# zpoline: a system call hook mechanism based on binary rewriting

Kenichi Yasukata[1], Hajime Tazaki[1], Pierre-Louis Aublin[1], Kenta Ishiguro[2]

[1] IIJ Research Laboratory
[2] Hosei University

# 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
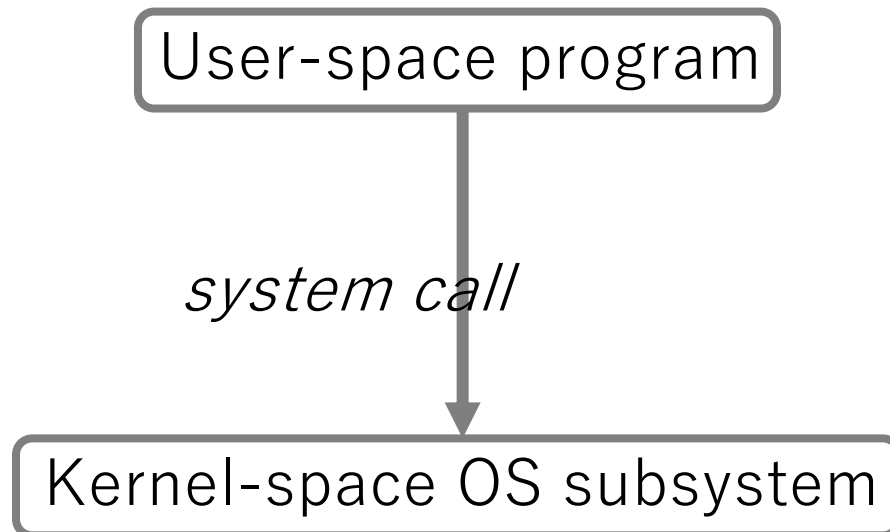
User-space program

# System Call

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

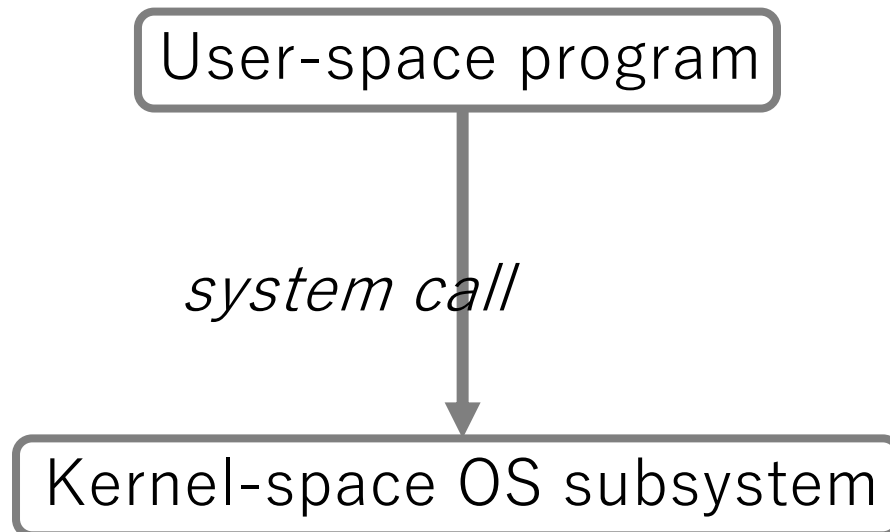User-space program

Kernel-space OS subsystem

# System Call

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

User-space program

*system call*

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

User-space program

*system call*
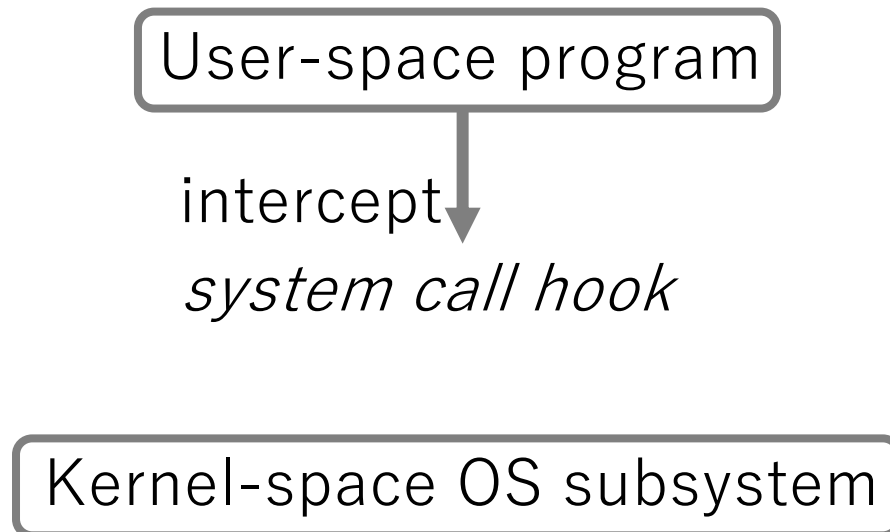
Kernel-space OS subsystem

# 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

```
┌─────────────────────────┐
│   User-space program    │
└─────────────────────────┘
           │
  intercept │
           ▼
   system call hook

┌─────────────────────────────┐
│  Kernel-space OS subsystem   │
└─────────────────────────────┘
```
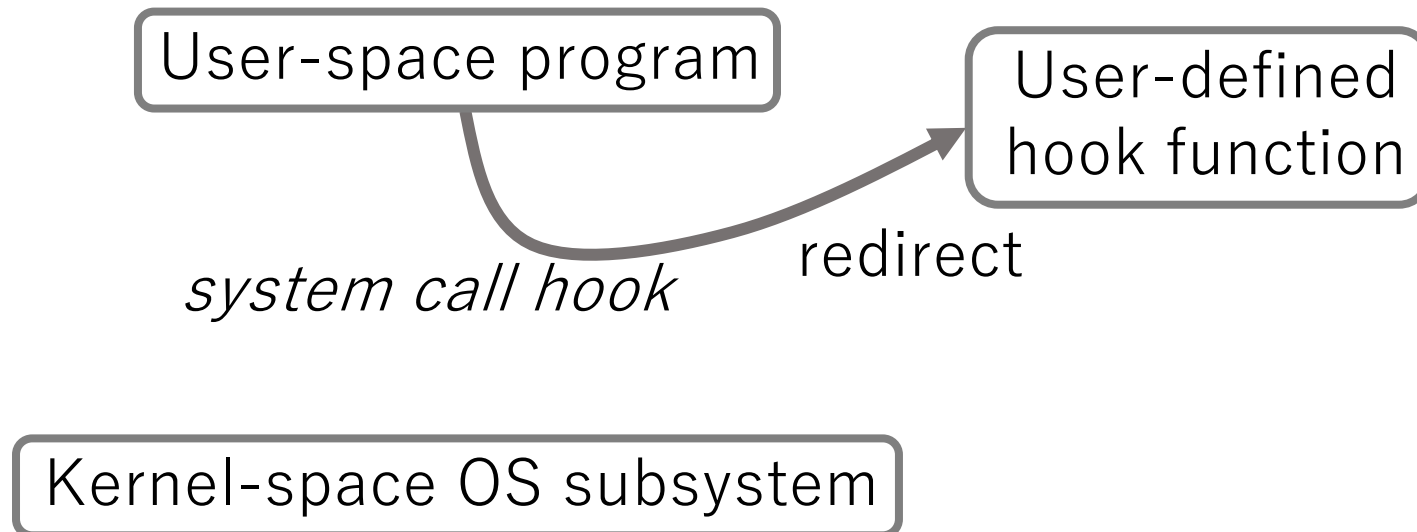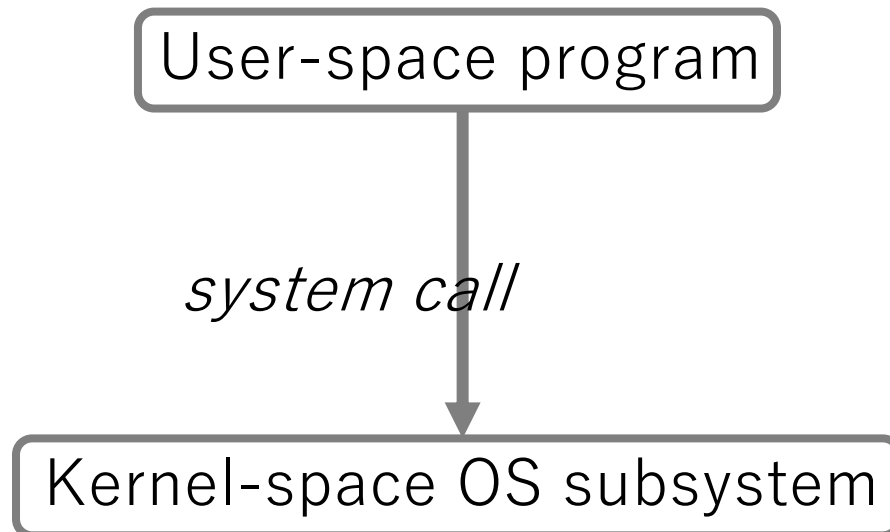
# System Call Hook

- 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
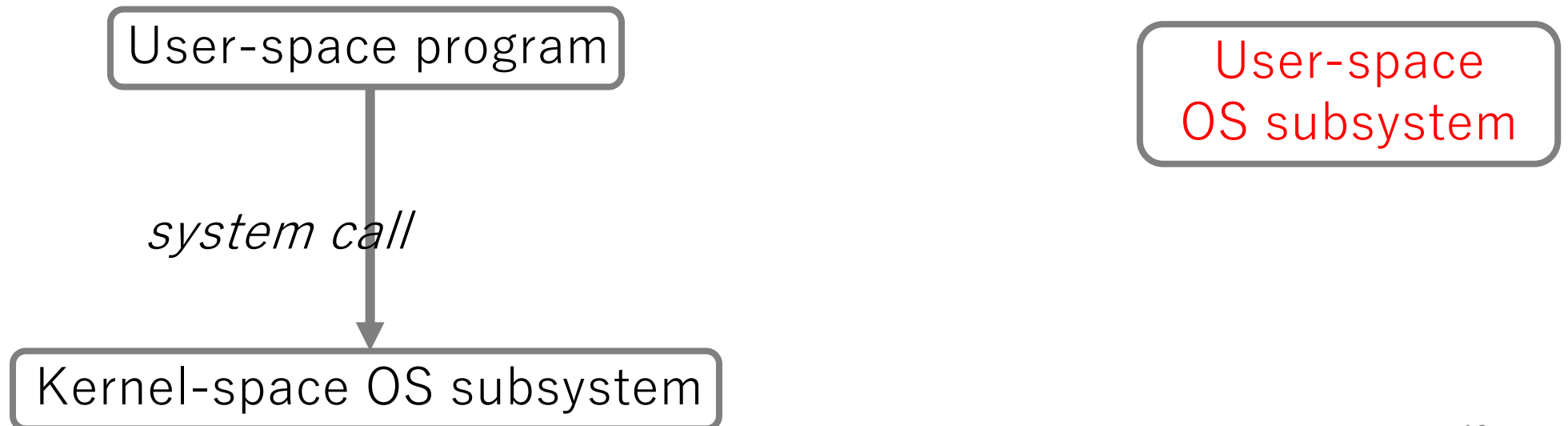
User-space program

*system call hook*

redirect

User-defined
hook function

Kernel-space OS subsystem

# Motivating Use Case

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

User-space program

*system call*

Kernel-space OS subsystem

# Motivating Use Case

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

User-space program

*system call*

Kernel-space OS subsystem
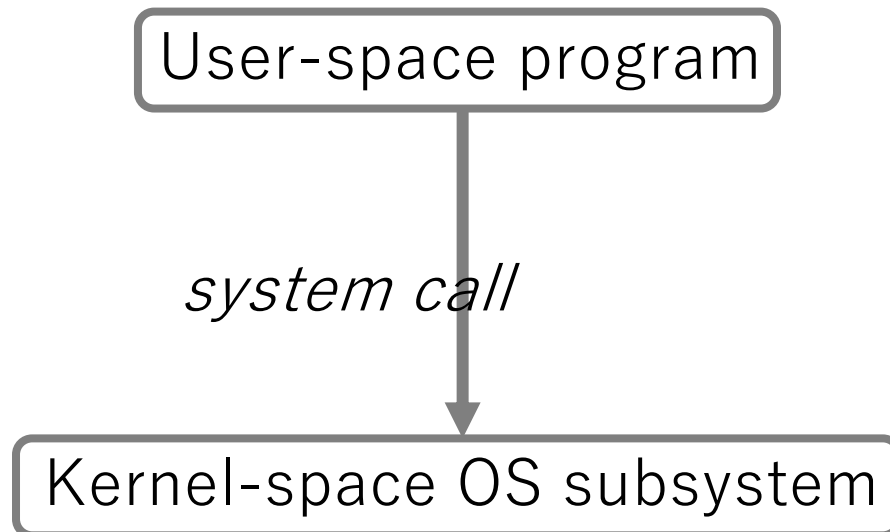
User-space OS subsystem

# Motivating Use Case

- System call hook mechanisms allow us to <u>transparently</u> apply <u style="color:red">user-space OS subsystems</u> to existing applications

User-space program

*system call*

Kernel-space OS subsystem

User-space
OS subsystem

**Highly performant**

# Motivating Use Case



TCP ping-pong performance

**5.3 times faster**

…o transparently apply
…plications

User-space
OS subsystem

**Highly performant**

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

# Motivating Use Case

- System call hook mechanisms allow us to <u>transparently</u> apply <u>user-space OS subsystems</u> to existing applications
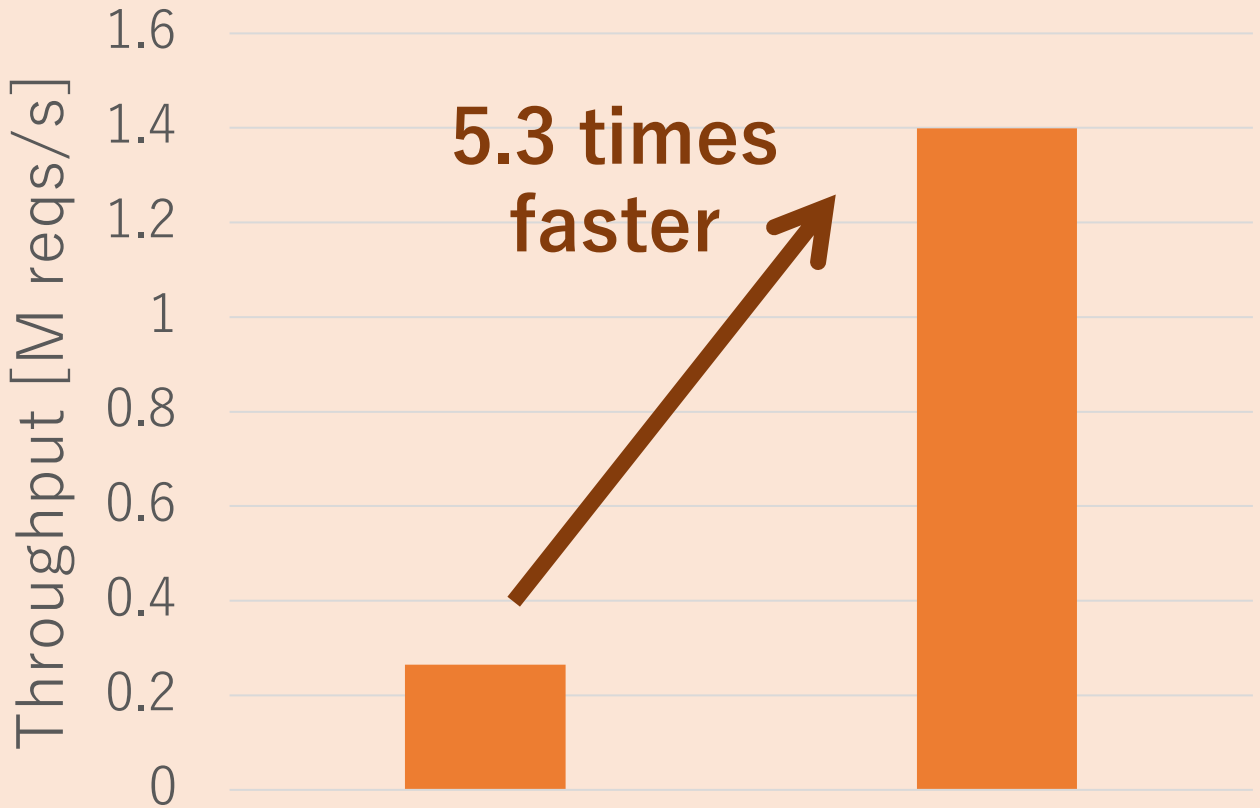
User-space program

*system call*

Kernel-space OS subsystem
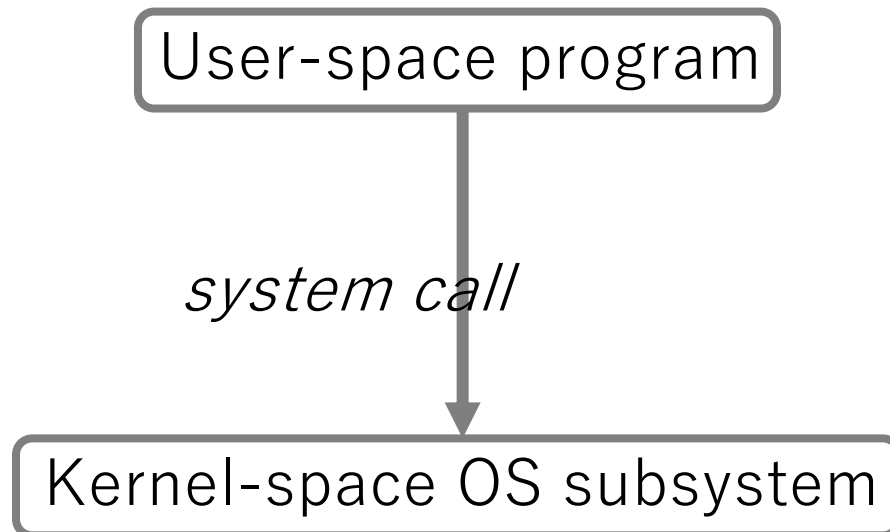
User-space OS subsystem

**Highly performant**

# Motivating Use Case

- System call hook mechanisms allow us to <u>transparently</u> apply <u style="color:red">user-space OS subsystems</u> to existing applications

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

User-space program

*system call*

Kernel-space OS subsystem

User-space OS subsystem

**Highly performant**

# Motivating Use Case

- System call hook mechanisms allow us to <u>transparently</u> apply <u>user-space OS subsystems</u> to existing applications
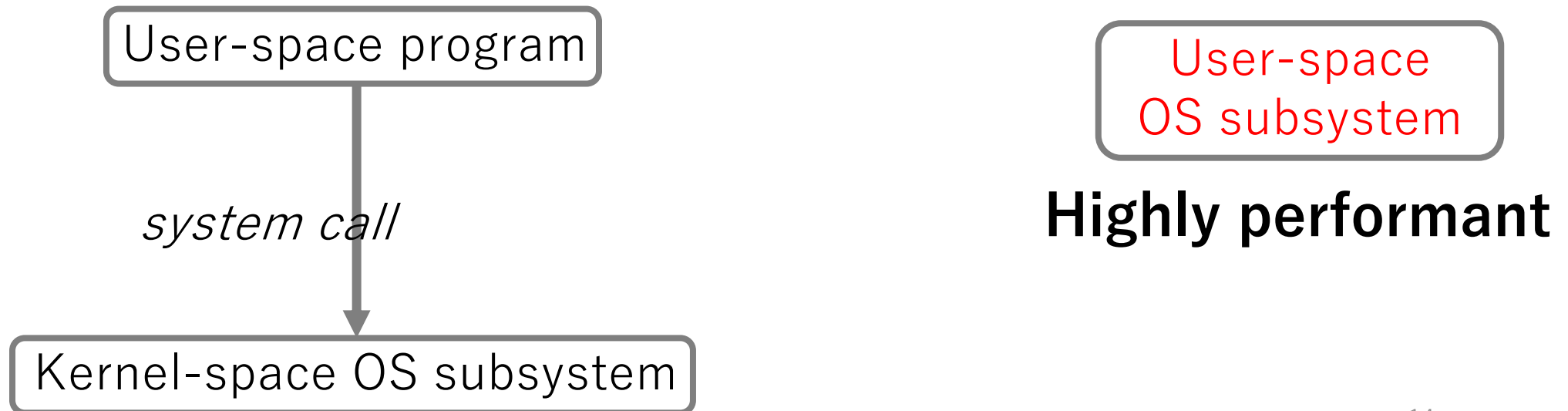
**Normally, adaptation requires changes of a user-space program to apply a specific API of a user-space OS subsystem**

**Modified**

| User-space program | Specific API → | User-space OS subsystem |

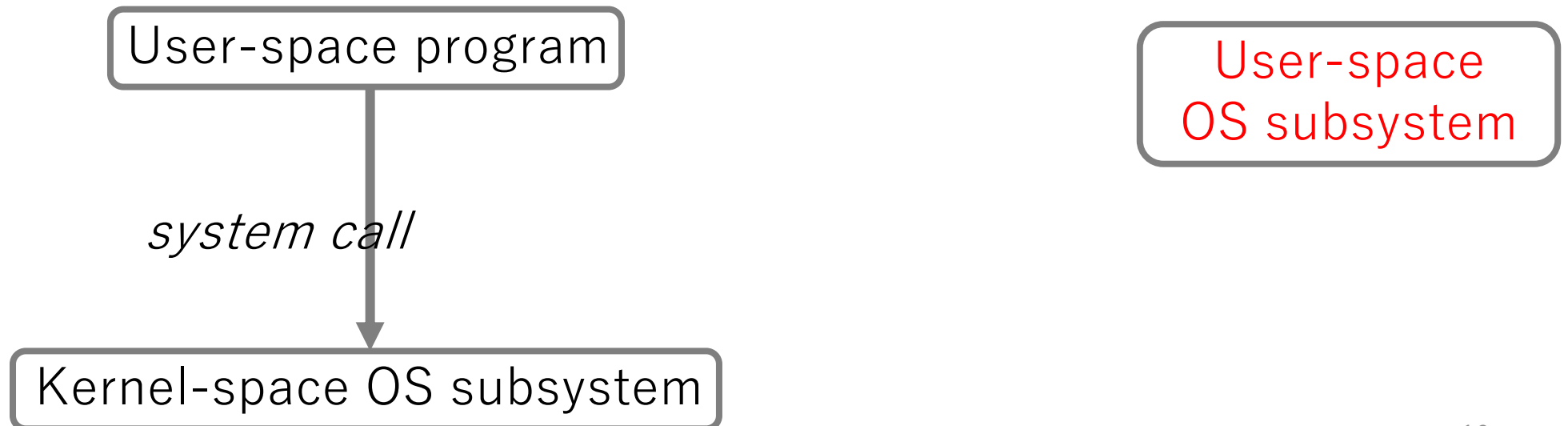*We need to change the program*

**Highly performant**

Kernel-space OS subsystem

# Motivating Use Case

- System call hook mechanisms allow us to <u>transparently</u> apply <span style="color:red"><u>user-space OS subsystems</u></span> to existing applications

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

User-space program

<span style="color:red">User-space
OS subsystem</span>

*system call*

Kernel-space OS subsystem

# Motivating Use Case

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

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

User-space program

*system call hook*

User-defined hook function | User-space OS subsystem

Kernel-space OS subsystem

# Motivating Use Case

- System call hook mechanisms allow us to transparently apply user-space OS subsystems to existing applications
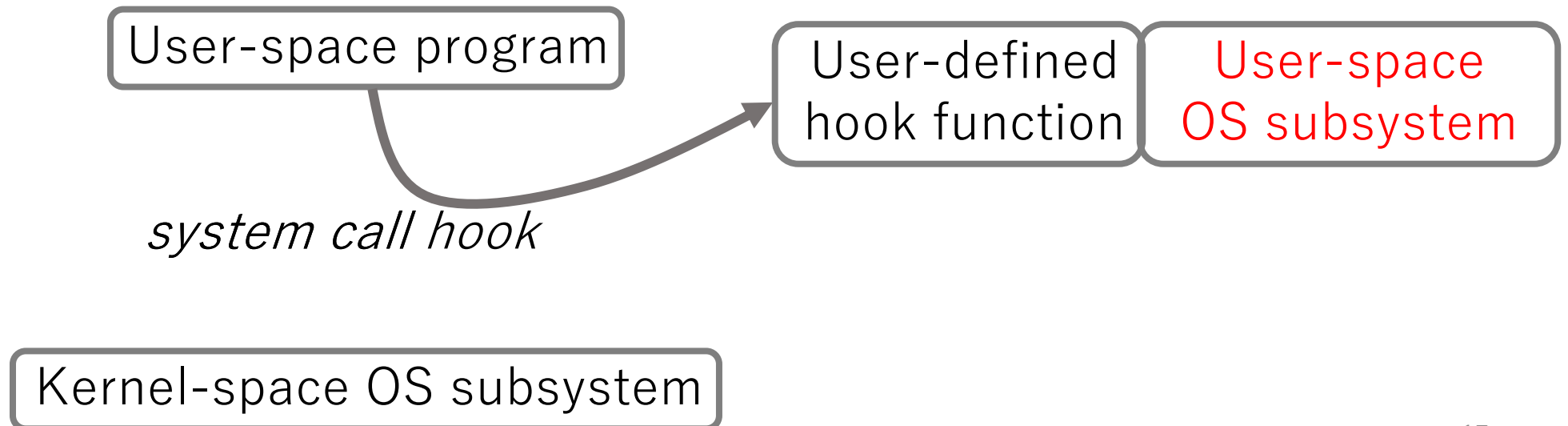
**If we use a system call hook mechanism, ...**

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

😄 | User-space program | → | User-defined hook function | User-space OS subsystem

*system call hook*
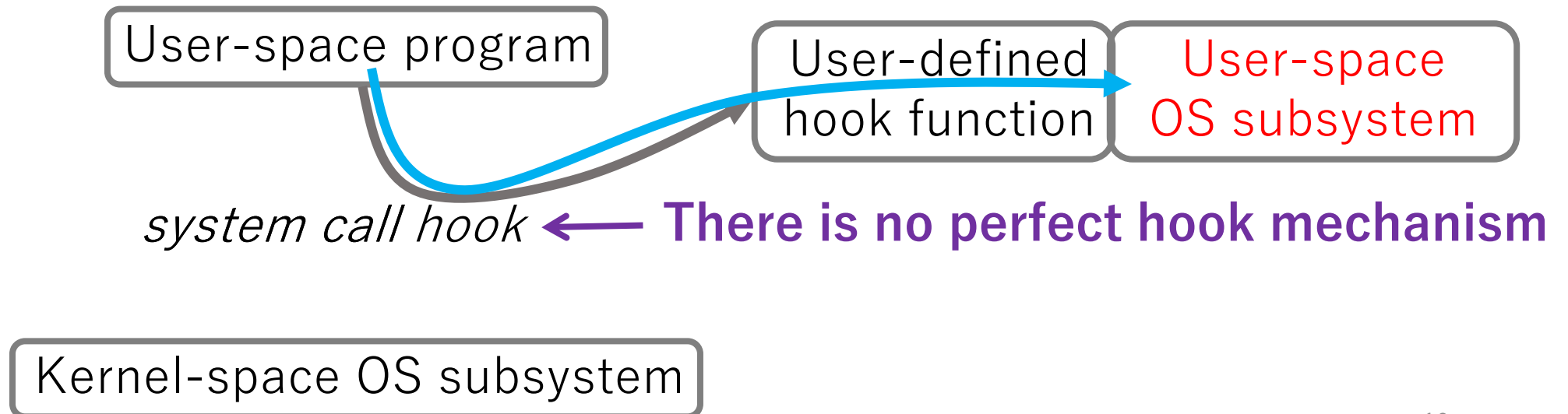
Kernel-space OS subsystem

# Problem

- System call hook mechanisms allow us to <u>transparently</u> apply <u><span style="color:red">user-space OS subsystems</span></u> to existing applications

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

<span style="color:deepskyblue">**no modification of the user-space program is necessary**</span>

User-space program

User-defined hook function

<span style="color:red">User-space OS subsystem</span>

*system call hook* ← <span style="color:purple">**There is no perfect hook mechanism**</span>

Kernel-space OS subsystem

# Problem

...ms allow us to <u>transparently</u> apply
... to existing applications

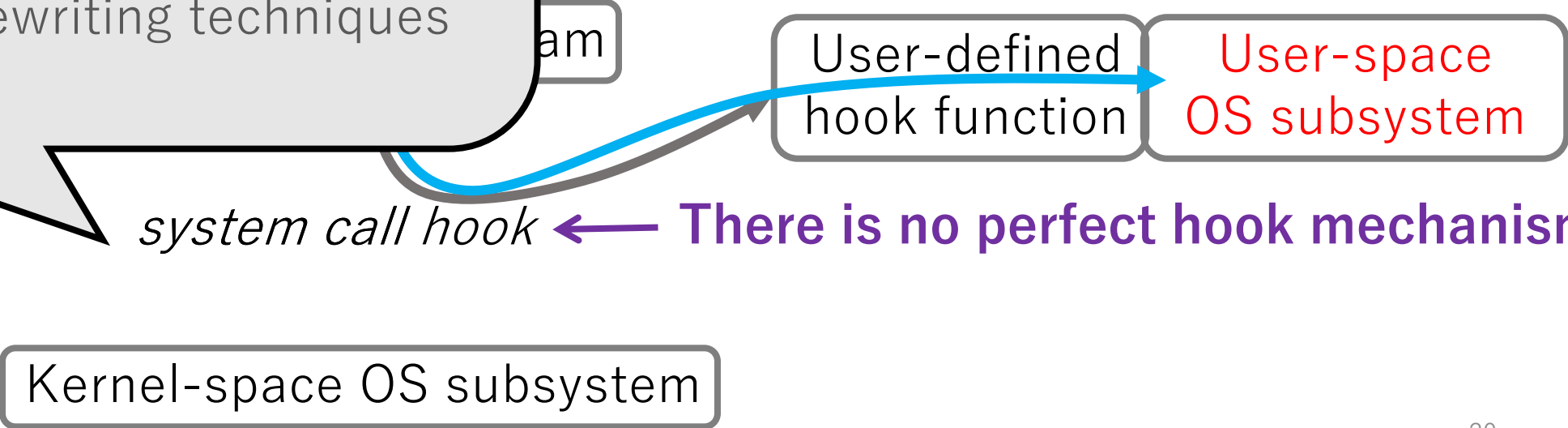**...ok mechanism, ...**

**...e user-space program is necessary**

**Existing Mechanisms**

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

...am

| User-defined hook function | User-space OS subsystem |

*system call hook* ← **There is no perfect hook mechanism**
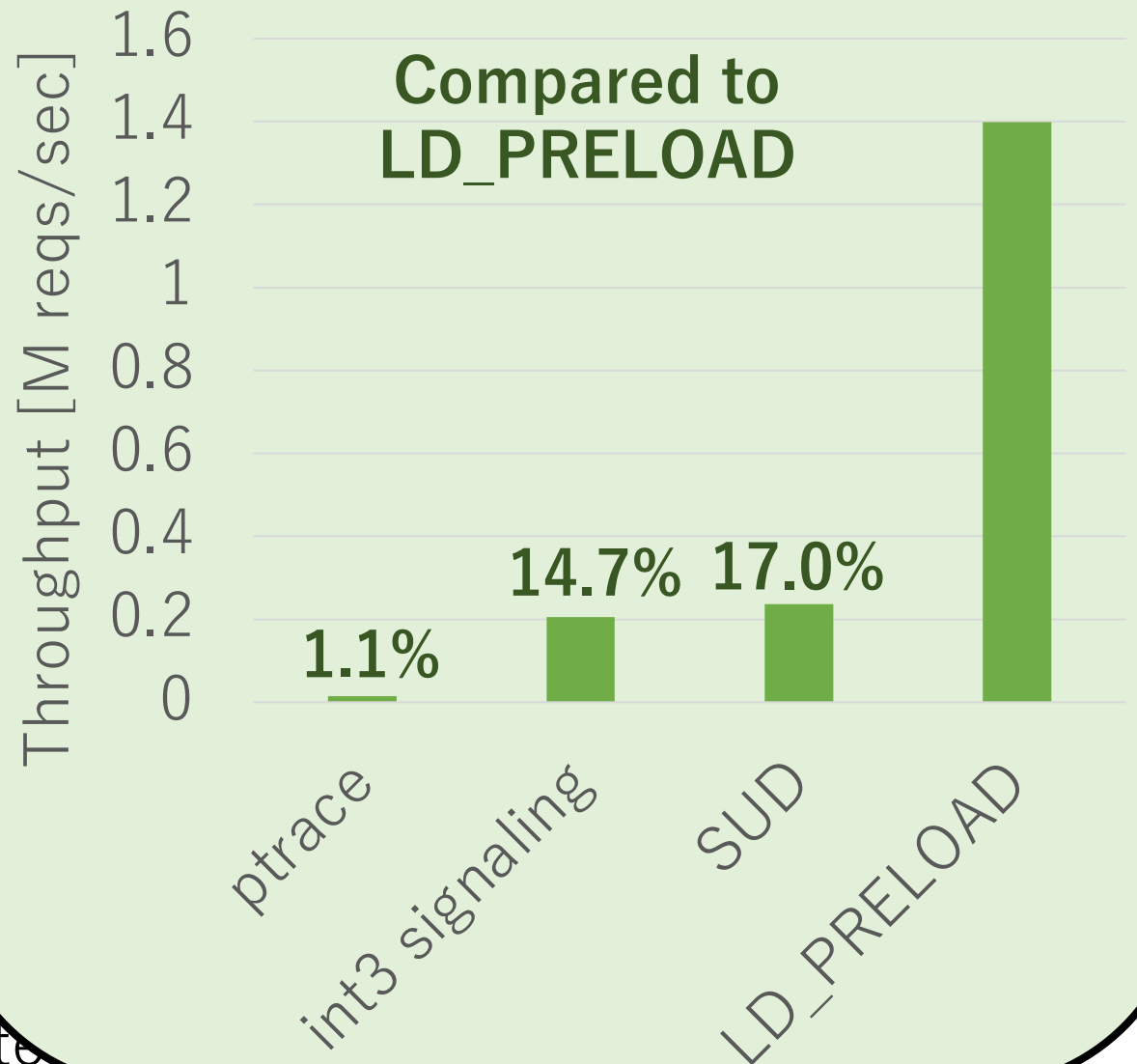
Kernel-space OS subsystem

# Problem

**Existing Mechanisms**

- ptrace
- int3 signaling technique
- Syscall User Dispatch (SUD)
- LD_PRELOAD trick
- Binary rewriting techniques
- …

*system call hook*

Kernel-space OS subsystem

lwIP on DPDK : TCP ping-pong

**Compared to LD_PRELOAD**

Throughput [M reqs/sec]

1.6
1.4
1.2
1
0.8
0.6
0.4
0.2
0

**1.1%** — ptrace
**14.7%** — int3 signaling
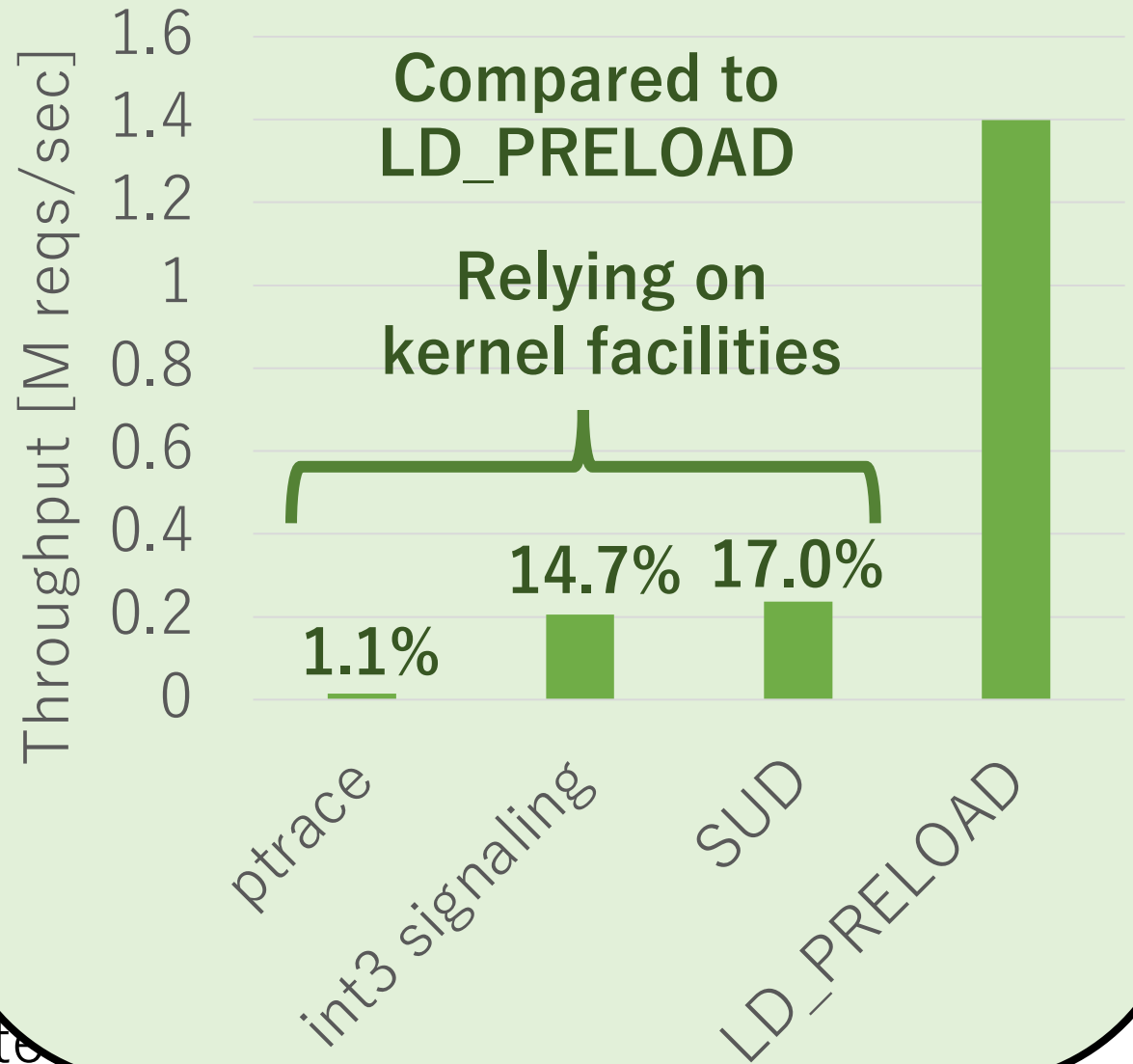**17.0%** — SUD
LD_PRELOAD

# Problem

**Existing Mechanisms**

- ptrace
- int3 signaling technique
- Syscall User Dispatch (SUD)
- LD_PRELOAD trick
- Binary rewriting techniques
- …

*system call hook*

Kernel-space OS subsystem
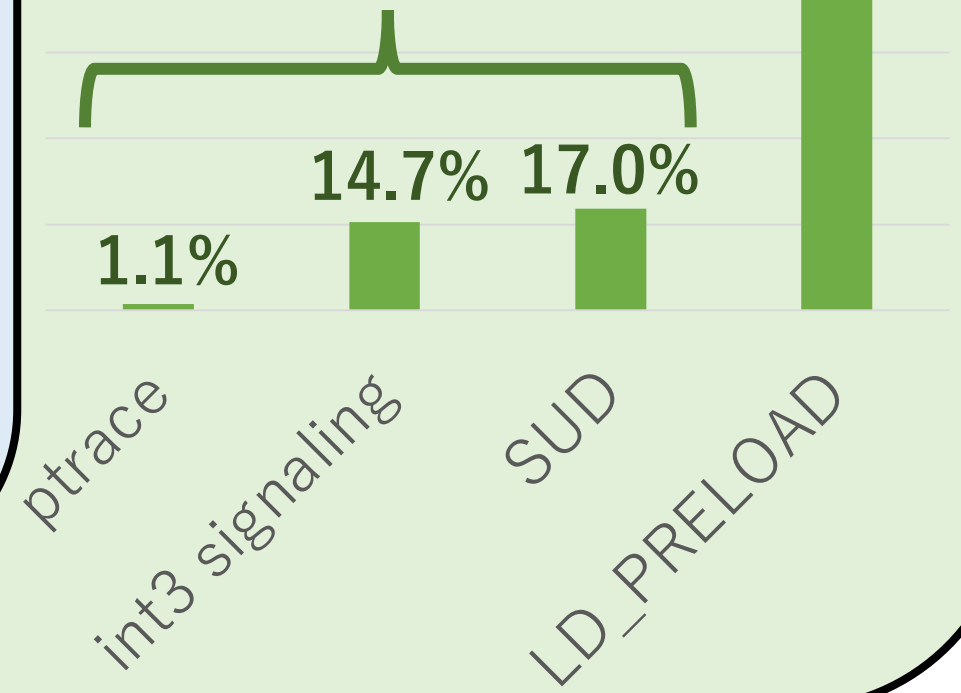
lwIP on DPDK : TCP ping-pong

**Compared to LD_PRELOAD**

**Relying on kernel facilities**

Throughput [M reqs/sec]

1.1%   14.7%   17.0%

ptrace   int3 signaling   SUD   LD_PRELOAD

# Problem

lwIP on DPDK : TCP ping-pong

- ptrace
  - overhead: process scheduling between the tracer and tracee processes

**Compared to LD_PRELOAD**

**Relying on kernel facilities**

1.1%   14.7%   17.0%

ptrace   int3 signaling   SUD   LD_PRELOAD
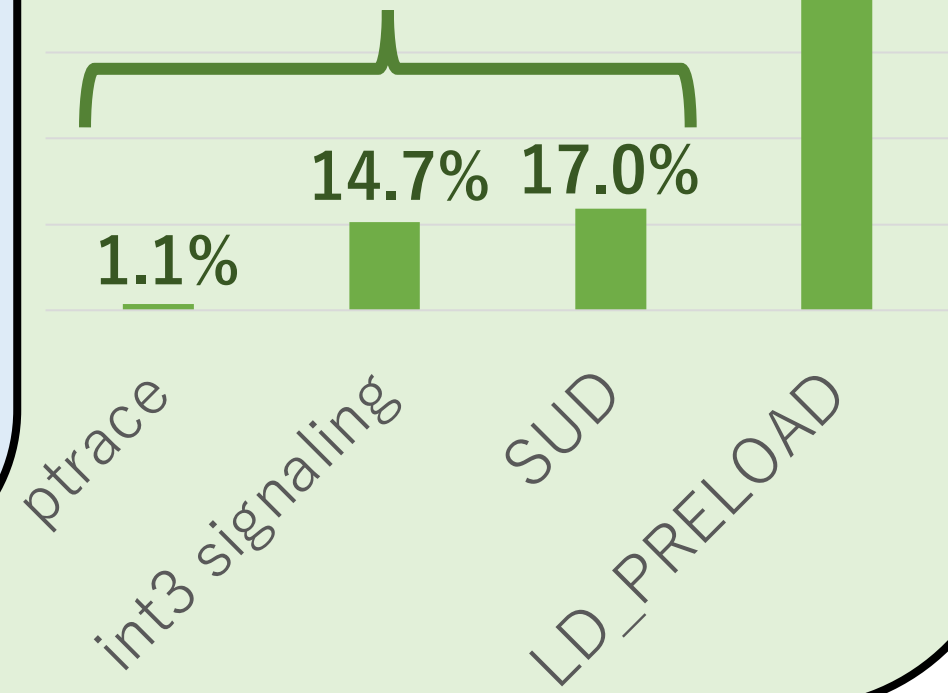
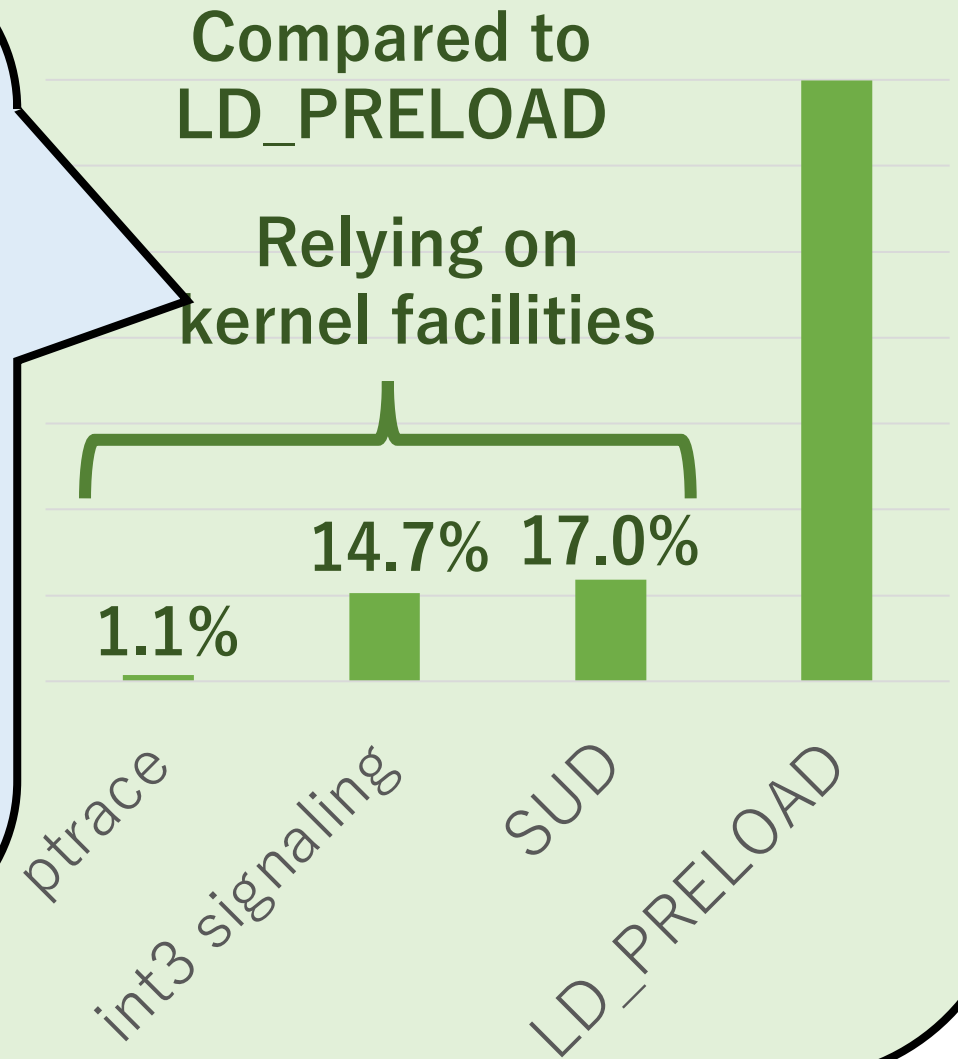Kernel-space OS subsyste

# Problem

- ptrace
  - overhead: process scheduling between the tracer and tracee processes
- int3 signaling / SUD
  - overhead: context manipulation for a signal() handler (SIGINT/SIGSYS)

Kernel-space OS subsyste

lwIP on DPDK : TCP ping-pong

**Compared to LD_PRELOAD**

**Relying on kernel facilities**

1.1%    14.7%    17.0%

ptrace    int3 signaling    SUD    LD_PRELOAD

# Problem

- ptrace
  - overhead: process scheduling between the tracer and tracee processes
- int3 signaling / SUD
  - overhead: context manipulation for a signal() handler (SIGINT/SIGSYS)

- LD_PRELOAD just replaces function calls, therefore, it is fast

Kernel-space OS subsyste

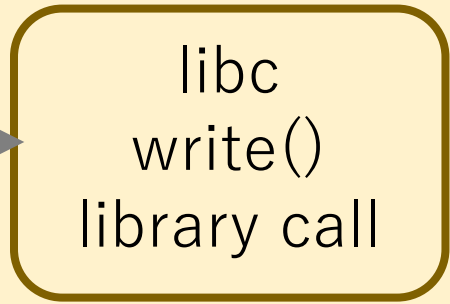lwIP on DPDK : TCP ping-pong

**Compared to LD_PRELOAD**

**Relying on kernel facilities**

1.1%   14.7%   17.0%

ptrace   int3 signaling   SUD   LD_PRELOAD

Prob...

- ptrac...
  - ove...
    the...
- int3 s...
  - ove...
    a s...

- LD_P...
  funct...

```
app_function(...)
{

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

libc
write()
library call

```
special_write(...)
{

    asm volatile (
        trigger
        write syscall
    )
}
```

User-defined
write()
library call

**LD_PRELOAD**

...OK : TCP ping-pong

...pared to
...PRELOAD

...elying on
...el facilities

14.7%  17.0%

...UD

LD_PRELOAD

# Problem

~~ms allow us to [transparently] apply~~
~~to existing applications~~

**~~ook~~ mechanism, ...**

**~~e~~ user-space program is necessary**

**Existing Mechanisms**

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

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

| User-defined hook function | User-space OS subsystem |
| --- | --- |

Kernel-space OS subsystem

# Problem

...ms allow us to <u>transparently</u> apply
... to existing applications

**...ok mechanism, ...**

**...e user-space program is necessary**

**Existing Mechanisms**

- int3 signal technique
- Syscall User Dispatch (SUD)
- LD_PRELOAD trick
- Binary rewriting techniques
- ...

**High performance penalty**

User-defined hook function → User-space OS subsystem

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

Kernel-space OS subsystem

# Problem

...ms allow us to <u>transparently</u> apply
... to existing applications

**...ok mechanism, ...**

**...e user-space program is necessary**

**Existing Mechanisms**

- int3 si... ...hnique
- Syscall User Dis... ...D)
- ...RELOAD trick
- Binary ... ...hniques
- ...

*High performance penalty*

*Sometimes fail to hook*

...am

| User-defined hook function | User-space OS subsystem |

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

| Kernel-space OS subsystem |

36

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

...ms allow us to <u>transparently</u> apply
...to existing applications

**...ook mechanism, ...**

**...e user-space program is necessary**

**Existing Mechanisms**

High performance penalty

Sometimes fail to hook

- int3 ~~signal~~ ~~technique~~
- Syscall User Dispatch (SUD)
- ~~LD_P~~RELOAD trick
- Binary ~~rewriting~~ ~~techniques~~
- ...

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

User-defined hook function → User-space OS subsystem

Kernel-space OS subsystem

# Contribution

- zpoline: a system call hook mechanism for x86-64 CPUs
  - based on binary rewriting
  - free from the drawbacks of the pervious mechanisms

- This work addresses a challenge that is specific to binary rewriting approaches

# Binary Rewriting Approach

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



0x0000
0x0001
0x0002

. . .

. . .

. . .

**syscall** 0x0f 0x05

. . .

# Binary Rewriting Approach

0x0000
0x0001
0x0002

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

**syscall** 0x0f 0x05

# Binary Rewriting Approach

0x0000
0x0001
0x0002

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

- What we wish to achieve

**syscall** 0x0f 0x05

# Binary Rewriting Approach

0x0000
0x0001
0x0002

. . .

. . .

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

**user-defined hook function**

. . .

- What we wish to achieve

**syscall** 0x0f 0x05

. . .

42

# 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**

Virtual Memory

0x0000
0x0001
0x0002

. . .

. . .

**user-defined hook function**

. . .

**syscall** 0x0f
0x05

. . .

# 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**

**user-defined hook function**

**????**

44

# Binary Rewriting Approach

0x0000
0x0001
0x0002

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

*jump*

**user-defined hook function**

- What we wish to achieve
  - replace syscall/sysenter instruction with **something**
  - to jump to a user-defined hook function

**????**

# Binary Rewriting Approach

0x0000
0x0001
0x0002

. . .

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

*jump*

. . .

**user-defined hook function**

- 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?**

**????**

46

# Challenge

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

0x0000
0x0001
0x0002

. . .

. . .

**user-defined hook function**

. . .

**syscall** 0x0f 0x05

. . . .

# Challenge

0x0000
0x0001
0x0002

. . .

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

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

**user-defined
hook function**

. . .

**syscall** 0x0f
0x05

. . .

48

# Challenge

0x0000
0x0001
0x0002

- 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**

. . .

**user-defined hook function**

. . .

syscall 0x0f 0x05

. . .

# 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**

Virtual Memory

0x0000
0x0001
0x0002

. . .

. . .

*ADDR*

**user-defined hook function**

. . .

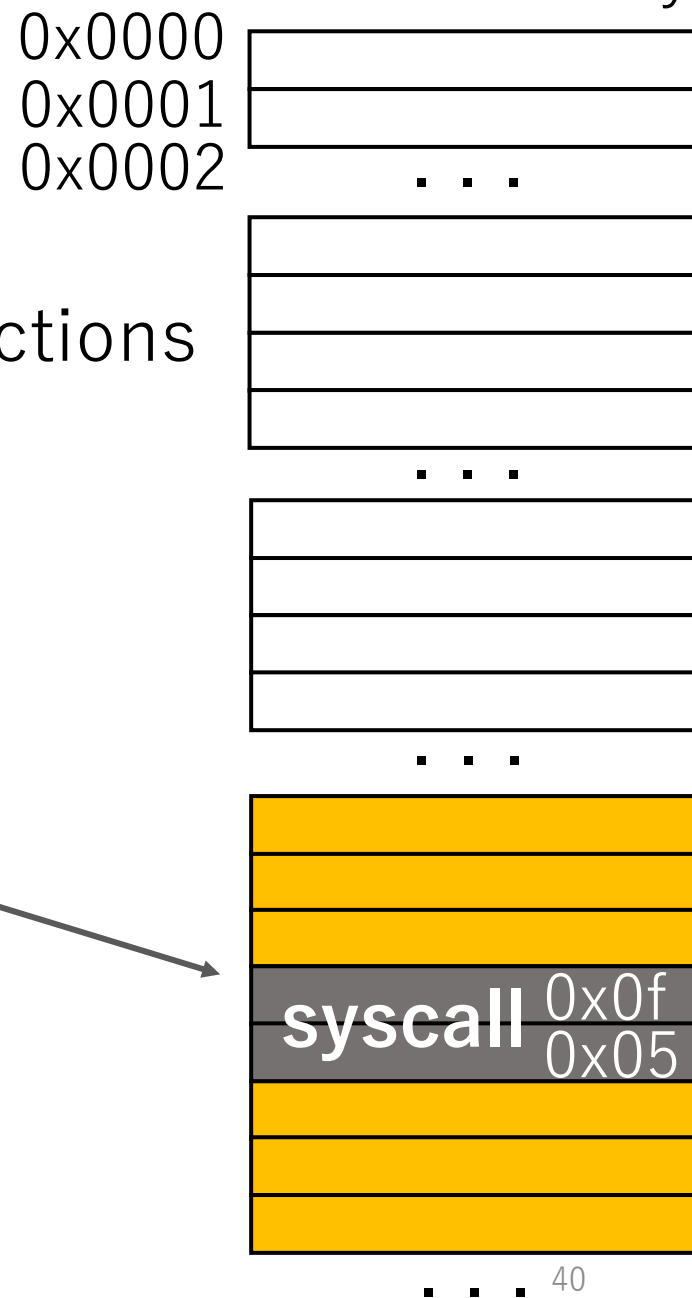**syscall** 0x0f 0x05

. . . 50

# Challenge

0x0000
0x0001
0x0002

. . .

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

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

**ADDR**

. . .

**user-defined hook function**

. . .

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

*ADDR*

ADDR is bigger than 2 bytes

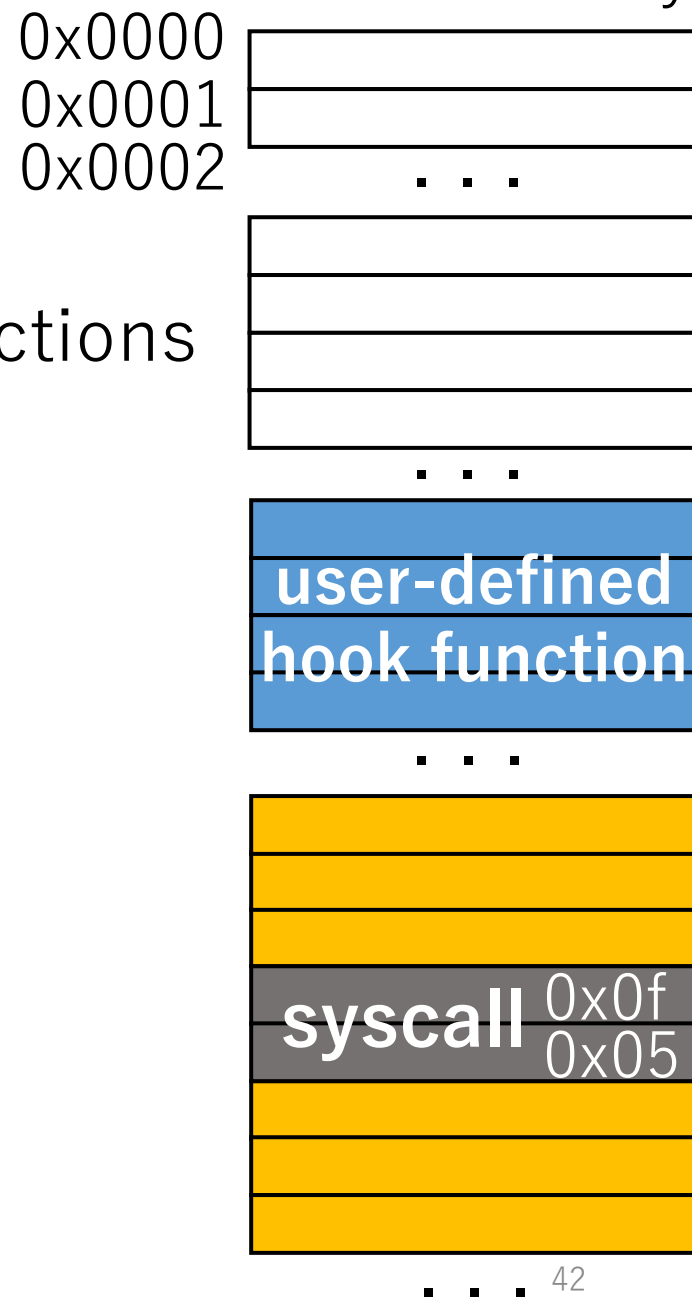If we put ADDR, subsequent instructions are overwritten

51

# Challenge

0x0000
0x0001
0x0002

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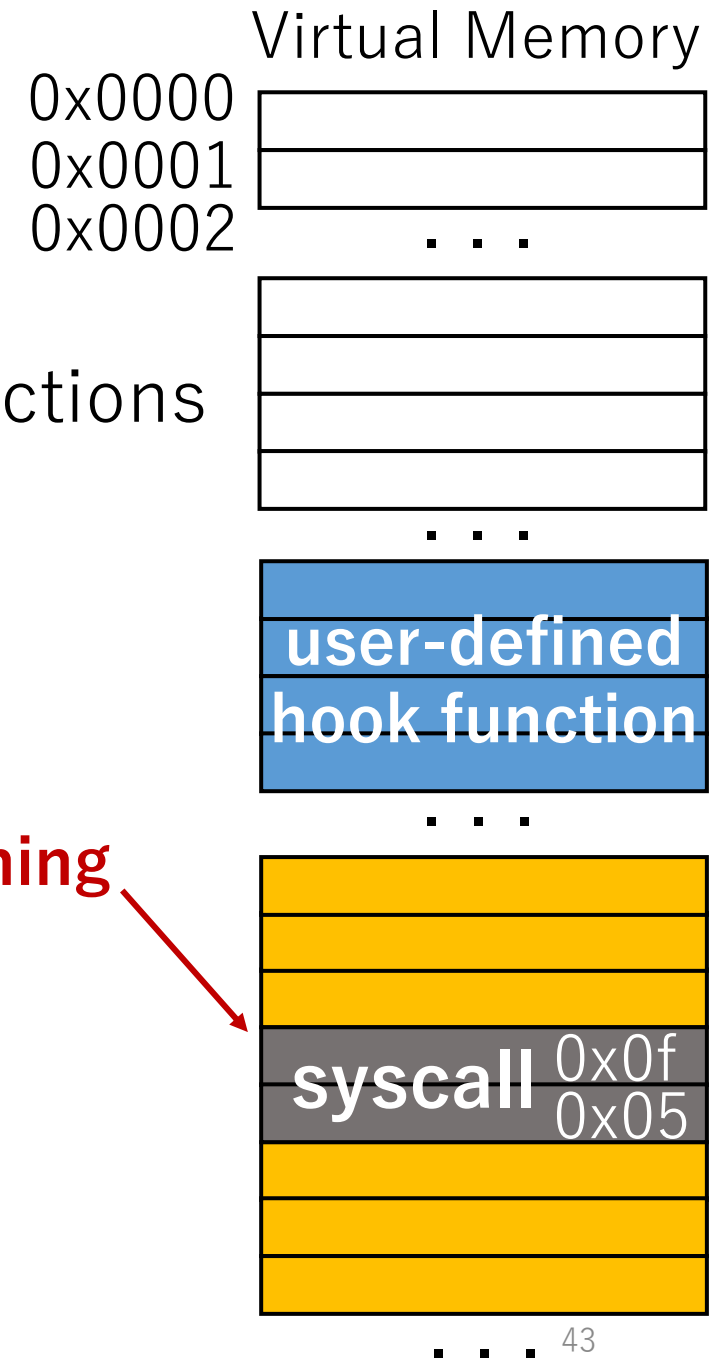  - syscall: <u>0x0f 0x05</u>, sysenter: <u>0x0f 0x34</u>

*ADDR*

**some program**

user-defined hook function

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

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

*ADDR*

ADDR is bigger than 2 bytes

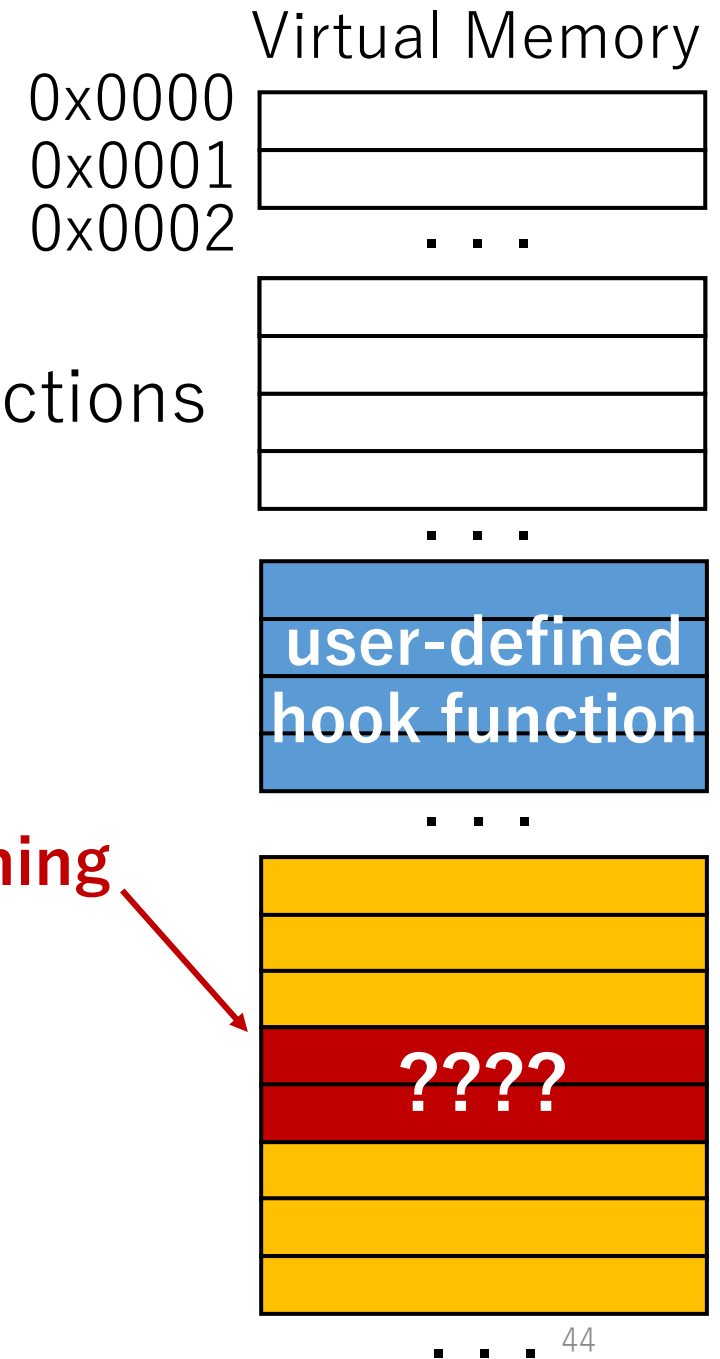If we put ADDR, subsequent instructions are overwritten

52

# Challenge

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

- 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

Virtual Memory

0x0000
0x0001
0x0002

. . .

~~some program~~

. . .

*ADDR*

*jump*

**user-defined hook function**

. . .

ADDR is bigger than 2 bytes

. . .

# Challenge

0x0000
0x0001
0x0002

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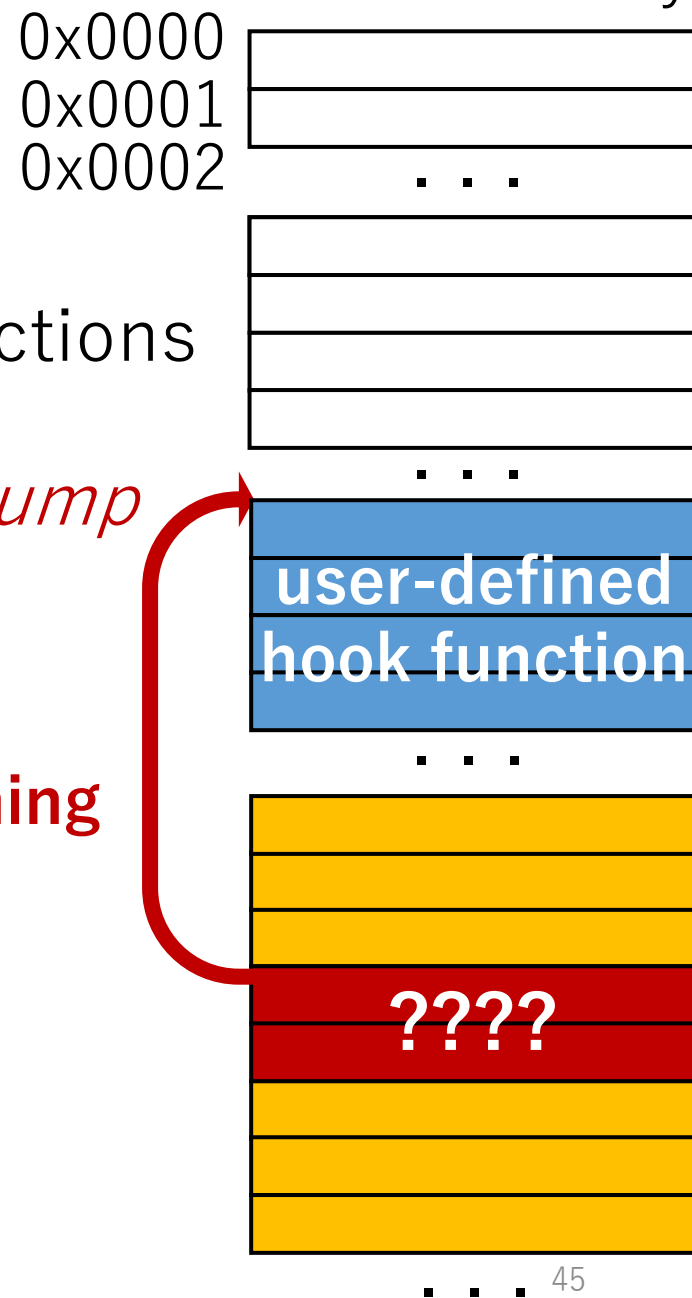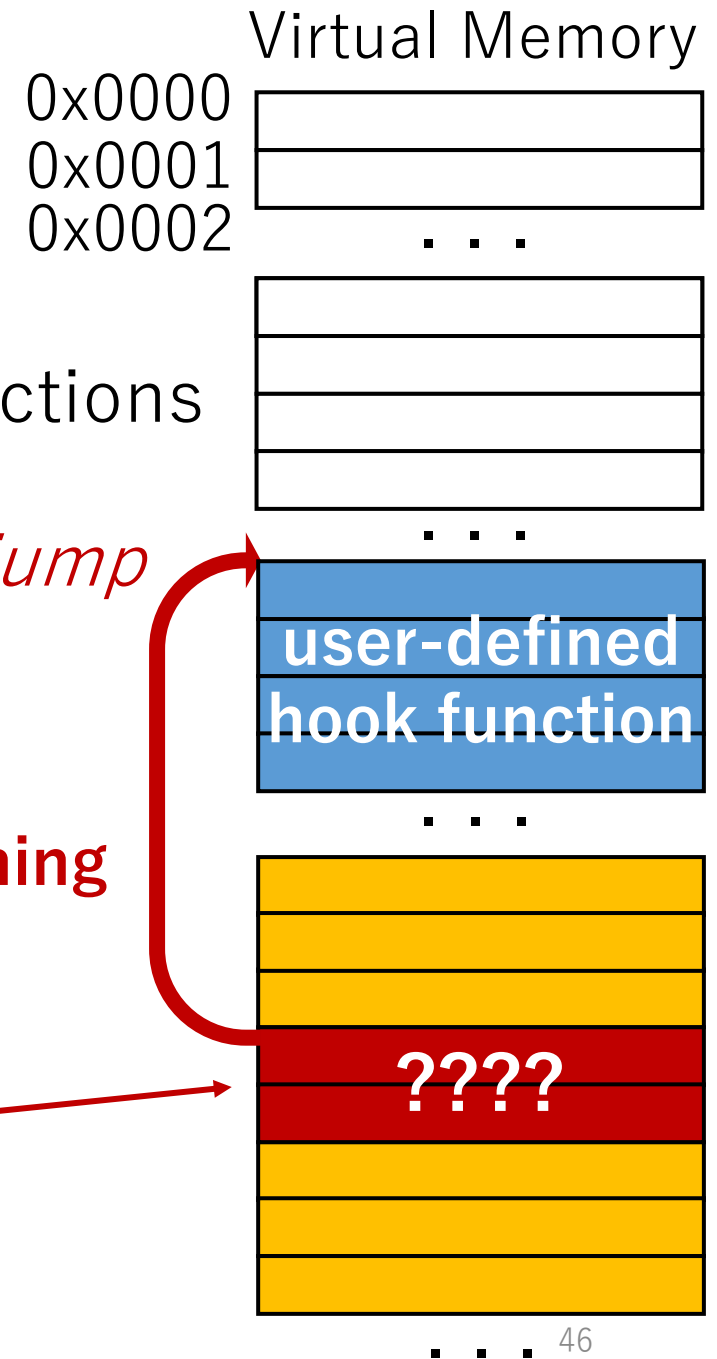
- On x86-64 CPUs, syscall and sysenter instructions trigger a system call
  - syscall: 0x0f 0x05, sysenter: 0x0f 0x34

**some program**

*ADDR*

Because of this issue, previous binary rewriting techniques
- could not ensure exhaustive hooking
- or, overwrite neighbour 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

ADDR is bigger than 2 bytes

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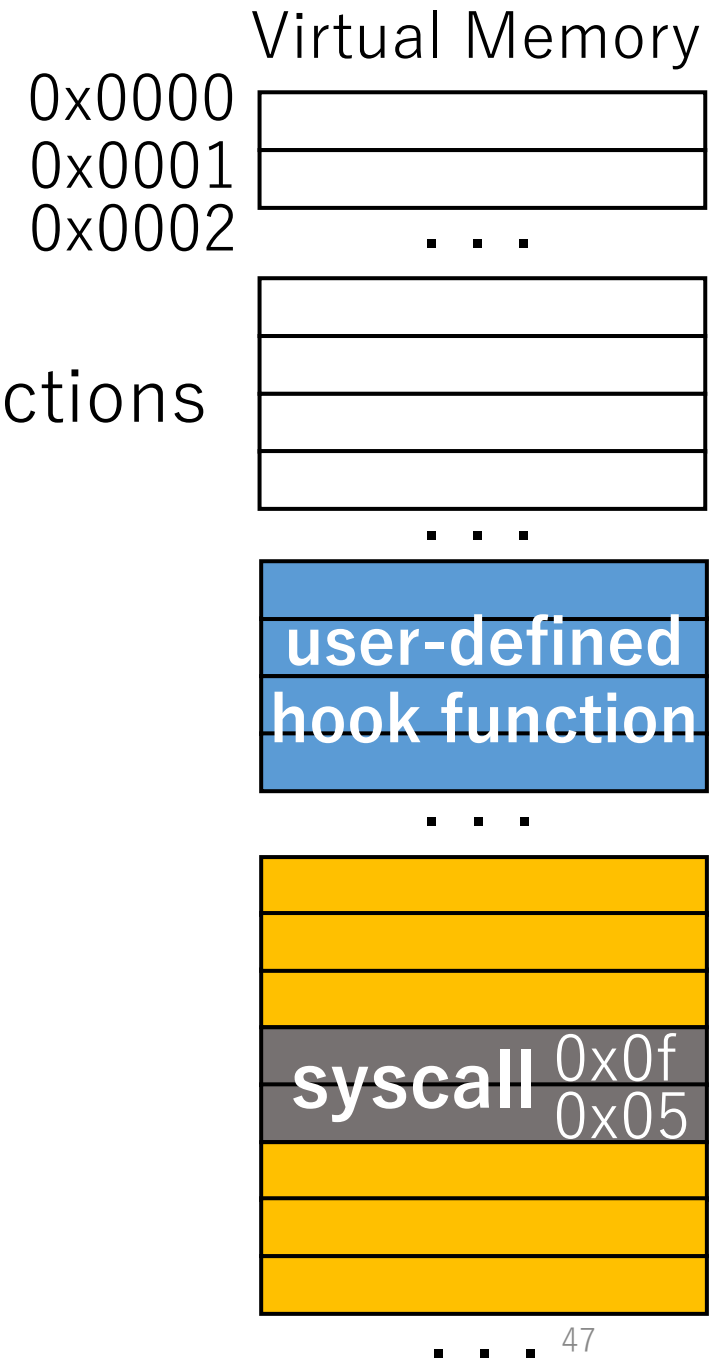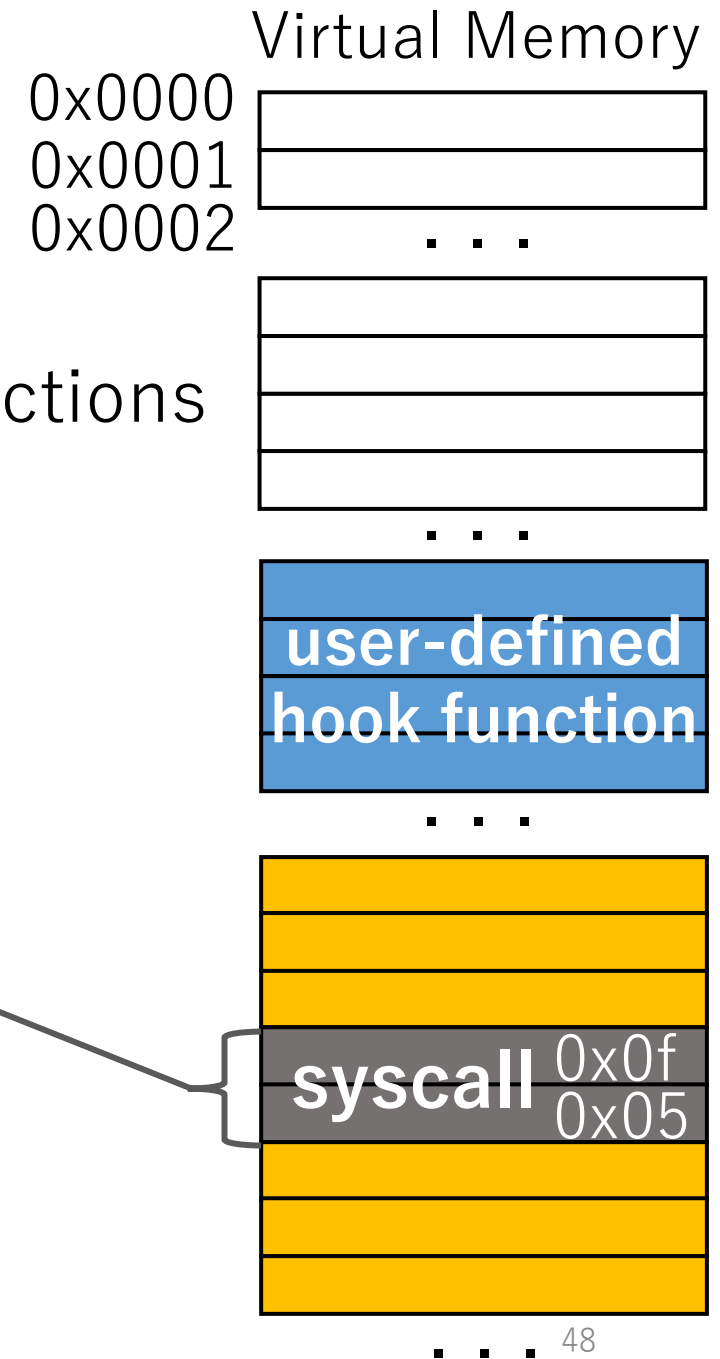
Virtual Memory

0x0000
0x0001
0x0002

. . .

some ~~program~~

. . .

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

*ADDR*

Because of this issue, previous binary rewriting techniques
- could not ensure exhaustive hooking
- or, overwrite neighbour instructions

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

ADDR is bigger 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

0x0000
0x0001
0x0002

. . .

. . .

**user-defined hook function**

. . .

**syscall** 0x0f
0x05

. . .

56

# Calling Convention

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

0x0000
0x0001
0x0002

. . .

. . .

**user-defined hook function**

. . .

**set syscall num to rax register**

**syscall** 0x0f
0x05

. . .

# Calling Convention

0x0000
0x0001
0x0002

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

**user-defined hook function**

**set syscall num to rax register**

**syscall** 0x0f
0x05

# Calling Convention

0x0000
0x0001
0x0002

. . .

. . .

- How to invoke a system call
  - A user-space program sets a system call number, predefined by the kenel, 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 ----

**user-defined hook function**

**context switch to kernel**

um

05

. . .

# Calling Convention

0x0000
0x0001
0x0002

- How to invoke a system call
  - A user-space program sets a system call number, predefined by the kenel, 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**

. . .

. . .

**user-defined hook function**

. . .

**set syscall num to rax register**

**syscall** 0x0f
0x05

. . .

# Calling Convention

0x0000
0x0001
0x0002

- How to invoke a system call
  - A user-space program sets a system call number, predefined by the kenel, 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, the kernel executes read()
    - if the rax register has 1, the kernel executes write()
    - if the rax register has 2, the kernel executes open()

**user-defined hook function**

**set syscall num to rax register**

**syscall** 0x0f
0x05

# Calling Convention

- How to invoke a system call
  - A user-space program sets a system call number, predefined by the kenel, 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,**

Virtual Memory

0x0000
0x0001
0x0002

. . .

. . .

**user-defined**
**hook function**

. . .

**set syscall num**
**to rax register**

**syscall** 0x0f
0x05

. . . .

# Calling Convention

0x0000
0x0001
0x0002

- How to invoke a system call
  - A user-space program sets a system call number, predefined by the kenel, 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)**

**user-defined
hook function**

**set syscall num
to rax register**

**syscall** 0x0f
0x05

63

# zpoline

• 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)**

Virtual Memory

0x0000
0x0001
0x0002

. . .

. . .

user-defined
hook function

. . .

set syscall num
to rax register

syscall 0x0f
0x05

. . . .

# zpoline

Virtual Memory

0x0000
0x0001
0x0002

. . .

. . .

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

**user-defined hook function**

. . .

**set syscall num to 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)**

**callq *%rax**

. . .

# zpoline

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

0x0000
0x0001
0x0002

. . .

. . .

**user-defined hook function**

. . .

**set syscall num to 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)**

**callq *%rax**

66

# zpoline

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

**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)**

. . .

**user-defined
hook function**

. . .

**set syscall num
to rax register**

**callq *%rax**

. . .

# zpoline

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

. . .

**user-defined
hook function**

. . .

**set syscall num
to rax register**

**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

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

**After the binary rewriting**

**<u>Point: Calling Convention</u>**

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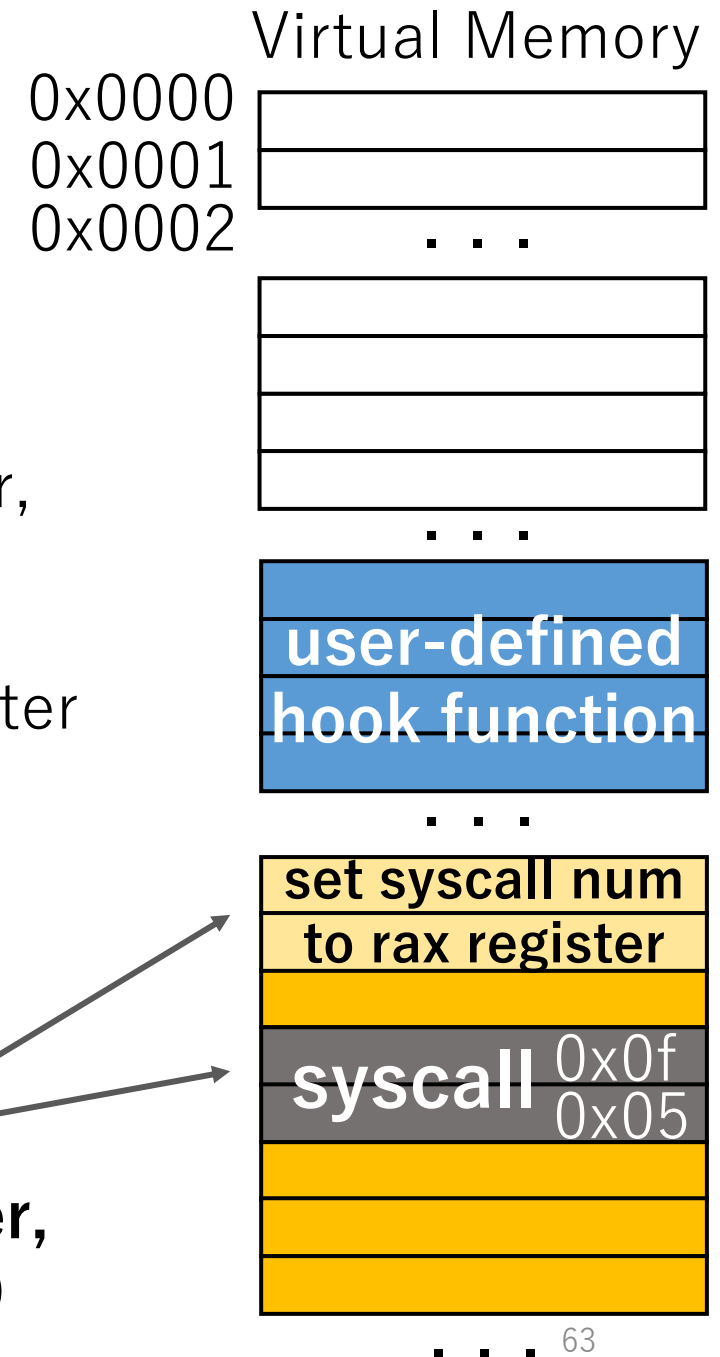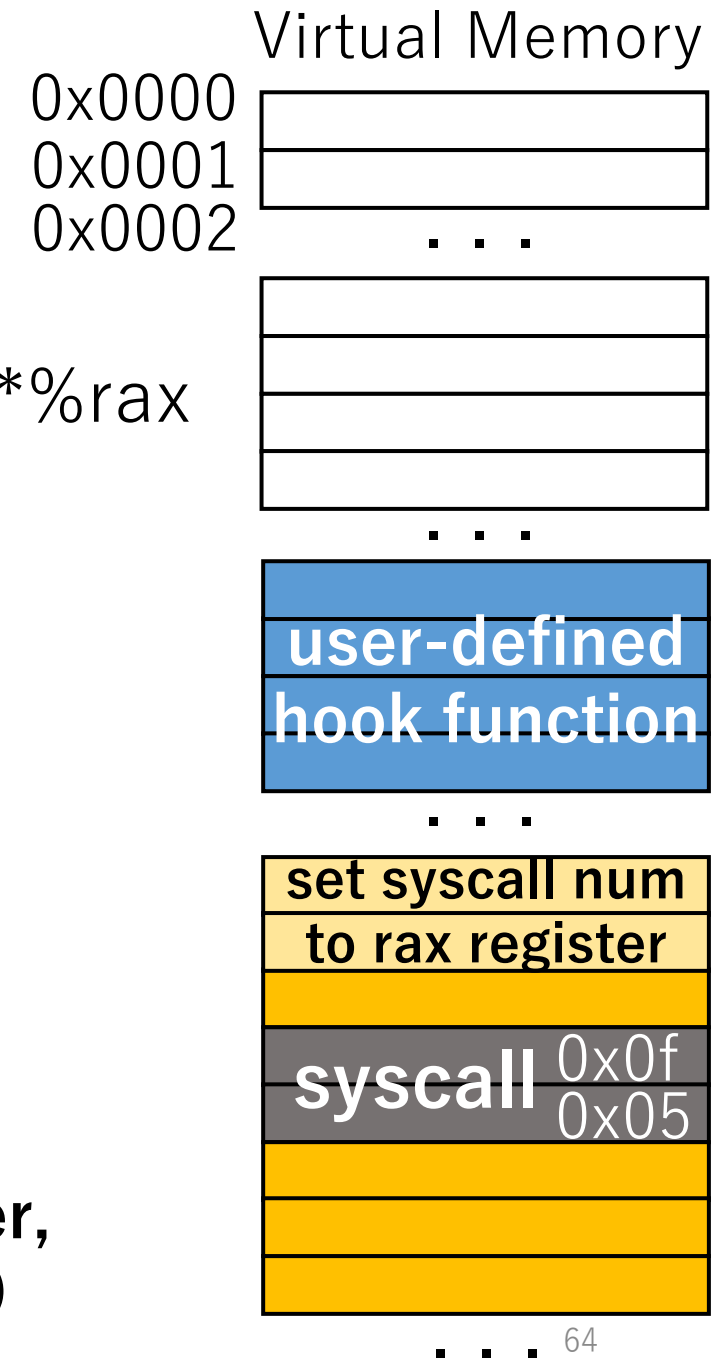which is 0 ~ around 500 (defined in the kernel)**

. . .

. . .

**user-defined
hook function**

. . .

**set syscall num
to rax register**

**callq *%rax**

. . .

# 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)**

Virtual Memory

0x0000
0x0001
0x0002

. . .

. . .

user-defined
hook function

. . .

set syscall num
to rax register

callq *%rax

. . .

# zpoline

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

. . .

**user-defined hook function**

. . .

**After the binary rewriting**

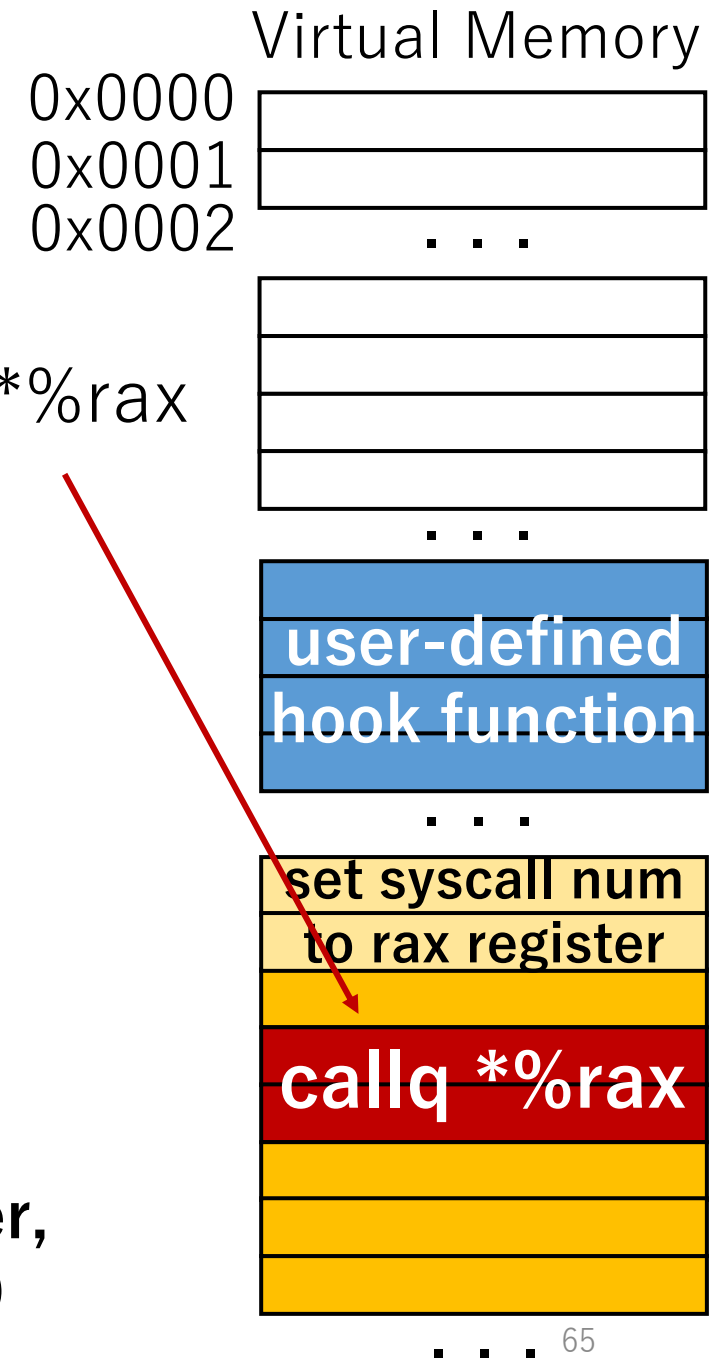**<u>Point: Calling Convention</u>**

**When ~~syscall/sysenter~~ <span style="color:red">callq *%rax</span> is executed, the rax register always has a system call number, which is 0 ~ around 500 (defined in the kernel)**

**set syscall num to rax register**

**callq *%rax**

. . . .

# zpoline

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

**After the binary rewriting**

**<u>Point: Calling Convention</u>**

**When ~~syscall/sysenter~~ <span style="color:red">callq *%rax</span> is executed, the rax register always has a system call number, which is 0 ~ around 500 (defined in the kernel)**

**user-defined hook function**

**set syscall num to rax register**
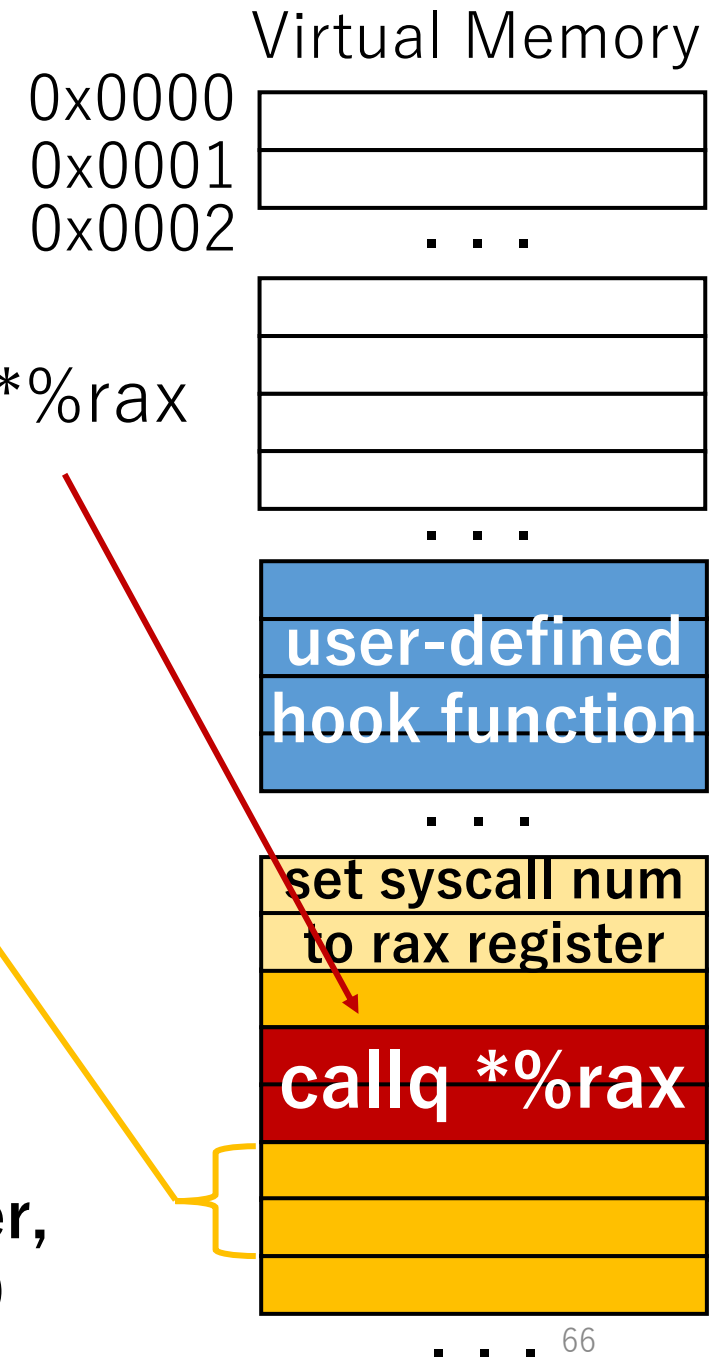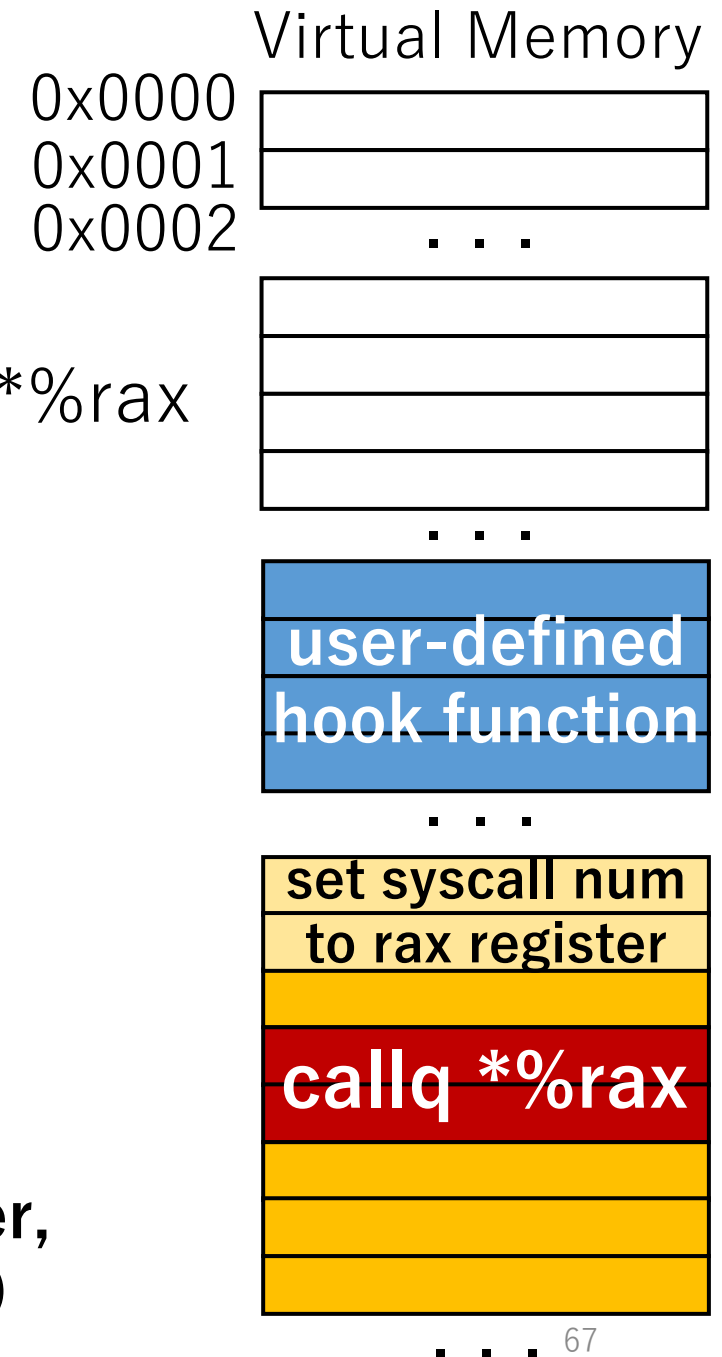
**callq *%rax**

72

# 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)**

Virtual Memory

0x0000
0x0001
0x0002

. . .

. . .

**user-defined hook function**

. . .

**set syscall num to rax register**

**callq *%rax**

. . .

73
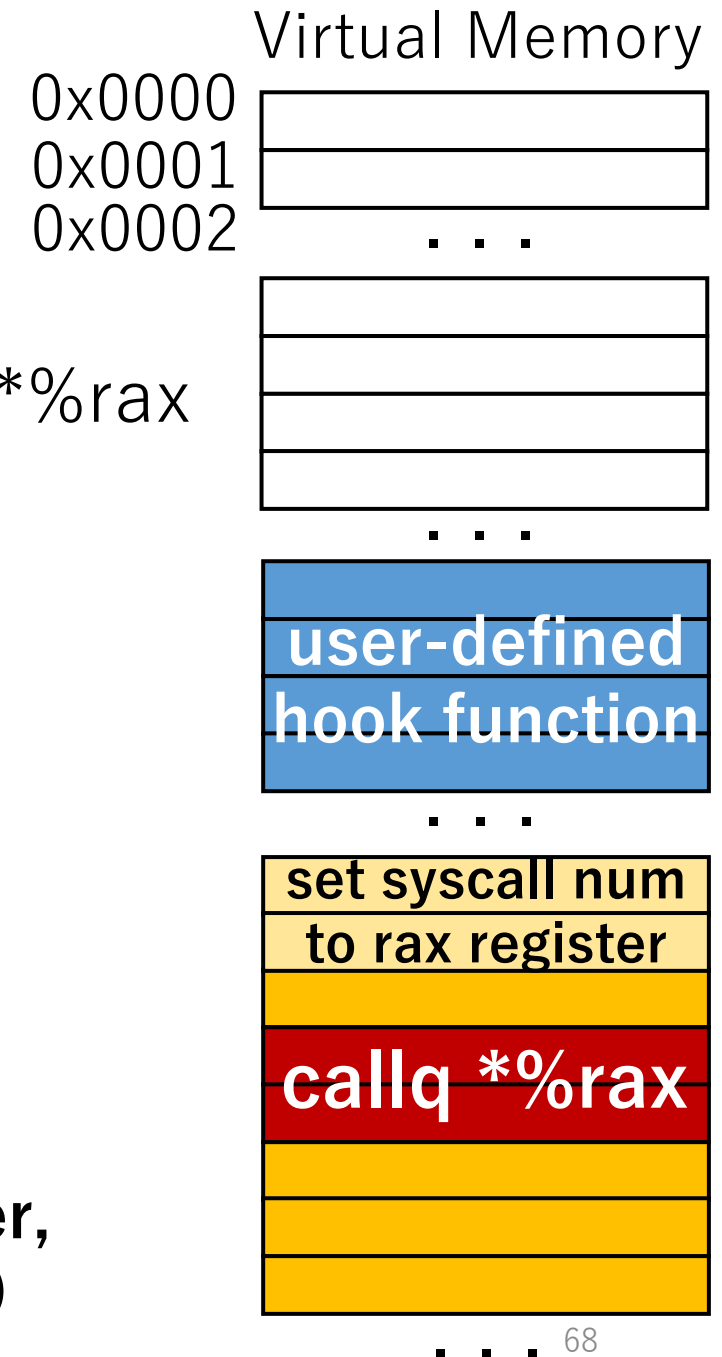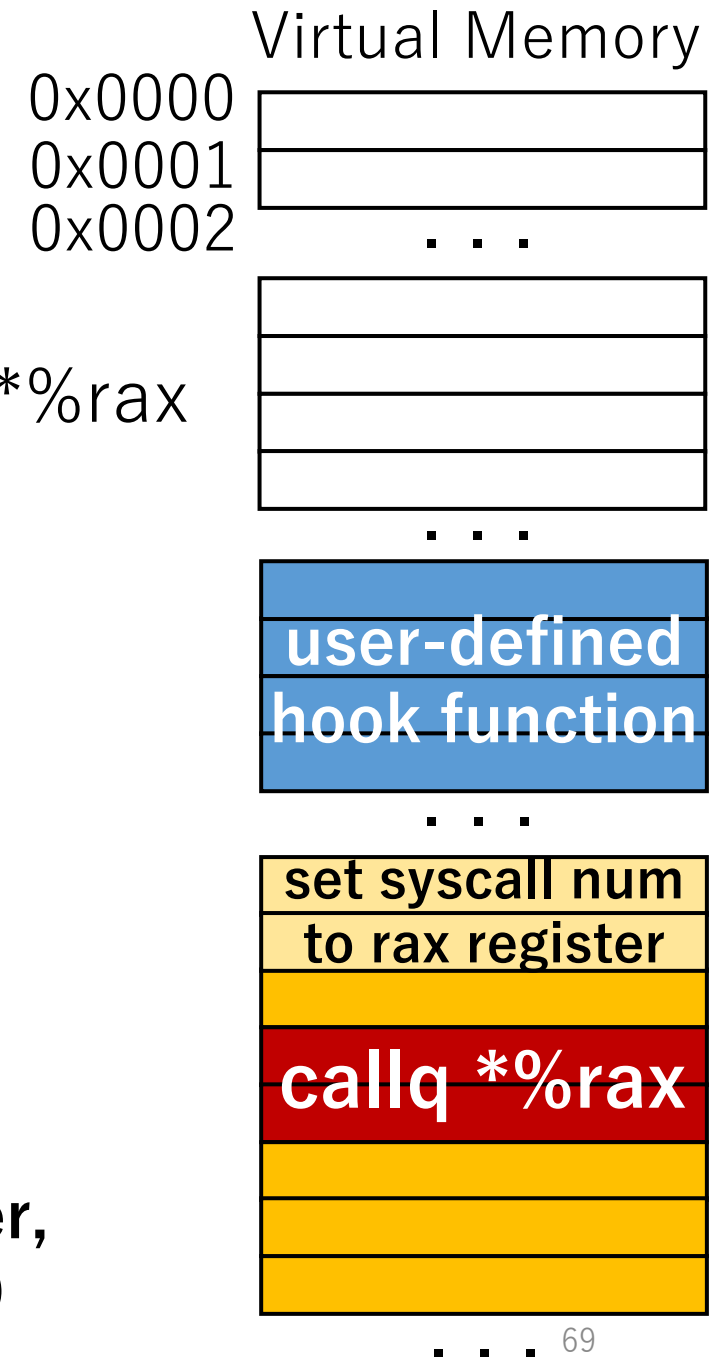
# zpoline

Virtual Memory

0x0000
0x0001
0x0002
around 500

- 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
  - <u>replaced callq *%rax jumps to address 0~around 500</u>

**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)**

user-defined hook function

. . .

set syscall num to rax register

**callq *%rax**

# zpoline

**address range, potentially replaced "callq *%rax" jumps to**
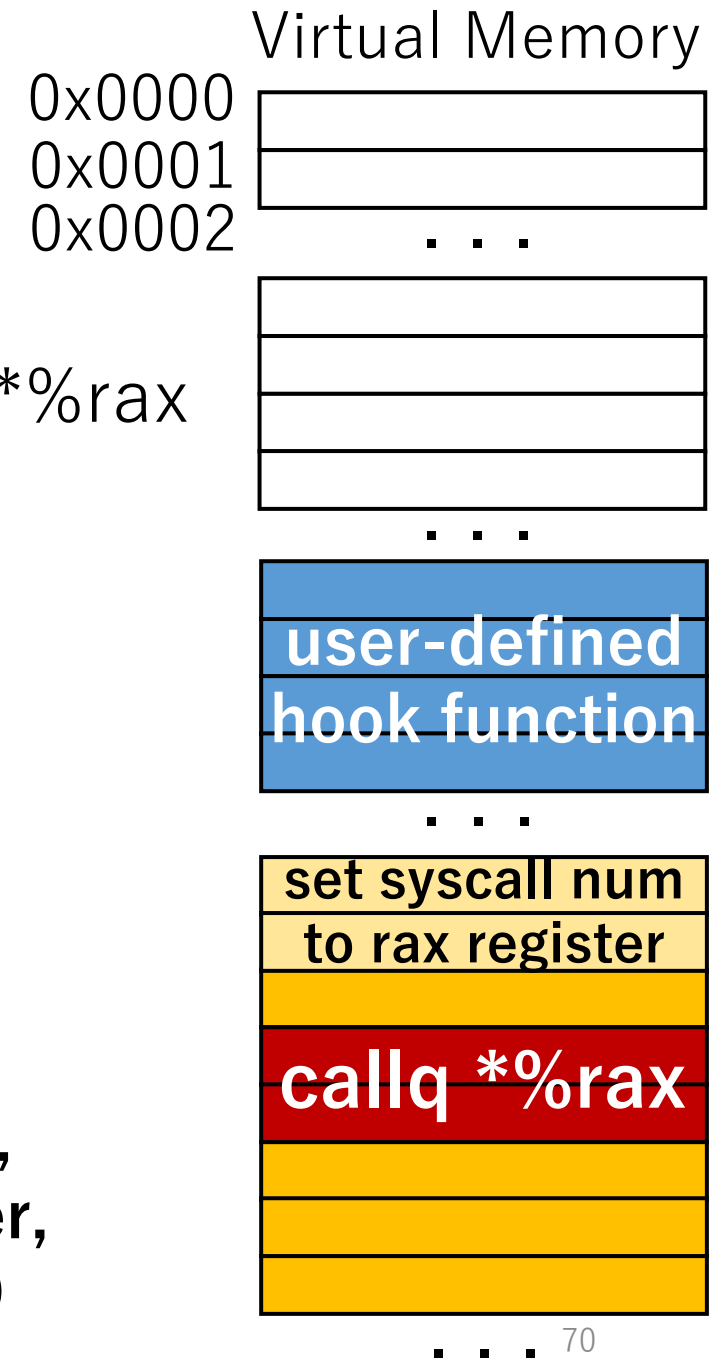( *N* is the max syscall number )

0x0000
0x0001
0x0002
*N*

- 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)**
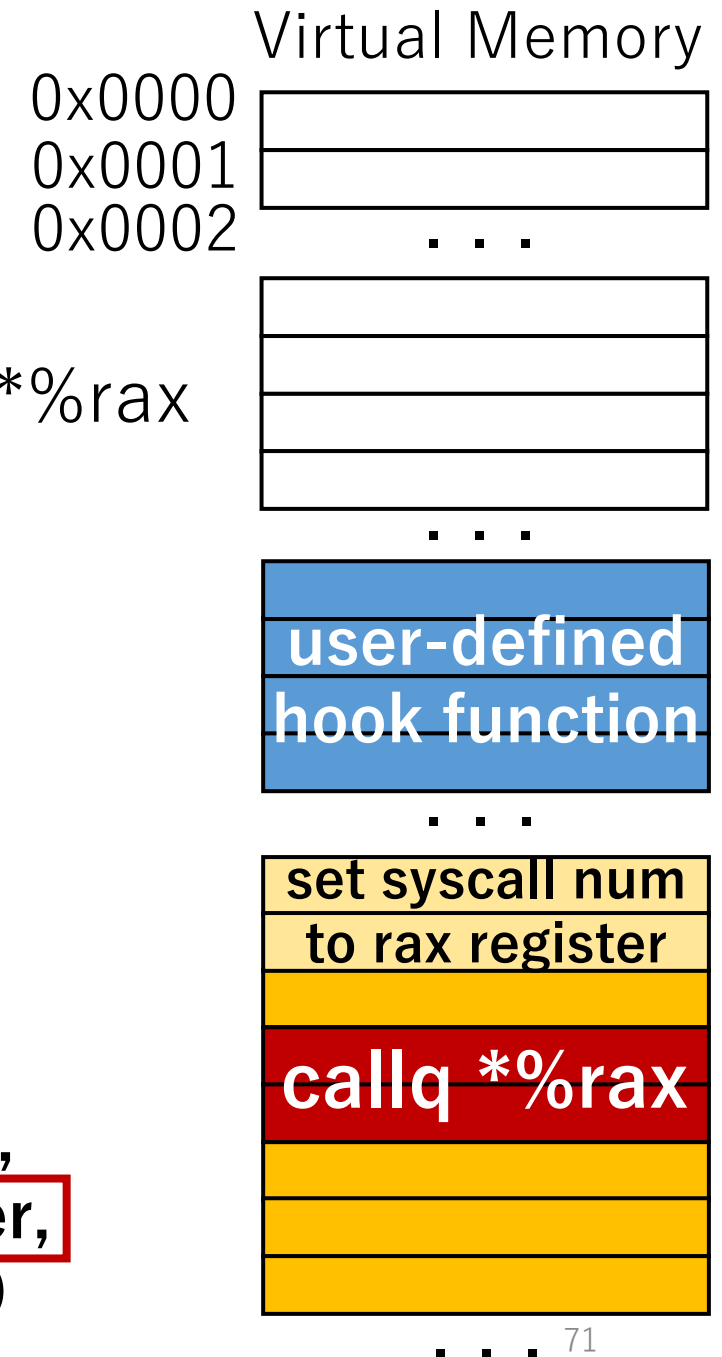
. . .

**user-defined hook function**

. . .

**set syscall num to rax register**

**callq *%rax**

. . .

75

# zpoline

**address range, potentially replaced "callq *%rax" jumps to**
( *N* is the max syscall number )

0x0000
0x0001
0x0002
*N*

- 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

. . .

**user-defined hook function**

. . .

**set syscall num to rax register**

**callq *%rax**

. . .

76

# zpoline

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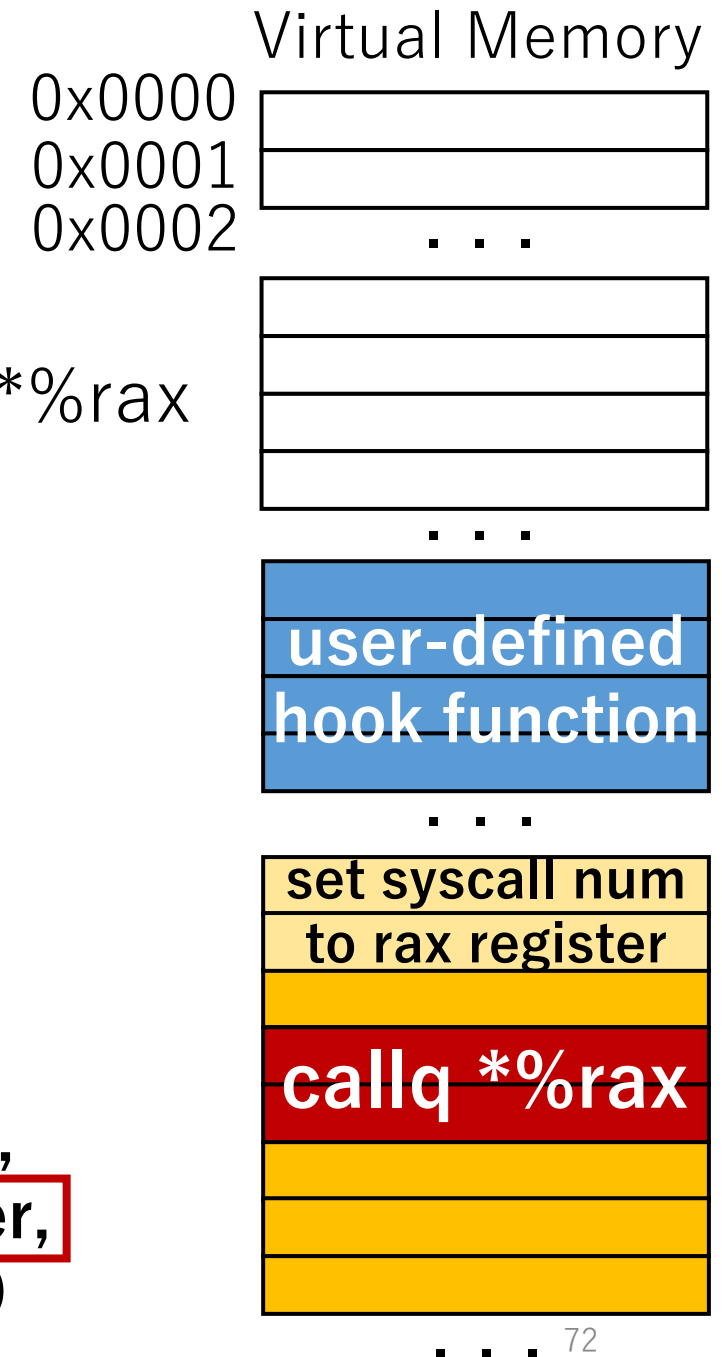( $N$ is the max syscall number )

Virtual Memory

0x0000
0x0001
0x0002
$N$

- 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?**

?

user-defined
hook function

. . .

set syscall num
to rax register

callq *%rax

# zpoline

**address range, potentially replaced "callq *%rax" jumps to**
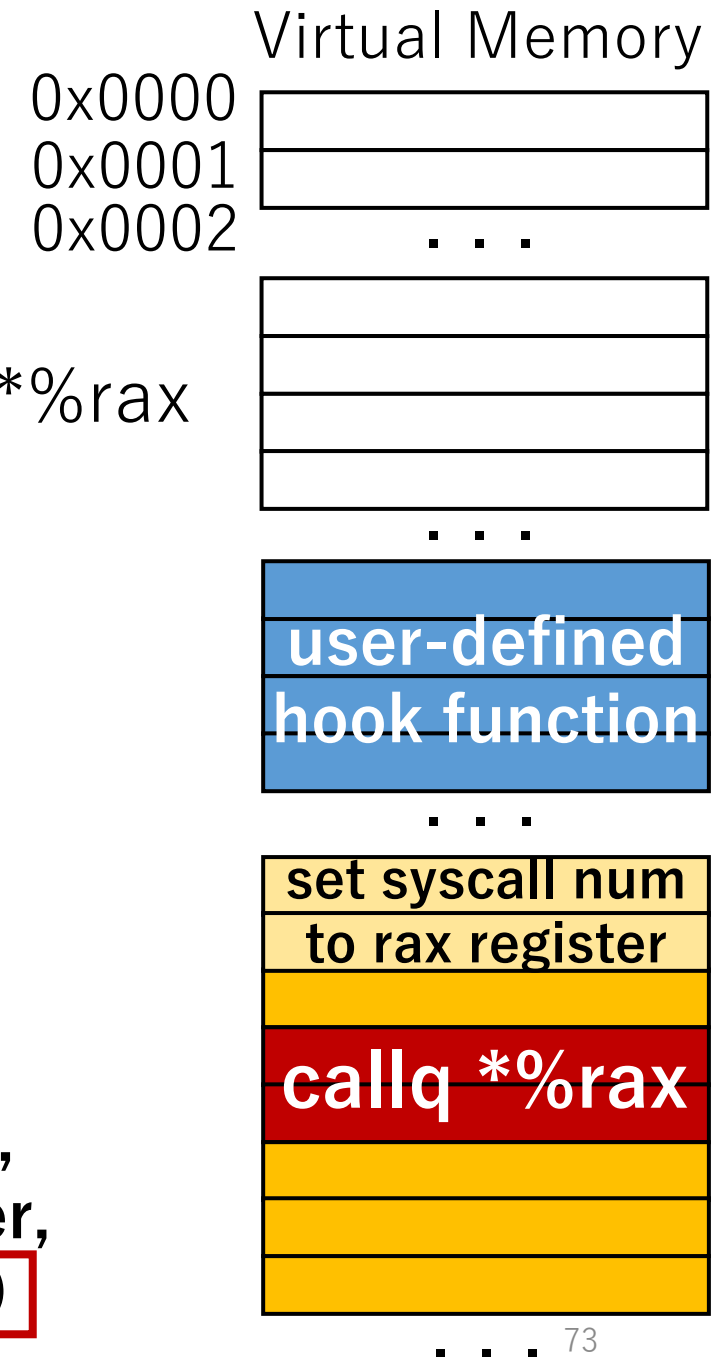( *N* is the max syscall number )

Virtual Memory

0x0000
0x0001
0x0002
*N*

- 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

trampoline code

user-defined hook function

set syscall num to rax register

callq *%rax

# zpoline

**address range, potentially replaced "callq *%rax" jumps to**
( *N* is the max syscall number )

0x0000
0x0001
0x0002
*N*

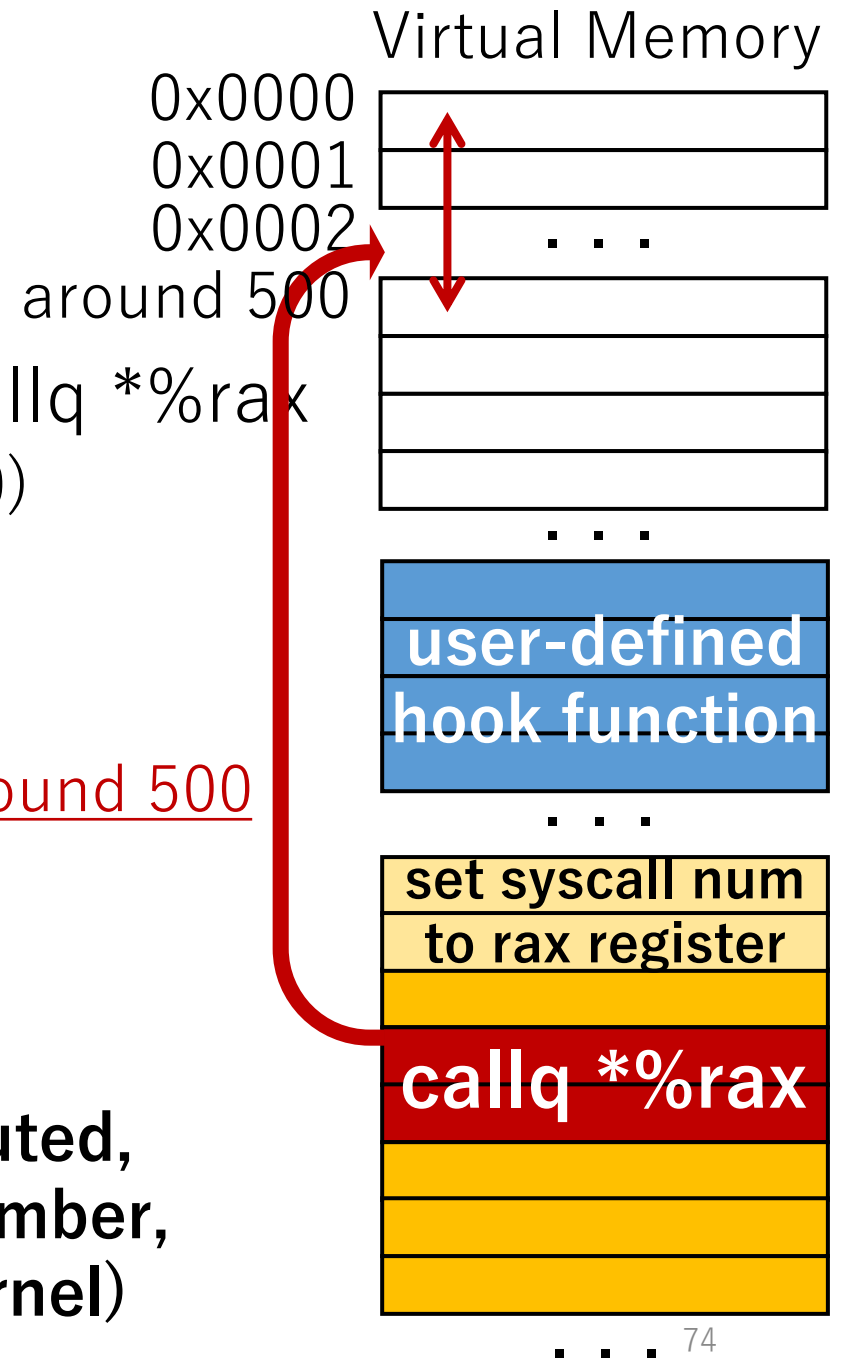| nop |
| --- |
| nop |

. . .

| nop |
| --- |
| |
| |
| |

. . .

- 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

| **user-defined hook function** |
| --- |

. . .

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

| **set syscall num to rax register** |
| --- |
| |

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

| **callq *%rax** |
| --- |
| |
| |
| |

. . .

# zpoline

**address range, potentially replaced "callq *%rax" jumps to** ( *N* is the max syscall number )

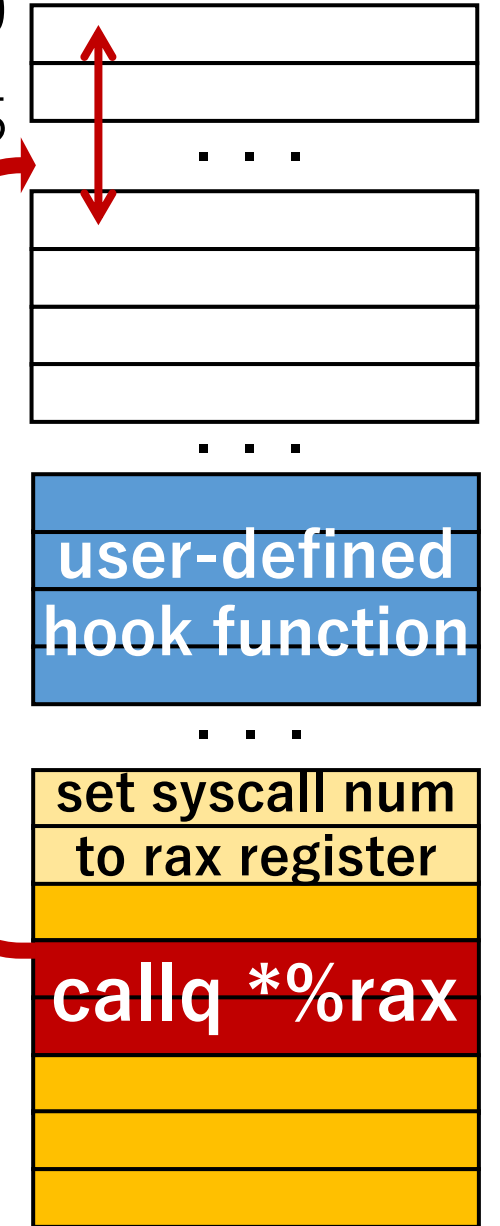0x0000
0x0001
0x0002
*N*

- 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)

Virtual Memory

fall through

nop
nop
. . .
nop

user-defined hook function
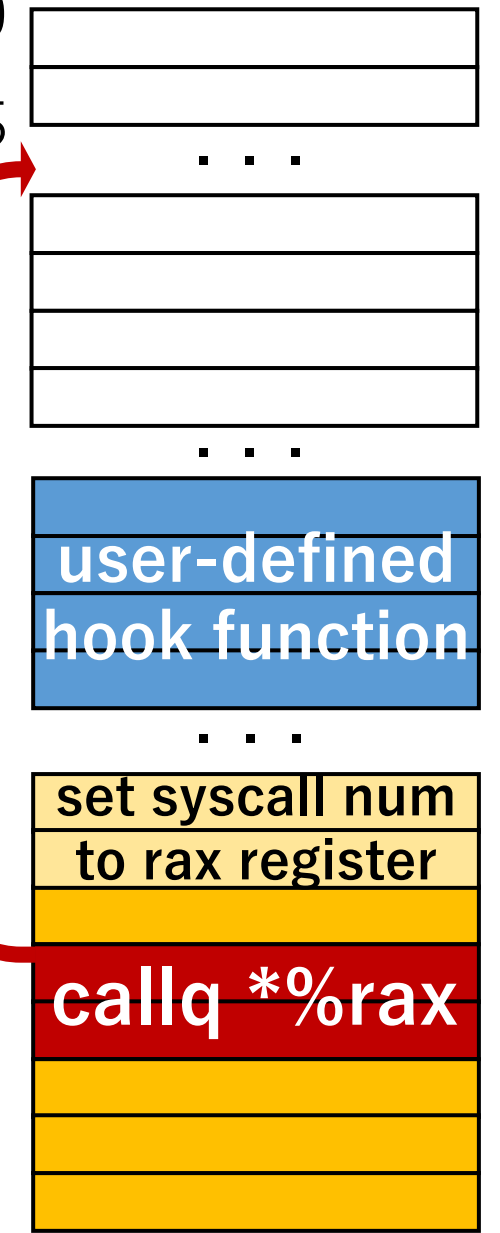
. . .

set syscall num to rax register

callq *%rax

# zpoline

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

```
0x0000
0x0001
0x0002
N
```
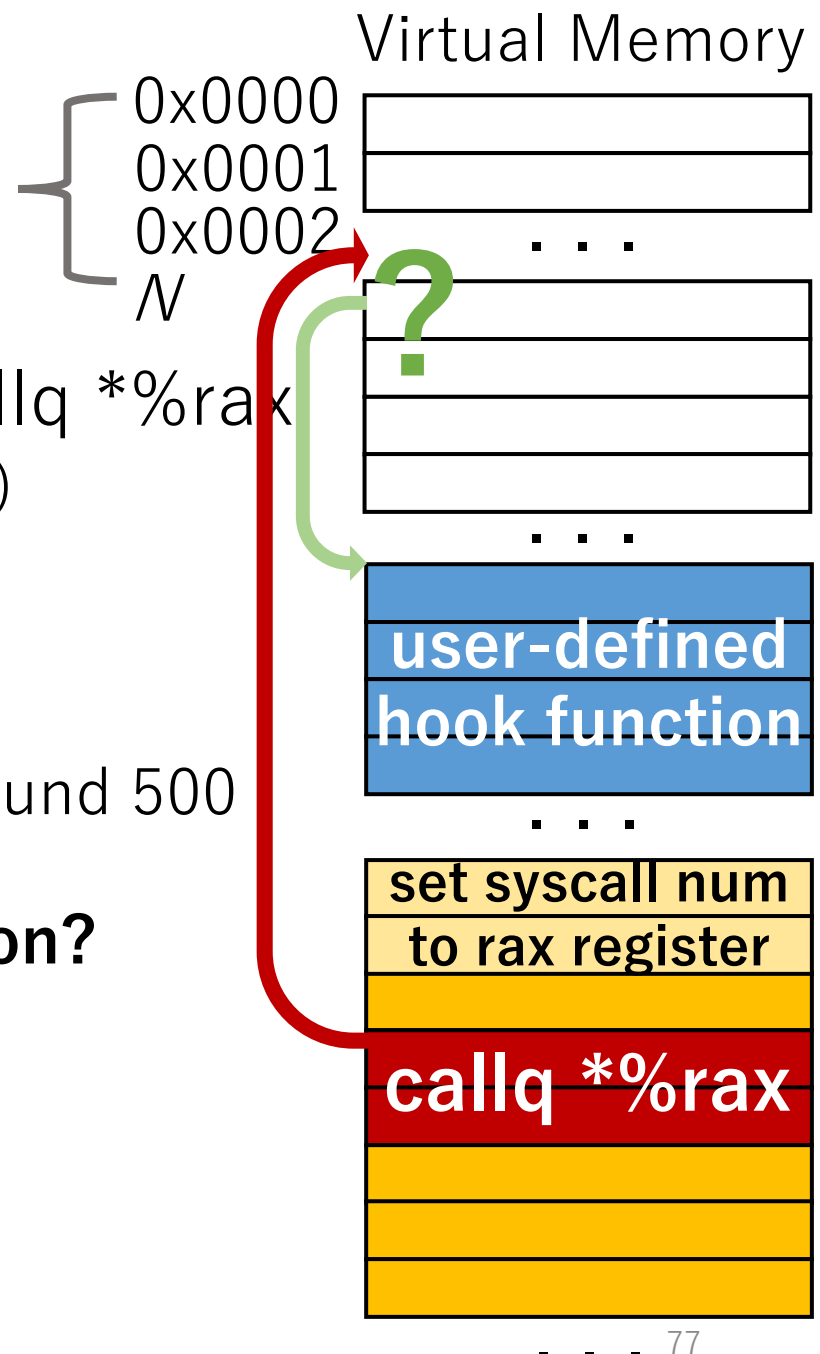
- zpoline replaces syscall/sysenter with callq *%rax
  - callq *%rax is a **2-byte** instruction (0xff 0xd0)
    - Neighbour instructions are not overwritten
  - callq *%rax is an instruction to jump to the address stored in the rax register
  - replaced callq *%rax jumps to address 0~around 500

**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

fall through

nop
nop
. . .
nop
**jump to hook function**

. . .

**user-defined hook function**

. . .

**set syscall num to rax register**

**callq *%rax**

. . .

# zpoline

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

0x0000
0x0001
0x0002
$N$

fall through

nop
nop
. . .
nop
jump to hook function

. . .

user-defined hook function

. . .

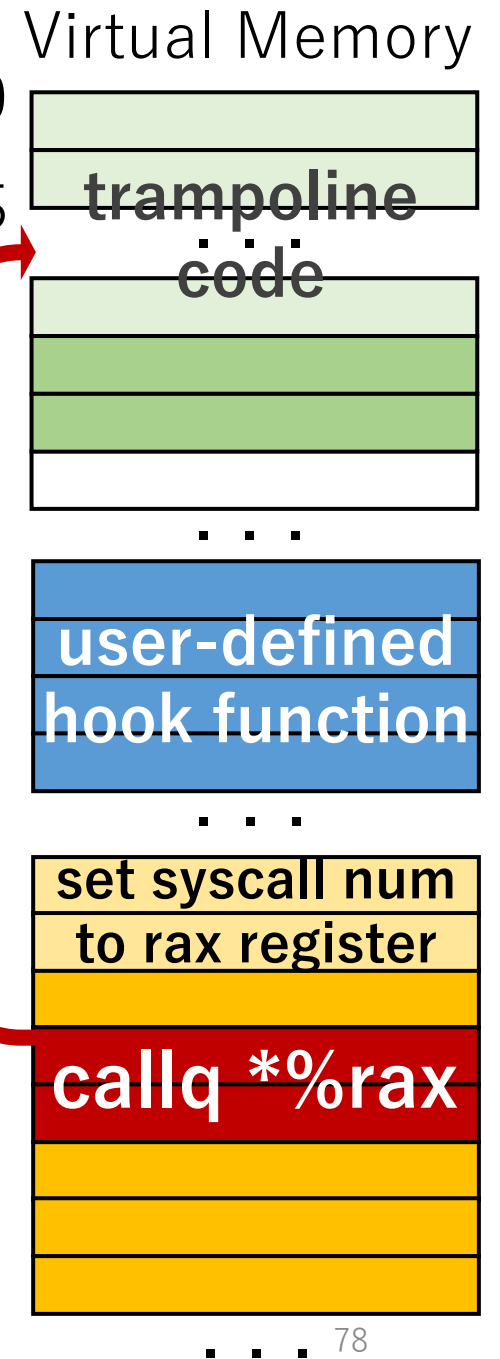set syscall num to rax register

callq *%rax

. . .

- 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
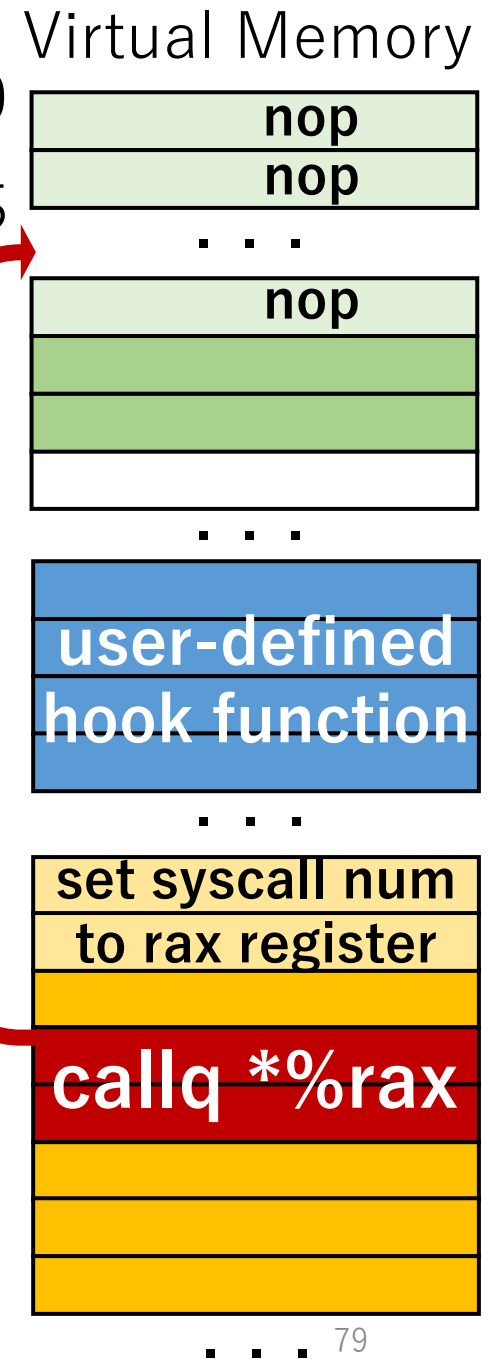  - replaced callq *%rax jumps to address 0~around 500

**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!**

# zpoline

← tram**poline** code at address 0 (**z**ero)

( *N* is the max syscall number )

**nop**
**nop**
. . .
**nop**

fall through

**jump to hook function**

- 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
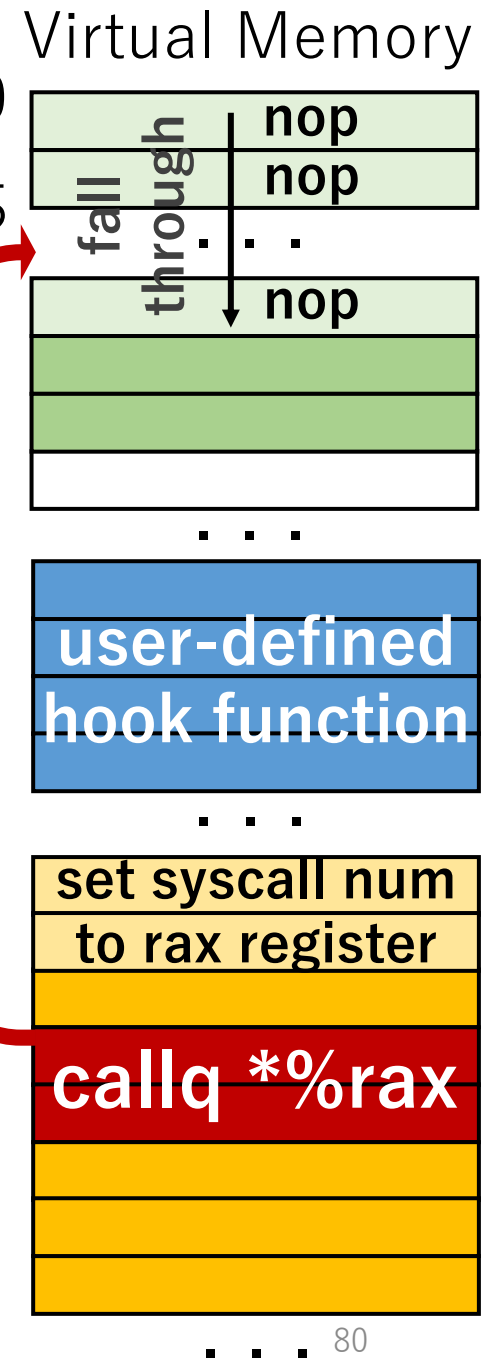  - replaced callq *%rax jumps to address 0~around 500

. . .

**user-defined hook function**

. . .

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

**set syscall num to rax register**

- 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

**callq *%rax**

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

# NULL Access Termination

- A buggy program may access NULL (address 0)

0x0000
0x0001
0x0002
N

fall through

nop
nop
. . .
nop
jump to
hook function

. . .

user-defined
hook function

. . .

set syscall num
to rax register

callq *%rax

. . .

# NULL Access Termination

- A buggy program may access NULL (address 0)

| | |
|---|---|
| 0x0000 | nop |
| 0x0001 | nop |
| 0x0002 | . . . |
| N | nop |

fall through

jump to
hook function

. . .

user-defined
hook function

. . .

set syscall num
to rax register

callq *%rax

Bug
Access NULL

. . .

# NULL Access Termination

- A buggy program may access NULL (address 0)

0x0000
0x0001
0x0002
*N*

| fall through | nop |
|---|---|
| | nop |
| | . . . |
| | nop |
| jump to | |
| hook function | |

. . .

**user-defined hook function**

. . .

**set syscall num to rax register**

**callq *%rax**

Bug
Access NULL

. . .

# NULL Access Termination

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

0x0000
0x0001
0x0002
*N*

nop
nop
. . .
nop
jump to
hook function

fall through

. . .

user-defined
hook function

. . .

set syscall num
to rax register

callq *%rax

Bug
Access NULL

. . .

# NULL Access Termination

Virtual Memory

0x0000
0x0001
0x0002
*N*

- 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

page fault!
no physical
memory
mapping!

. . .

user-defined
hook function

. . .

set syscall num
to rax register

callq *%rax

Bug
Access NULL

. . .

# NULL Access Termination

0x0000
0x0001
0x0002
N

fall
through

nop
nop
. . .
nop
jump to
hook function

. . .

user-defined
hook function

. . .

set syscall num
to rax register

callq *%rax

Bug
Access NULL

. . .

- 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

89

# 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

0x0000
0x0001
0x0002
*N*

fall through

nop
nop

· · ·

nop

**jump to hook function**

· · ·

**user-defined hook function**

· · ·

**set syscall num to rax register**

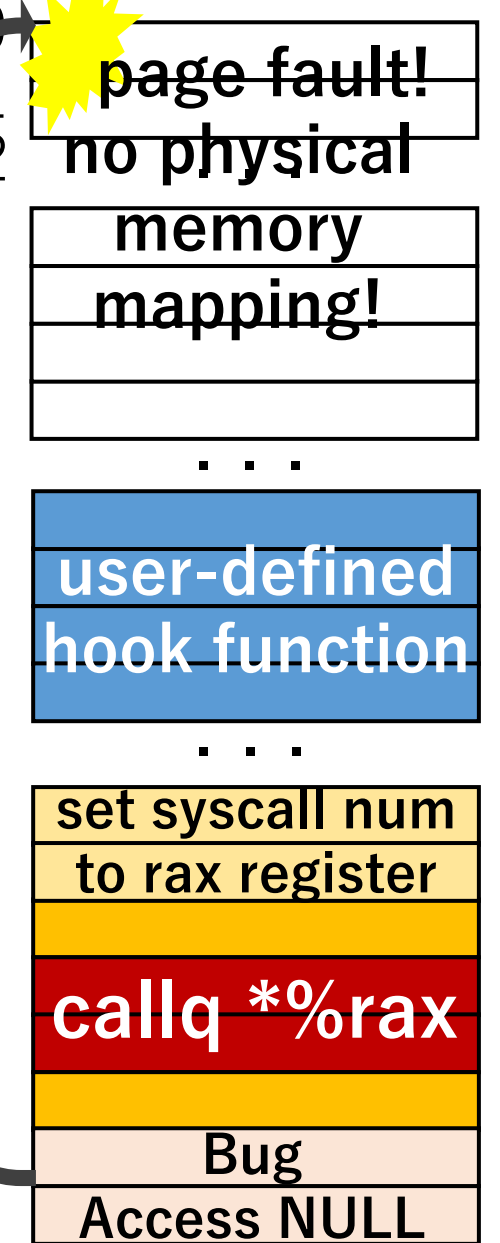**callq *%rax**

Bug
**Access NULL**

· · ·

# 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?**

Virtual Memory

0x0000
0x0001
0x0002
$N$

fall through

nop
nop
. . .
nop
jump to
hook function

. . .

user-defined
hook function

. . .

set syscall num
to rax register

callq *%rax

Bug
Access NULL

. . .

# NULL Access Termination

- Memory access: read / write / execute



Virtual Memory

| 0x0000 | nop |
| 0x0001 | nop |
| 0x0002 | . . . |
| N | nop |
| | jump to |
| | hook function |

. . .

user-defined
hook function

. . .

| | set syscall num |
| | to rax register |
| A | callq *%rax |
| B | |
| | Bug |
| | Access NULL |

. . .

# NULL Access Termination

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

0x0000
0x0001
0x0002
*N*

eXecute
Only
Memory
(XOM)

. . .

user-defined
hook function

. . .

set syscall num
to rax register

*A*

**callq *%rax**

*B*

Bug

Access NULL

. . .

93

# 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

Virtual Memory

0x0000
0x0001
0x0002
*N*

eXecute
Only
Memory
(XOM)

nop
nop
hook function

user-defined
hook function

set syscall num
to rax register

*A*

callq *%rax

*B*

Bug
Access NULL

# 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

0x0000
0x0001
0x0002
$N$

eXecute Only Memory (XOM)

user-defined hook function

. . .

set syscall num to rax register

$A$

callq *%rax

$B$

Bug
Access NULL

. . .

95

# 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. during the binary rewriting phase,
       we collect the addresses of replaced syscall/sysenter

Virtual Memory

0x0000
0x0001
0x0002
$N$

eXecute Only Memory (XOM)

. . .

user-defined hook function

. . .

set syscall num to rax register

$A$

callq *%rax

$B$

Bug
Access NULL

. . .

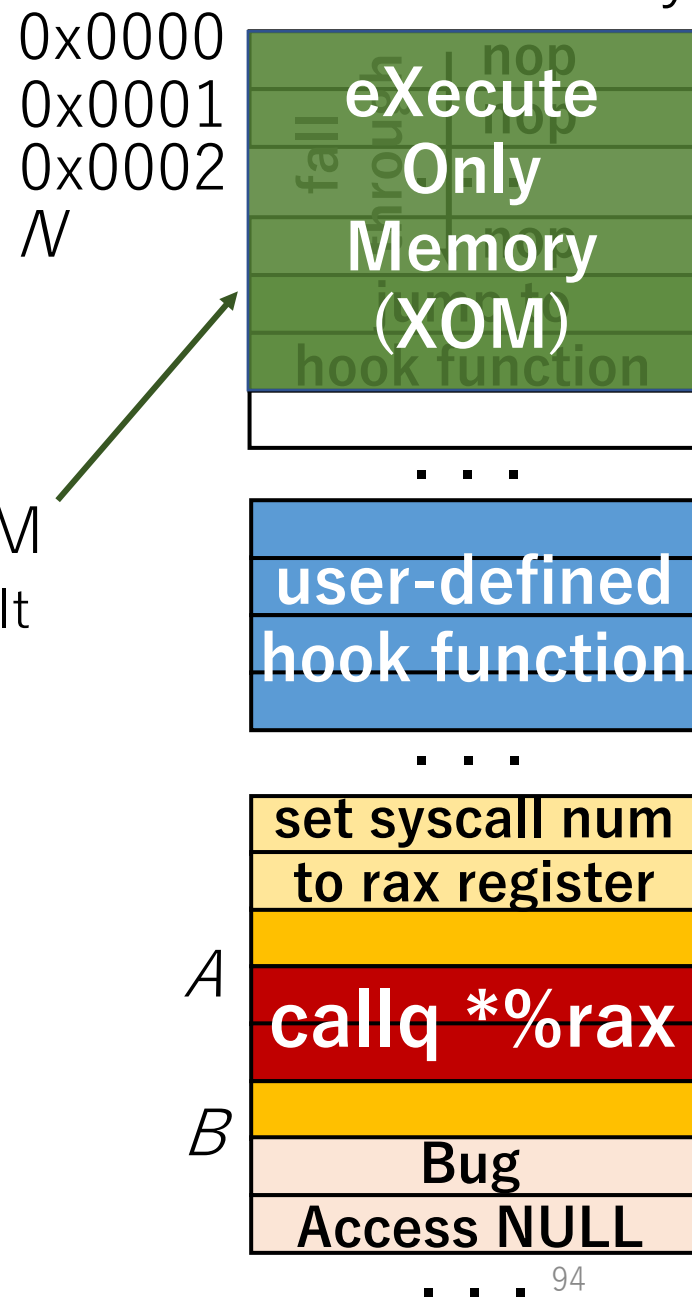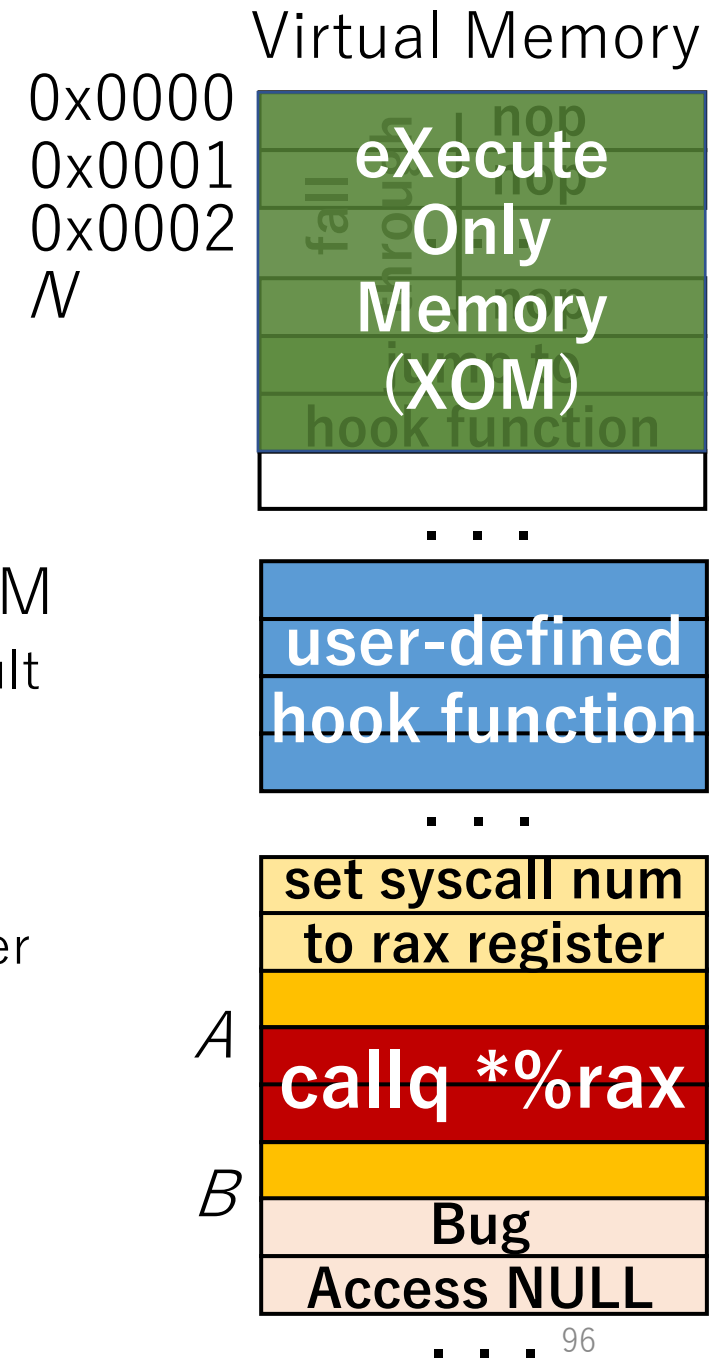# NULL Access Termination

- 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. during the binary rewriting phase,
           we collect the addresses of replaced syscall/sysenter

**During binary rewriting phase ...**

Virtual Memory

0x0000
0x0001
0x0002
$N$

eXecute
Only
Memory
(XOM)

fall through
hook function

. . .

**user-defined
hook function**

. . .

**set syscall num
to rax register**

$A$

**callq *%rax**

$B$

**Bug**

**Access NULL**

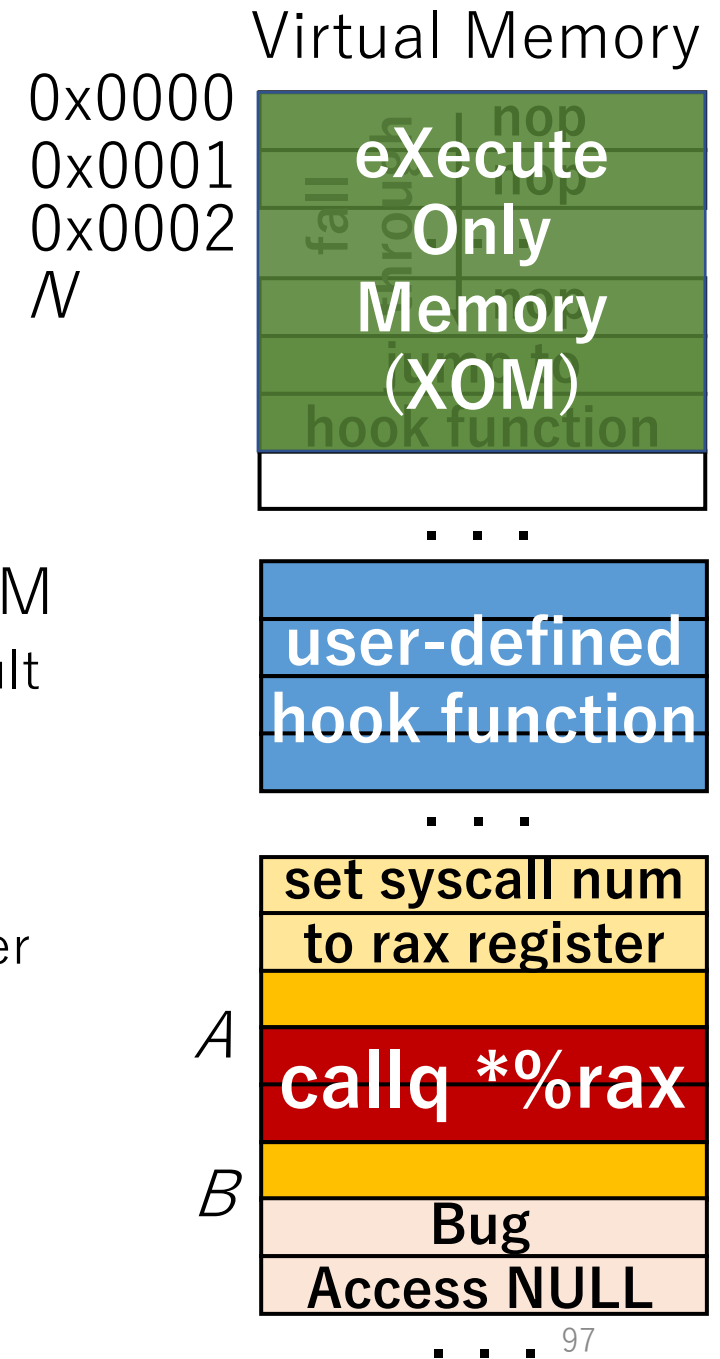. . .

# 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. during the binary rewriting phase,
       we collect the addresses of replaced syscall/sysenter

**During binary rewriting phase ...**

List of replaced addresses : [...]



Virtual Memory

0x0000
0x0001
0x0002
$N$

eXecute Only Memory (XOM)

user-defined hook function

set syscall num to rax register

$A$
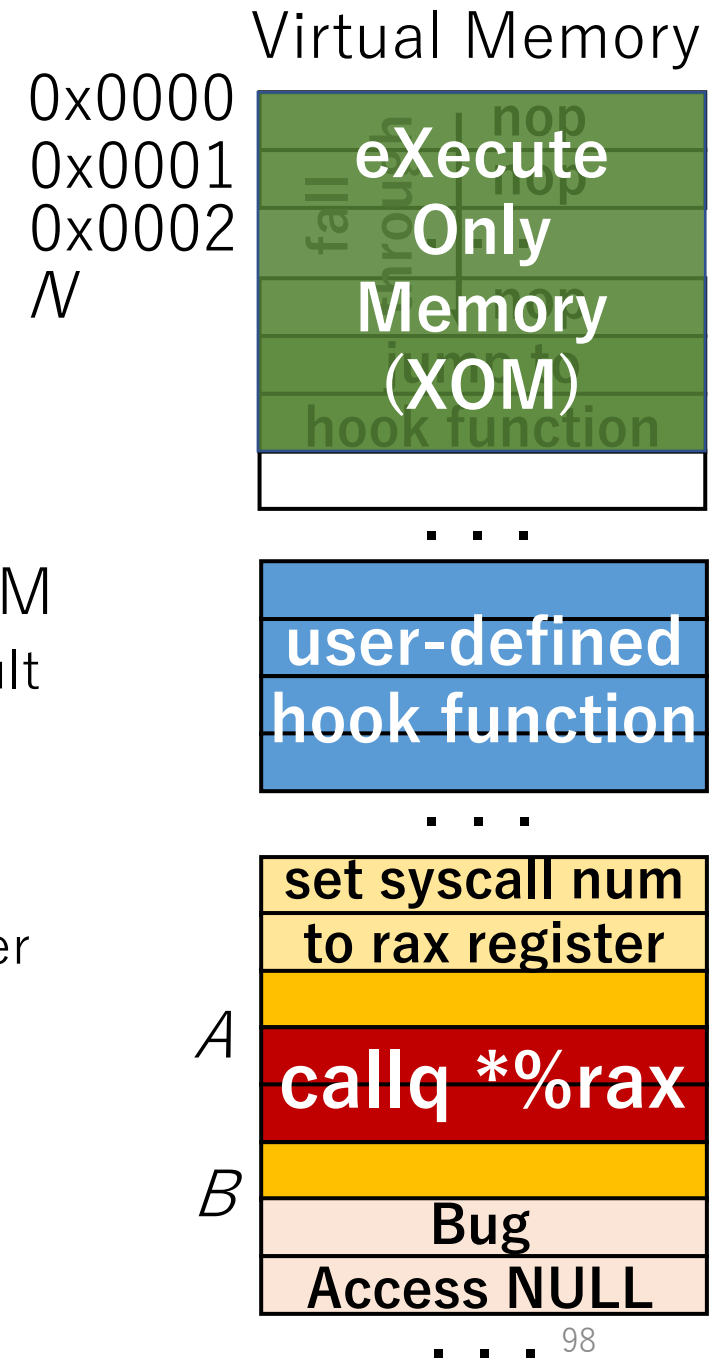
callq *%rax

$B$

Bug

Access NULL

# 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. during the binary rewriting phase,
       we collect the addresses of replaced syscall/sysenter

**During binary rewriting phase ...**

List of replaced addresses : [...]

0x0000
0x0001
0x0002
$N$

eXecute
Only
Memory
(XOM)

user-defined
hook function

set syscall num
to rax register

$A$

callq *%rax
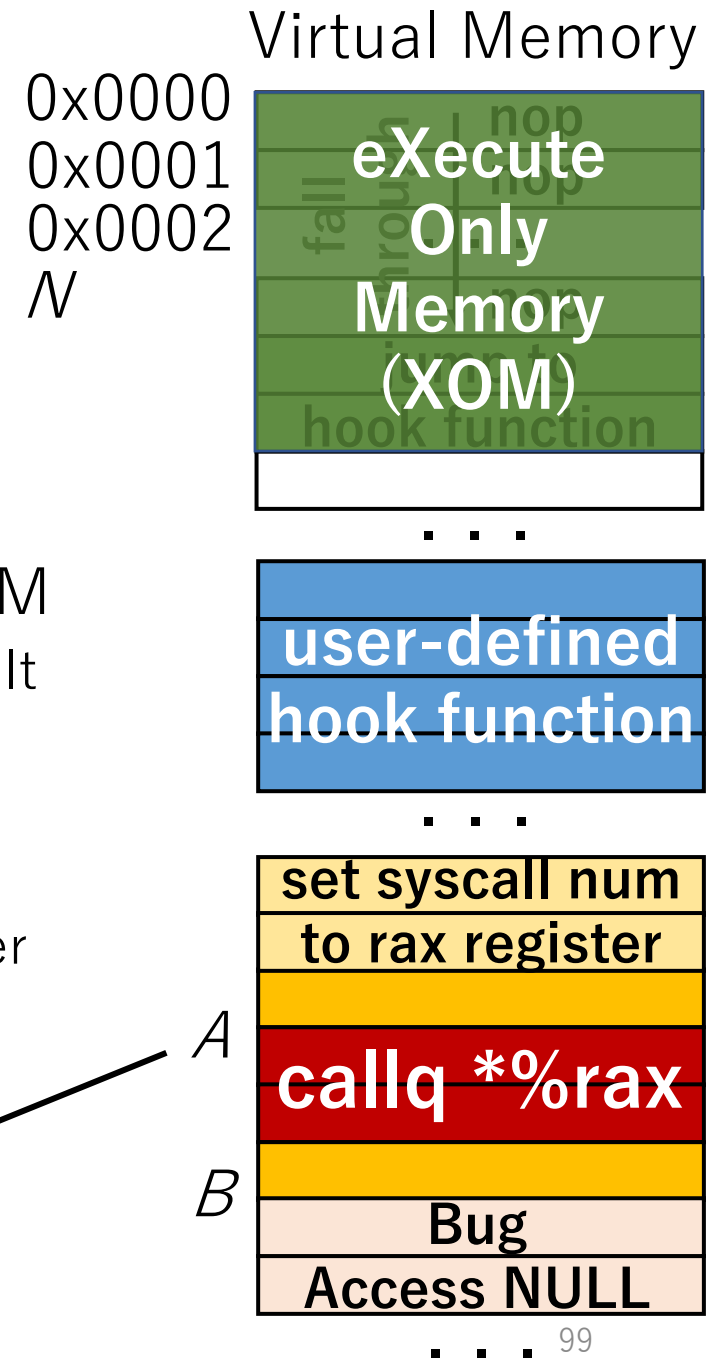
$B$

Bug
Access NULL

99

# 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. during the binary rewriting phase,
       we collect the addresses of replaced syscall/sysenter

**During binary rewriting phase ...**

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

Virtual Memory

0x0000
0x0001
0x0002
$N$

eXecute Only Memory (XOM)

nop
nop
fall through
hook function

· · ·

user-defined hook function

· · ·

set syscall num to rax register

$A$
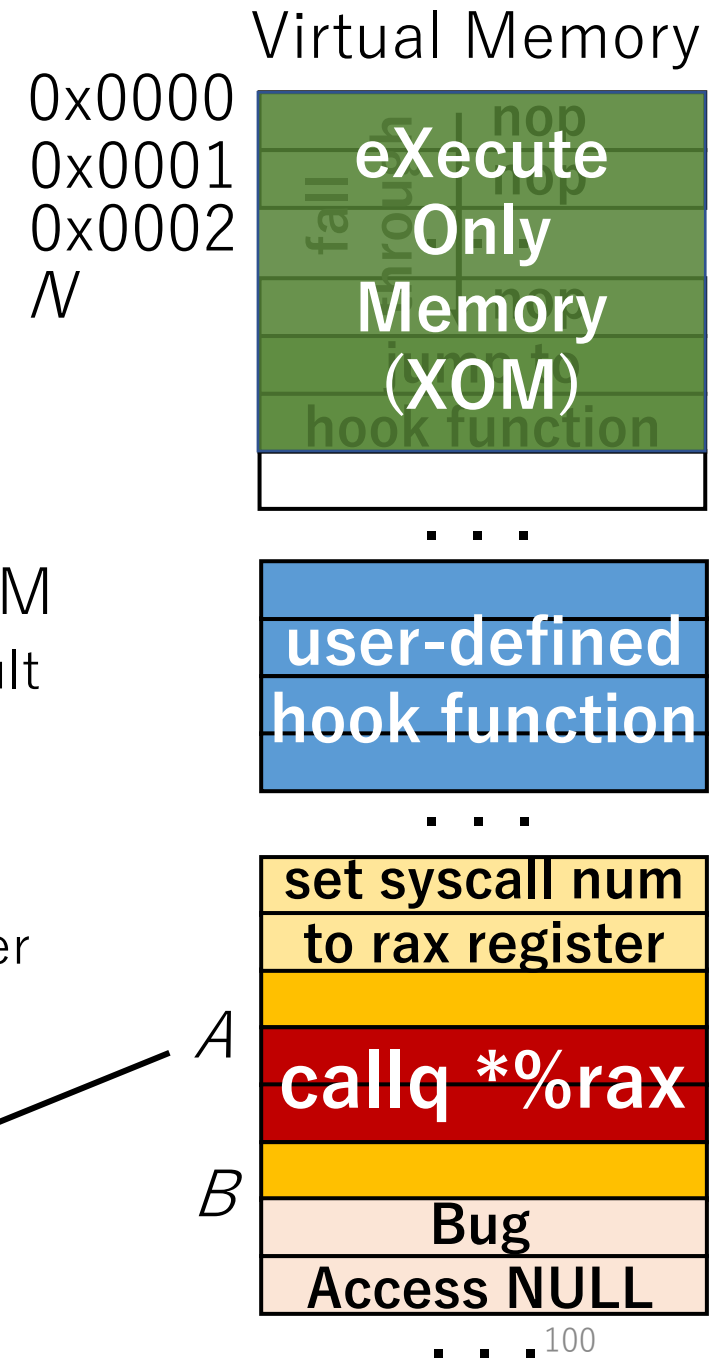
callq *%rax

$B$

Bug
Access NULL

· · ·

# 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. during the binary rewriting phase,
       we collect the addresses of replaced syscall/sysenter
    2. at runtime, in the hook function,
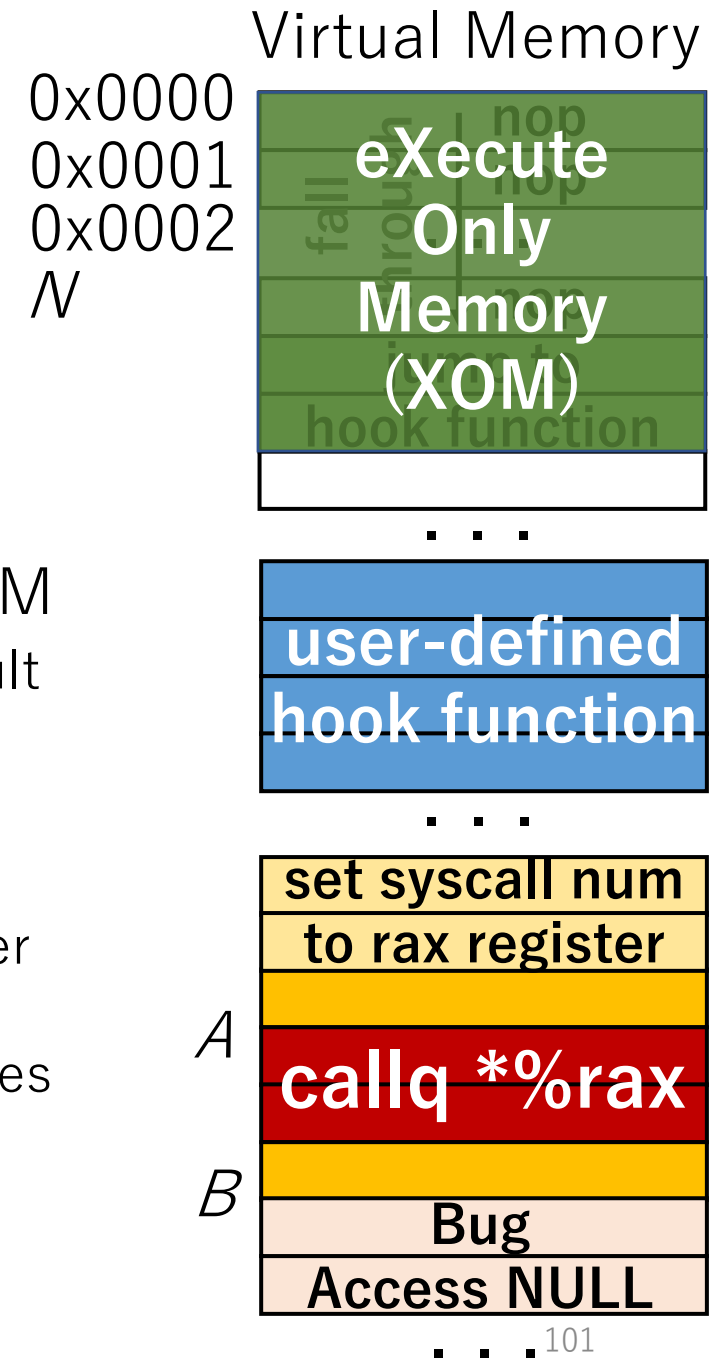       we check if the caller is one of the replaced addresses

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

Virtual Memory

0x0000
0x0001
0x0002
$N$

eXecute Only Memory (XOM)

user-defined hook function

set syscall num to rax register

$A$

callq *%rax

$B$

Bug
Access NULL

# 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. during the binary rewriting phase,
       we collect the addresses of replaced syscall/sysenter

    2. at runtime, in the hook function,
       we check if the caller is one of the replaced addresses
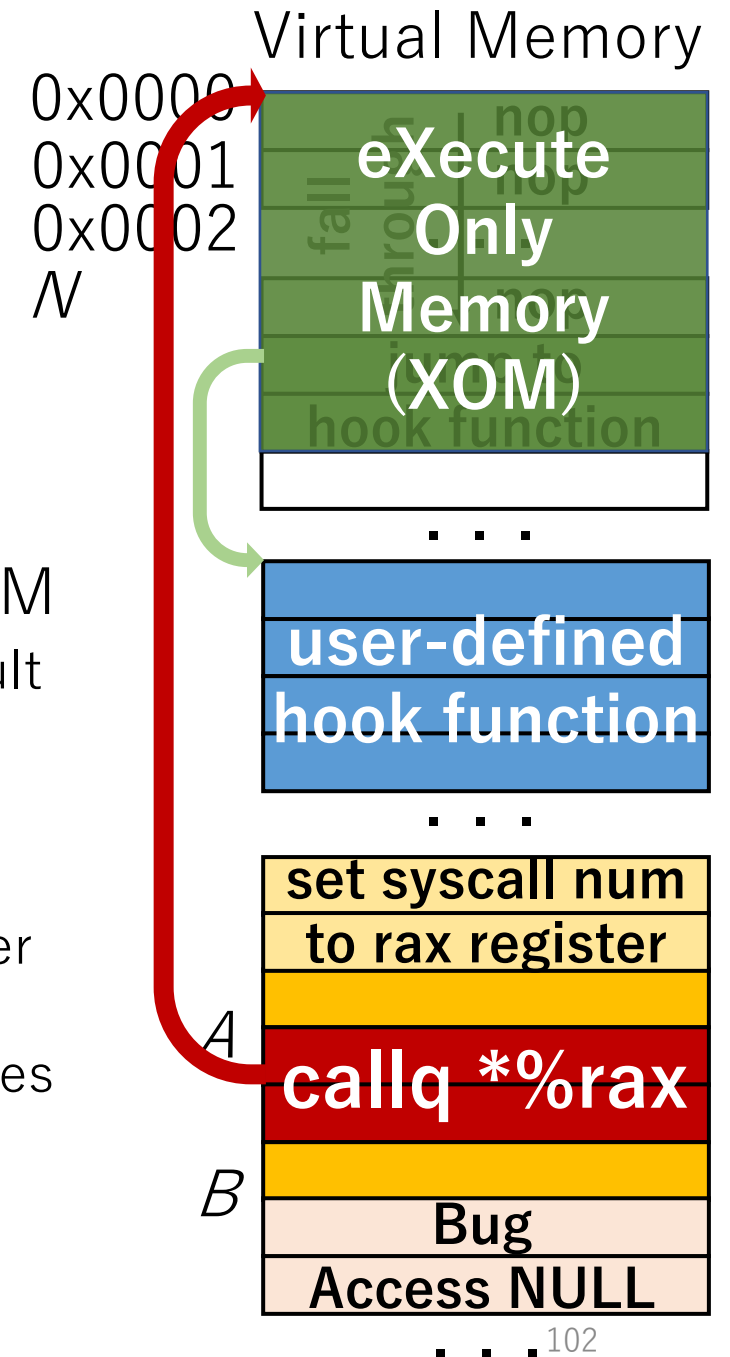
    List of replaced addresses : $[A, ...]$



Virtual Memory

0x0000
0x0001
0x0002
$N$

eXecute Only Memory (XOM)

nop
nop
fall through
hook function

user-defined hook function

set syscall num to rax register

$A$

callq *%rax

$B$

Bug
Access NULL

# NULL Access Termination

## At runtime ...

- Memory acces[s]

- Solution
  - read/write: [...] XOM
    - read/write access to the tra[...] code causes a fault
    - This can be done by mprotect() system call
  - execute: check the caller address
    1. during the binary rewriting phase,
       we collect the addresses of replaced syscall/sysenter
    2. at runtime, in the hook function,
       we check if the caller is one of the replaced addresses
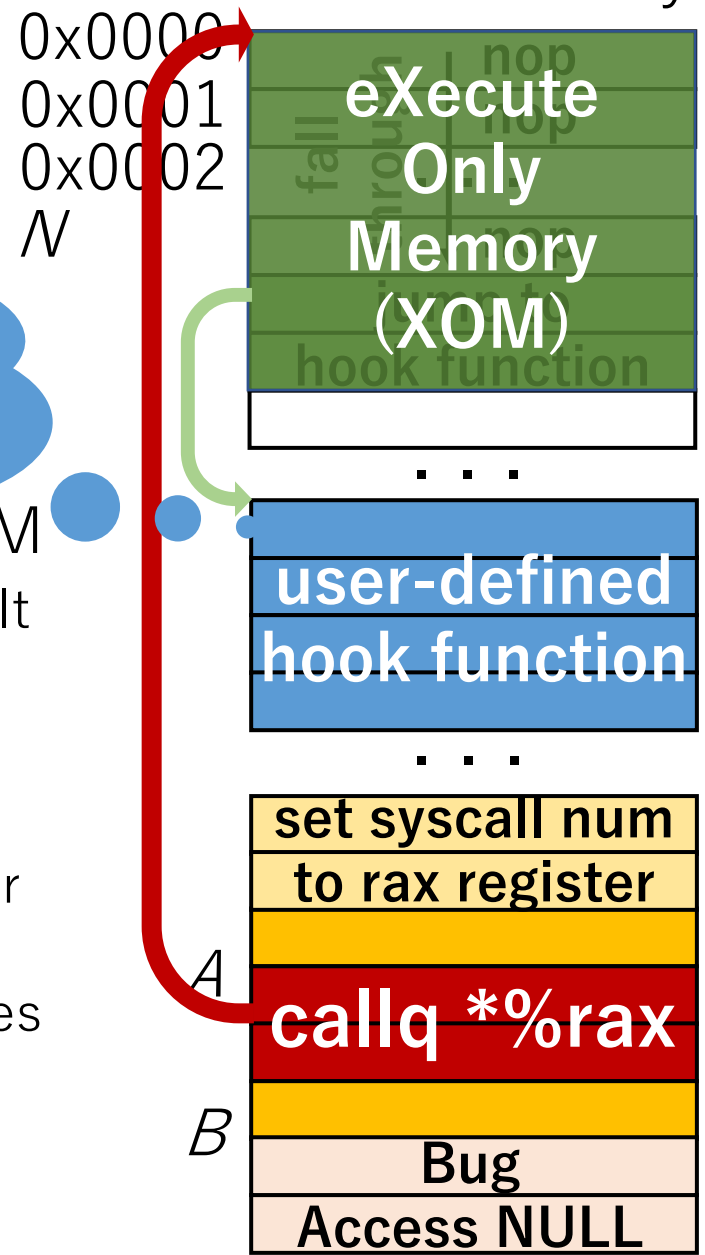
List of replaced addresses : [$A$ , ...]

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

Virtual Memory

0x0000
0x0001
0x0002
$N$

eXecute
Only
Memory
(XOM)

nop
nop
fall through

hook function

user-defined
hook function

. . .

set syscall num
to rax register

$A$

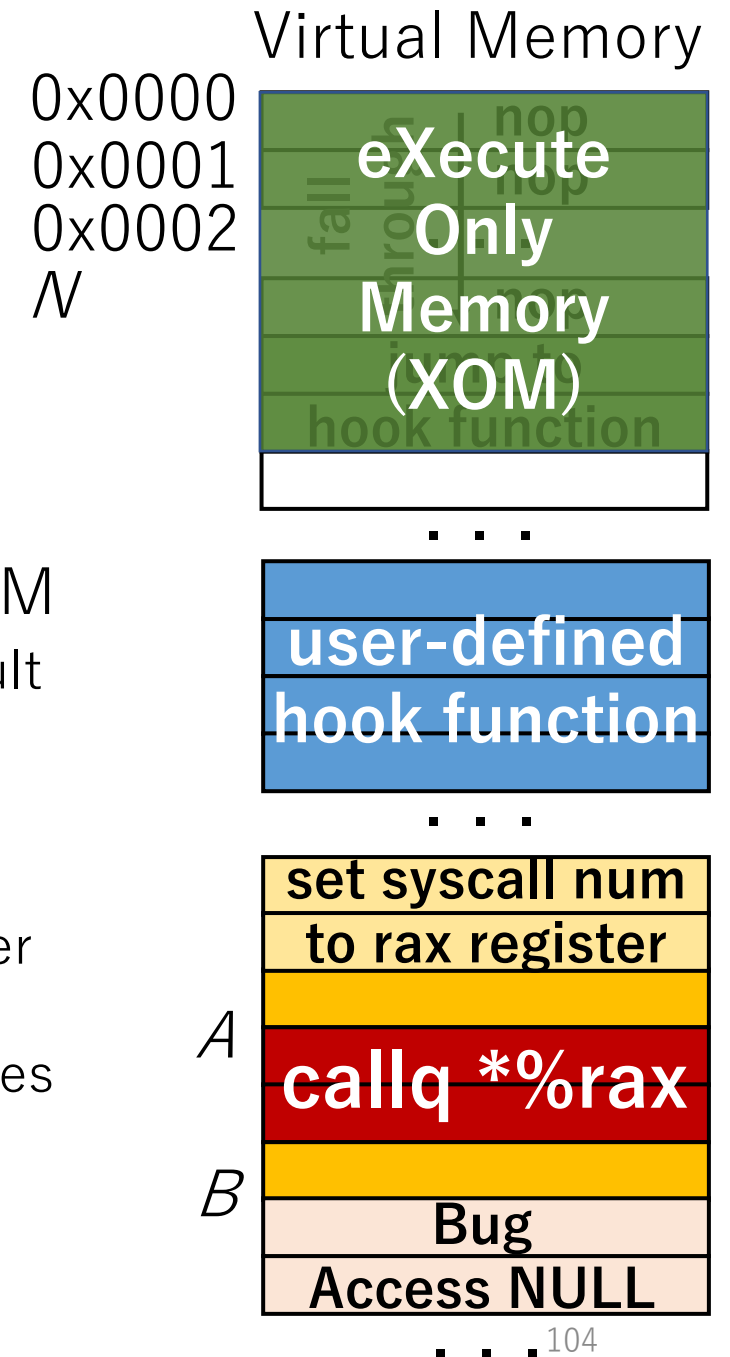callq *%rax

$B$

Bug
Access NULL

. . .

# 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. during the binary rewriting phase,
       we collect the addresses of replaced syscall/sysenter

    2. at runtime, in the hook function,
       we check if the caller is one of the replaced addresses
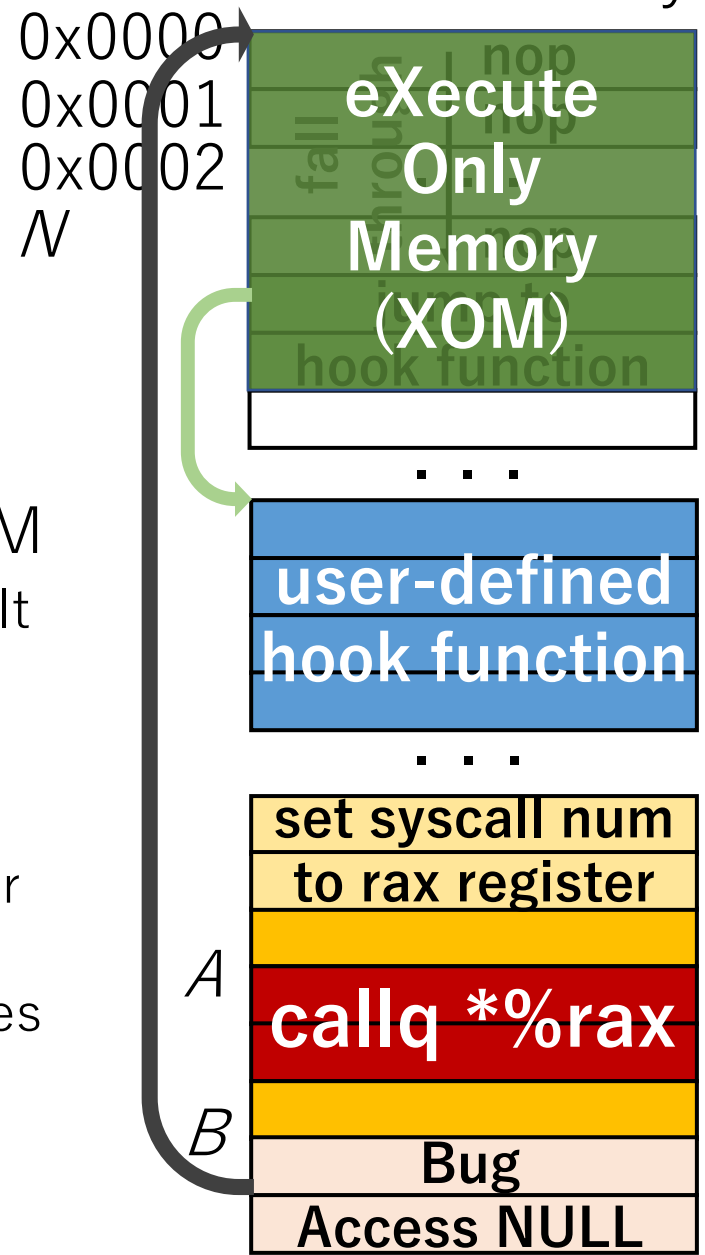
  List of replaced addresses : $[A, ...]$

0x0000
0x0001
0x0002
$N$

eXecute
Only
Memory
(XOM)

hook function

. . .

user-defined
hook function

. . .

set syscall num
to rax register

$A$

**callq *%rax**

$B$

Bug

Access NULL

. . .

# 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. during the binary rewriting phase,
       we collect the addresses of replaced syscall/sysenter
    2. at runtime, in the hook function,
       we check if the caller is one of the replaced addresses
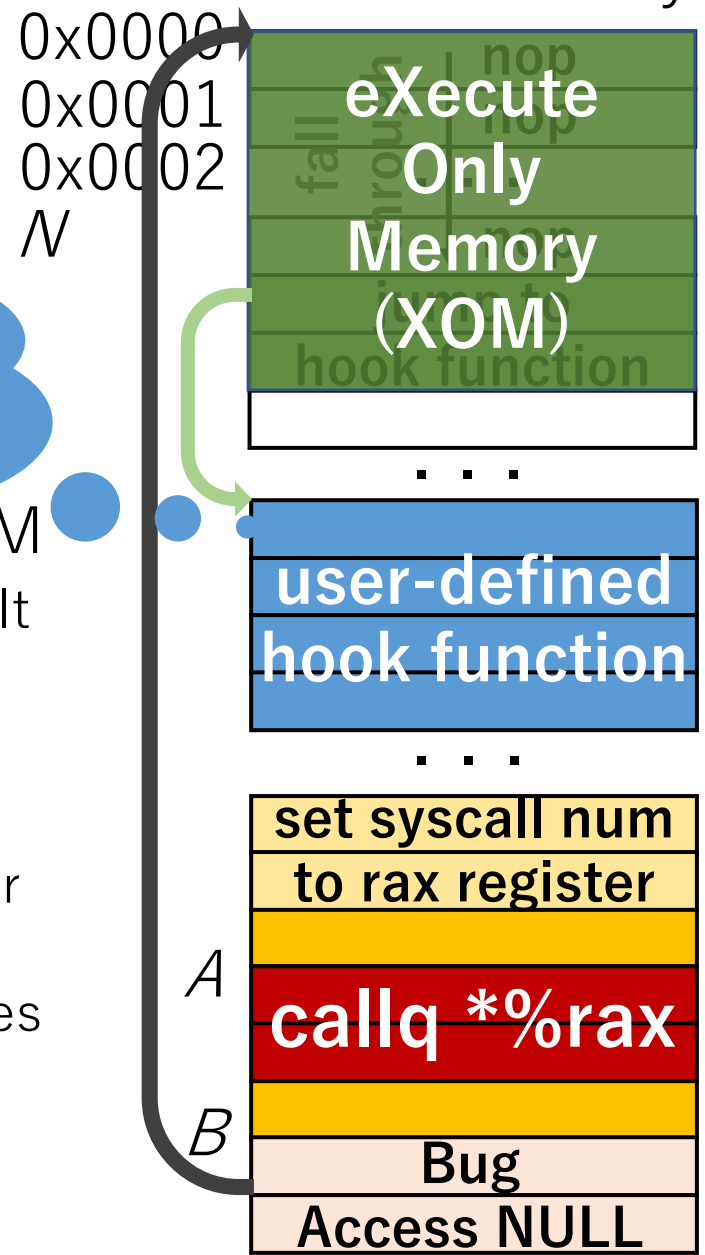
List of replaced addresses : [$A$ , ...]

Virtual Memory

0x0000
0x0001
0x0002
$N$

nop
eXecute Only Memory (XOM)
nop
fall through
hook function

user-defined hook function

. . .

set syscall num to rax register

$A$
callq *%rax

$B$
Bug
Access NULL

. . .

# NULL Access Termination

## At runtime ...

- Memory acces... 

- Solution

  - read/write... XOM
    - read/write access to the trans... code causes a fault
      - This can be done by mprotect() system call

  - execute: check the caller address
    1. during the binary rewriting phase,
       we collect the addresses of replaced syscall/sysenter
    2. at runtime, in the hook function,
       we check if the caller is one of the replaced addresses

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

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

Virtual Memory

0x0000
0x0001
0x0002
$N$

eXecute
Only
Memory
(XOM)

fall through, then
hook function

user-defined
hook function

. . .

set syscall num
to rax register

$A$

**callq *%rax**

$B$

Bug
Access NULL
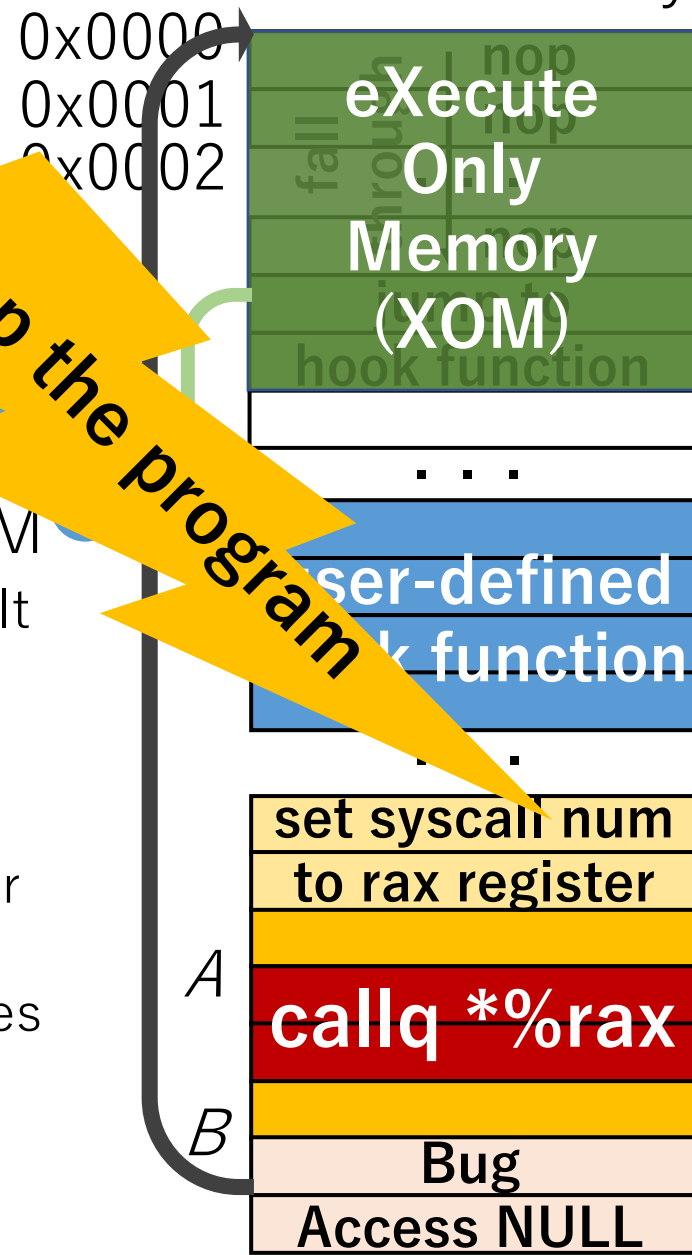
. . .

# NULL Access Termination

## At runtime ...

- Memory acces

- Solution

  - read/write:

    - read/write access to the tra~~nsp~~ code causes a fault

      - This can be done by mprotect() system call

  - execute: check the caller address

    1. during the binary rewriting phase,
       we collect the addresses of replaced syscall/sysenter

    2. at runtime, in the hook function,
       we check if the caller is one of the replaced addresses

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

**The caller address is $B$
$B$ is NOT in the list, so
this is an <u>invalid</u> access**

**stop the program**

0x0000
0x0001
0x0002

**eXecute
Only
Memory
(XOM)**

fall through
hook function

ser-defined
k function

**set syscall num
to rax register**

$A$

**callq *%rax**

$B$

**Bug
Access NULL**

# NULL Access Termination

**At runtime ...**

- Memory acces...

- Solution

  - read/write... XOM
    - read/write access to the tra... code causes a fault
      - This can be done by mprotect() system call

  - execute: check the caller address
    1. during the binary rewriting phase,
       we collect the addresses of replaced syscall/sysenter
    2. at runtime, in the hook function,
       we check if the caller is one of the replaced addresses
       - Current prototype uses bitmap to implement this check

  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**

Virtual Memory

0x0000
0x0001
0x0002

eXecute
Only
Memory
(XOM)

nop
nop

hook function

user-defined
hook function

. . .

set syscall num
to rax register
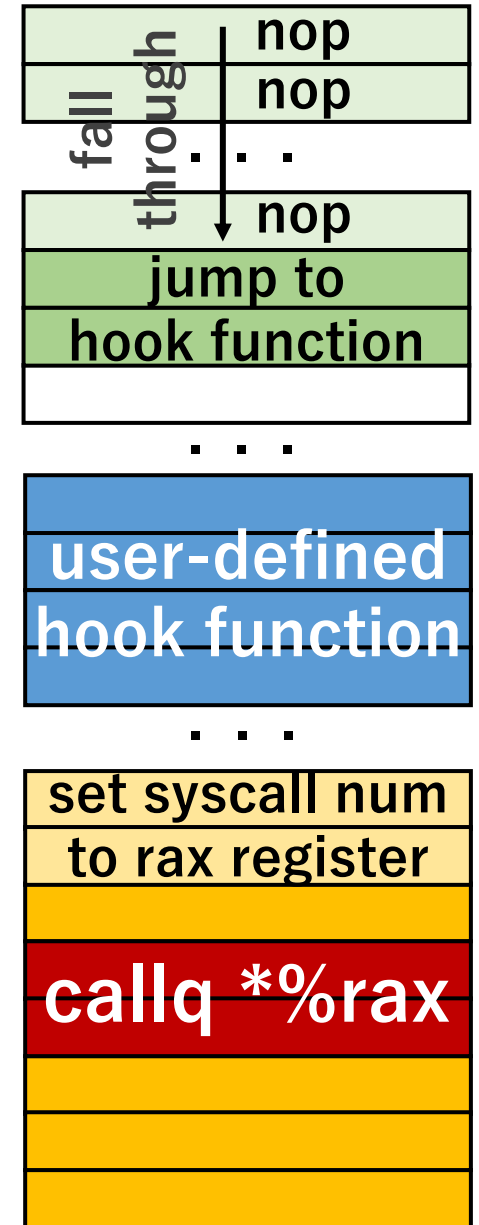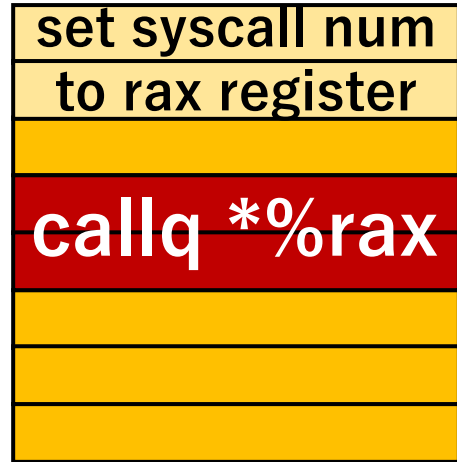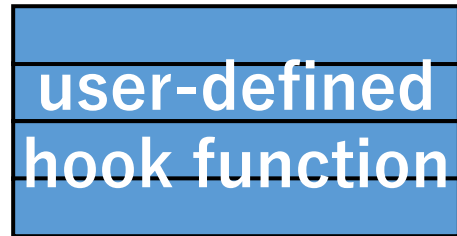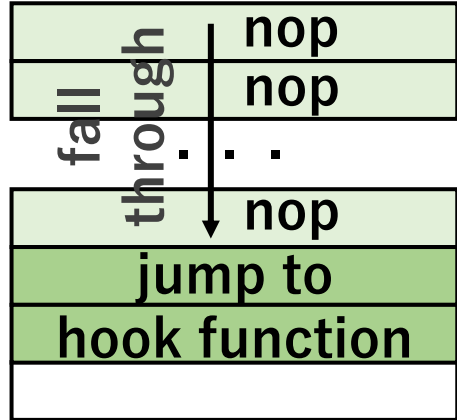
$A$

**callq *%rax**

$B$

**Bug**
**Access NULL**

. . .

# 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 |
| int3 signaling | 1342 |
| SUD | 1156 |
| zpoline | 41 |
| LD_PRELOAD | 6 |

716x

32.7x

28.1x

improvement

Virtual Memory

fall through

**nop**
**nop**
. . .
**nop**
**jump to hook function**

. . .

**user-defined hook function**

. . .

**set syscall num to rax register**

**callq *%rax**

# System Call Hook Overhead

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

| Mechanism | Time [ns] |
|-----------|-----------|
| ptrace | 31201 |
| int3 signaling | 1342 |
| SUD | 1156 |
| zpoline    NULL exec check: 1 ns out of | 41 |
| LD_PRELOAD | 6 |

Virtual Memory



fall through

nop
nop
. . .
nop
jump to
hook function

. . .

user-defined
hook function

. . .

set syscall num
to rax register

callq *%rax

. . . .

# 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 |

zpoline: +35ns overhead

fall through

nop
nop
. . .
nop
jump to hook function

. . .

user-defined hook function

. . .

set syscall num to rax register

**callq *%rax**

. . .

# 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 |

+35ns overhead

additional overhead

Virtual Memory

fall through

nop
nop
. . .
nop
jump to hook function

. . .

user-defined hook function

. . .

set syscall num to rax register

callq *%rax

. . .

# Application Performance

- We **<u>transparently</u>** apply lwIP + DPDK to an application using different system call hook mechanisms
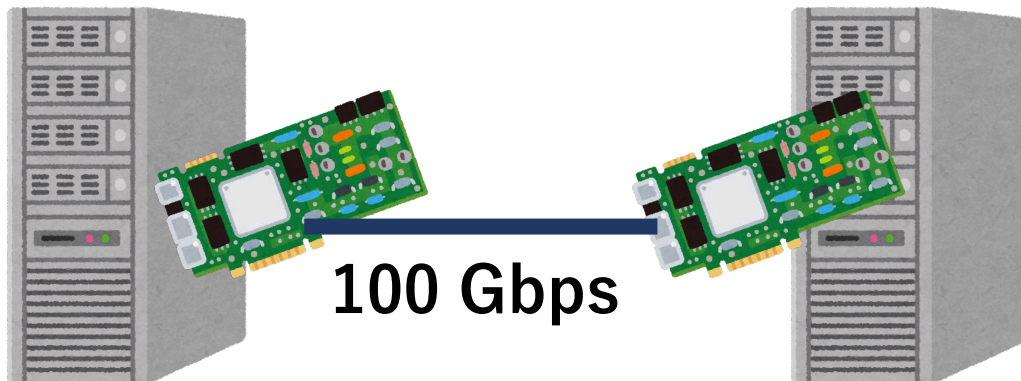
**100 Gbps**

# Application Performance

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

**Simple HTTP server**

lwIP + DPDK

**100 Gbps**
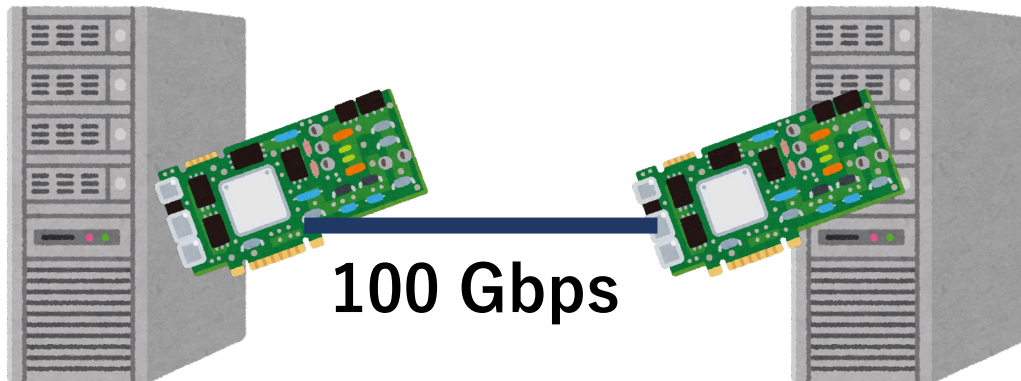
# Application Performance

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

**Simple HTTP server**

ptrace, int3, SUD, zpoline, LD_PRELOAD

*syscall hook*

lwIP + DPDK

**100 Gbps**

# Application Performance

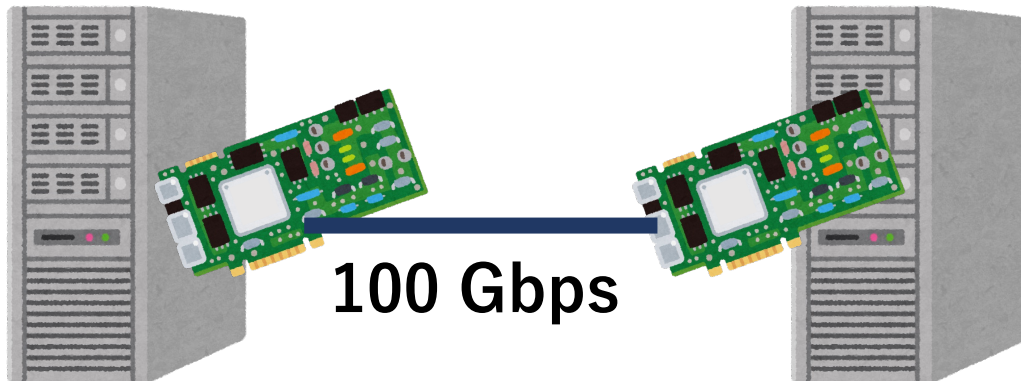- We **<u>transparently</u>** apply lwIP + DPDK to an application using different system call hook mechanisms

**Simple HTTP server**

*syscall hook*

ptrace, int3, SUD, zpoline, LD_PRELOAD

wrk: benchmark client
fetch 64B content

lwIP + DPDK

**100 Gbps**

# Application Performance

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

**Simple HTTP server**
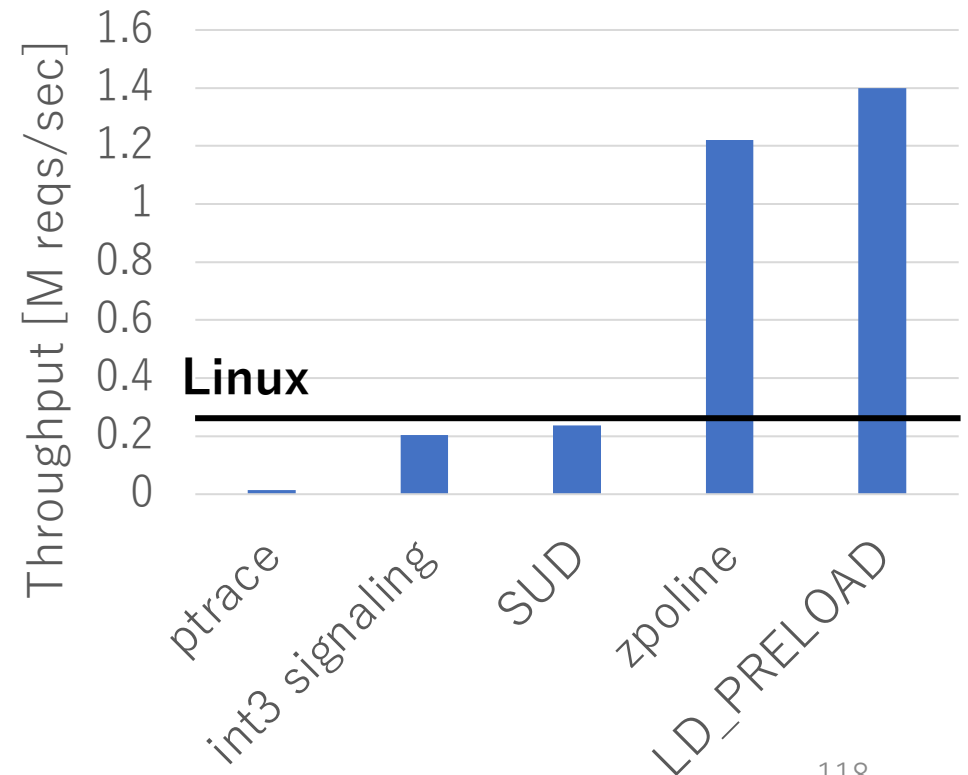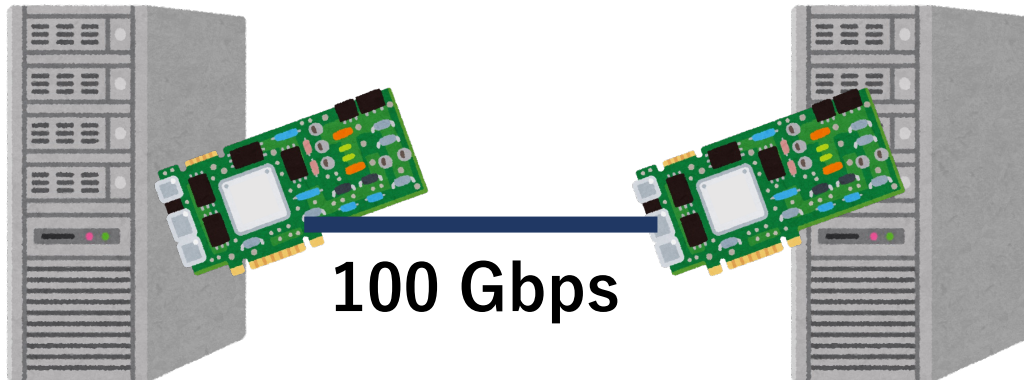
ptrace, int3, SUD, zpoline, LD_PRELOAD

*syscall hook*

lwIP + DPDK

wrk: benchmark client
fetch 64B content

**100 Gbps**

Throughput [M reqs/sec]

1.6
1.4
1.2
1
0.8
0.6
0.4
0.2
0

**Linux**

ptrace
int3 signaling
SUD
zpoline
LD_PRELOAD

# Application Performance

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

**Simple HTTP server**

*syscall hook*

ptrace, int3, SUD, zpoline, LD_PRELOAD

lwIP + DPDK

wrk: benchmark client
fetch 64B content

**100 Gbps**



Throughput [M reqs/sec]

**Compared to LD_PRELOAD**

87.3%

**Linux** 14.7% 17.0%

1.1%

ptrace
int3 signaling
SUD
zpoline
LD_PRELOAD

# Application Performance

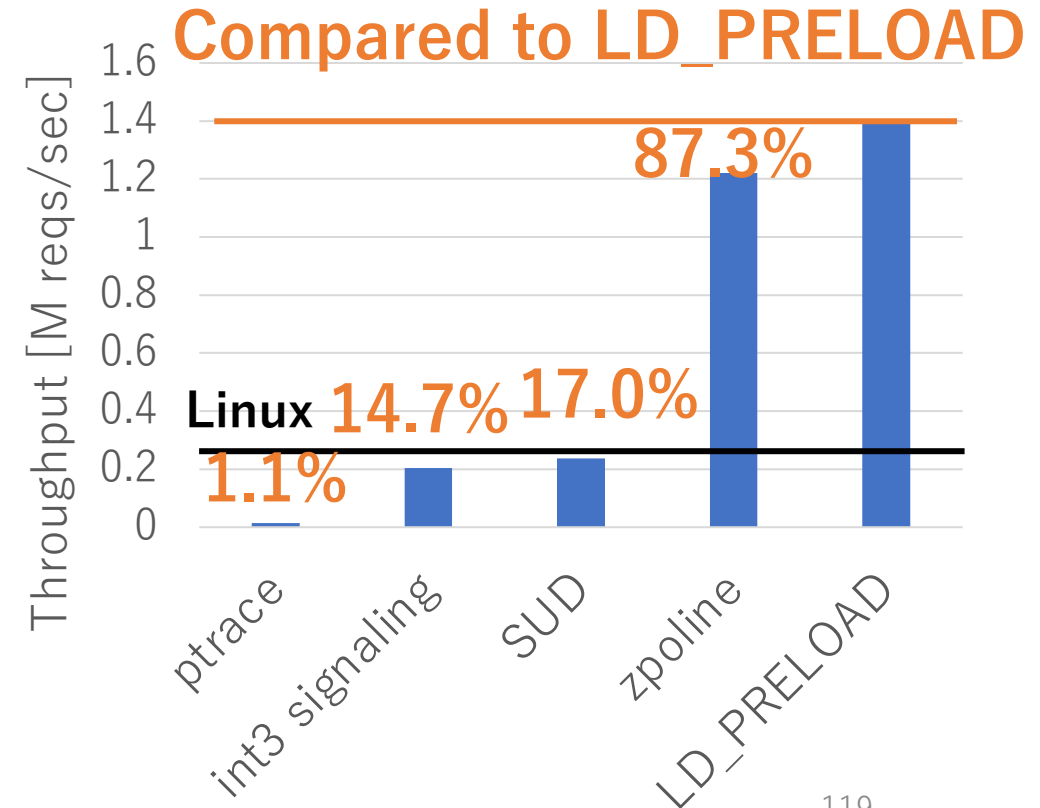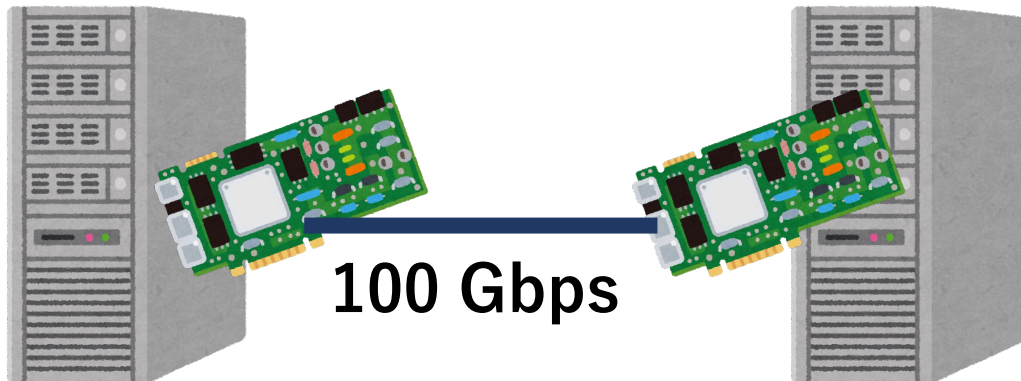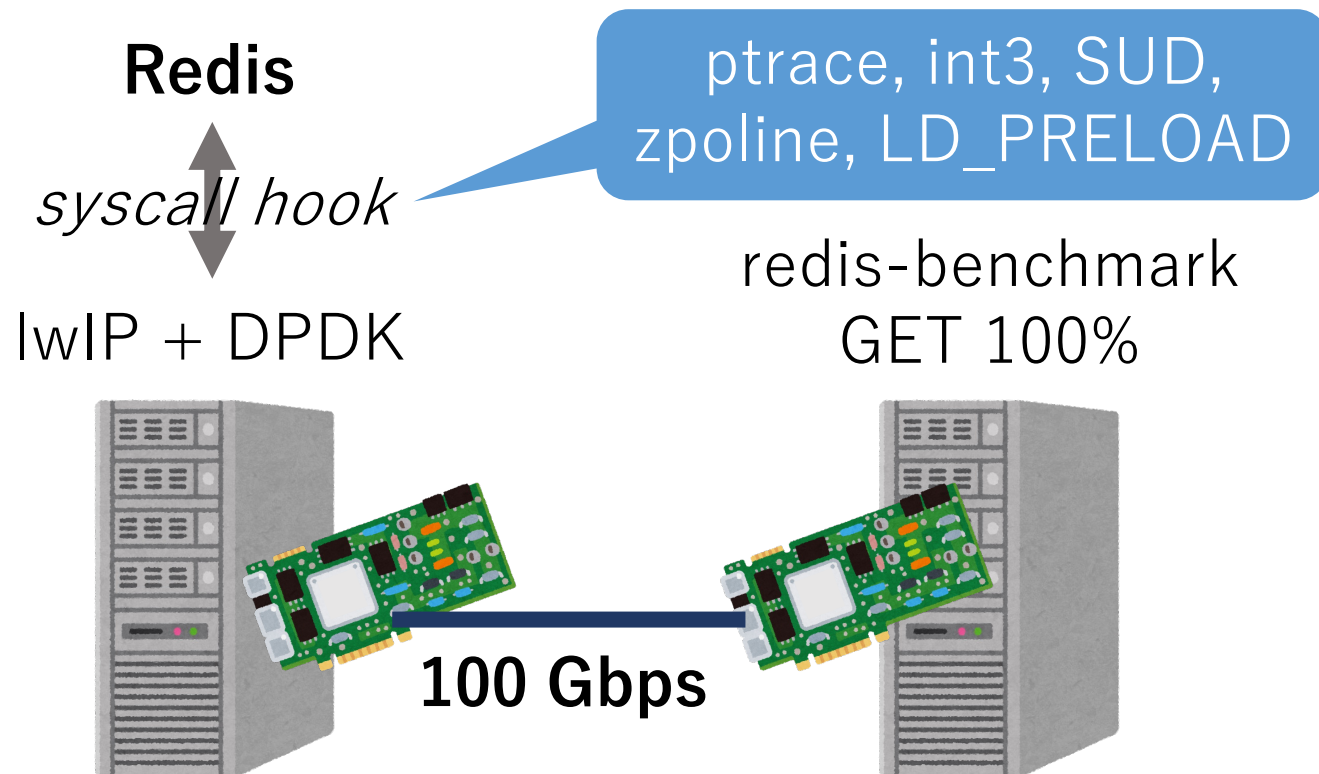- We **<u>transparently</u>** apply lwIP + DPDK to an application using different system call hook mechanisms

**Redis**

*syscall hook*

ptrace, int3, SUD, zpoline, LD_PRELOAD

redis-benchmark
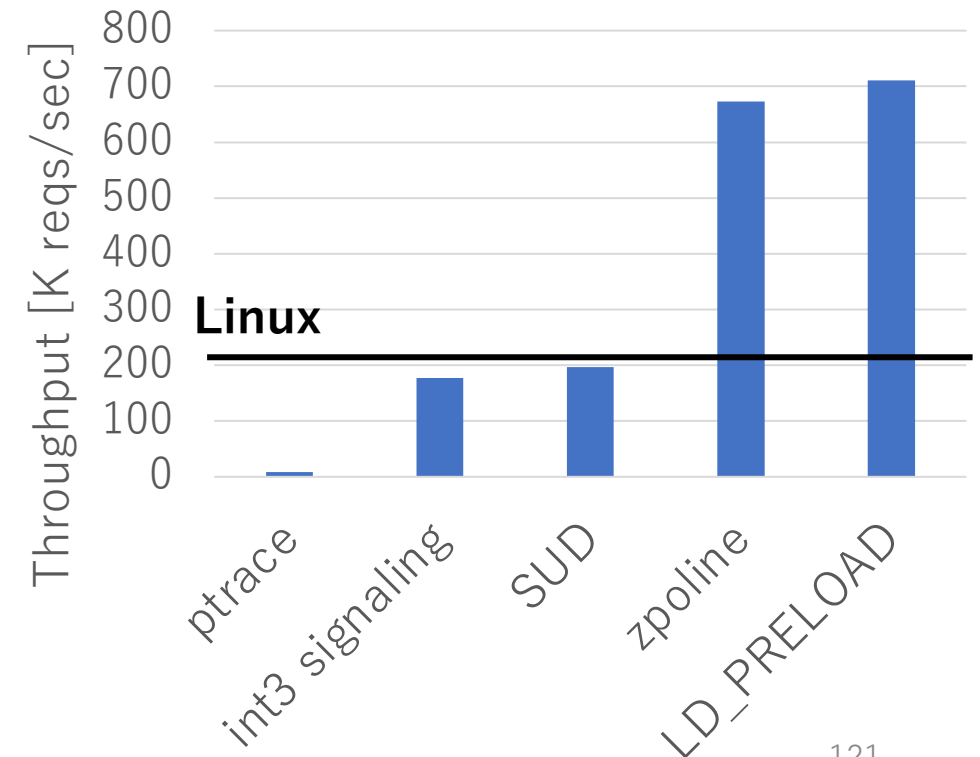GET 100%

lwIP + DPDK

**100 Gbps**

# Application Performance

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

**Redis**

*syscall hook*

ptrace, int3, SUD, zpoline, LD_PRELOAD

lwIP + DPDK

redis-benchmark
GET 100%

**100 Gbps**

Throughput [K reqs/sec]

800
700
600
500
400
300
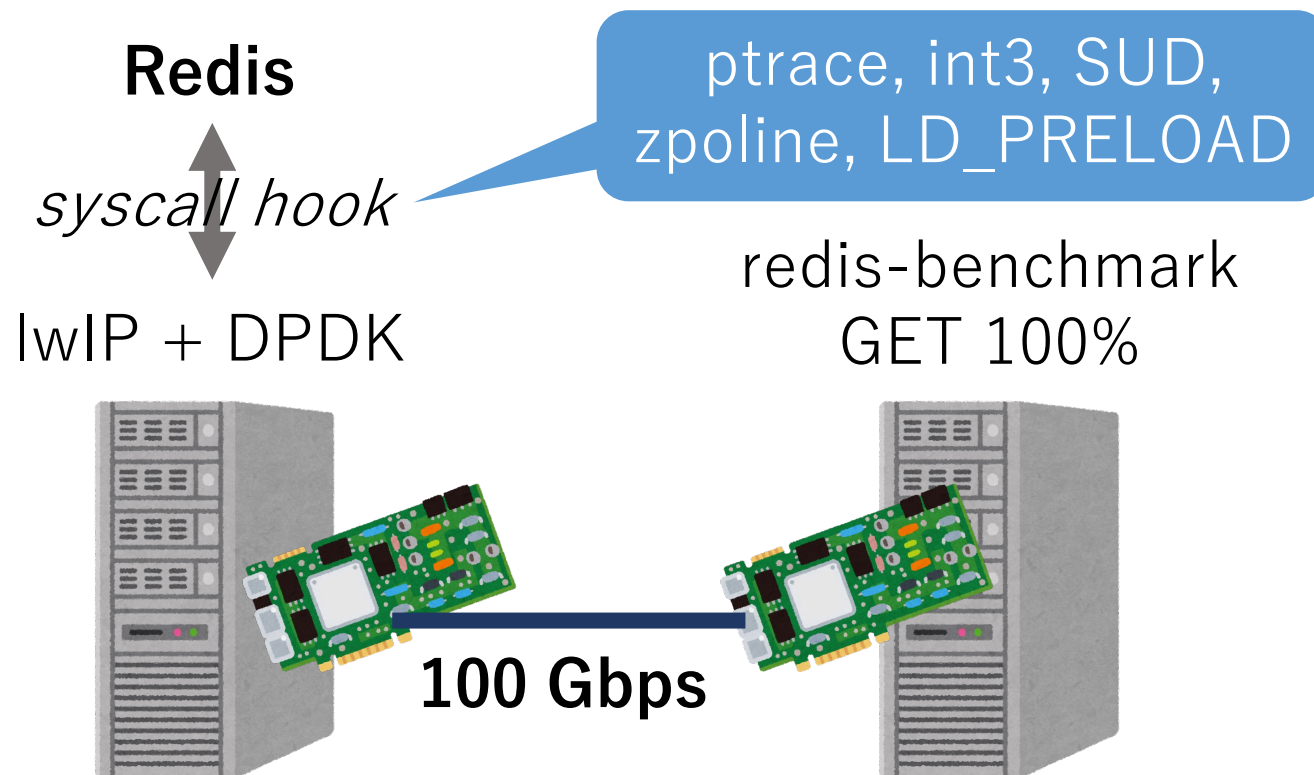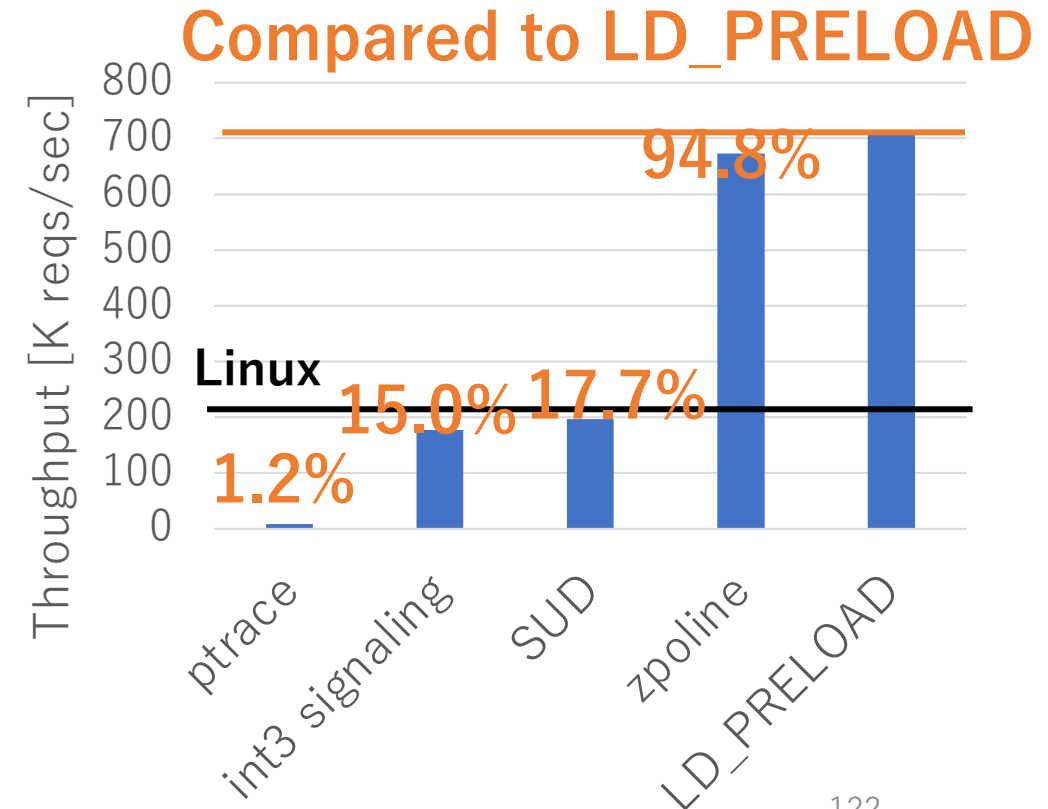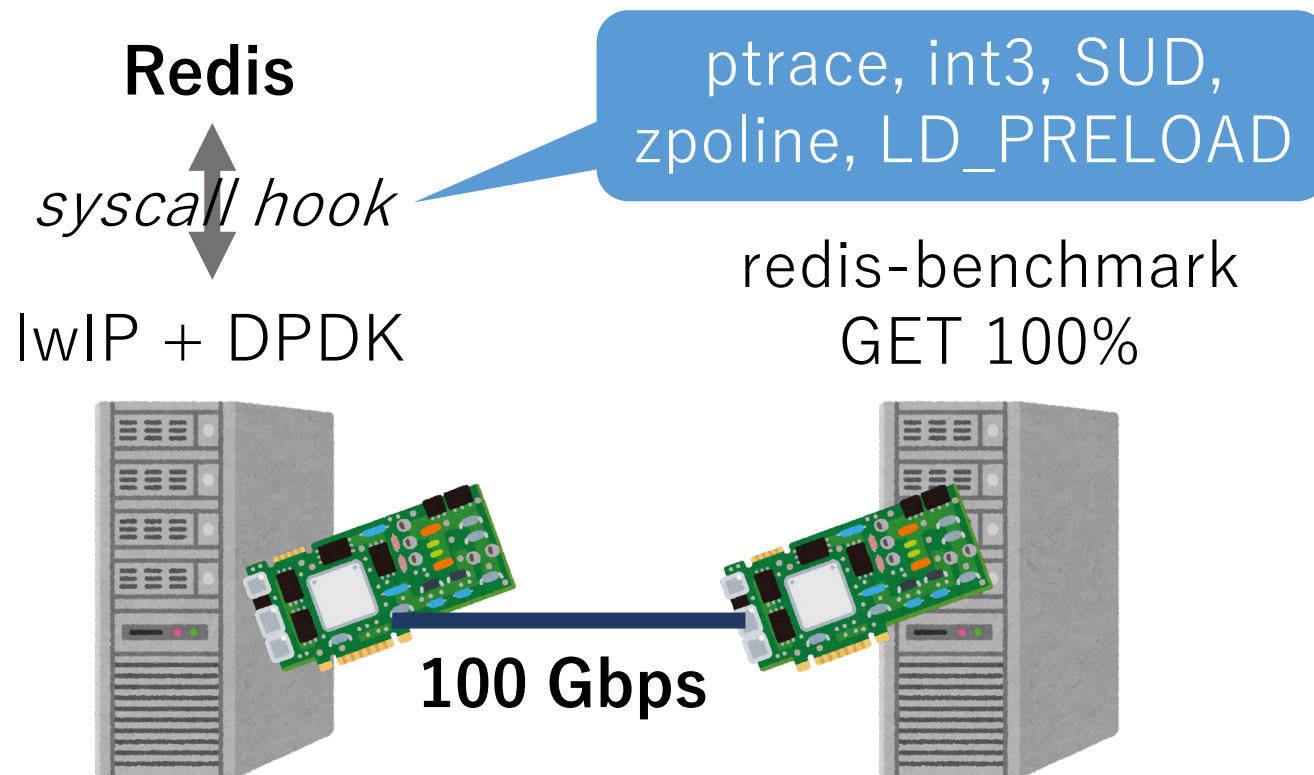200
100
0

**Linux**

ptrace
int3 signaling
SUD
zpoline
LD_PRELOAD

# Application Performance

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

**Redis**

*syscall hook*

lwIP + DPDK

ptrace, int3, SUD, zpoline, LD_PRELOAD

redis-benchmark
GET 100%

**100 Gbps**

**Compared to LD_PRELOAD**

Throughput [K reqs/sec]

800
700
600
500
400
300
200
100
0

**Linux**

**1.2%**   **15.0%**   **17.7%**   **94.8%**

ptrace   int3 signaling   SUD   zpoline   LD_PRELOAD

# Summary

- zpoline: a system call hook mechanism for x86-64 CPUs
  - based on binary rewriting
    - replaces syscall/sysenter with callq *%rax
    - instantiates the trampoline code at virtual address 0 (zero)
  - free from the drawbacks of the pervious mechanisms
  - keeps the performance benefit of user-space OS subsystems

- Source code: https://github.com/yasukata/zpoline
  - since October 2021

123

# Speeding up the Trampoline Code

- Inspired from USENIX ATC'23 reviewers who suggested to employ a one-byte short jump instruction for speeding up
  - Put it on the addresses corresponding to obsolete system calls

- Optimization: repeat **0xeb 0x6a 0x90** instead of nops
  - Hook overhead reduction from 41 ns to 10 ns

Syscall number:          3 x n + 0                3 x n + 1                3 x n + 2

```
                    jmp 0x6a          push 0x90          nop
                    nop               jmp 0x6a           jmp 0x6a
We pop 0x90 in      jmp 0x6a          nop                nop
the hook function   nop               jmp 0x6a           jmp 0x6a
```