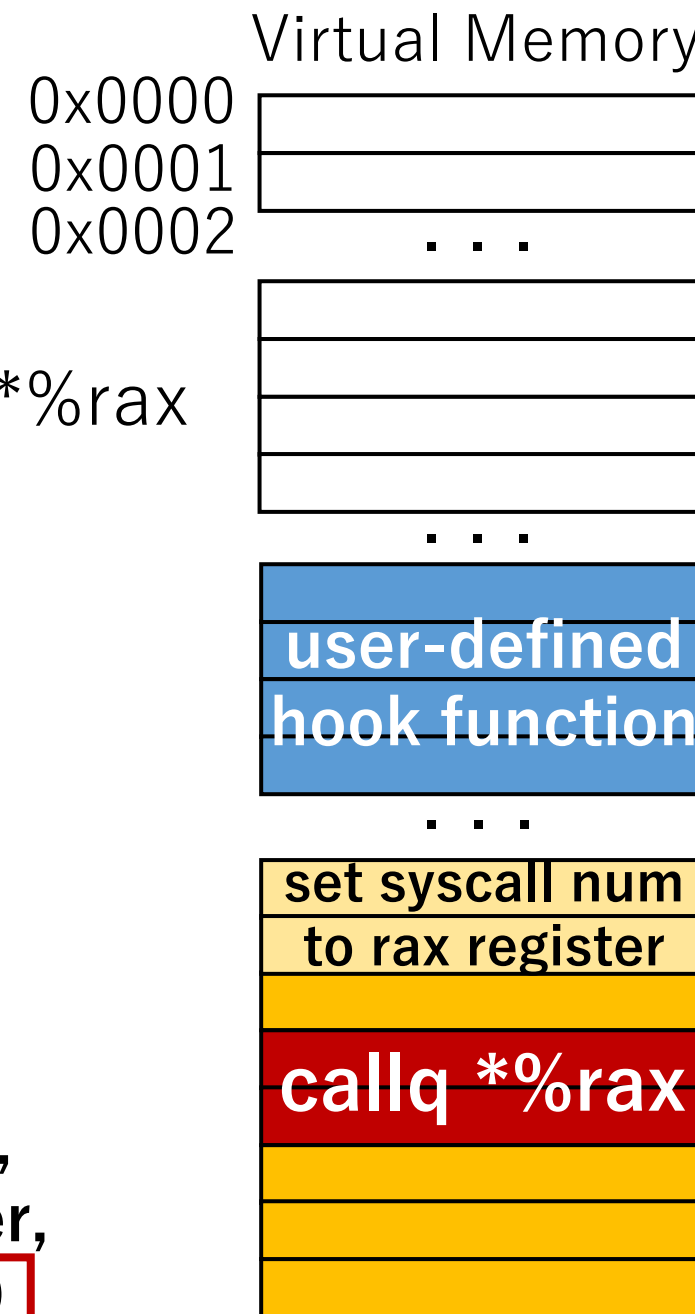
# zpoline

- zpoline replaces syscall/sysenter with callq \*%rax
  - callq \*%rax is a **2-byte** instruction (0xff 0xd0)
    - Neighbour instructions are not overwritten
  - callq \*%rax is an instruction to jump to the address stored in the rax register

**After the binary rewriting**

**Point: Calling Convention**

When ~~syscall/sysenter~~ **callq \*%rax** is executed, the rax register always has a system call number, **which is 0 ~ around 500 (defined in the kernel)**



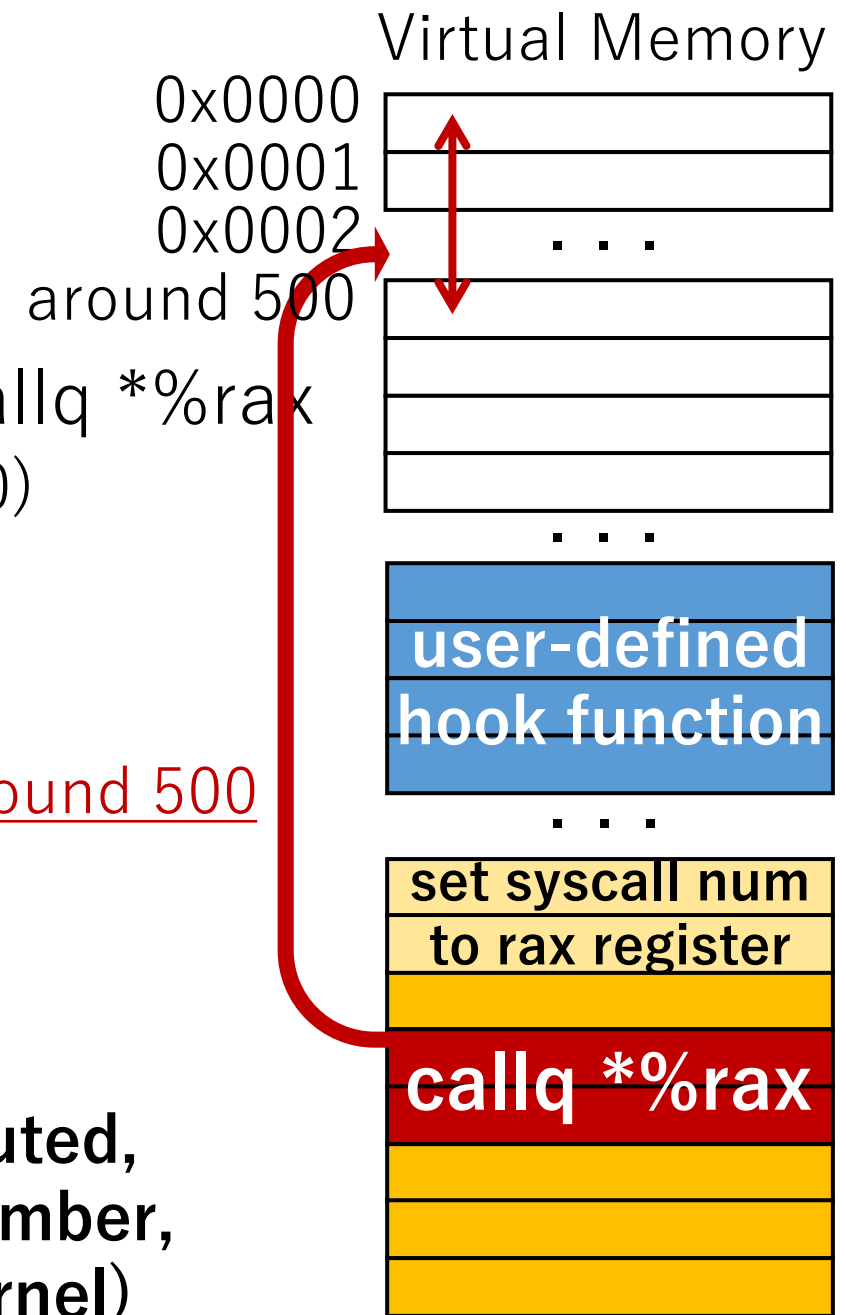
# zpoline

- zpoline replaces syscall/sysenter with callq \*%rax
  - callq \*%rax is a **2-byte** instruction (0xff 0xd0)
    - Neighbour instructions are not overwritten
  - callq \*%rax is an instruction to jump to the address stored in the rax register
  - replaced callq \*%rax jumps to address 0~around 500

**After the binary rewriting**

**Point: Calling Convention**

**When ~~syscall/sysenter~~ **callq \*%rax** is executed, the rax register always has a system call number, which is 0 ~ around 500 (defined in the kernel)**

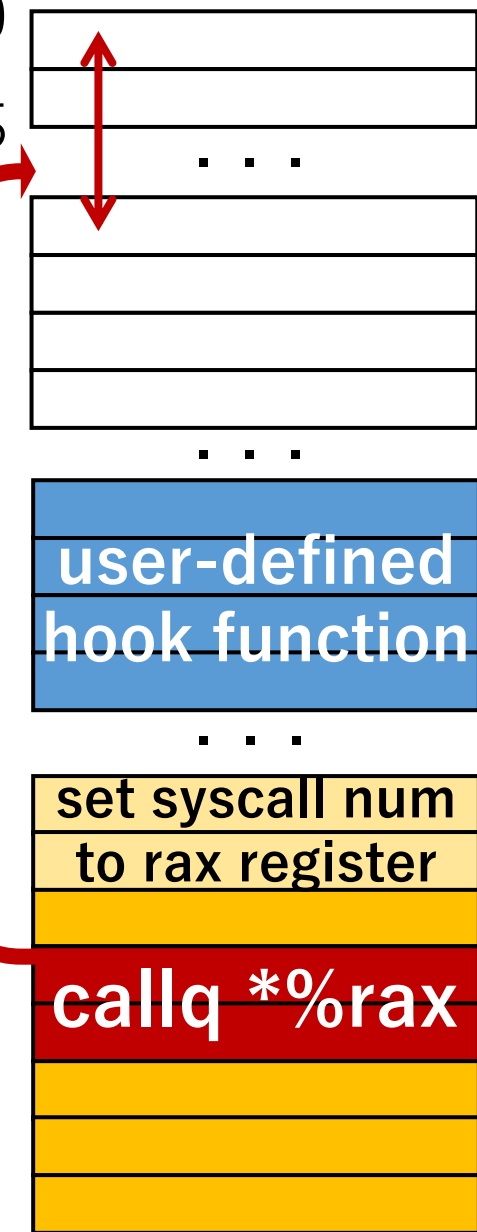


# zpoline

address range, potentially  
replaced “`callq *%rax`” jumps to  
(  $N$  is the max syscall number )

0x0000  
0x0001  
0x0002  
...

Virtual Memory



- zpoline replaces `syscall/sysenter` with `callq *%rax`
  - `callq *%rax` is a **2-byte** instruction (0xff 0xd0)
    - Neighbour instructions are not overwritten
  - `callq *%rax` is an instruction to jump to the address stored in the `rax` register
  - replaced `callq *%rax` jumps to address 0~around 500

**After the binary rewriting**

**Point: Calling Convention**

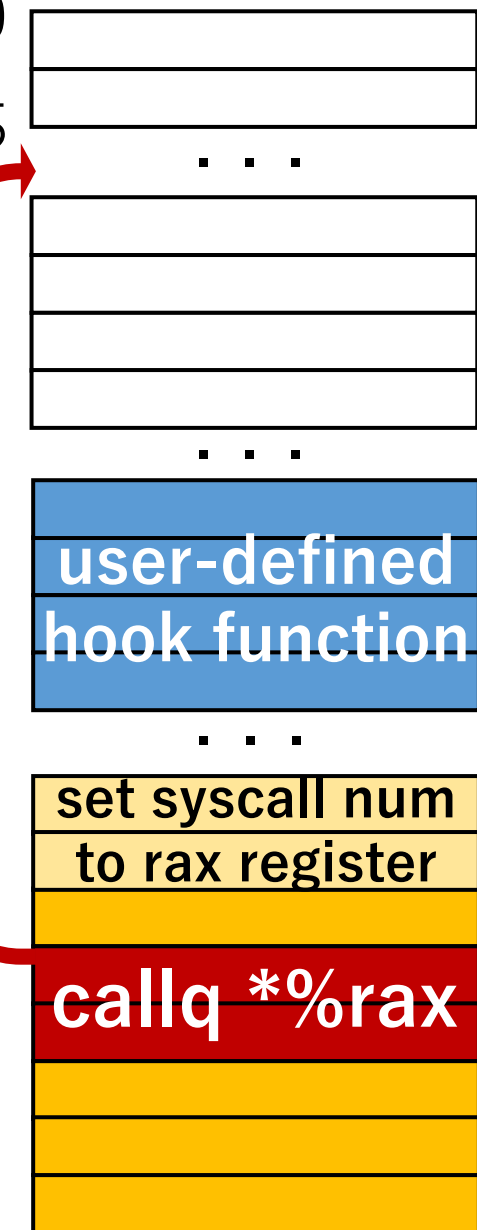
When ~~`syscall/sysenter`~~ **`callq *%rax`** is executed,  
the `rax` register always has a system call number,  
which is 0 ~ around 500 (defined in the kernel)

# zpoline

address range, potentially  
replaced “callq \*%rax” jumps to  
(  $N$  is the max syscall number )

0x0000  
0x0001  
0x0002  
...

Virtual Memory



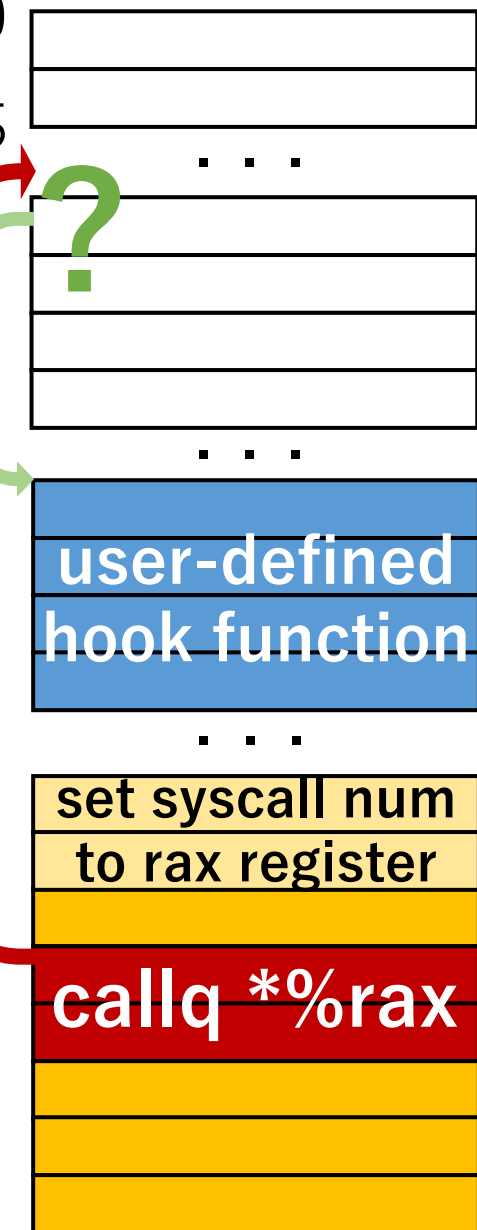
- zpoline replaces syscall/sysenter with callq \*%rax
  - callq \*%rax is a **2-byte** instruction (0xff 0xd0)
    - Neighbour instructions are not overwritten
  - callq \*%rax is an instruction to jump to the address stored in the rax register
  - replaced callq \*%rax jumps to address 0~around 500

# zpoline

address range, potentially  
replaced “callq \*%rax” jumps to  
(  $N$  is the max syscall number )

0x0000  
0x0001  
0x0002  
...

Virtual Memory



- zpoline replaces syscall/sysenter with callq \*%rax
  - callq \*%rax is a **2-byte** instruction (0xff 0xd0)
    - Neighbour instructions are not overwritten
  - callq \*%rax is an instruction to jump to the address stored in the rax register
  - replaced callq \*%rax jumps to address 0~around 500

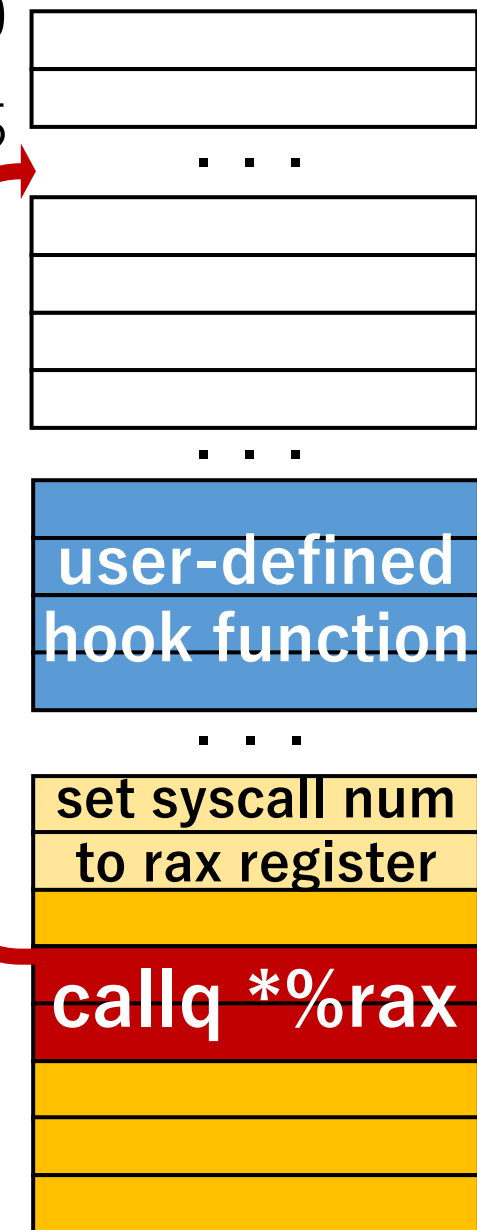
**How to redirect to the user-defined hook function?**

# zpoline

address range, potentially  
replaced “callq \*%rax” jumps to  
(  $N$  is the max syscall number )

0x0000  
0x0001  
0x0002  
...

Virtual Memory



- zpoline replaces syscall/sysenter with callq \*%rax
  - callq \*%rax is a **2-byte** instruction (0xff 0xd0)
    - Neighbour instructions are not overwritten
  - callq \*%rax is an instruction to jump to the address stored in the rax register
  - replaced callq \*%rax jumps to address 0~around 500

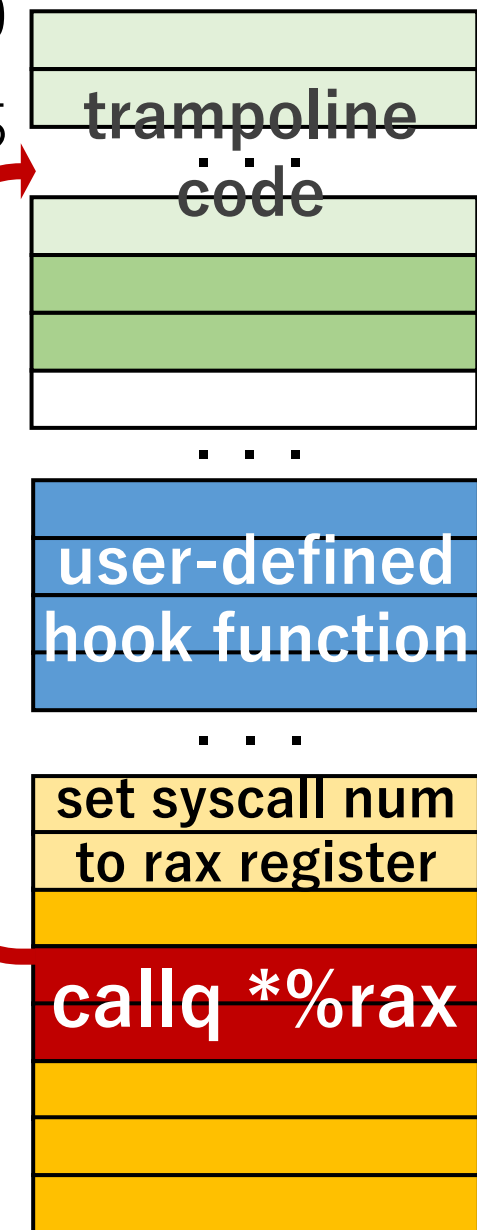
How to redirect to the user-defined hook function?

# zpoline

address range, potentially  
replaced “callq \*%rax” jumps to  
(  $N$  is the max syscall number )

0x0000  
0x0001  
0x0002  
 $N$

Virtual Memory



- zpoline replaces syscall/sysenter with callq \*%rax
  - callq \*%rax is a **2-byte** instruction (0xff 0xd0)
    - Neighbour instructions are not overwritten
  - callq \*%rax is an instruction to jump to the address stored in the rax register
  - replaced callq \*%rax jumps to address 0~around 500

## How to redirect to the user-defined hook function?

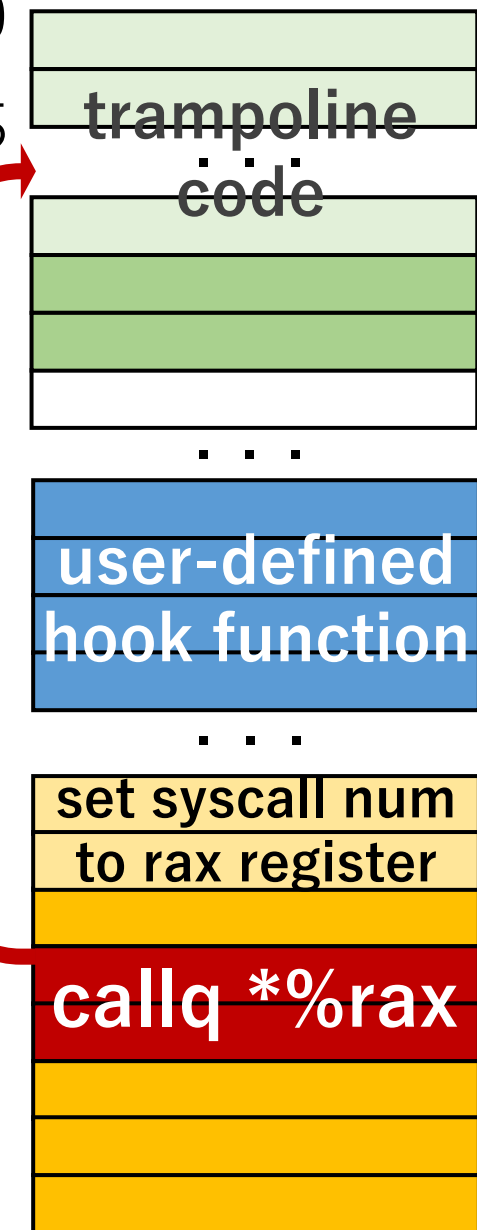
- zpoline instantiates trampoline code at address 0

# zpoline

address range, potentially  
replaced “`callq *%rax`” jumps to  
(  $N$  is the max syscall number )

0x0000  
0x0001  
0x0002  
 $N$

Virtual Memory



- zpoline replaces `syscall/sysenter` with `callq *%rax`
  - `callq *%rax` is a **2-byte** instruction (0xff 0xd0)
    - Neighbour instructions are not overwritten
  - `callq *%rax` is an instruction to jump to the address stored in the `rax` register
  - replaced `callq *%rax` jumps to address 0~around 500

## How to redirect to the user-defined hook function?

- zpoline instantiates trampoline code at address 0

Tram**poline** code at address 0 (**zero**) → **zpoline**

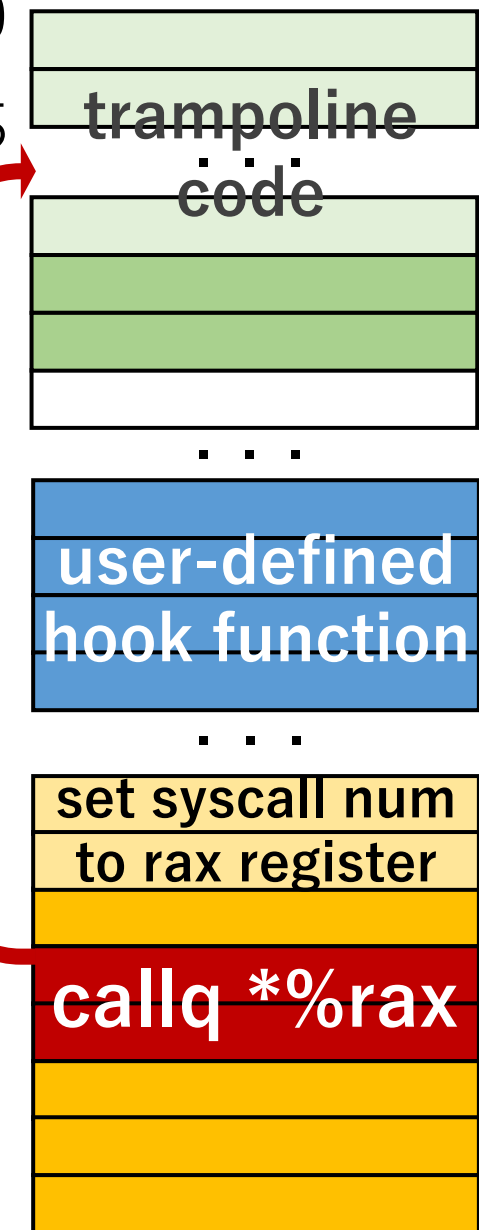


# zpoline

address range, potentially  
replaced “callq \*%rax” jumps to  
(  $N$  is the max syscall number )

0x0000  
0x0001  
0x0002  
 $N$

Virtual Memory



- zpoline replaces syscall/sysenter with callq \*%rax
  - callq \*%rax is a **2-byte** instruction (0xff 0xd0)
    - Neighbour instructions are not overwritten
  - callq \*%rax is an instruction to jump to the address stored in the rax register
  - replaced callq \*%rax jumps to address 0~around 500

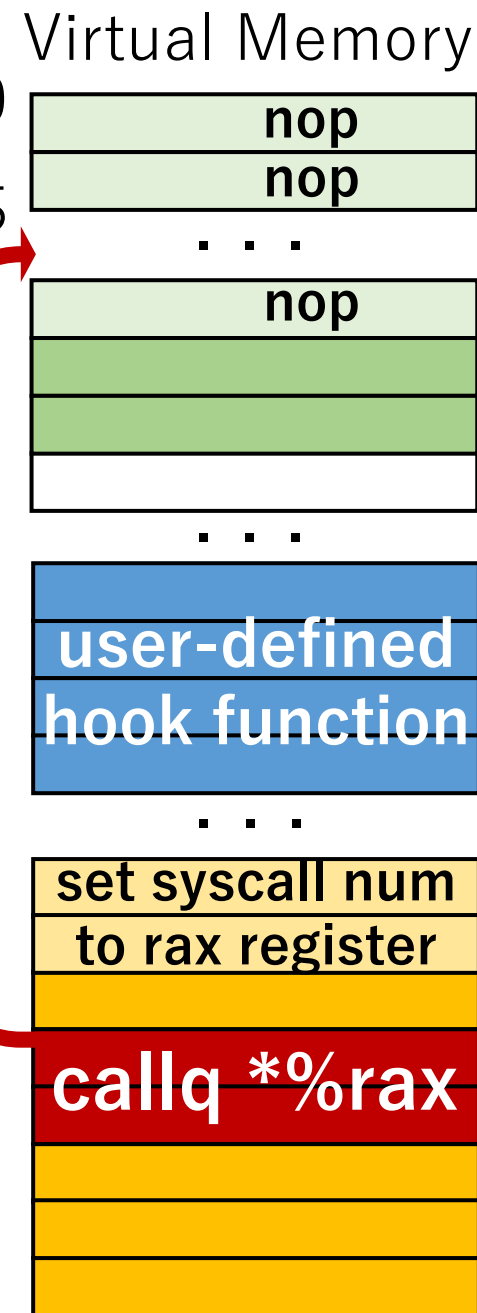
## How to redirect to the user-defined hook function?

- zpoline instantiates trampoline code at address 0

# zpoline

address range, potentially  
replaced “callq \*%rax” jumps to  
(  $N$  is the max syscall number )

0x0000  
0x0001  
0x0002  
...



- zpoline replaces syscall/sysenter with callq \*%rax
  - callq \*%rax is a **2-byte** instruction (0xff 0xd0)
    - Neighbour instructions are not overwritten
  - callq \*%rax is an instruction to jump to the address stored in the rax register
  - replaced callq \*%rax jumps to address 0~around 500

## How to redirect to the user-defined hook function?

- zpoline instantiates trampoline code at address 0
  - fills address range 0 to  $N$  with nop (0x90)

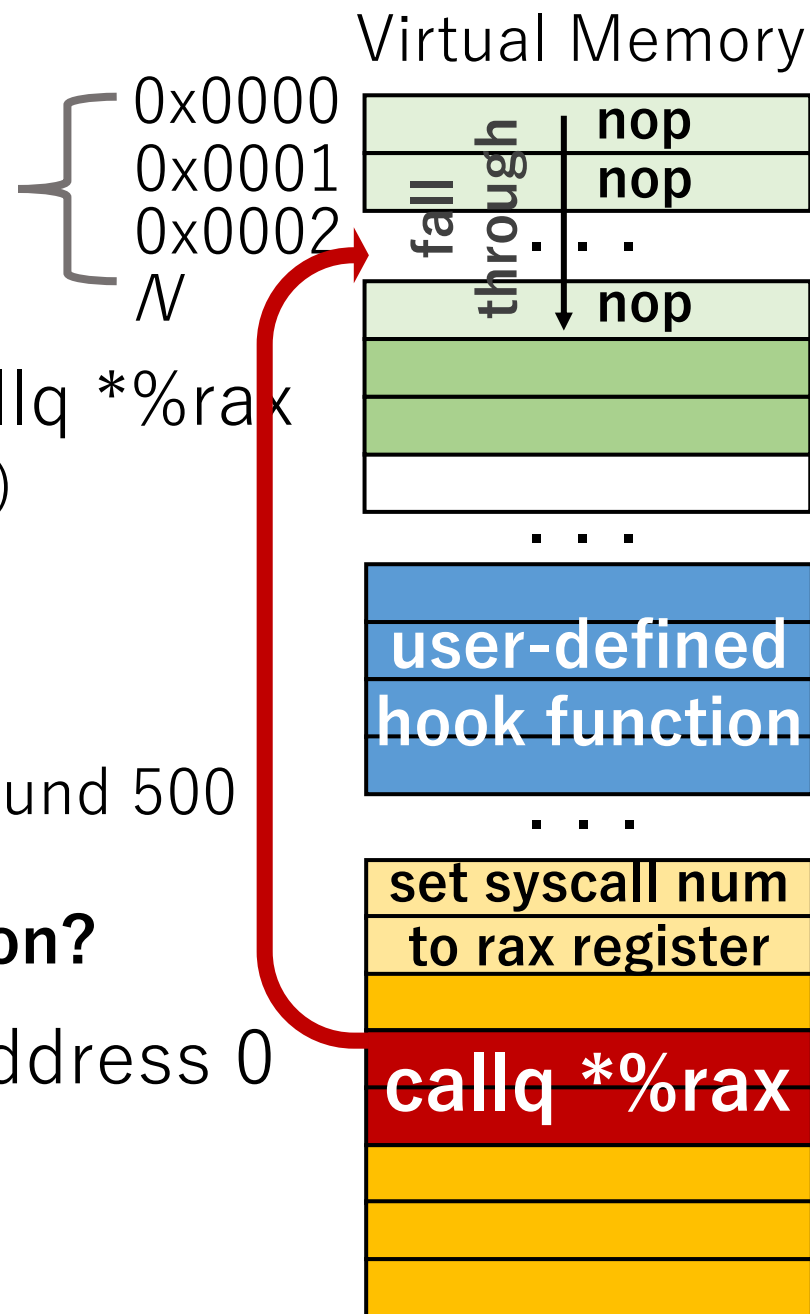
zpline

address range, potentially  
replaced “callq \*%rax” jumps to  
(  $N$  is the max syscall number )

- zpline replaces syscall/sysenter with callq \*%rax
  - callq \*%rax is a **2-byte** instruction (0xff 0xd0)
    - Neighbour instructions are not overwritten
  - callq \*%rax is an instruction to jump to the address stored in the rax register
  - replaced callq \*%rax jumps to address 0~around 500

## How to redirect to the user-defined hook function?

- zpline instantiates trampoline code at address 0
  - fills address range 0 to  $N$  with nop (0x90)

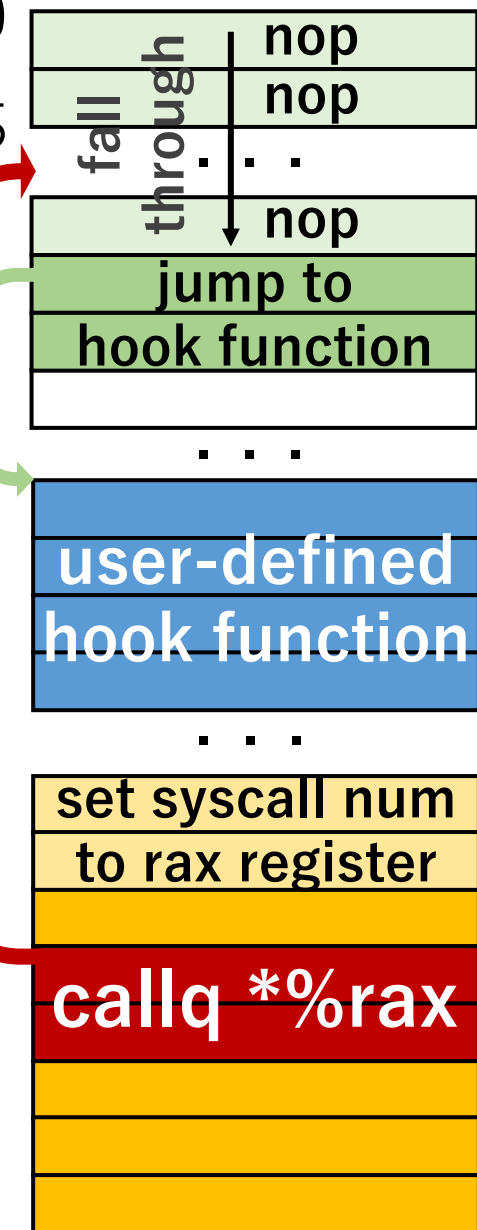


# zpoline

address range, potentially  
replaced “callq \*%rax” jumps to  
(  $N$  is the max syscall number )

0x0000  
0x0001  
0x0002  
...

Virtual Memory



- zpoline replaces syscall/sysenter with callq \*%rax
  - callq \*%rax is a **2-byte** instruction (0xff 0xd0)
    - Neighbour instructions are not overwritten
  - callq \*%rax is an instruction to jump to the address stored in the rax register
  - replaced callq \*%rax jumps to address 0~around 500

## How to redirect to the user-defined hook function?

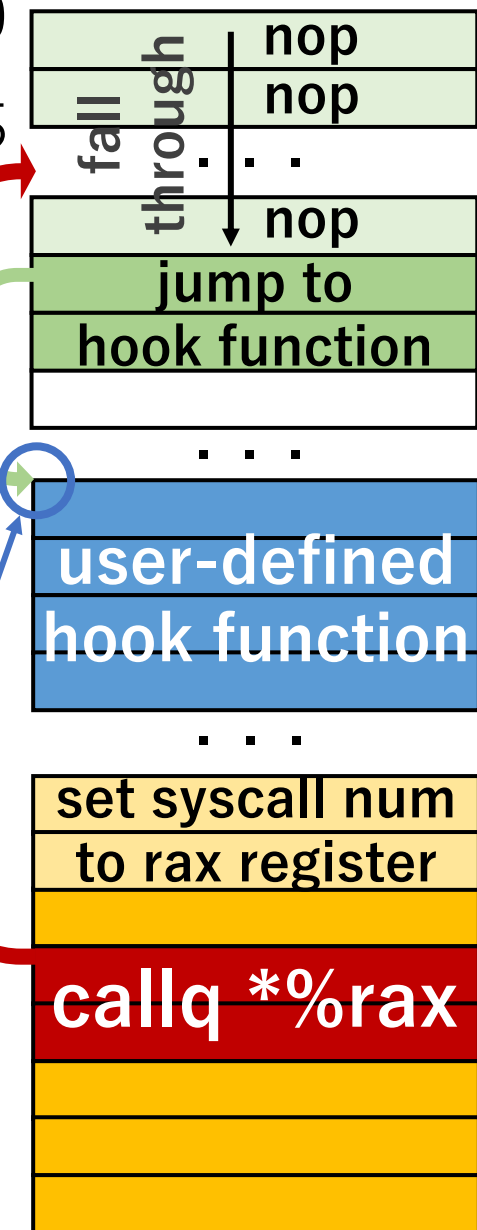
- zpoline instantiates trampoline code at address 0
  - fills address range 0 to  $N$  with nop (0x90)
    - puts code to jump to the hook function next to the last nop

# zpoline

address range, potentially  
replaced “callq \*%rax” jumps to  
(  $N$  is the max syscall number )

0x0000  
0x0001  
0x0002  
...

Virtual Memory



- zpoline replaces syscall/sysenter with callq \*%rax
  - callq \*%rax is a **2-byte** instruction (0xff 0xd0)
    - Neighbour instructions are not overwritten
  - callq \*%rax is an instruction to jump to the address stored in the rax register
  - replaced callq \*%rax jumps to address 0~around 500

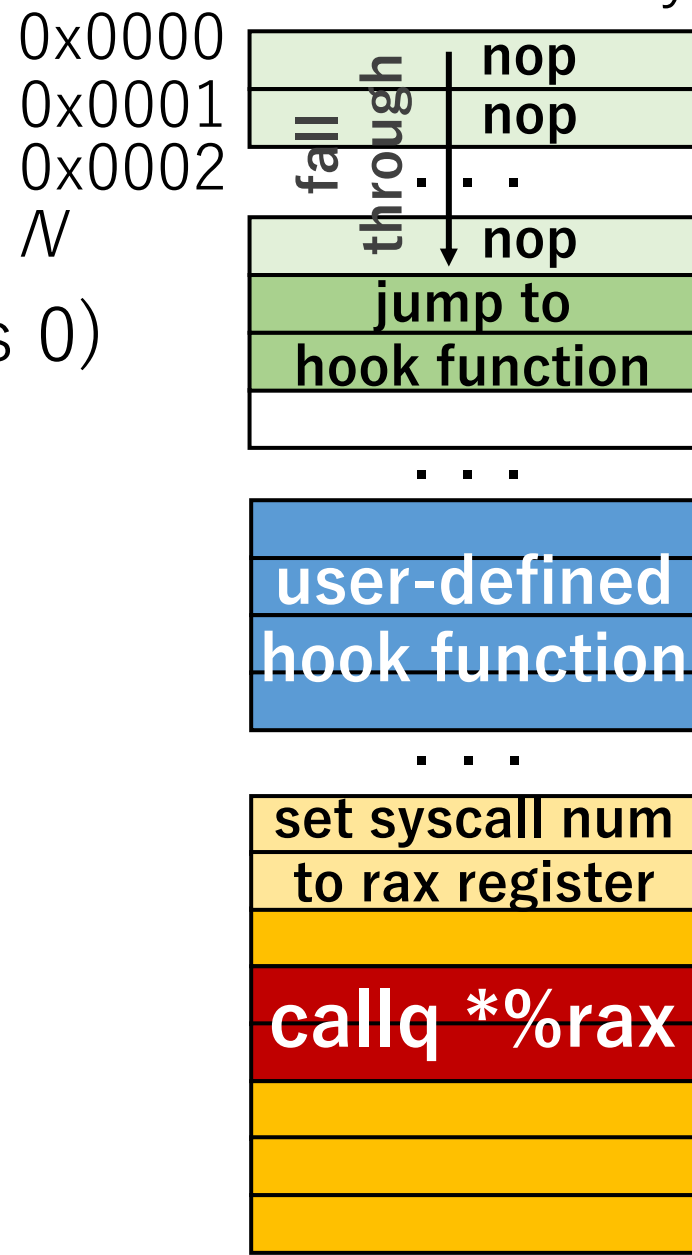
## How to redirect to the user-defined hook function?

- zpoline instantiates trampoline code at address 0
  - fills address range 0 to  $N$  with nop (0x90)
    - puts code to jump to the hook function next to the last nop

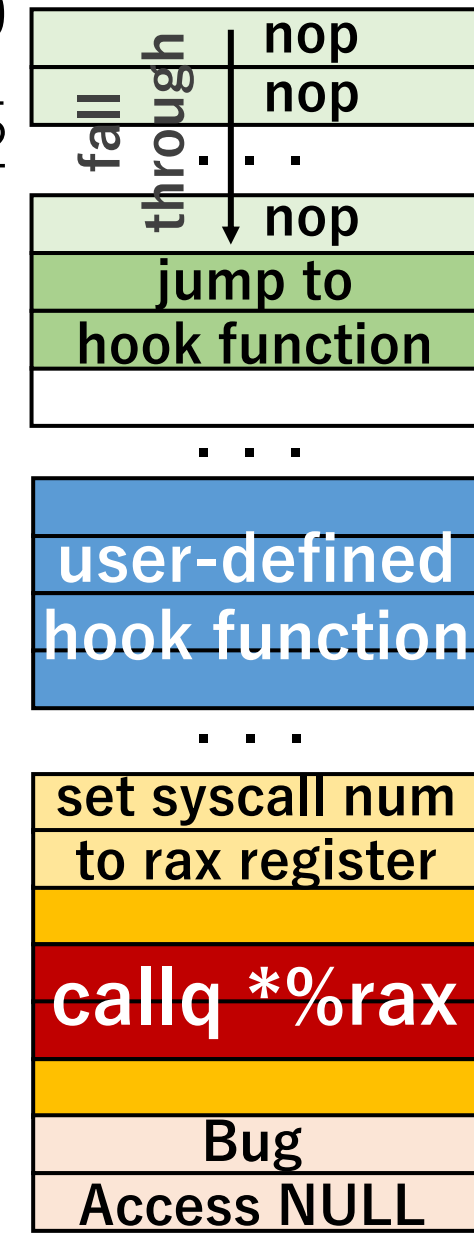
**We could reach the user-defined hook function!**

# NULL Access Termination

- A buggy program may access NULL (address 0)



# Virtual Memory

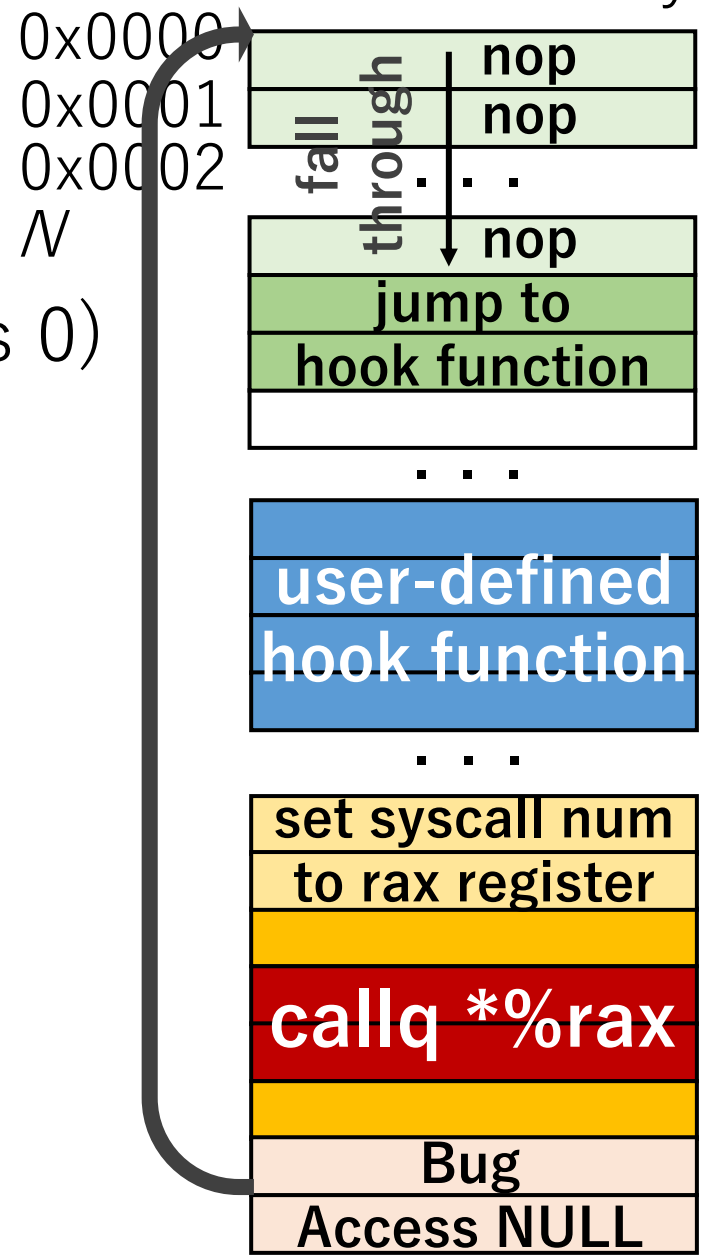


## NULL Access Termination

- A buggy program may access NULL (address 0)

# NULL Access Termination

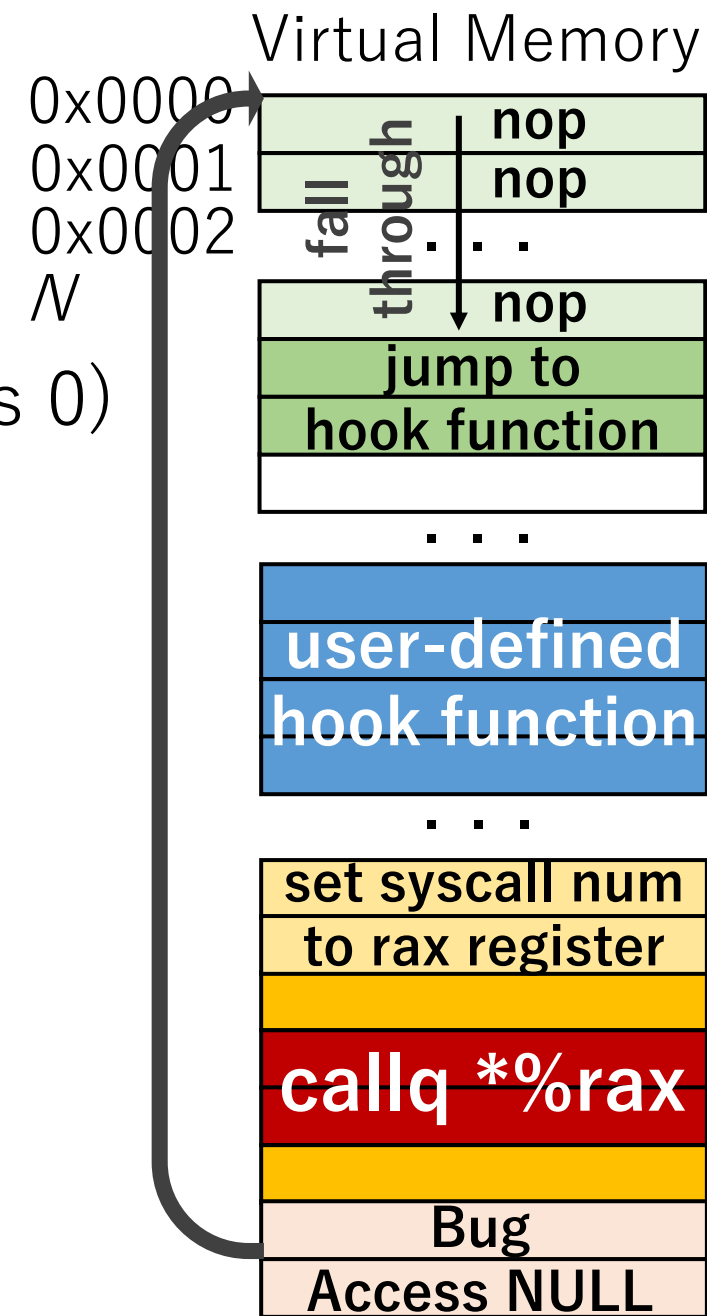
- A buggy program may access NULL (address 0)





# NULL Access Termination

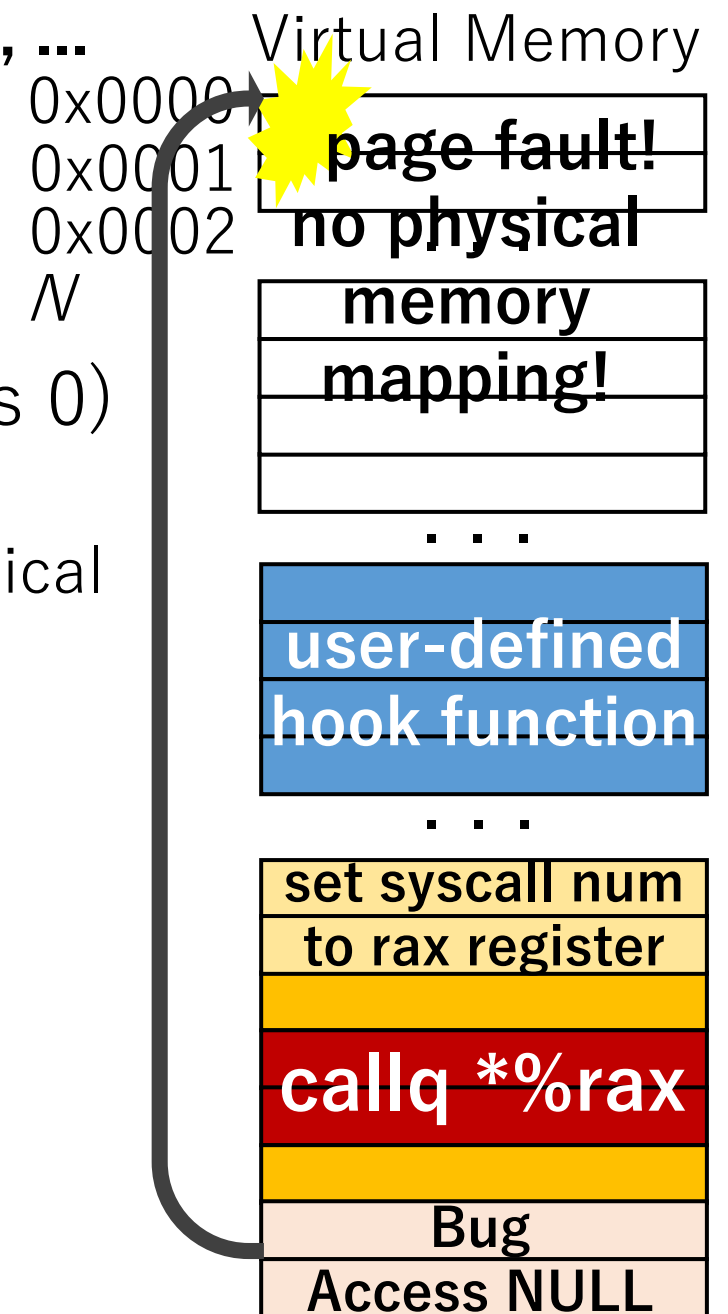
- A buggy program may access NULL (address 0)
  - In principle, NULL access has to be terminated



# NULL Access Termination

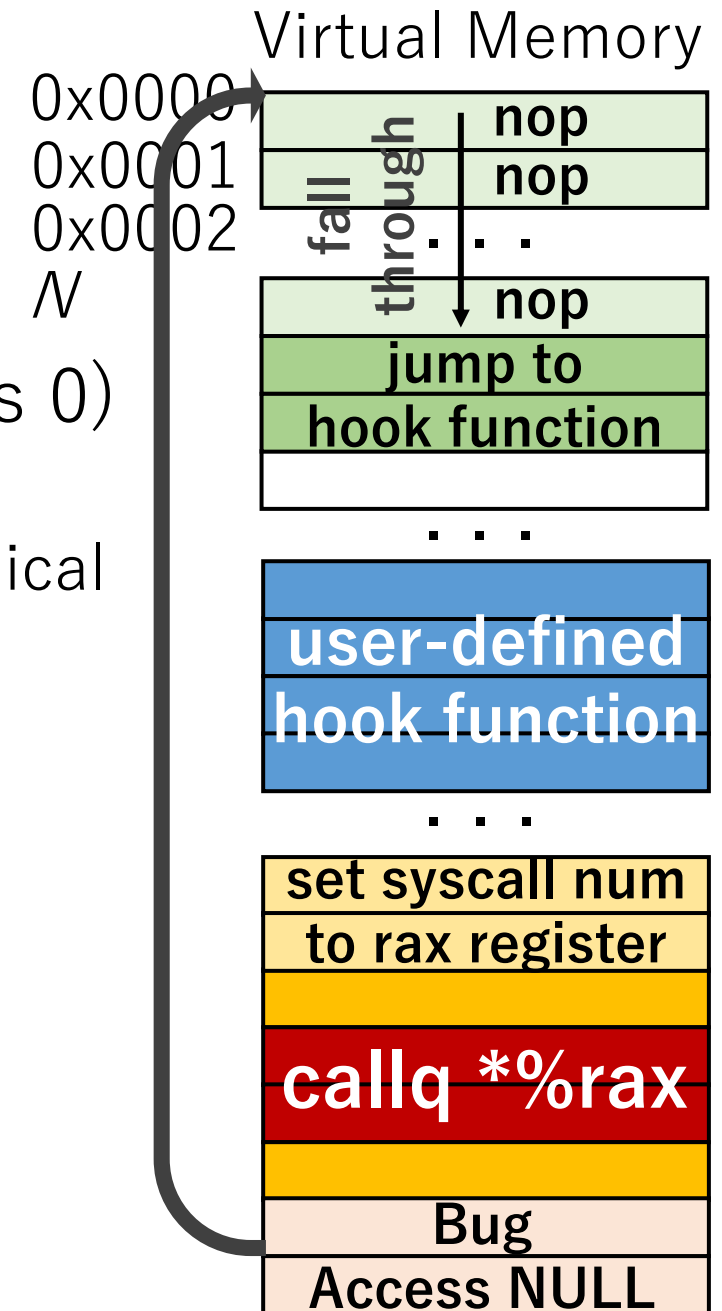
- A buggy program may access NULL (address 0)
  - In principle, NULL access has to be terminated
  - Normally, a page fault happens because no physical memory is mapped to virtual address 0

Normally, ...



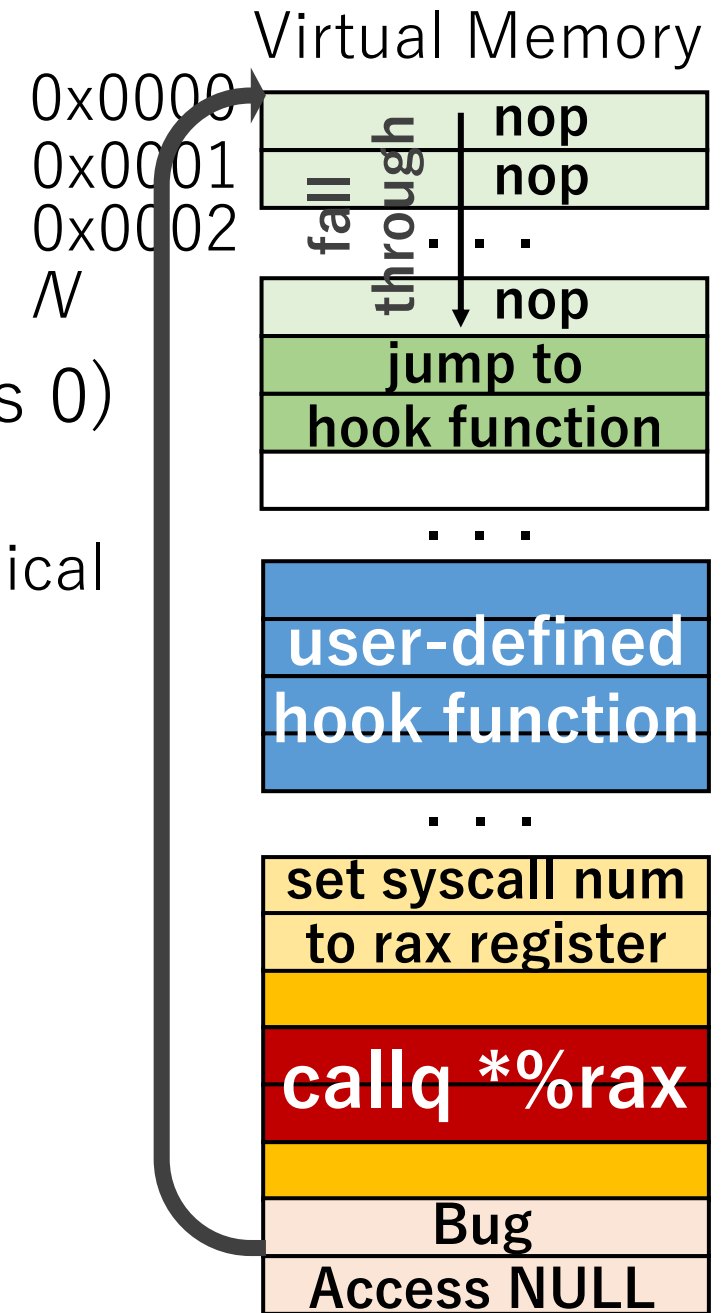
# NULL Access Termination

- A buggy program may access NULL (address 0)
  - In principle, NULL access has to be terminated
  - Normally, a page fault happens because no physical memory is mapped to virtual address 0
  - zpoline uses virtual address 0, therefore, the page fault does not happen



# NULL Access Termination

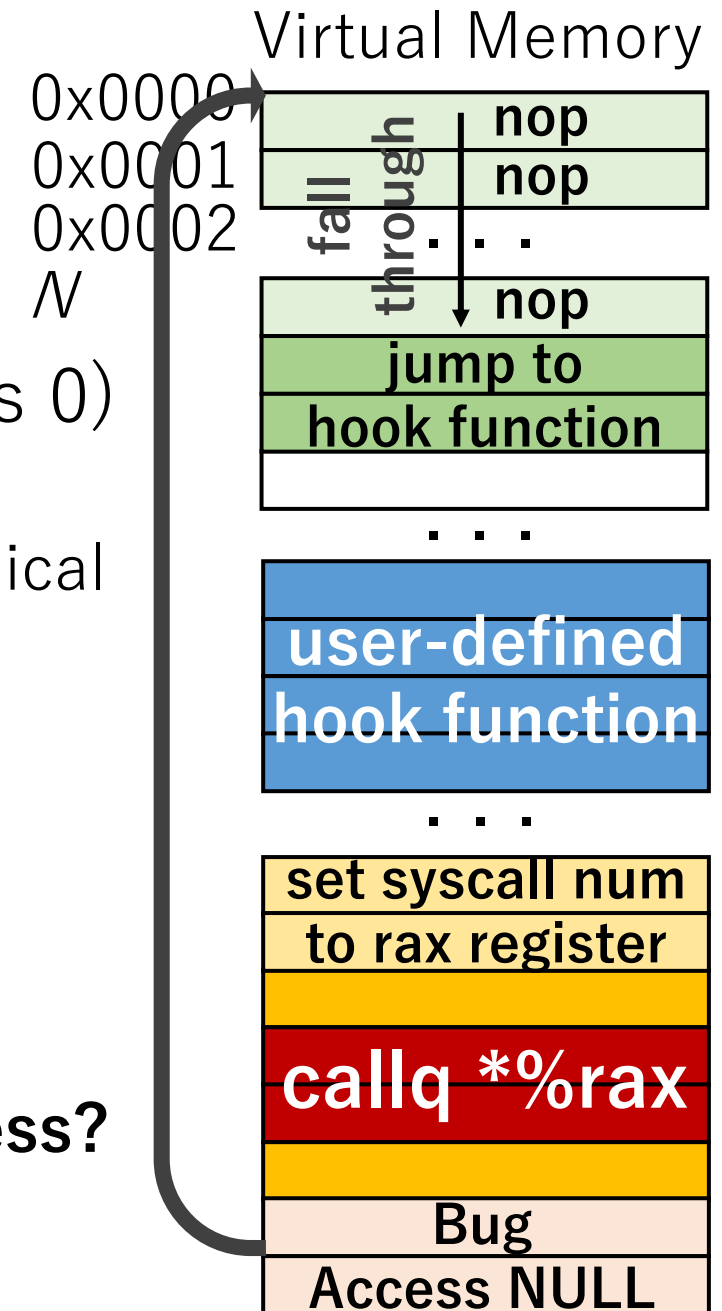
- A buggy program may access NULL (address 0)
    - In principle, NULL access has to be terminated
    - Normally, a page fault happens because no physical memory is mapped to virtual address 0
    - zpoline uses virtual address 0, therefore, the page fault does not happen
- ↓
- The buggy program continues to run



# NULL Access Termination

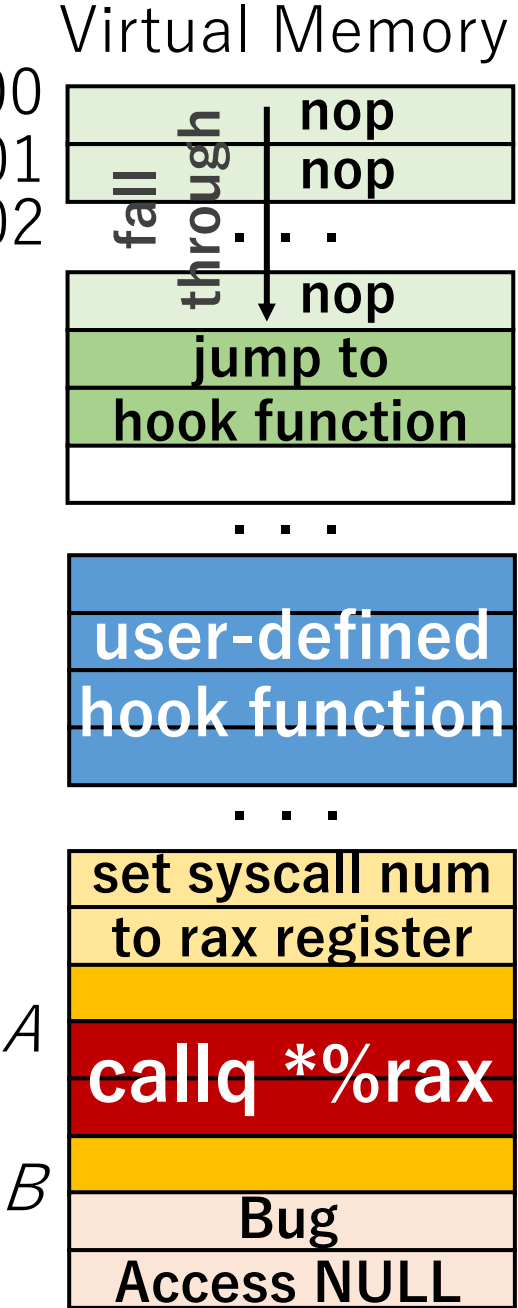
- A buggy program may access NULL (address 0)
    - In principle, NULL access has to be terminated
    - Normally, a page fault happens because no physical memory is mapped to virtual address 0
    - zpoline uses virtual address 0, therefore, the page fault does not happen
- ↓
- The buggy program continues to run

**How can we detect and terminate a buggy NULL access?**



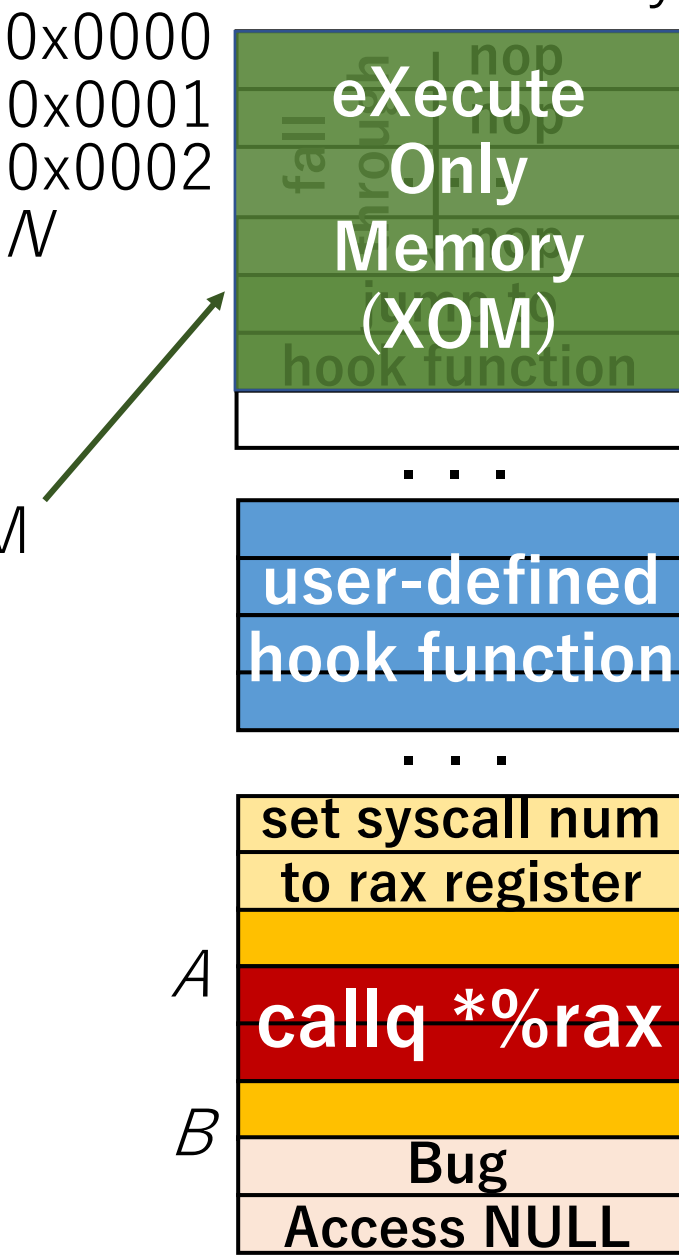
# NULL Access Termination

- Memory access: read / write / execute



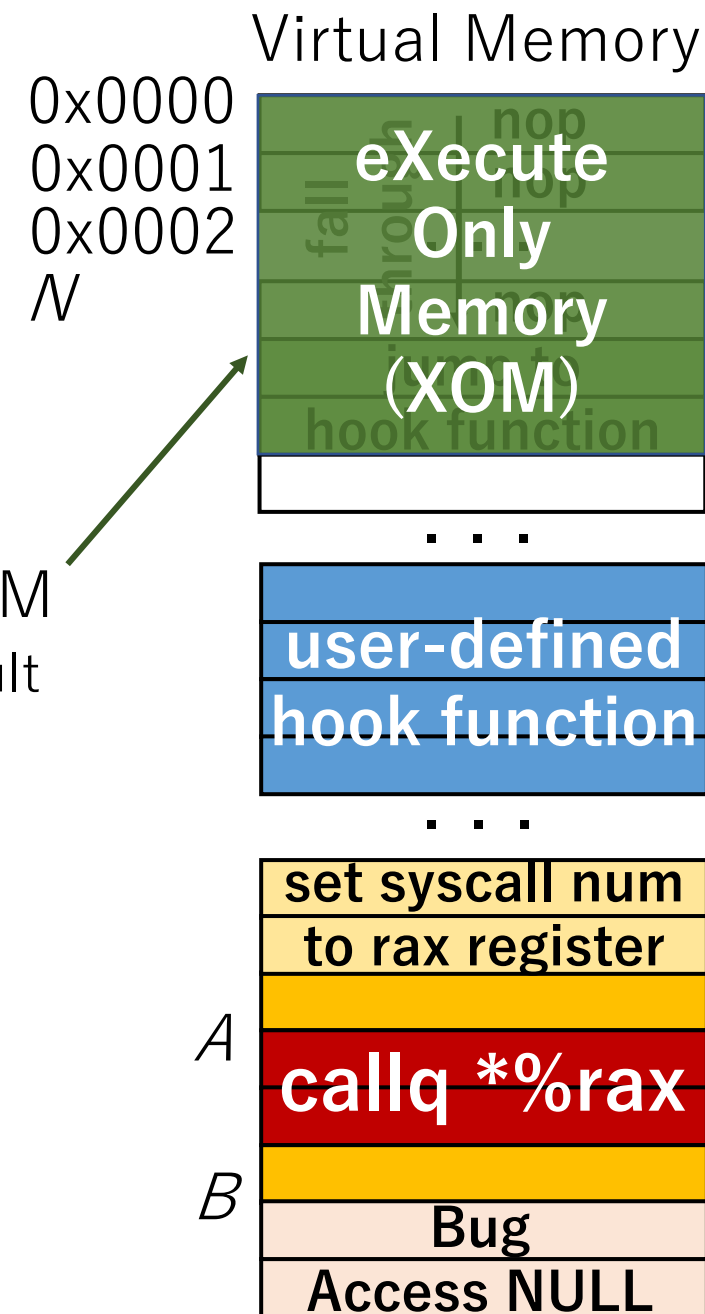
# NULL Access Termination

- Memory access: read / write / execute
- Solution
  - read/write: configure the trampoline code as XOM



# NULL Access Termination

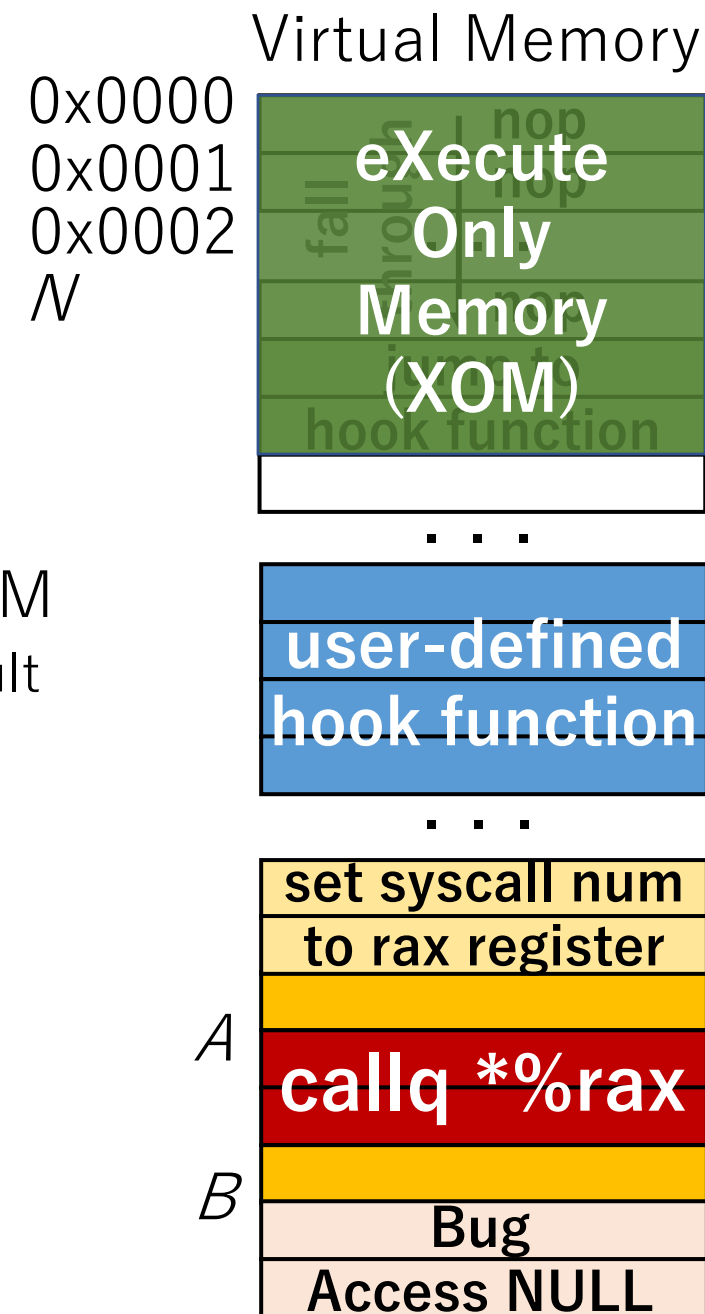
- Memory access: read / write / execute
- Solution
  - read/write: configure the trampoline code as XOM
    - read/write access to the trampoline code causes a fault
      - This can be done by mprotect() system call





# NULL Access Termination

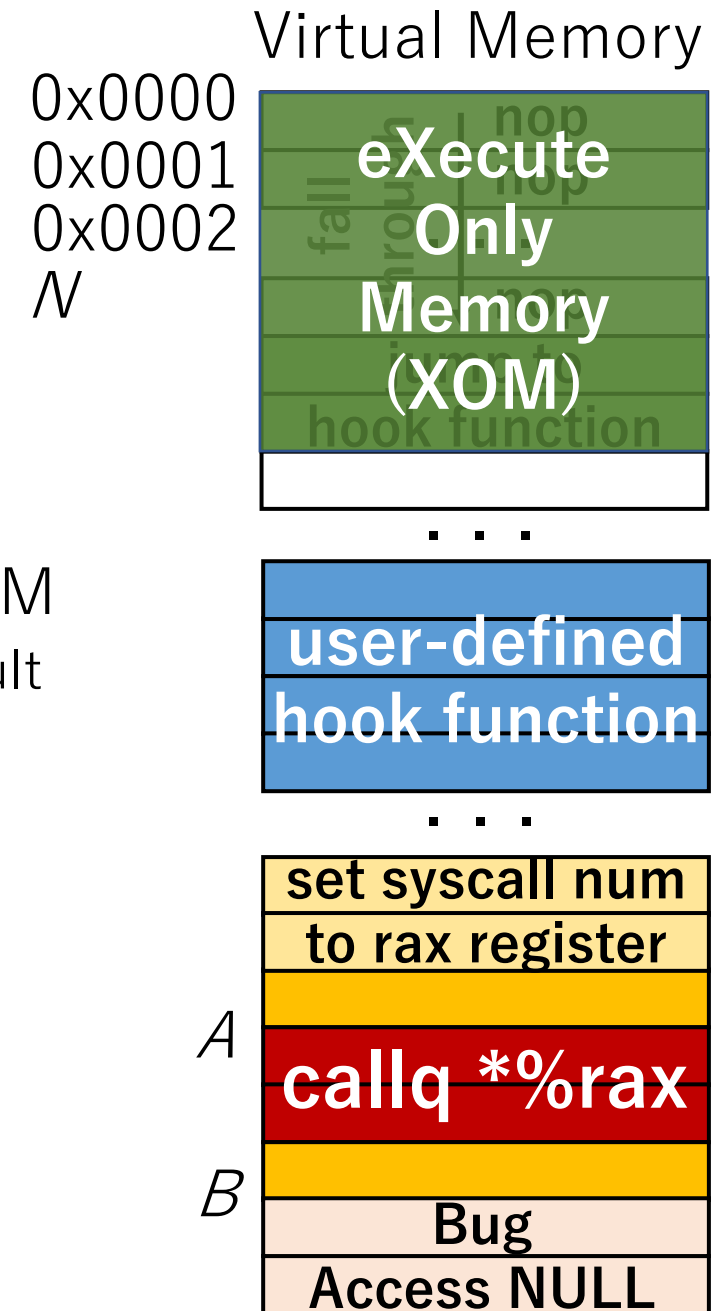
- Memory access: read / write / execute
- Solution
  - read/write: configure the trampoline code as XOM
    - read/write access to the trampoline code causes a fault
      - This can be done by mprotect() system call
  - execute: check the caller address



# NULL Access Termination

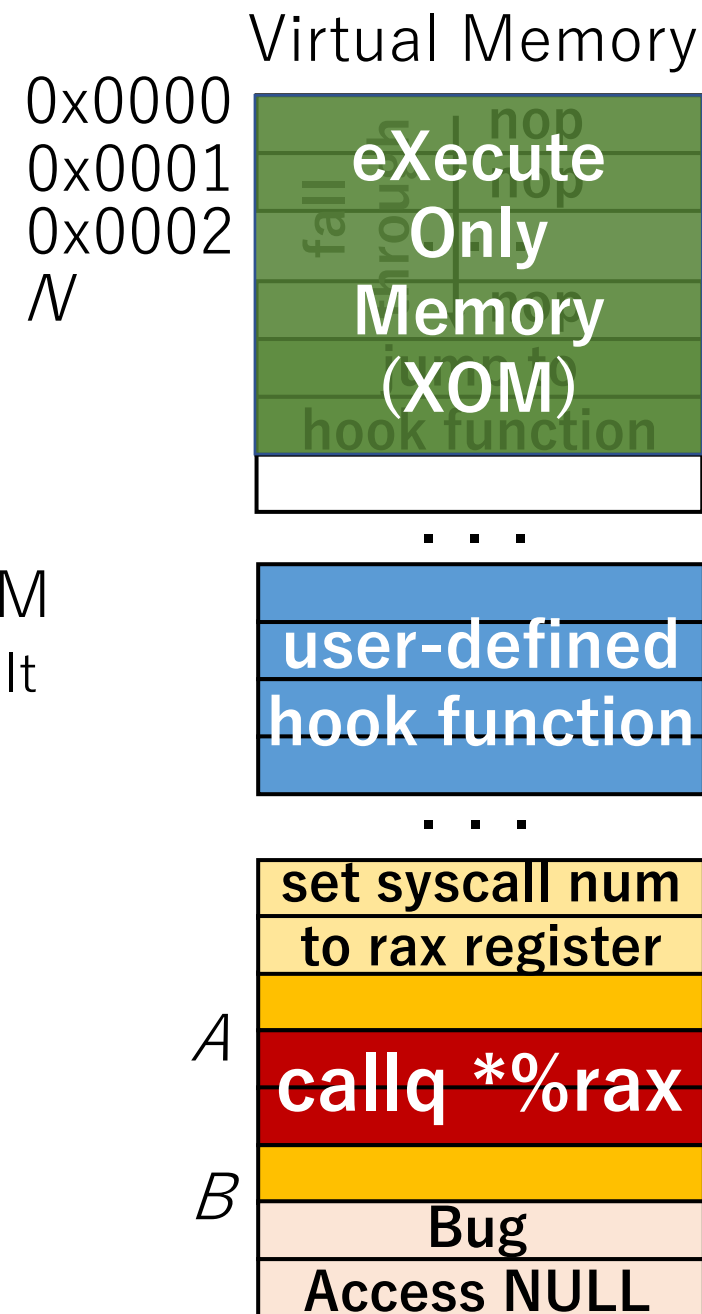
- Memory access: read / write / execute
- Solution
  - read/write: configure the trampoline code as XOM
    - read/write access to the trampoline code causes a fault
      - This can be done by mprotect() system call
  - execute: check the caller address

List of replaced addresses : [...]



# NULL Access Termination

- Memory access: read / write / execute
- Solution
  - read/write: configure the trampoline code as XOM
    - read/write access to the trampoline code causes a fault
      - This can be done by mprotect() system call
  - execute: check the caller address
    1. collect the addresses of replaced syscall/sysenter during the binary rewriting phase

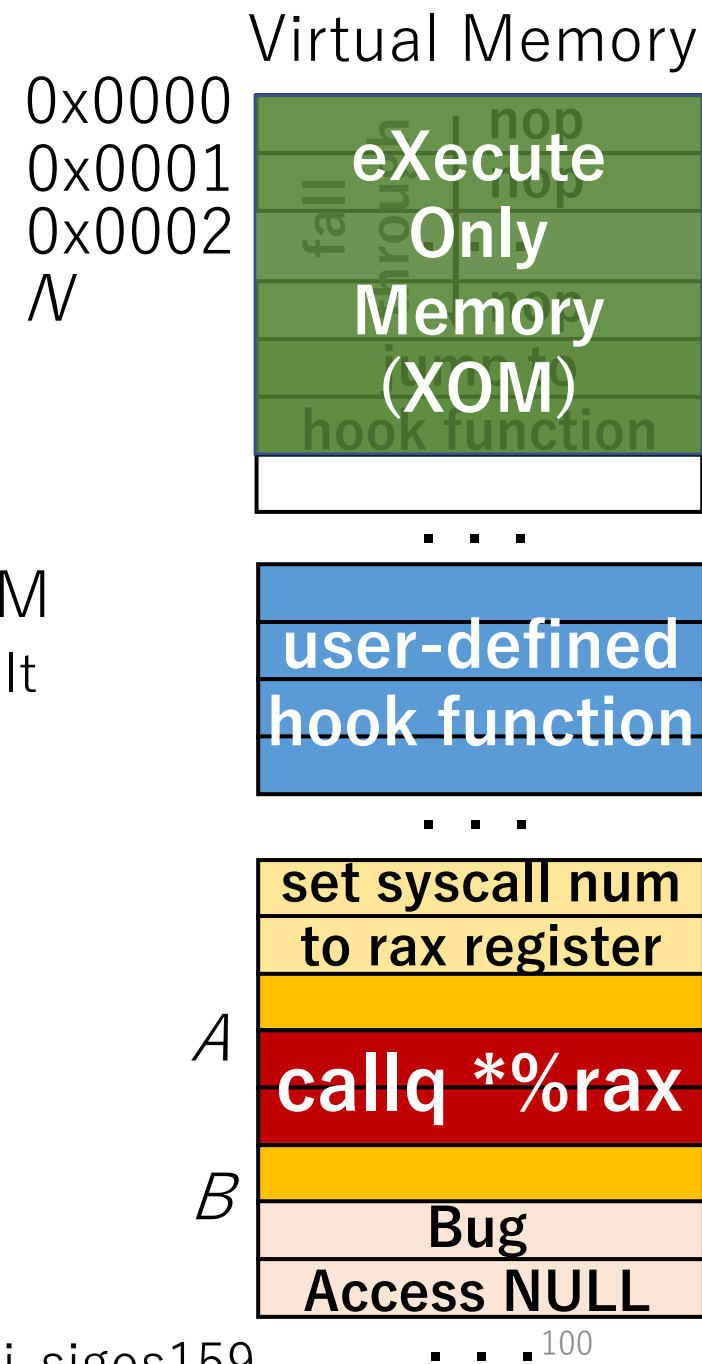


# NULL Access Termination

- Memory access: read / write / execute
- Solution
  - read/write: configure the trampoline code as XOM
    - read/write access to the trampoline code causes a fault
      - This can be done by mprotect() system call
  - execute: check the caller address
    1. collect the addresses of replaced syscall/sysenter during the binary rewriting phase

## During binary rewriting phase ...

List of replaced addresses : [...]

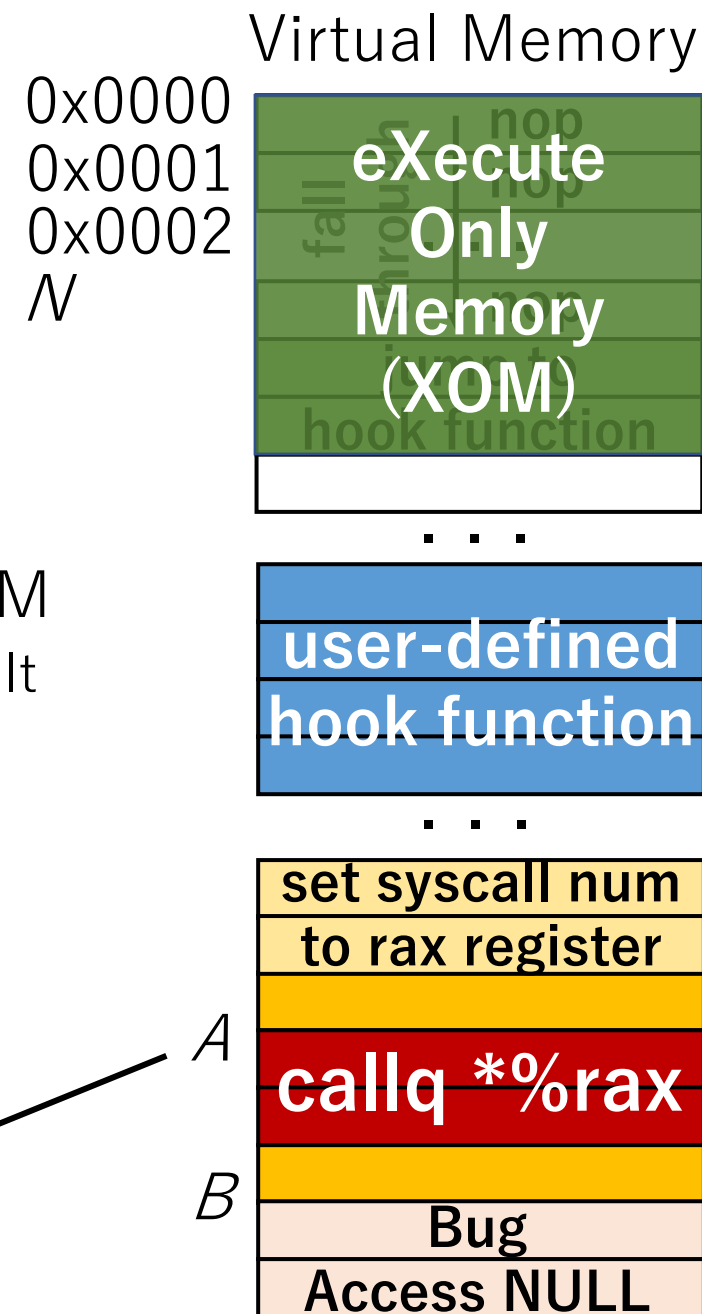


# NULL Access Termination

- Memory access: read / write / execute
- Solution
  - read/write: configure the trampoline code as XOM
    - read/write access to the trampoline code causes a fault
      - This can be done by mprotect() system call
  - execute: check the caller address
    1. collect the addresses of replaced syscall/sysenter during the binary rewriting phase

**During binary rewriting phase ...**

List of replaced addresses : [...]

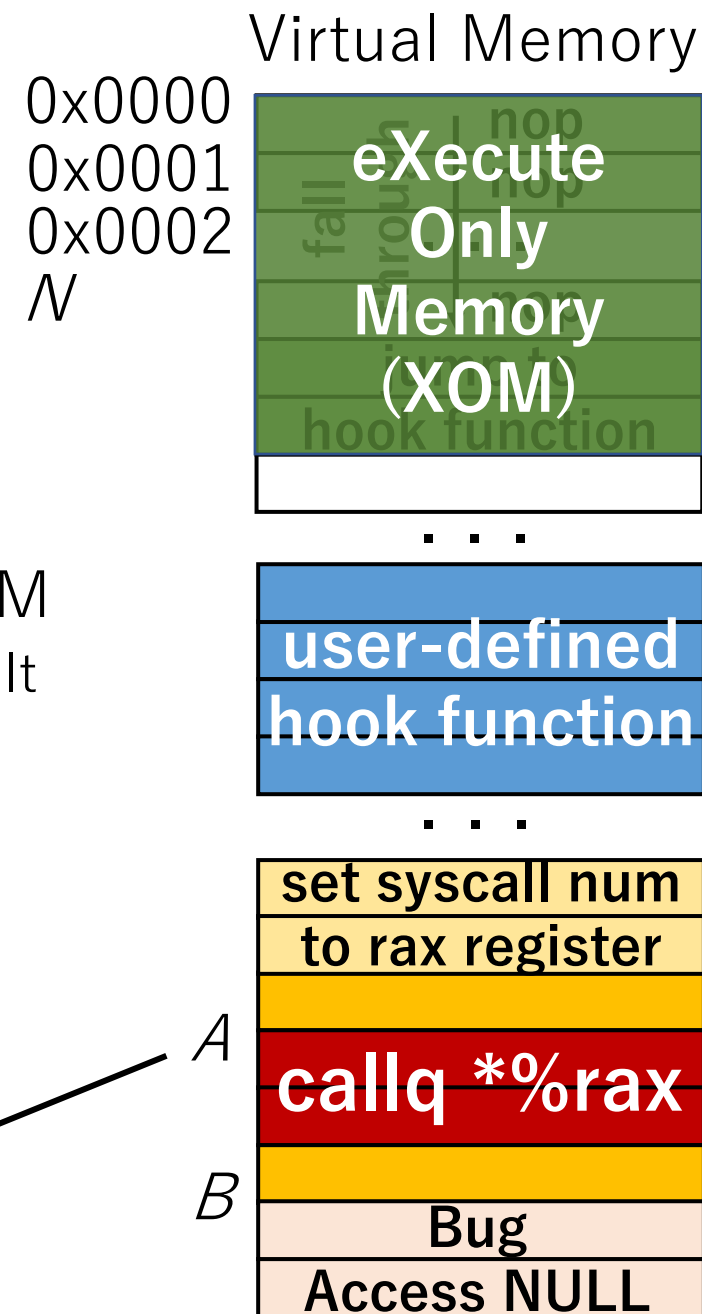


# NULL Access Termination

- Memory access: read / write / execute
- Solution
  - read/write: configure the trampoline code as XOM
    - read/write access to the trampoline code causes a fault
      - This can be done by mprotect() system call
  - execute: check the caller address
    1. collect the addresses of replaced syscall/sysenter during the binary rewriting phase

**During binary rewriting phase ...**

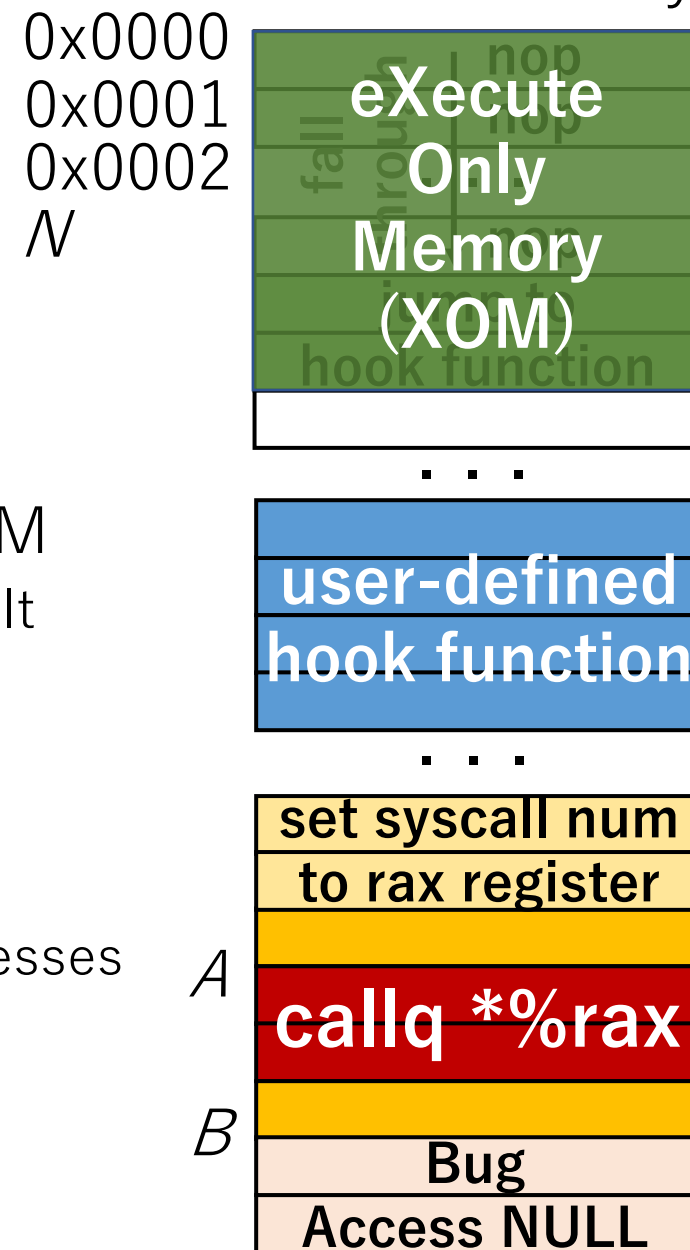
List of replaced addresses : [A , ...]



# NULL Access Termination

- Memory access: read / write / execute
- Solution
  - read/write: configure the trampoline code as XOM
    - read/write access to the trampoline code causes a fault
      - This can be done by mprotect() system call
  - execute: check the caller address
    1. collect the addresses of replaced syscall/sysenter during the binary rewriting phase
    2. check the caller address is one of the replaced addresses in the hook function

List of replaced addresses : [A , ...]

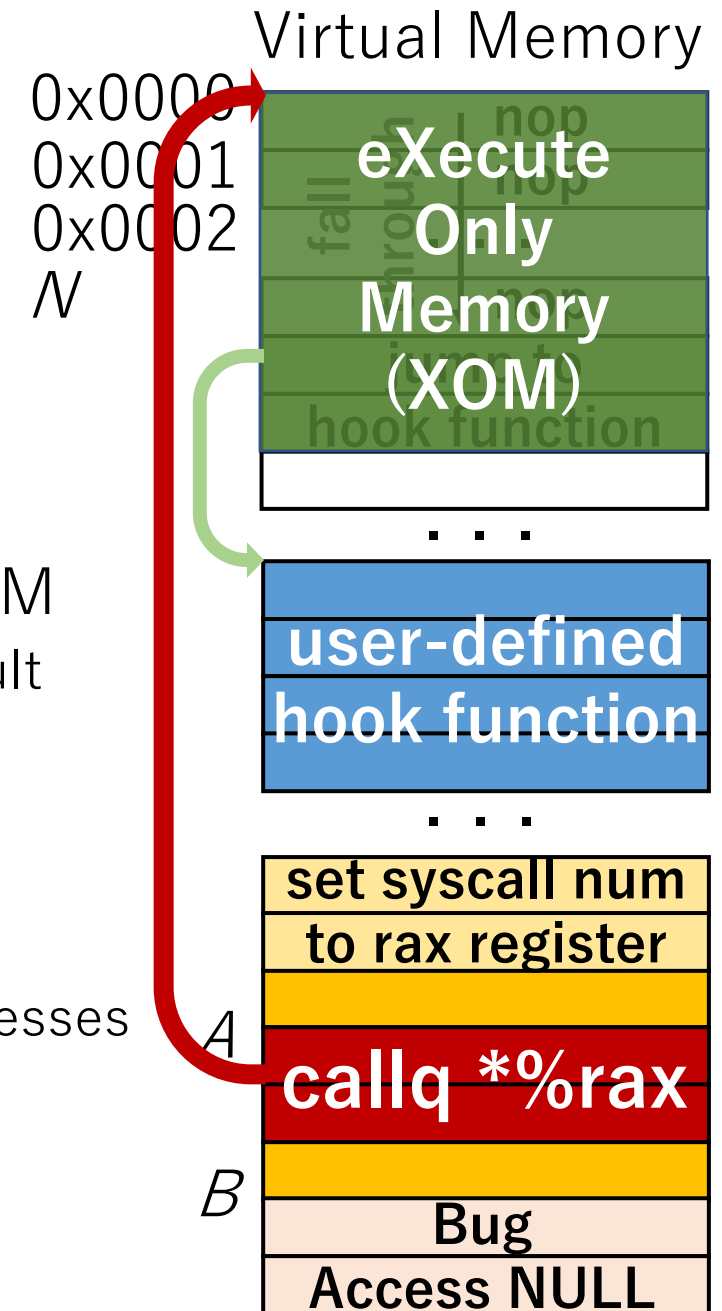


# NULL Access Termination

## At runtime ...

- Memory access: read / write / execute
- Solution
  - read/write: configure the trampoline code as XOM
    - read/write access to the trampoline code causes a fault
      - This can be done by mprotect() system call
  - execute: check the caller address
    1. collect the addresses of replaced syscall/sysenter during the binary rewriting phase
    2. check the caller address is one of the replaced addresses in the hook function

List of replaced addresses : [A , ...]





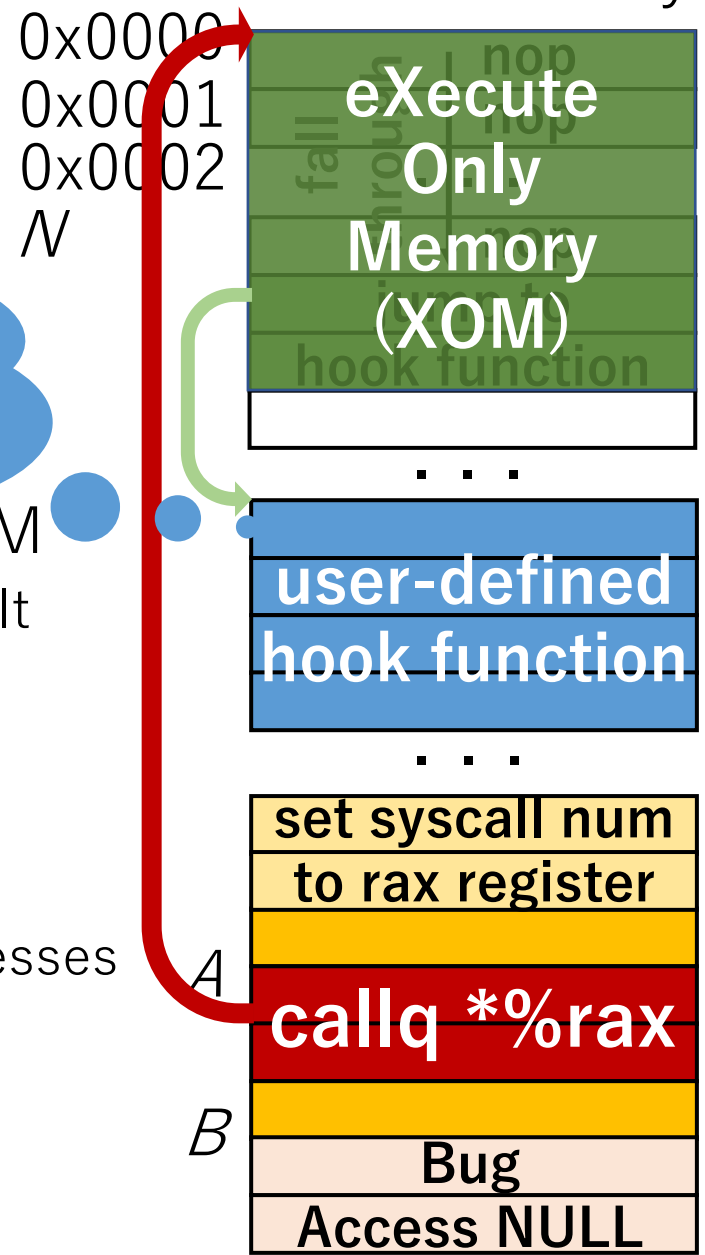
# NULL Access Termination

## At runtime ...

- Memory access
- Solution
  - read/write:
    - read/write access to the trap code causes a fault
      - This can be done by mprotect() system call
  - execute: check the caller address
    1. collect the addresses of replaced syscall/sysenter during the binary rewriting phase
    2. check the caller address is one of the replaced addresses in the hook function

The caller address is *A*  
*A* is in the list, so  
this is a valid access

List of replaced addresses : [*A* , ...]

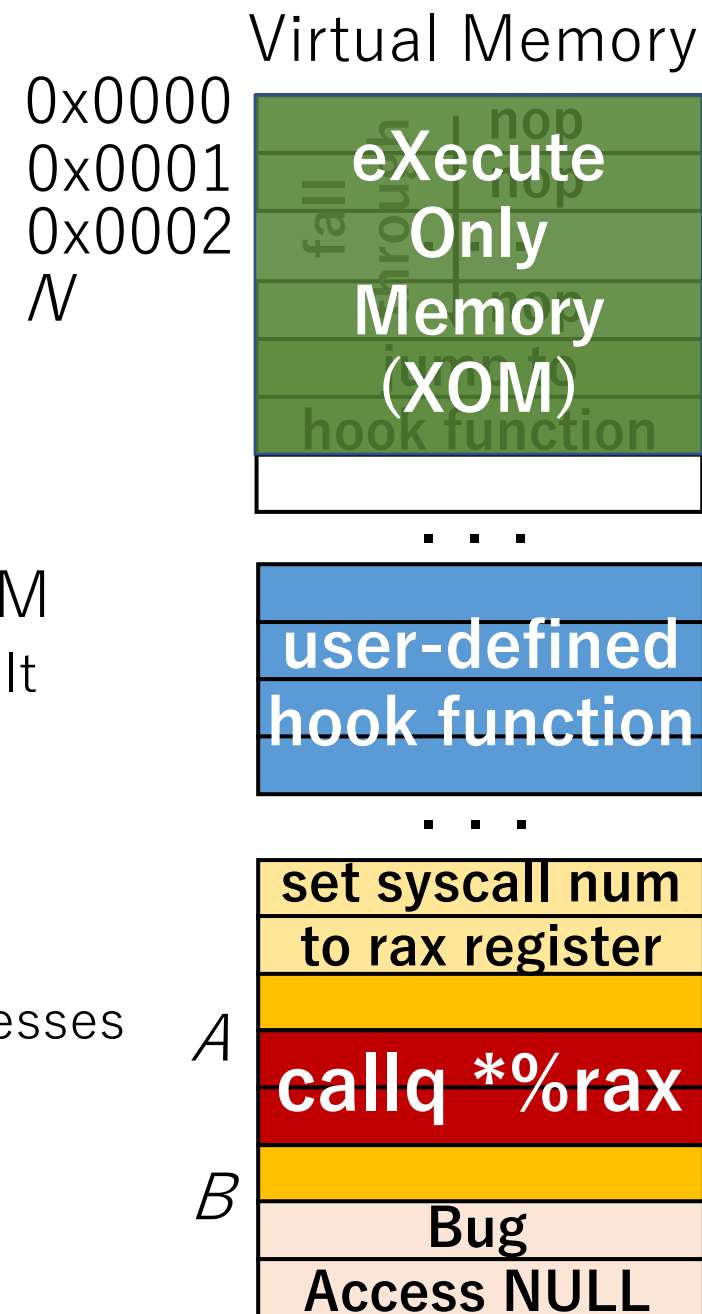


# NULL Access Termination

## At runtime ...

- Memory access: read / write / execute
- Solution
  - read/write: configure the trampoline code as XOM
    - read/write access to the trampoline code causes a fault
      - This can be done by mprotect() system call
  - execute: check the caller address
    1. collect the addresses of replaced syscall/sysenter during the binary rewriting phase
    2. check the caller address is one of the replaced addresses in the hook function

List of replaced addresses : [A , ...]

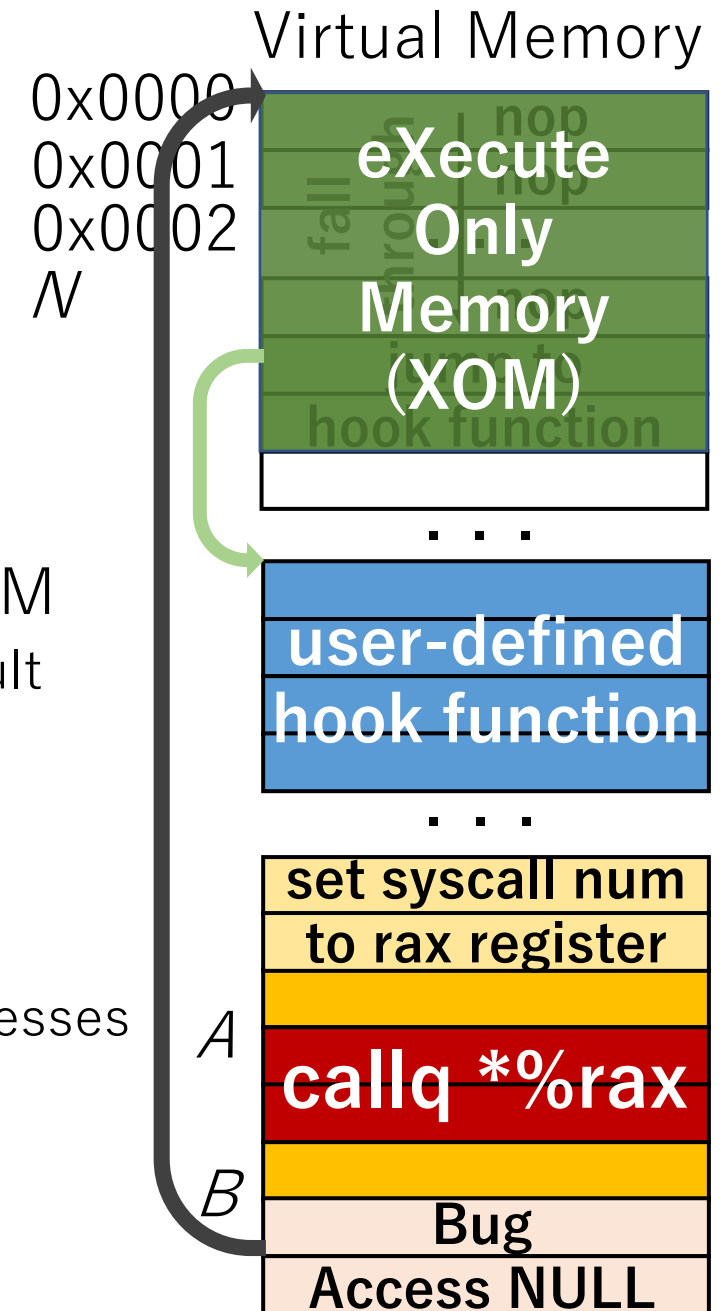


# NULL Access Termination

## At runtime ...

- Memory access: read / write / execute
- Solution
  - read/write: configure the trampoline code as XOM
    - read/write access to the trampoline code causes a fault
      - This can be done by mprotect() system call
  - execute: check the caller address
    1. collect the addresses of replaced syscall/sysenter during the binary rewriting phase
    2. check the caller address is one of the replaced addresses in the hook function

List of replaced addresses : [A , ...]

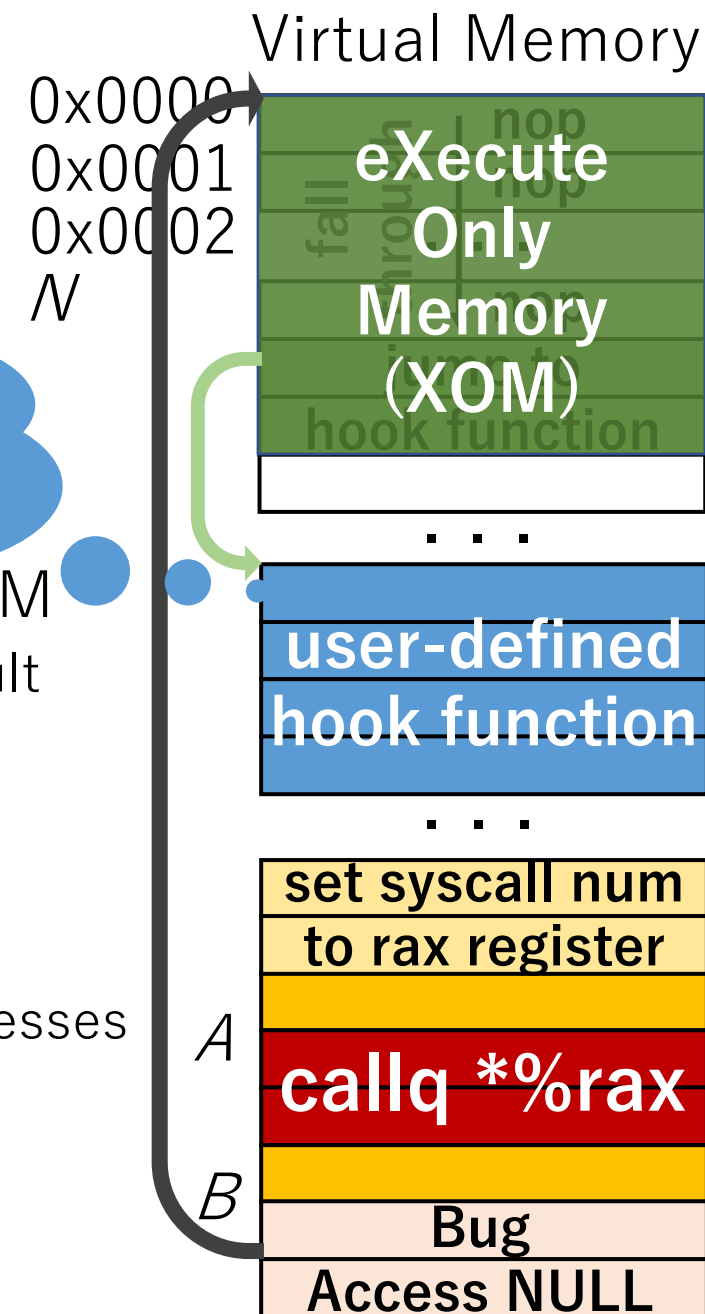


# NULL Access Termination

## At runtime ...

- Memory access
- Solution
  - read/write:
    - read/write access to the trap code causes a fault
      - This can be done by mprotect() system call
  - execute: check the caller address
    1. collect the addresses of replaced syscall/sysenter during the binary rewriting phase
    2. check the caller address is one of the replaced addresses in the hook function

List of replaced addresses : [A, ...]



# NULL Access Termination

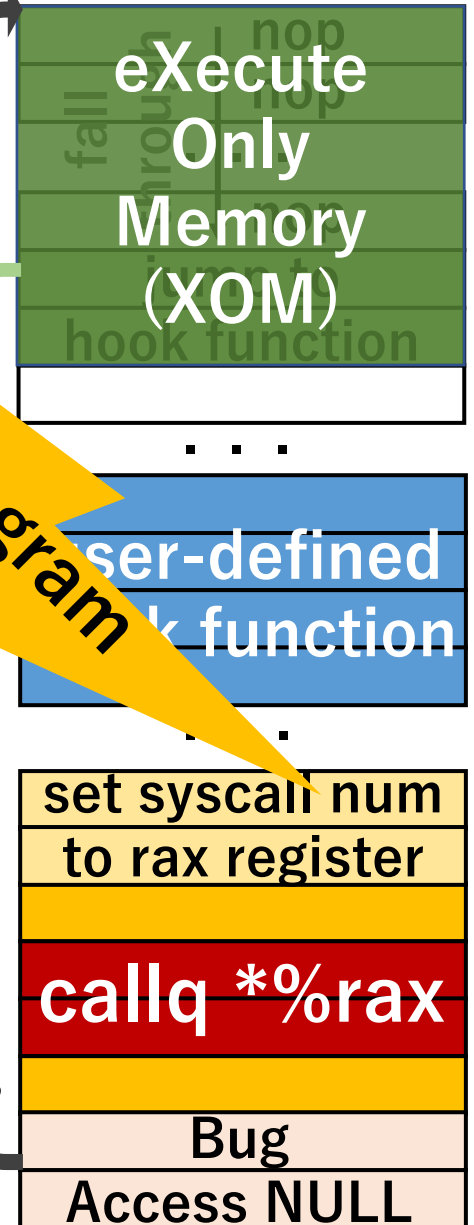
## At runtime ...

- Memory access
- Solution
  - read/write.
    - read/write access to the trap code causes a fault
      - This can be done by mprotect() system call
  - execute: check the caller address
    1. collect the addresses of replaced syscall/sysenter during the binary rewriting phase
    2. check the caller address is one of the replaced addresses in the hook function

List of replaced addresses : [A, ...]

The caller address is *B*  
*B* is NOT in the list, so  
 this is an invalid access

stop the program



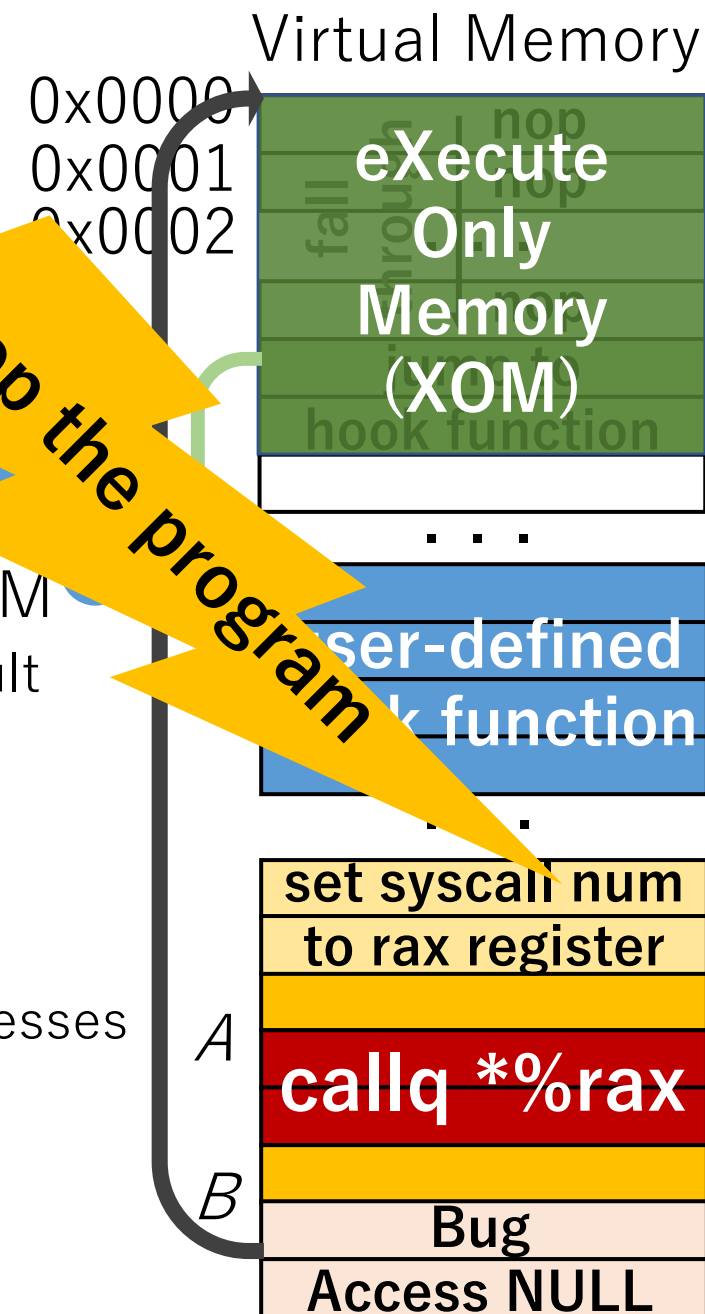
# NULL Access Termination

## At runtime ...

- Memory access
- Solution
  - read/write:
    - read/write access to the trap code causes a fault
      - This can be done by mprotect() system call
  - execute: check the caller address
    1. collect the addresses of replaced syscall/sysenter during the binary rewriting phase
    2. check the caller address is one of the replaced addresses in the hook function
      - Current prototype uses bitmap to implement this check

List of replaced addresses : [A , ...]

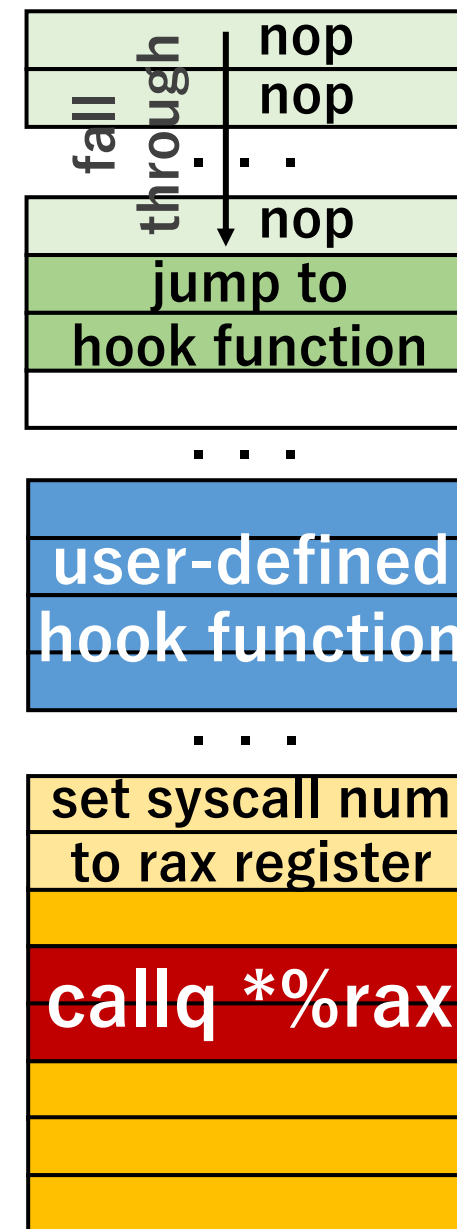
stop the program



# System Call Hook Overhead

- Time to hook getpid() and return a dummy value

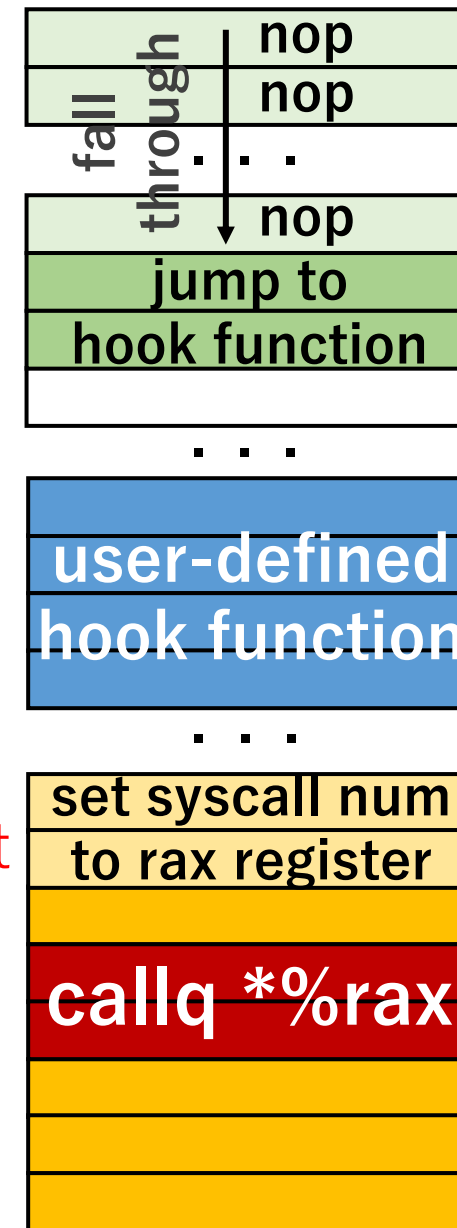
Mechanism	Time [ns]
ptrace	31201
int3 signaling	1342
SUD	1156
zpoline	41
LD_PRELOAD	6



# System Call Hook Overhead

- Time to hook getpid() and return a dummy value

Mechanism	Time [ns]	
ptrace	31201	716x
int3 signaling	1342	32.7x
SUD	1156	28.1x
zpoline	41	improvement
LD_PRELOAD	6	

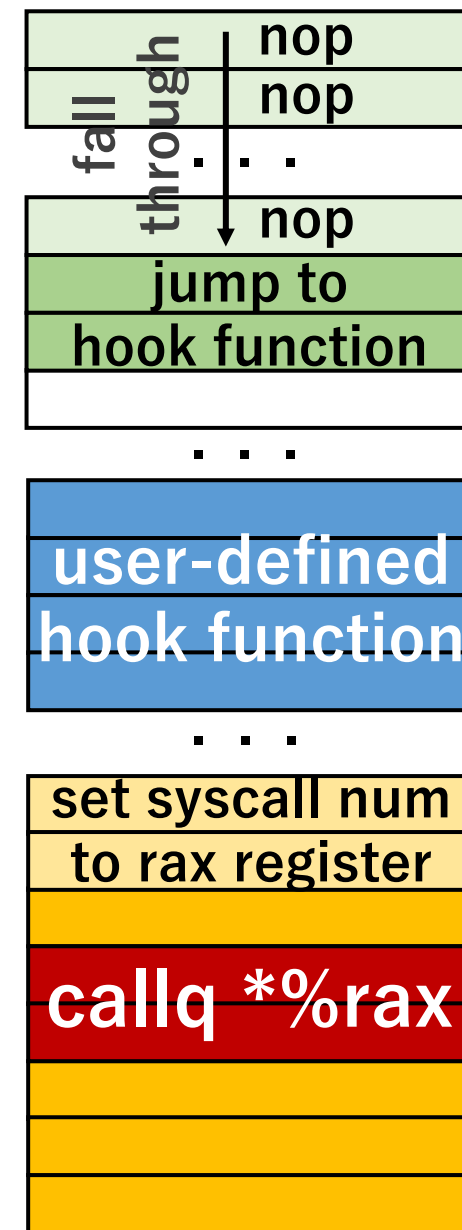




# System Call Hook Overhead

- Time to hook getpid() and return a dummy value

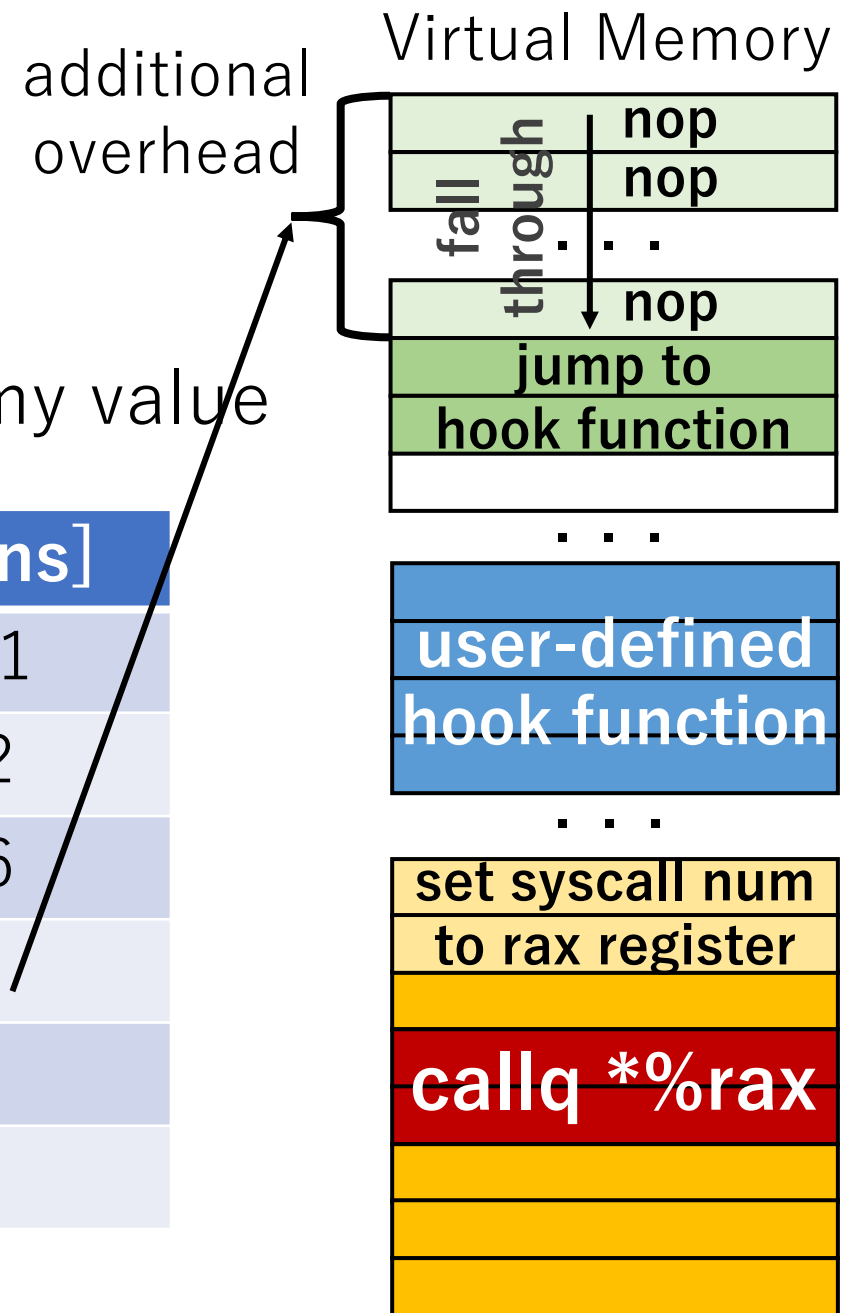
Mechanism	Time [ns]
ptrace	31201
int3 signaling	1342
SUD	1156
zpoline	41
zpoline (w/o NULL exec check)	40
LD_PRELOAD	6



# System Call Hook Overhead

- Time to hook getpid() and return a dummy value

Mechanism	Time [ns]
ptrace	31201
int3 signaling	1342
SUD	1156
zpoline	41
zpoline (w/o NULL exec check)	40
LD_PRELOAD	6



# Application Performance

- We **transparently** apply lwIP + DPDK to an application using different system call hook mechanisms

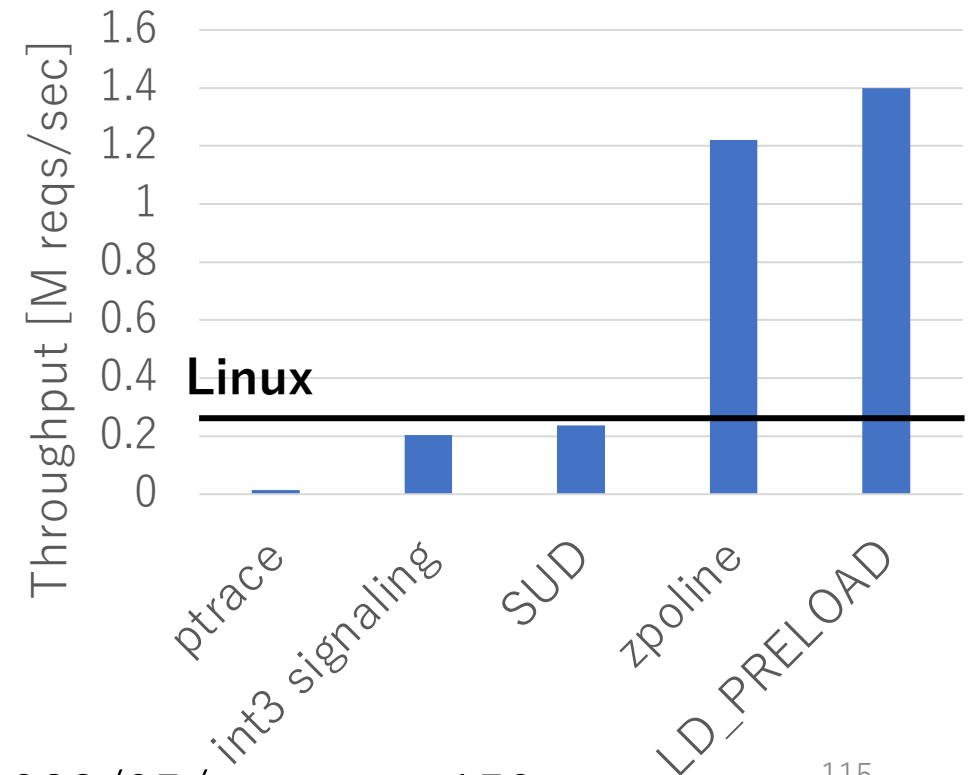
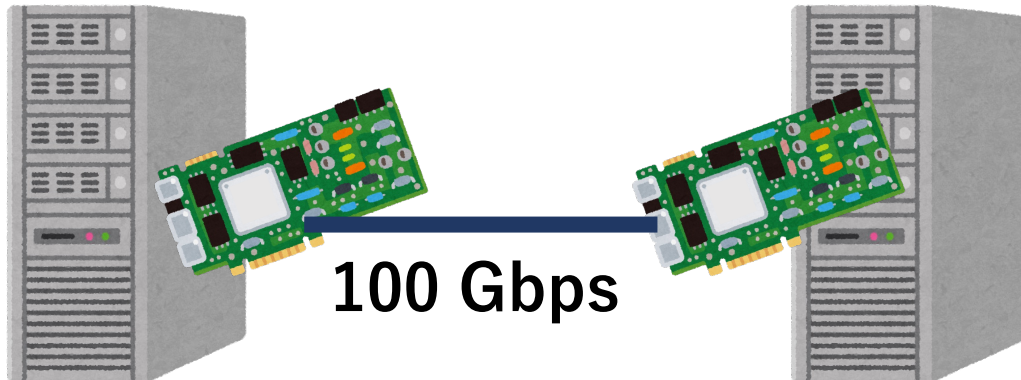
**Simple HTTP server**

*syscall hook*

lwIP + DPDK

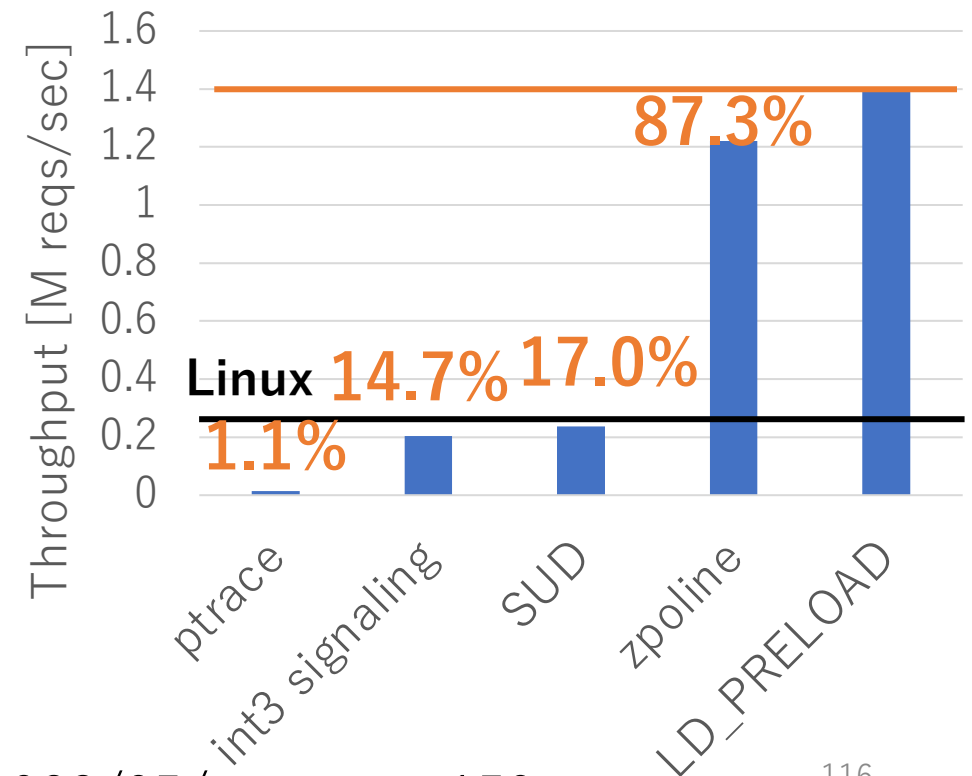
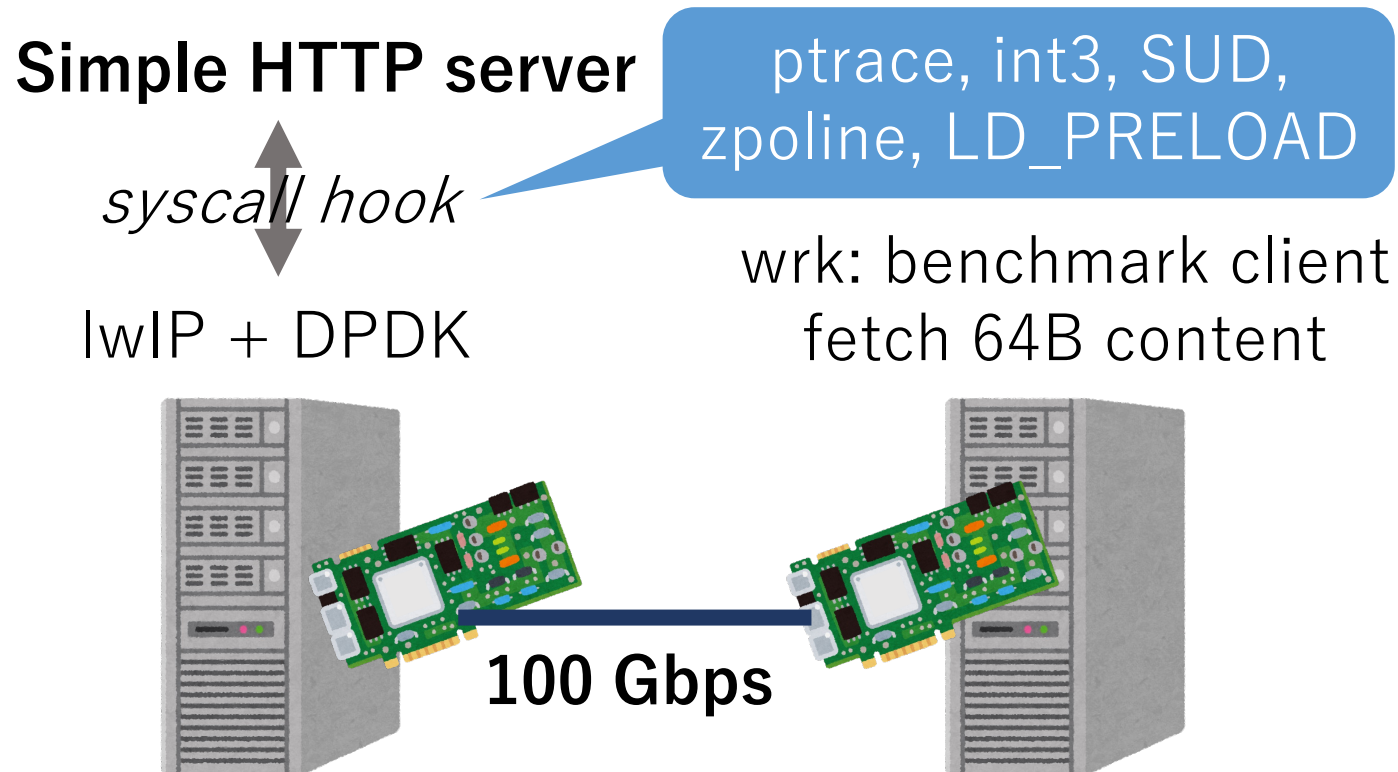
ptrace, int3, SUD,  
zpoline, LD\_PRELOAD

wrk: benchmark client  
fetch 64B content



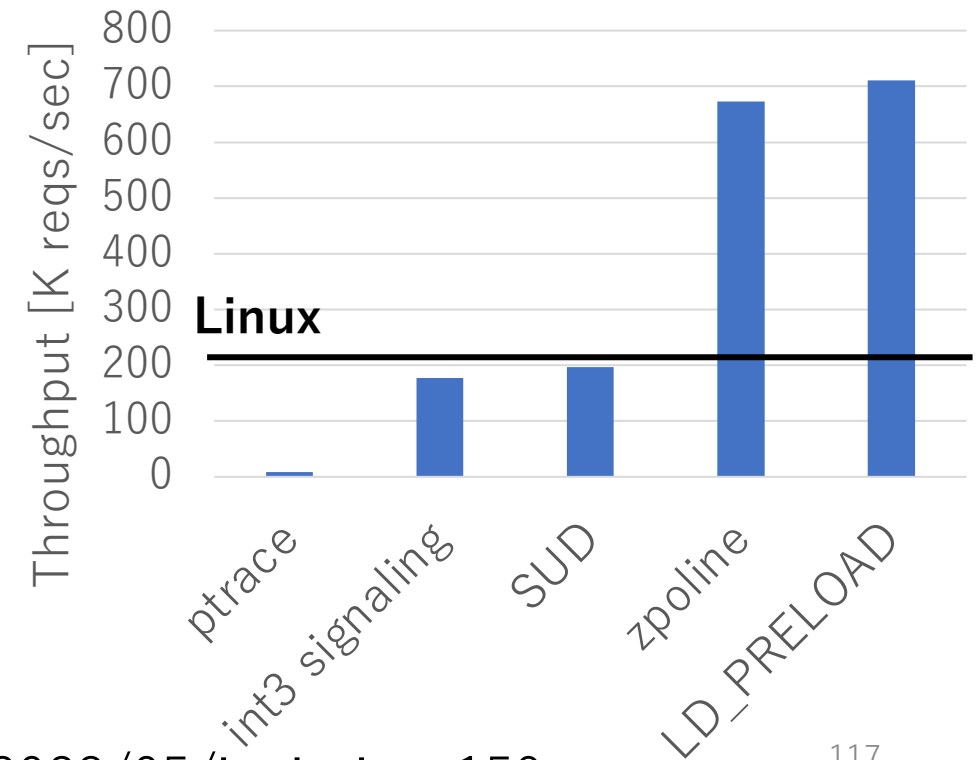
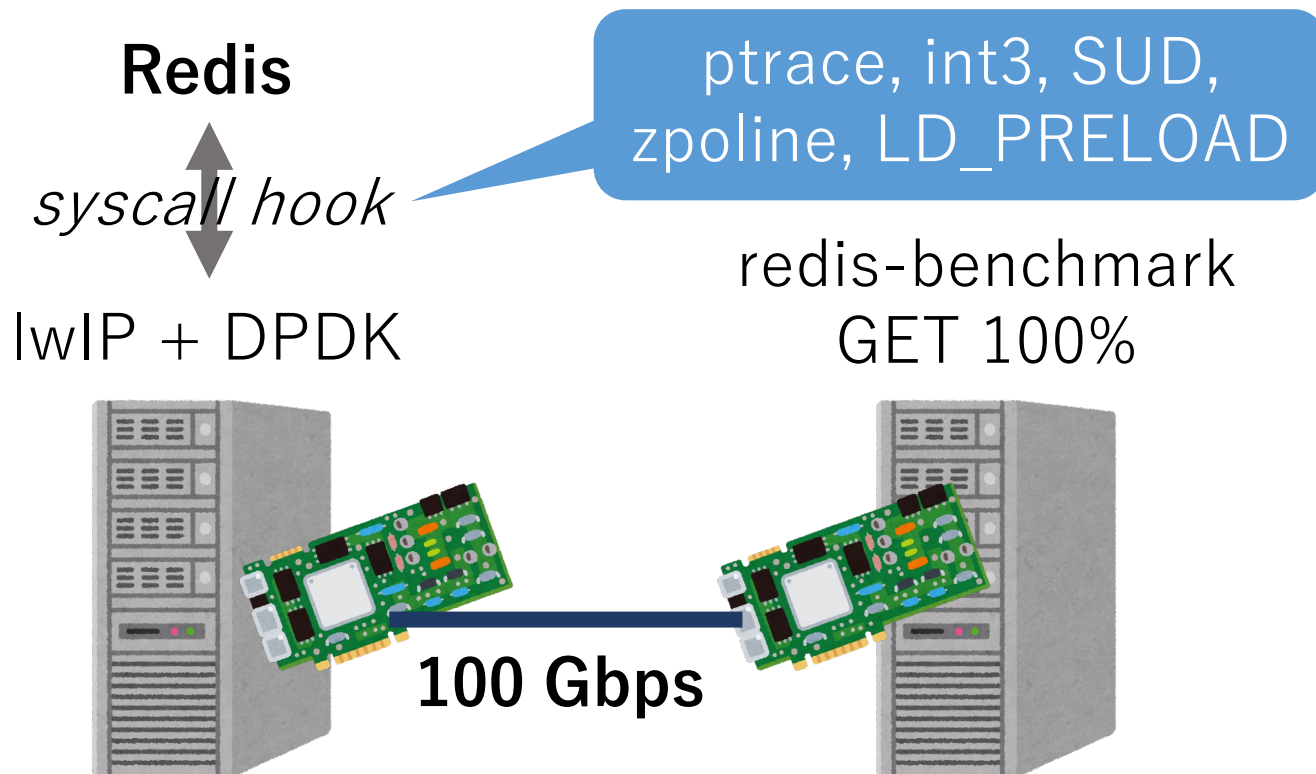
# Application Performance

- We **transparently** apply lwIP + DPDK to an application using different system call hook mechanisms



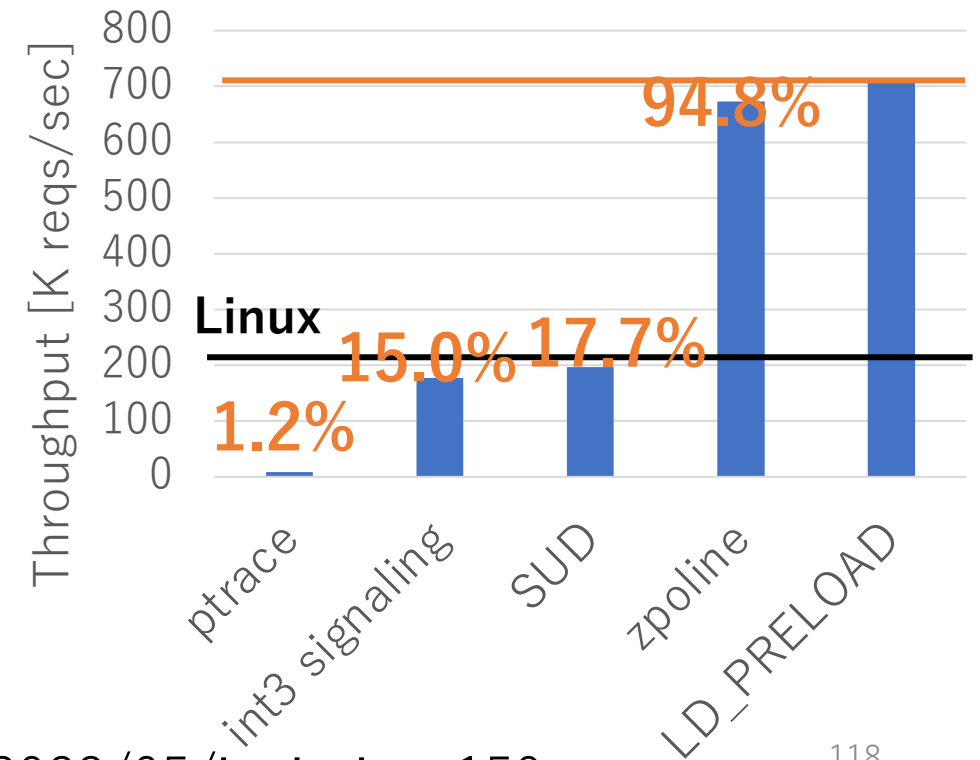
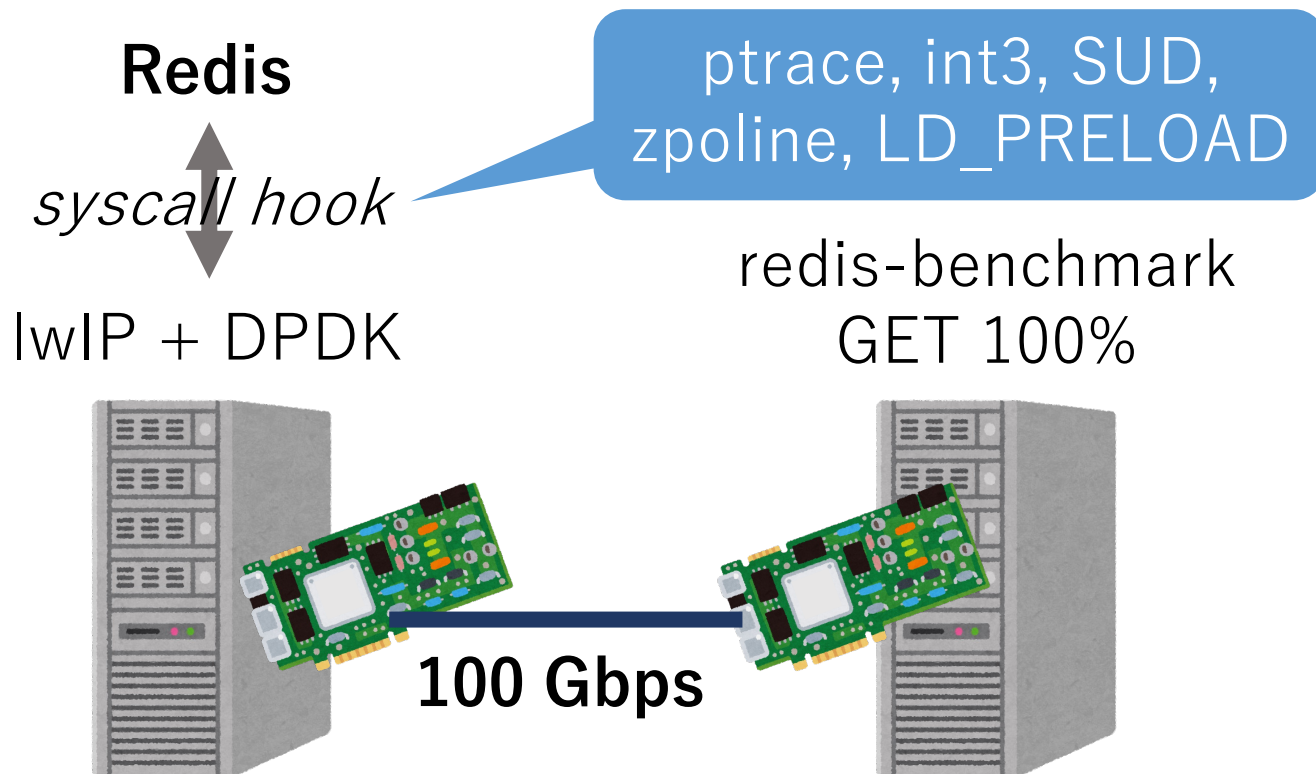
# Application Performance

- We **transparently** apply lwIP + DPDK to an application using different system call hook mechanisms



# Application Performance

- We **transparently** apply lwIP + DPDK to an application using different system call hook mechanisms



# Summary

- zpoline: a system call hook mechanism
  - replaces syscall/sysenter with callq \*%rax
  - instantiates the trampoline code at virtual address 0 (zero)
- 6 advantages: good for transparently applying user-space OS subsystems to existing applications
  1. Low hook overhead
  2. Exhaustive hooking
  3. No breakage of user-space program logic
  4. No kernel change and no additional kernel module are necessary
  5. No source code of a user-space program is needed
  6. It can be used for emulating system calls