# Security Isolation for System

Zeyu Mi

SHANGHAI JIAO TONG UNIVERSITY

IPADS
INSTITUTE OF PARALLEL
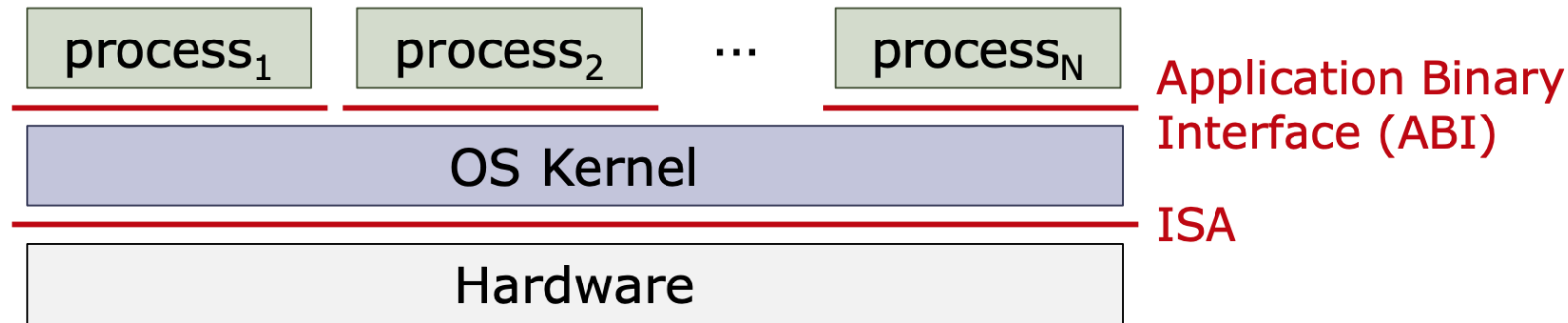AND DISTRIBUTED SYSTEMS

# What is system isolation?

- **An entity may encounter failures or even be malicious**

- **System isolation**
  - Contains effects of failures
  - Performance fairness

- **Unit of isolation**
  - Process
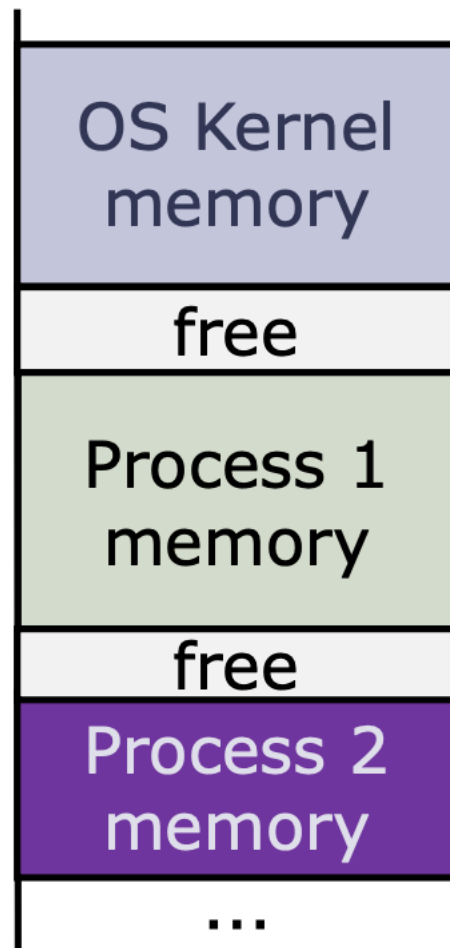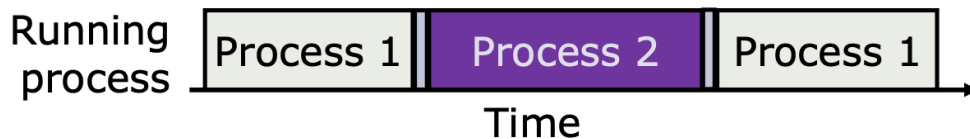  - Virtual Machine
  - Or smaller entities

# Operating system (OS) and its Goals

- **Resource Management**
  - OS controls how processes share hardware (CPU, memory, storage, network, etc.)

- **Abstraction**
  - Hide underline details
  - Provide usable interfaces

- **Protection and Privacy**
  - Process cannot access other process data

| process$_1$ | process$_2$ | $\cdots$ | process$_N$ |
|---|---|---|---|

Application Binary Interface (ABI)
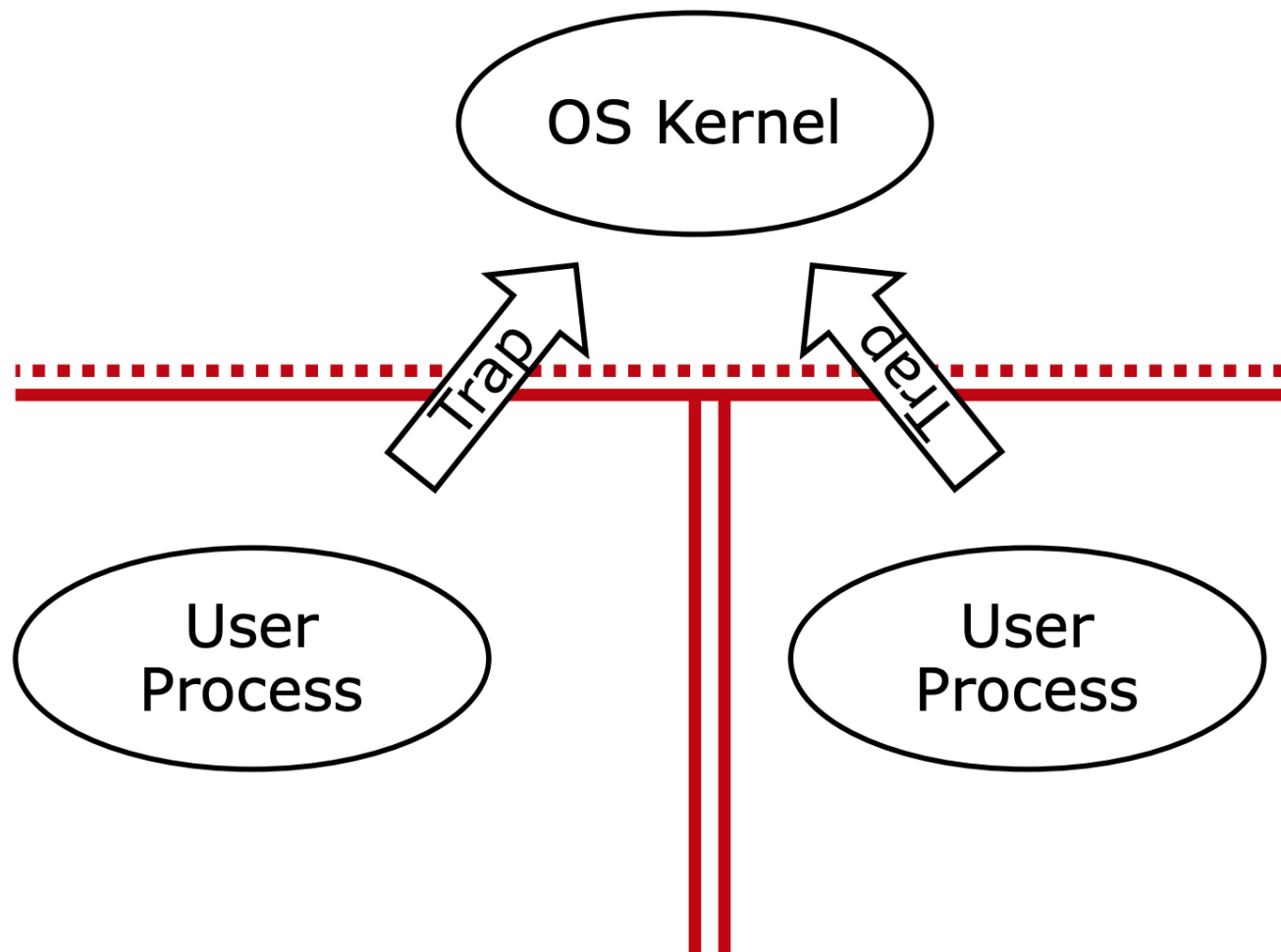
OS Kernel

ISA

Hardware

# Process Abstraction

- **Each process has a private address space**
  - Provided by OS
  - Cannot access other process spaces

- **The OS kernel schedules processes into cores**
  - Each process has a scheduling time slice

- **A process can invoke OS services via system calls**

| OS Kernel memory |
| free |
| Process 1 memory |
| free |
| Process 2 memory |
| ... |

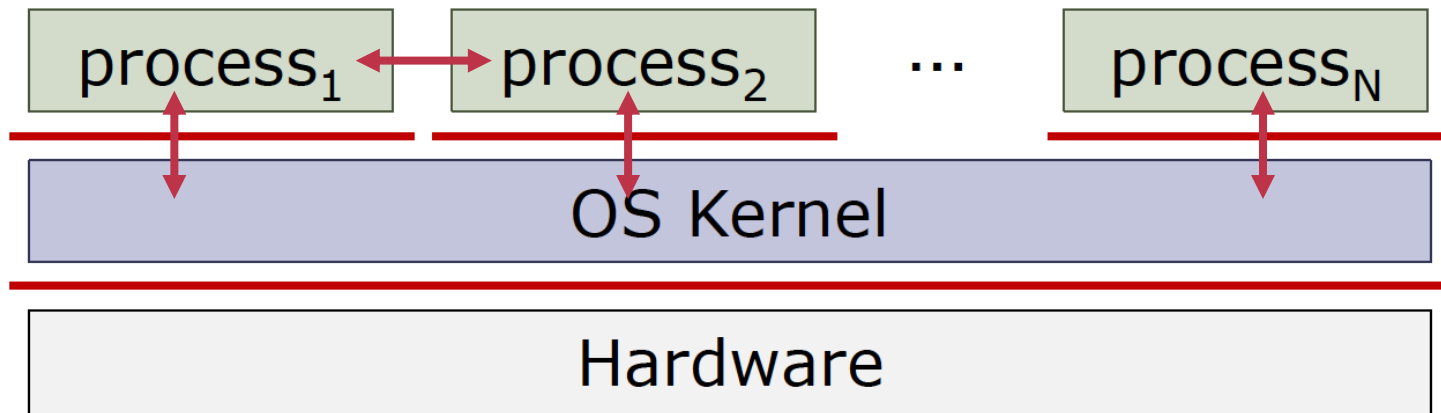Running process: | Process 1 | Process 2 | Process 1 |
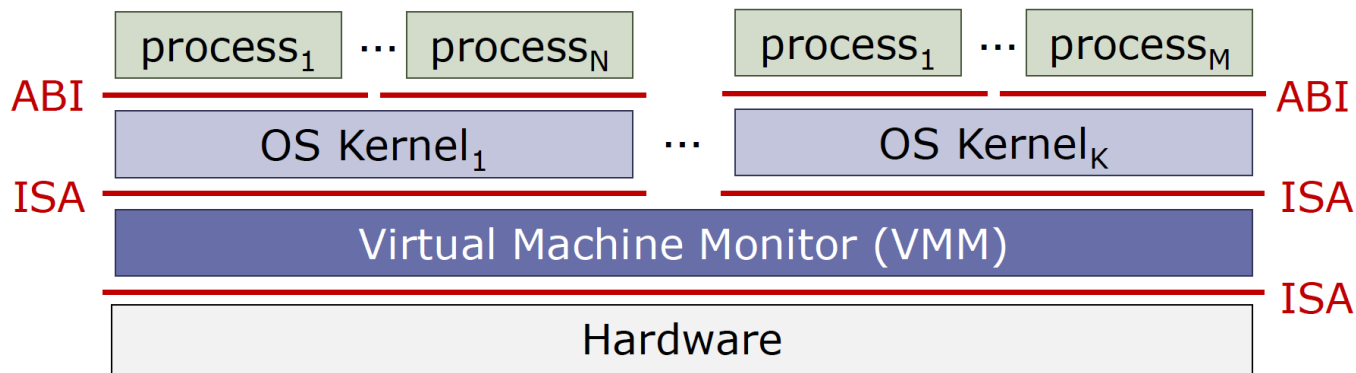
Time

# Protection for a Single OS

# System Isolation in OS

- **Isolation between processes (Horizontal Isolation)**
  - Enforced by OS kernel
  - Useful techniques: page tables, context switch, file abstraction…

- **Isolation between processes and the kernel (Vertical Isolation)**
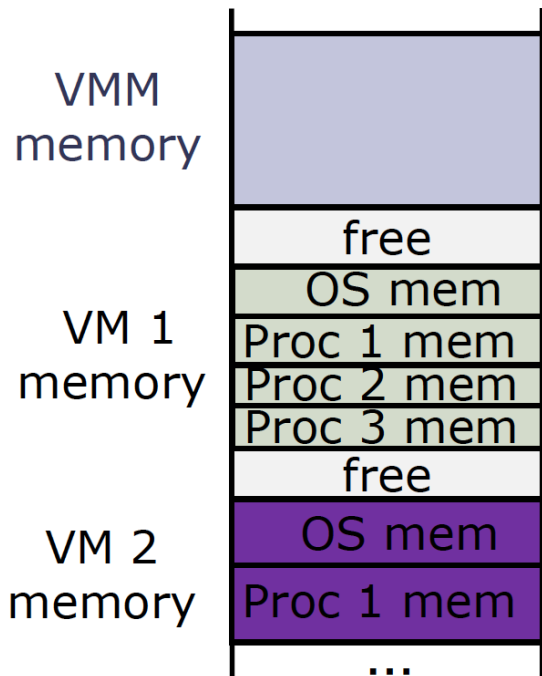  - CPU privilege and carefully designed interface

# Virtual Machine Monitor

- **A VMM (aka Hypervisor) provides a system virtual machine to each OS**

- **VMM can run directly on hardware (type-1) or on another OS (type 2)**

- **Hardware virtualization**
  - VT-x: root and non-root mode
  - EPT
  - IOMMU and SR-IOV



ABI

| process$_1$ | $\cdots$ | process$_N$ | | process$_1$ | $\cdots$ | process$_M$ |

OS Kernel$_1$ $\cdots$ OS Kernel$_K$

ISA

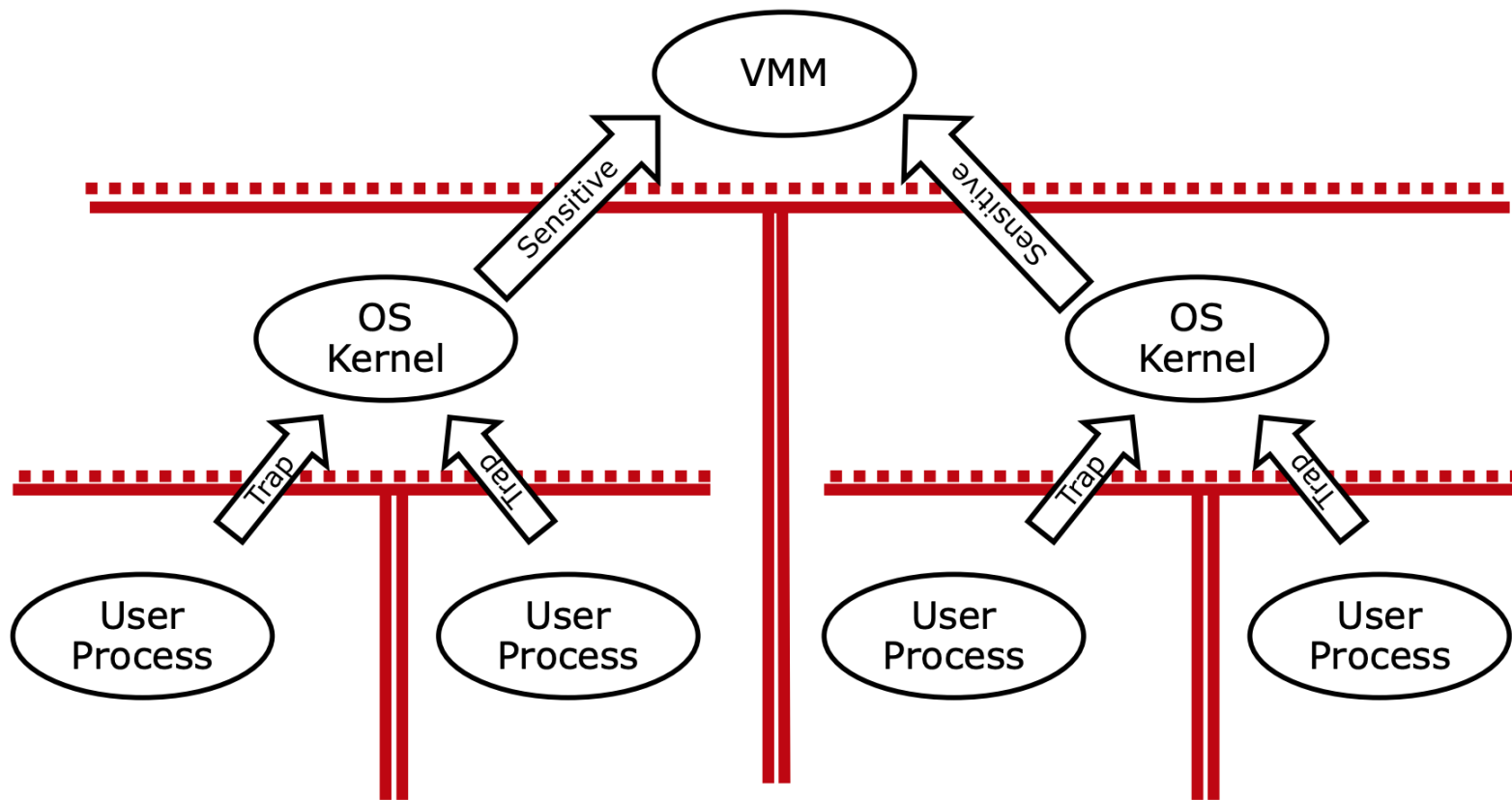Virtual Machine Monitor (VMM)

ISA

Hardware

# Virtual Machine (VM) Abstraction

- **Each VM has a private address space**
  - Provided by the hypervisor
  - Cannot access other VM spaces

- **The hypervisor schedules processes into cores**
  - Each process has a scheduling time slice

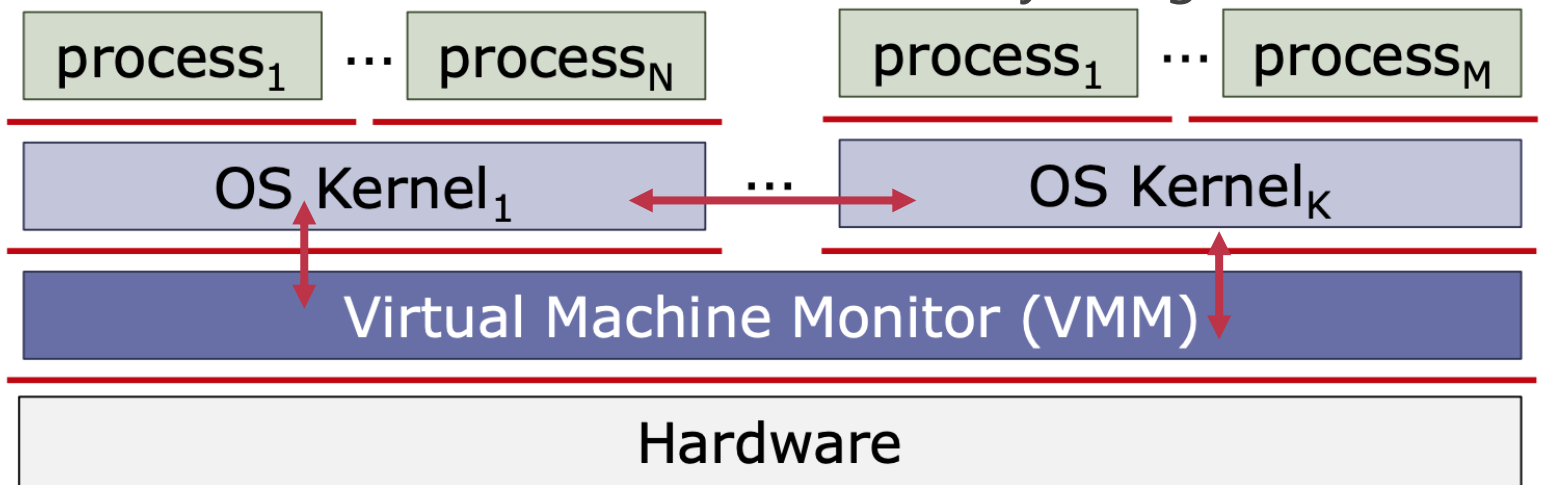- **A VM can invoke hypervisor services via hypercalls**

# Protection for Multiple OSes

# System Isolation in VMM

- **Isolation between OS kernels (Horizontal Isolation)**
  - Enforced by the VMM
  - Useful techniques: EPTs

- **Isolation between OS kernels and VMM (Vertical Isolation)**
  - Hardware virtualization and carefully designed interface

# Techniques for System Isolation

- **CPU privileges**

- **Segmentation**

- **Page table**

- **Extended page table**

- **Memory domain**

- **Secure hardware modules**

# Techniques for System Isolation

- **CPU privileges**

- **Segmentation**

- **Page table**

- **Extended page table**

- **Memory domain**

- **Secure hardware modules**

# CPU privileges

- **Privileged Mode**
  - Non-root Ring 0
  - Root Ring 0

- **Operations in privileged mode**
  - Configure address spaces
  - Access devices
  - Read/write special registers

| | | |
|---|---|---|
| | **Ring 3** | Processes |
| **Non-root Mode** | **Ring 2** | |
| | **Ring 1** | |
| | **Ring 0** | Guest Kernel |
| | **Ring 3** | |
| **Root Mode** | **Ring 2** | |
| | **Ring 1** | |
| | **Ring 0** | VMM |

# Techniques for System Isolation
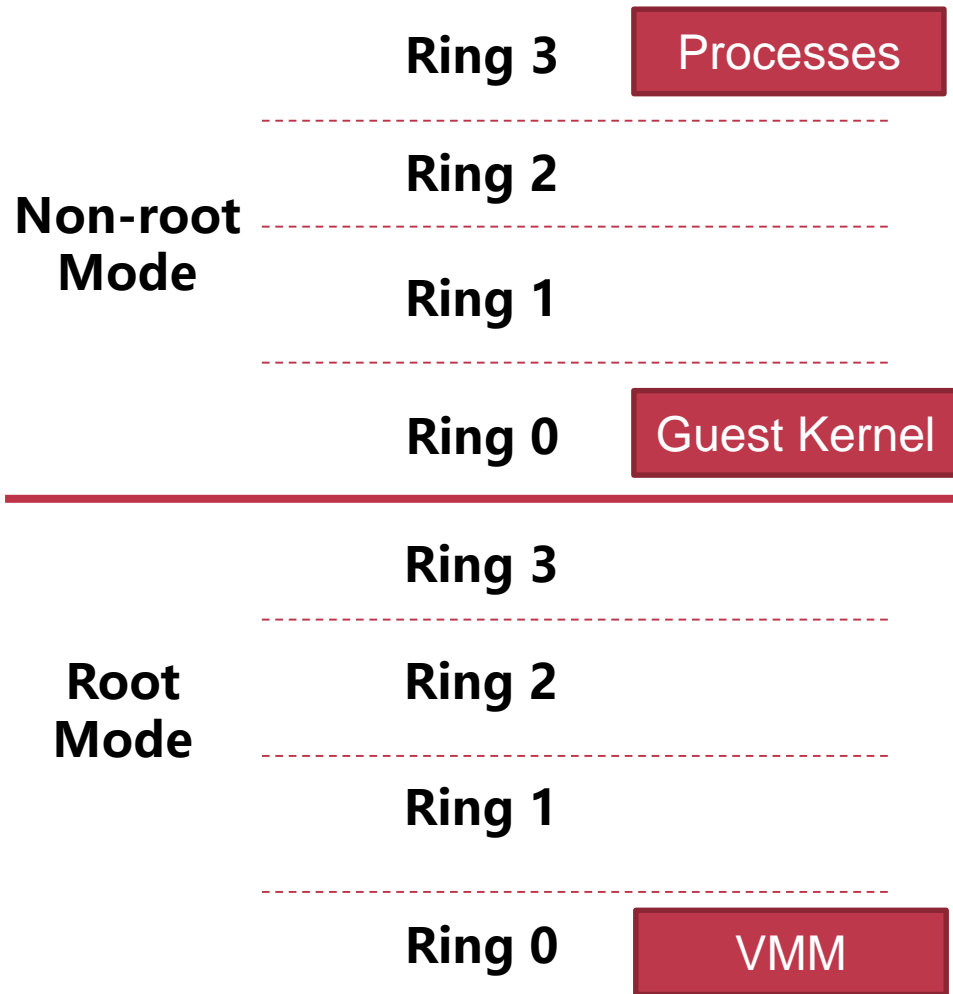
- CPU privileges

- **Segmentation**

- Page table

- Extended page table
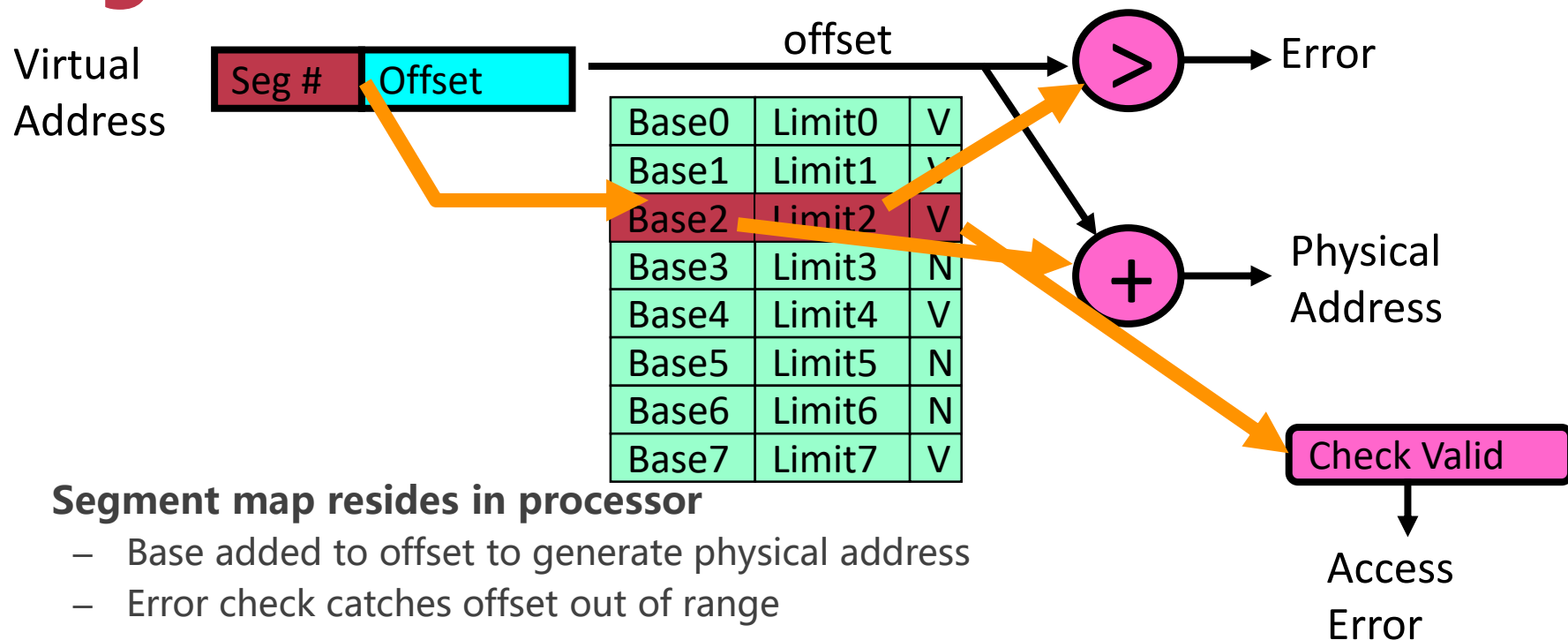
- Memory domain

- Secure hardware modules
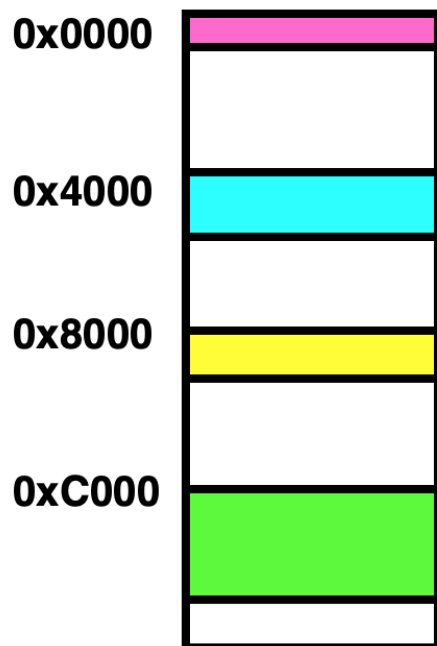
# Segmentation



- **Segment map resides in processor**
  - Base added to offset to generate physical address
  - Error check catches offset out of range
- **As many chunks of physical memory as entries**
  - Segment addressed by portion of virtual address
- **What is "V/N" (valid / not valid)?**
  - Can mark segments as invalid; requires check as well

20

| Seg ID # | Base | Limit |
|---|---|---|
| 0 (code) | 0x4000 | 0x0800 |
| 1 (data) | 0x4800 | 0x1400 |
| 2 (shared) | 0xF000 | 0x1000 |
| 3 (stack) | 0x0000 | 0x3000 |

**Virtual Address Format**

| Seg | Offset |
|---|---|

15  14  13                    0

**Virtual Address Space**

0x0000

0x4000

0x8000

0xC000

**Physical Address Space**

0x0000

| Seg ID # | Base | Limit |
|---|---|---|
| 0 (code) | 0x4000 | 0x0800 |
| 1 (data) | 0x4800 | 0x1400 |
| 2 (shared) | 0xF000 | 0x1000 |
| 3 (stack) | 0x0000 | 0x3000 |

**Virtual Address Format**

| Seg | Offset |
|---|---|

15  14 13                    0

**SegID = 0**

0x0000

0x4000

0x8000

0xC000

**Virtual Address Space**

0x0000

0x4000
0x4800

**Physical Address Space**

| Seg ID # | Base | Limit |
|----------|--------|--------|
| 0 (code) | 0x4000 | 0x0800 |
| 1 (data) | 0x4800 | 0x1400 |
| 2 (shared) | 0xF000 | 0x1000 |
| 3 (stack) | 0x0000 | 0x3000 |

**Virtual Address Format**

Seg | Offset

15  14  13  0

SegID = 0

SegID = 1

0x0000
0x4000
0x8000
0xC000

**Virtual
Address Space**

0x0000
0x4000
0x4800
0x5C00

Might
be shared

Space for
Other Apps

Shared with
Other Apps

**Physical
Address Space**

| Seg ID # | Base | Limit |
|----------|--------|--------|
| 0 (code) | 0x4000 | 0x0800 |
| 1 (data) | 0x4800 | 0x1400 |
| 2 (shared) | 0xF000 | 0x1000 |
| 3 (stack) | 0x0000 | 0x3000 |

**Virtual Address Format**

Seg | Offset

15  14  13                    0

SegID = 0

SegID = 1

0x0000
0x4000
0x8000
0xC000

**Virtual Address Space**

0x0000
0x4000
0x4800
0x5C00
0xF000

Might be shared

Space for Other Apps

Shared with Other Apps
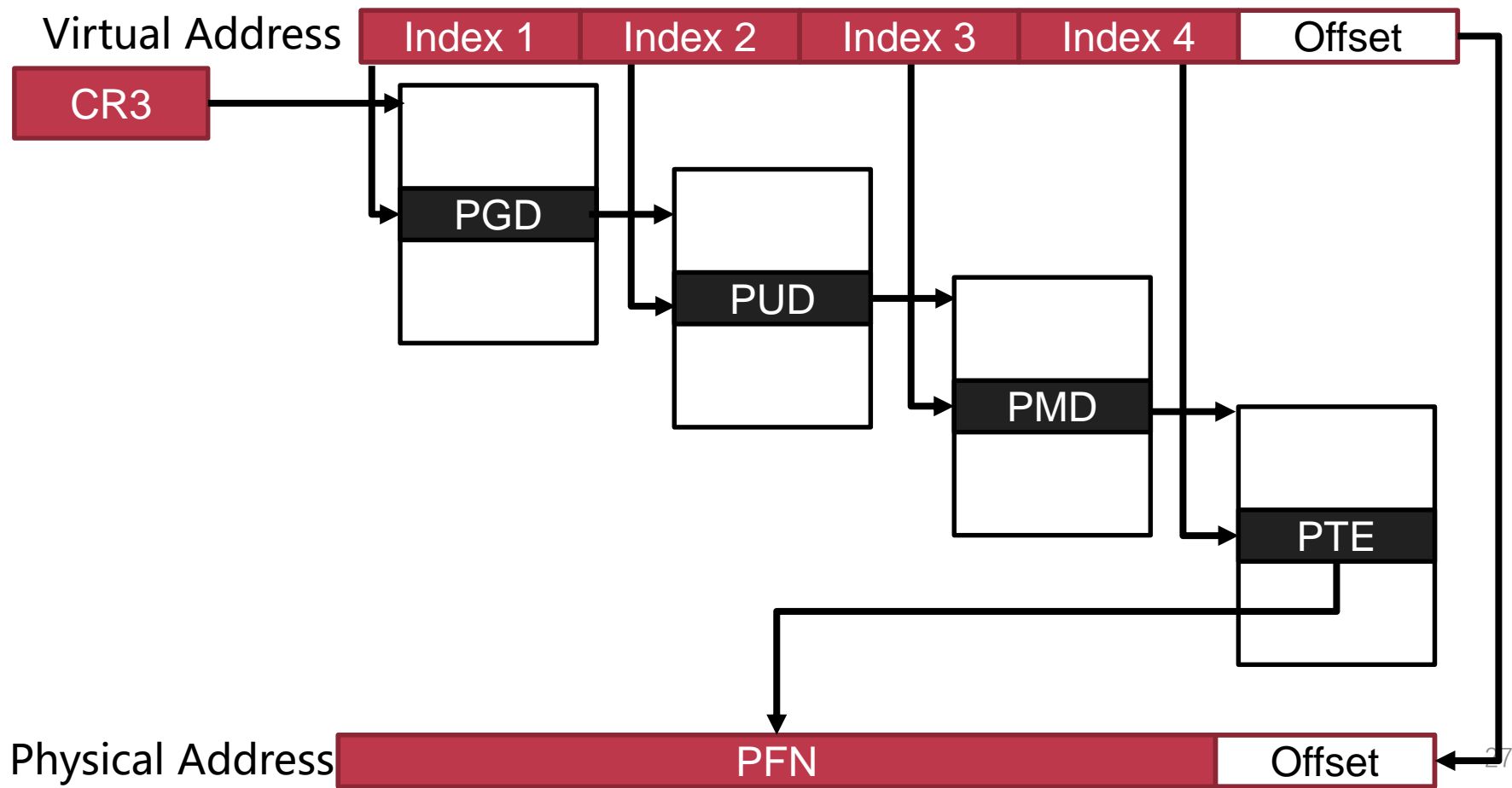
**Physical Address Space**

# Problems with Segmentation

- **Not supported in x86-64 mode**

- **Must fit variable-sized chunks into physical memory**

- **Fragmentation: wasted space**
  - External: free gaps between allocated chunks
  - Internal: do not need all memory within allocated chunks

# Techniques for System Isolation

- CPU privileges

- Segments

- **Page table**

- Extended page table

- Memory domain

- Secure hardware modules

# Hierarchical Page Table



Virtual Address | Index 1 | Index 2 | Index 3 | Index 4 | Offset

CR3

PGD

PUD

PMD

PTE

Physical Address | PFN | Offset

# Techniques for System Isolation

- CPU privileges

- Segments

- Page table

- **Extended page table**
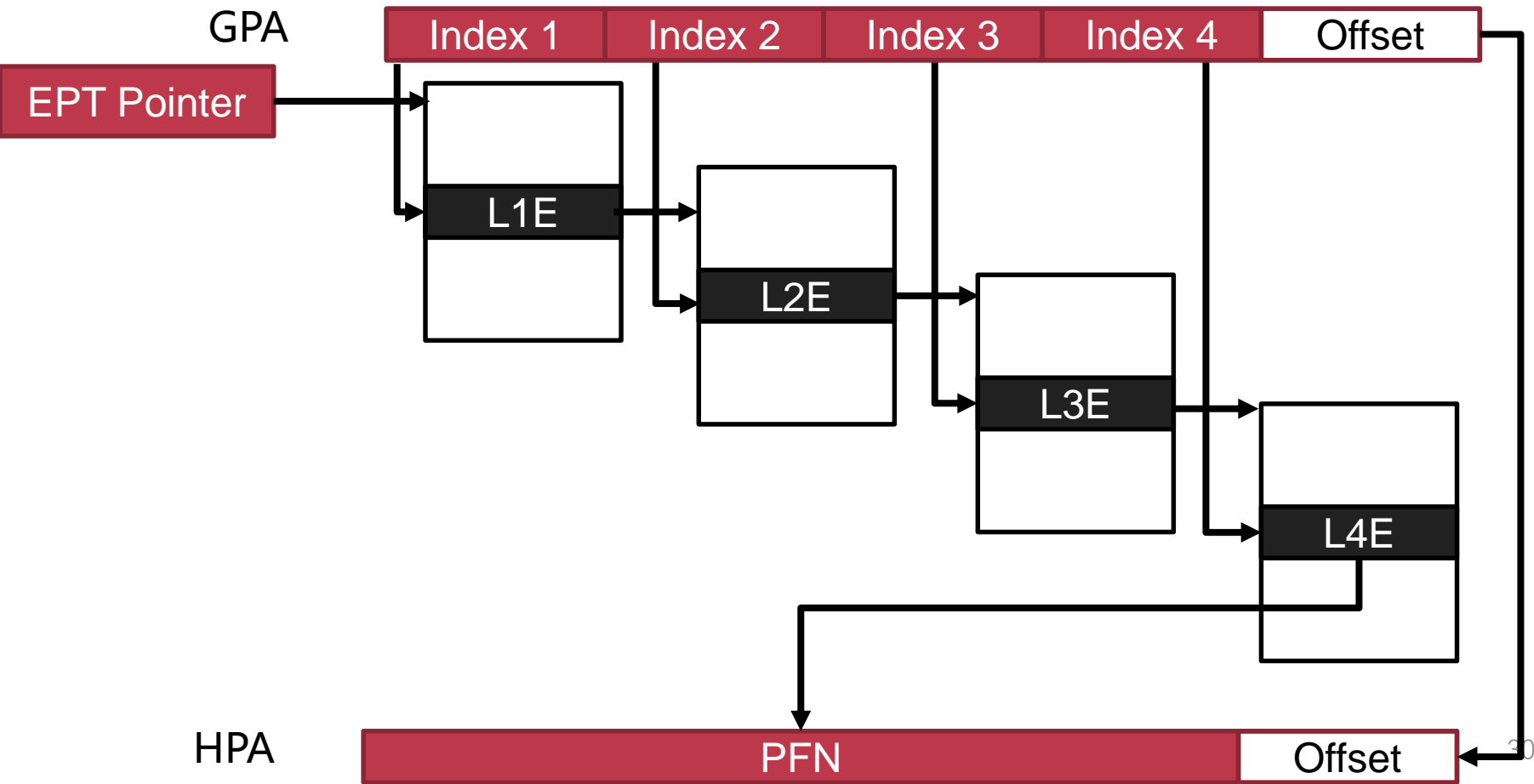
- Memory domain

- Secure hardware modules

# EPT

- **Translate guest physical addr to host physical addr**
  - The two-level translation are all done by hardware



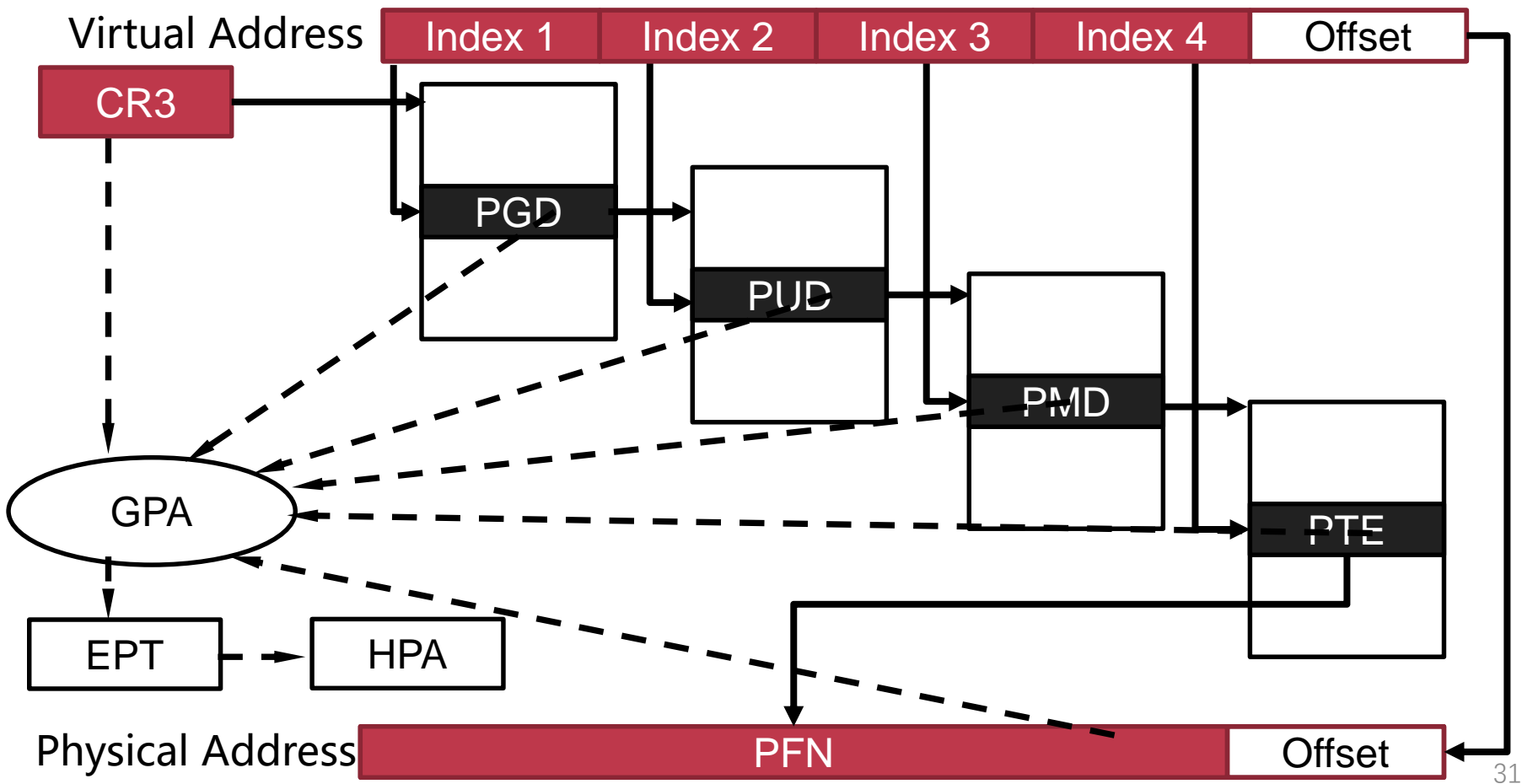| Guest Virtual Address (GVA) | →Guest Page table→ | Guest Physical Address (GPA) | →EPT→ | Host Physical Address (HPA) |

- **EPT is manipulated and maintained by hypervisor**
  - Hypervisor controls how guest accesses physical address
  - Any EPT violation triggers a VM Exit to hypervisor

# Hierarchical Extended Page Table



GPA: Index 1 | Index 2 | Index 3 | Index 4 | Offset

EPT Pointer

L1E

L2E

L3E
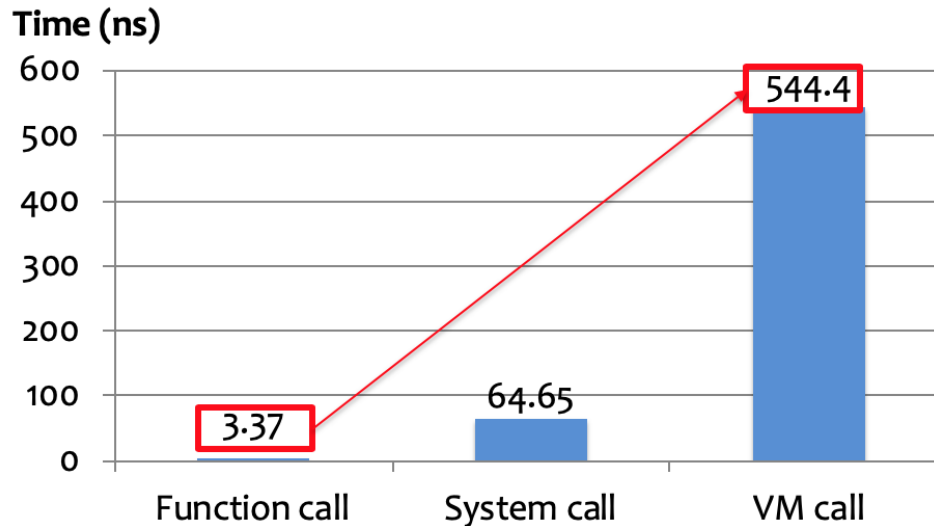
L4E

HPA: PFN | Offset

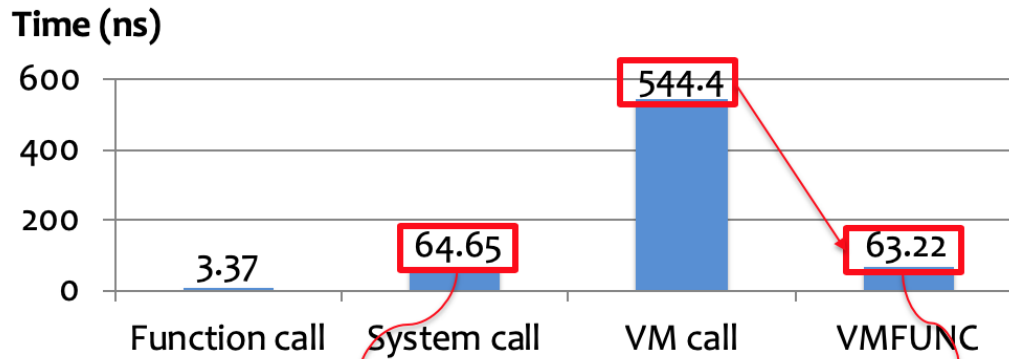# Any GPA is translated to HPA

# VMFUNC

- **Context switch introduces large overhead**
  - Every EPT switch is intervened by hypervisor
  - A VM Exit takes much more time than function call

# VMFUNC

- **VM Functions: Intel virtualization extension**
  - Non-root guest VMs can directly invoke some functions without VM Exit

- **VM Function 0: EPTP Switching**
  - Software in guest VM can directly load a new EPT pointer

**Time (ns)**

| | Function call | System call | VM call | VMFUNC |
|---|---|---|---|---|
| | 3.37 | 64.65 | 544.4 | 63.22 |

**VMFUNC can provide the hypervisor-level function at the cost of system calls**

# VMFUNC



1: **vmfunc**(0x0, 0x0); // Func ID is 0x0, switch to EPTP 0
2: **vmfunc**(0x0, 0x2); // Func ID is 0x0, switch to EPTP 2
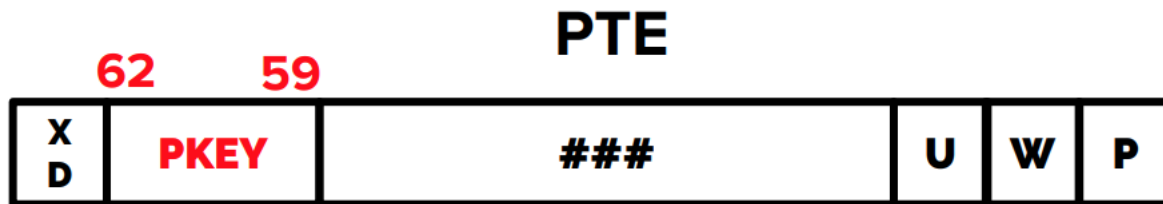3: **vmfunc**(0x0, 0x4); // Func ID is 0x0, Error

# Techniques for System Isolation

- CPU privileges

- Segments

- Page table

- Extended page table

- **Memory domain**

- Secure hardware modules

# Intel Memory Protection Keys (MPK)

- **32-bit PKRU register (Access/Write Disable)**

- **WRPKRU/RDPKRU**

## PKRU

| 31 | 30 | | 3 | 2 | 1 | 0 |
|---|---|---|---|---|---|---|
| WD 15 | AD 15 | ... | WD 1 | AD 1 | WD 0 | AD 0 |

PKEY 15       PKEY 1     PKEY 0

## PTE

62     59

| X D | PKEY | ### | U | W | P |
|---|---|---|---|---|---|

# Intel Memory Protection Keys (MPK)



OS-managed

**DTLB**

| page# | pkey | perm. |
|-------|------|-------|
| 120 | **8** | **r/w** |
| 232 | **1** | **r/o** |
| 456 | **8** | **r/o** |
| ... | | |

**ITLB**

| page# | perm. |
|-------|-------|
| 232 | **x** |
| ... | |

*(per-process, synchronized)*

Userspace process

*RDPKRU* ⇄ *WRPKRU*

**PKRU (core_a)**

| pkey: | 0 | **1** | | **8** | | 15 |
|-------|---|-------|-----|-------|-----|-----|
| perm: | r/w | **n/a** | ... | **r/w** | ... | n/a |

**(core_a)**

| page# | effective perm. |
|-------|------------------|
| 120 | **r/w** |
| 232 | **x** |
| 456 | **r** |

**PKRU (core_b)**

| pkey: | 0 | **1** | | **8** | | 15 |
|-------|---|-------|-----|-------|-----|-----|
| perm: | r/w | **r/w** | ... | **r/o** | ... | n/a |

**(core_b)**

| page# | effective perm. |
|-------|------------------|
| 120 | **r** |
| 232 | **r/x** |
| 456 | **r** |

*(per-core, asynchronous)*
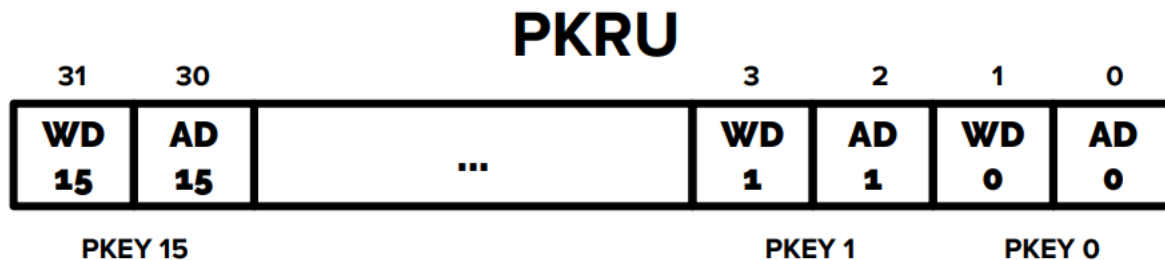
**red** : pkey = 1
**blue** : pkey = 8

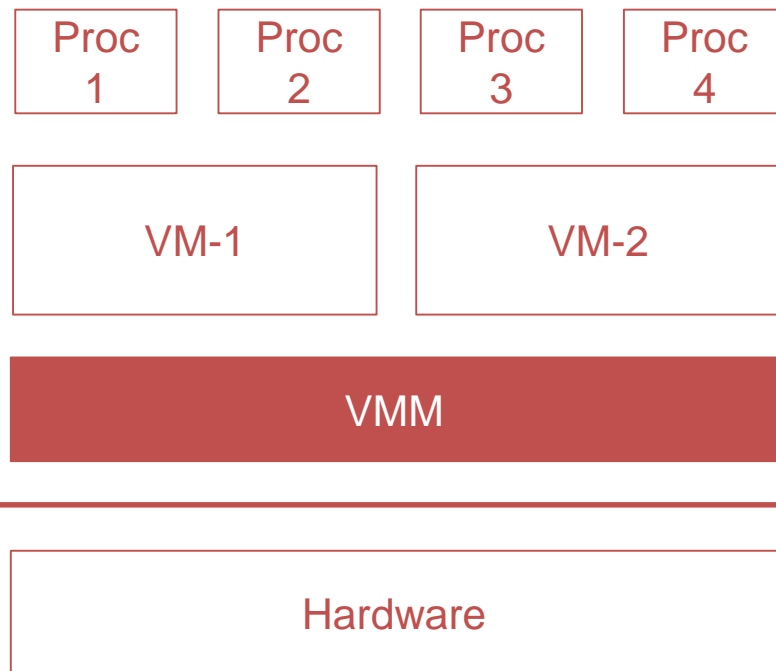37

# Techniques for System Isolation

- CPU privileges

- Segments

- Page table

- Extended page table

- Memory domain

- **Secure hardware modules**

# Secure Hardware Modules

- **Many companies has released their secure hardware modules**
  - Intel SGX
  - AMD SME/SEV

- **Encrypted memory for processes or VMs**
  - Privileged software cannot access these memory

- **Will be introduced in subsequent courses**

# Principles for System Isolation

- **Fine-grained isolation**

- **Reduced attack surface**

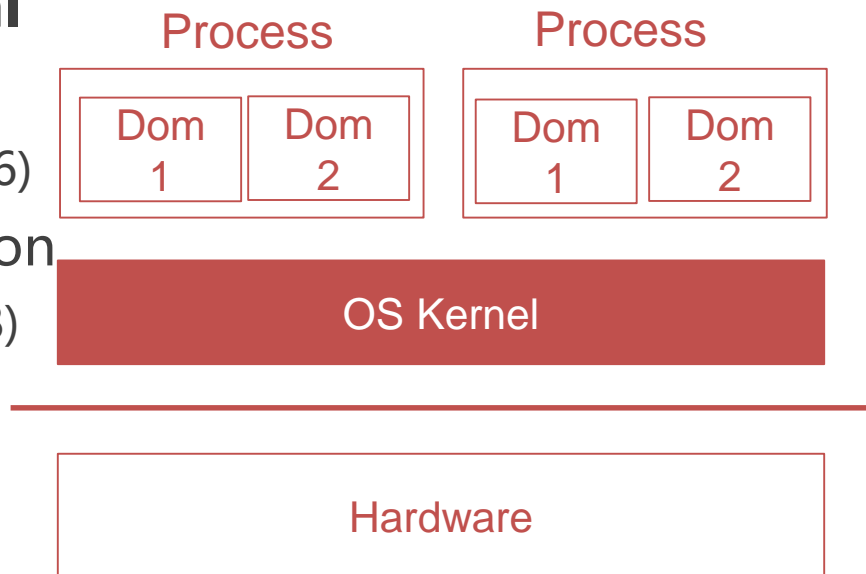- **Small TCB and Reference monitor**

- **Defense in depth**

| Proc 1 | Proc 2 | Proc 3 | Proc 4 |
|--------|--------|--------|--------|

| VM-1 | VM-2 |
|------|------|

**VMM**

Hardware

# Principles for System Isolation

- **Fine-grained isolation**

- **Reduced attack surface**

- **Small TCB and Reference monitor**

- **Defense in depth**

# Fine-grained isolation

- **A process/kernel usually contains different modules**
  - Vulnerable modules and confidential ones
  - Untrusted libraries and trusted libraries

- **These modules should have different privileges to the sensitive data**

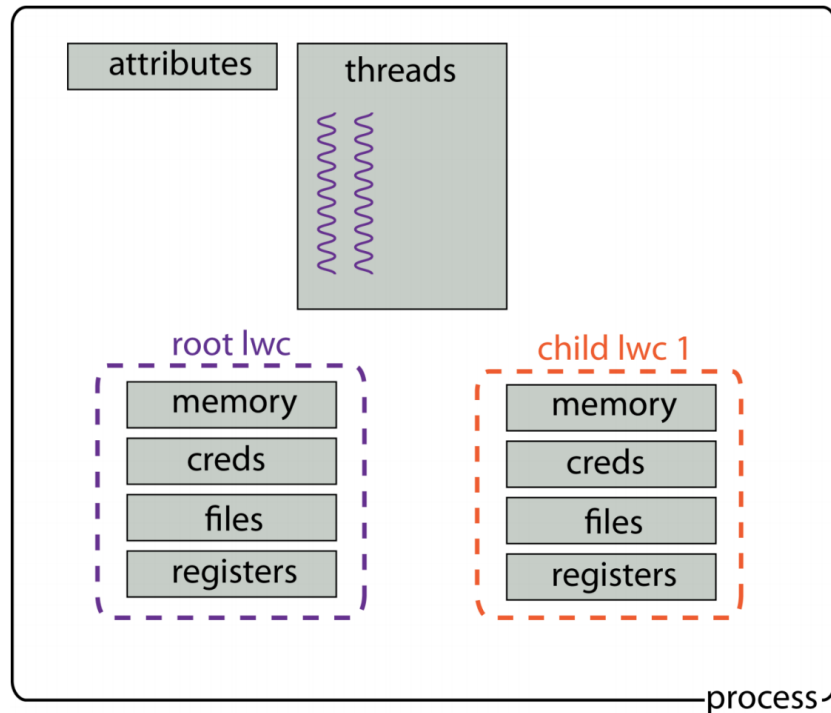- **There is a need to provide fine-grained isolation**

# Fine-grained isolation

- **Many techniques can be leveraged to achieve this goal**
  - Page table
    - Light-weight context (OSDI 2016)
  - CPU privilege and Segmentation
    - Lord of the x86 Rings (CCS 2018)
  - EPT and VMFUNC
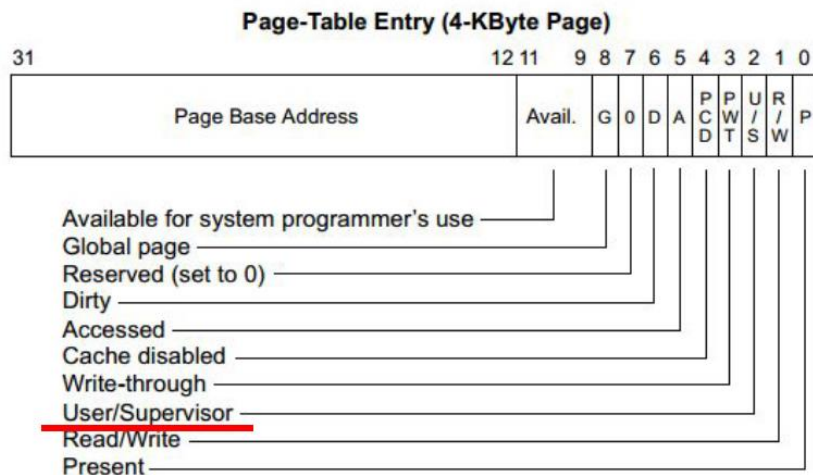    - SeCage (CCS 2015)
  - MPK
    - Hodor (ATC 2019)

| Process | | Process | |
|---|---|---|---|
| Dom 1 | Dom 2 | Dom 1 | Dom 2 |

OS Kernel

Hardware

# Light-Weight Context (LWC)

- **Multiple LWC for one process**

- **Each LWC has a private page table**

- **Switching LWC should get trapped into the kernel**
  - The kernel then changes related environment, including the page table



44

# Lord of the x86 Rings

- **Switching LWC is slow**
  - Cost 6050 cycles in a Skylake CPU

- **Why not using CPU rings?**
  - Ring 0-2 can access supervisor pages
  - Ring 3 cannot access these pages
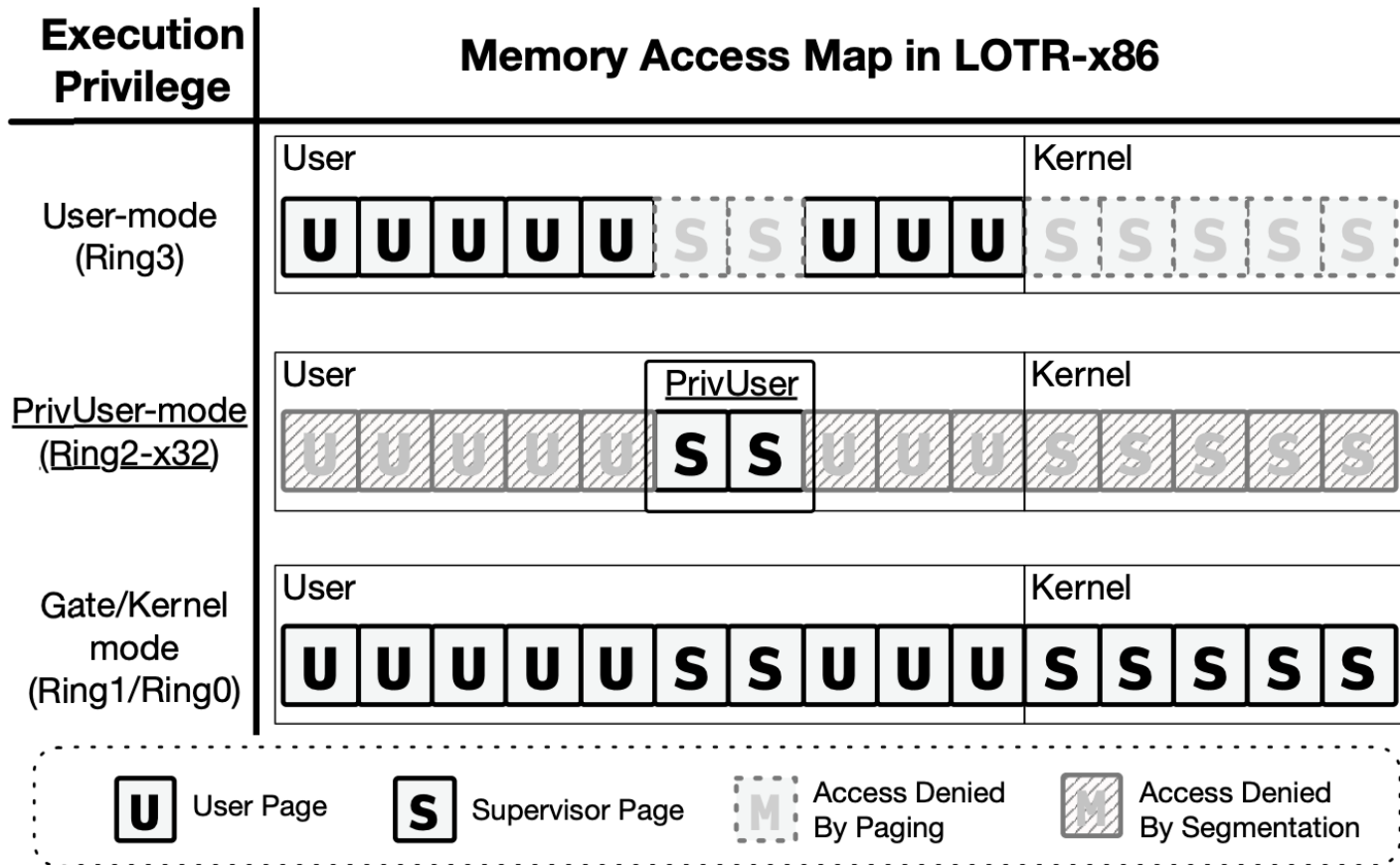  - Only Ring 0 can execute privilege instructions

**Page-Table Entry (4-KByte Page)**

```
31                          12 11  9 8 7 6 5 4 3 2 1 0
 ┌──────────────────────────┬──────┬─┬─┬─┬─┬─┬─┬─┬─┬─┬─┐
 │                          │      │ │ │ │ │ │P│P│U│R│ │
 │   Page Base Address      │Avail.│G│0│D│A│P│W│/│/│P│
 │                          │      │ │ │ │ │C│T│S│W│ │
 │                          │      │ │ │ │ │D│ │ │ │ │
 └──────────────────────────┴──────┴─┴─┴─┴─┴─┴─┴─┴─┴─┴─┘

Available for system programmer's use
Global page
Reserved (set to 0)
Dirty
Accessed
Cache disabled
Write-through
User/Supervisor
Read/Write
Present
```

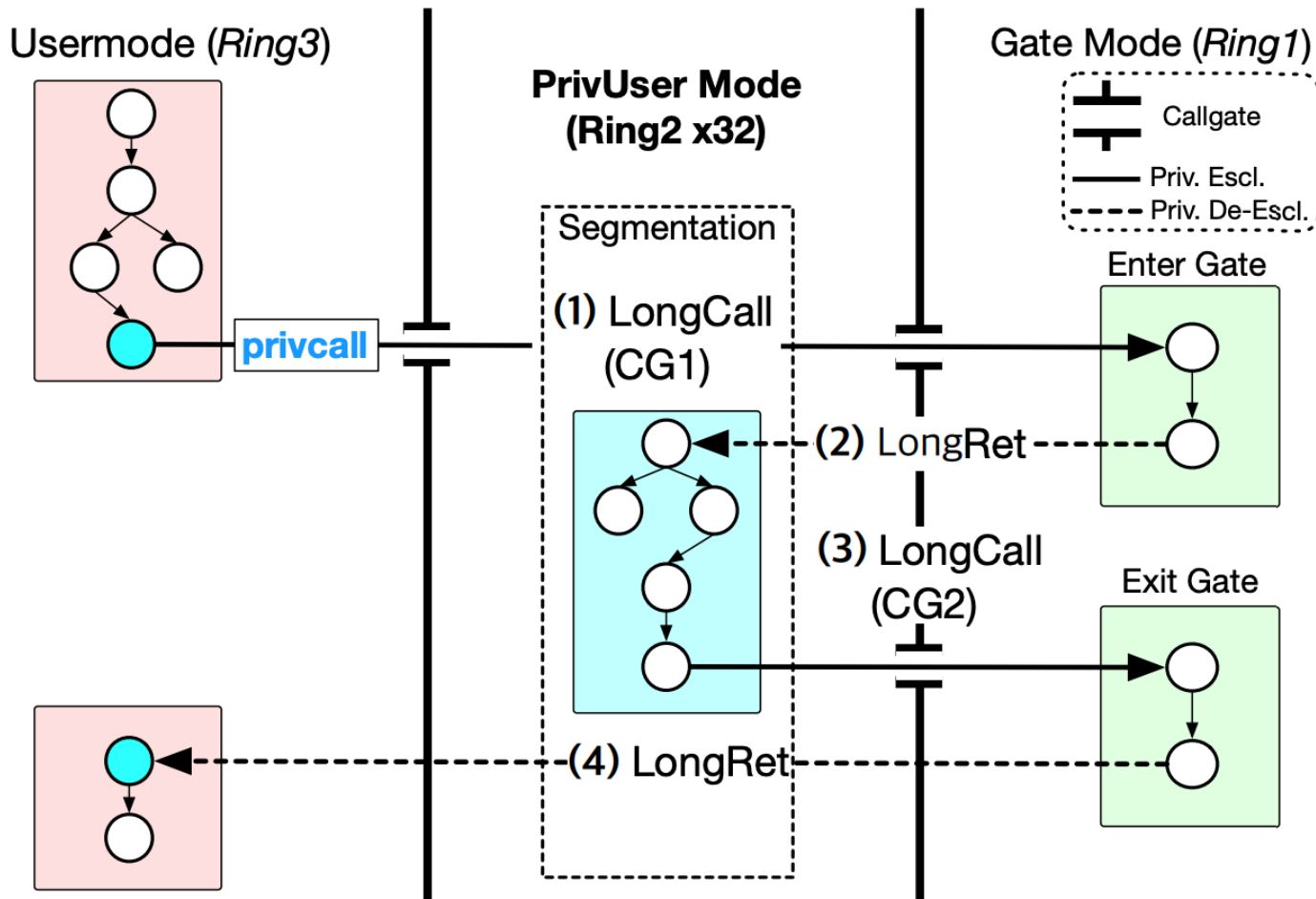|                         | Ring0 | **Ring1** | **Ring2** | Ring3 |
|-------------------------|-------|-----------|-----------|-------|
| Privileged instruction  | ✓     | ✗         | ✗         | ✗     |
| Supervisor page access  | ✓     | ✓         | ✓         | ✗     |

45

# Lord of the x86 Rings

- **Put privileged user to ring 1 or ring 2**

- **Map privileged user data and code to supervisor mode**

- **Question**
  - What if the program in Ring 1 or Ring 2 is malicious?
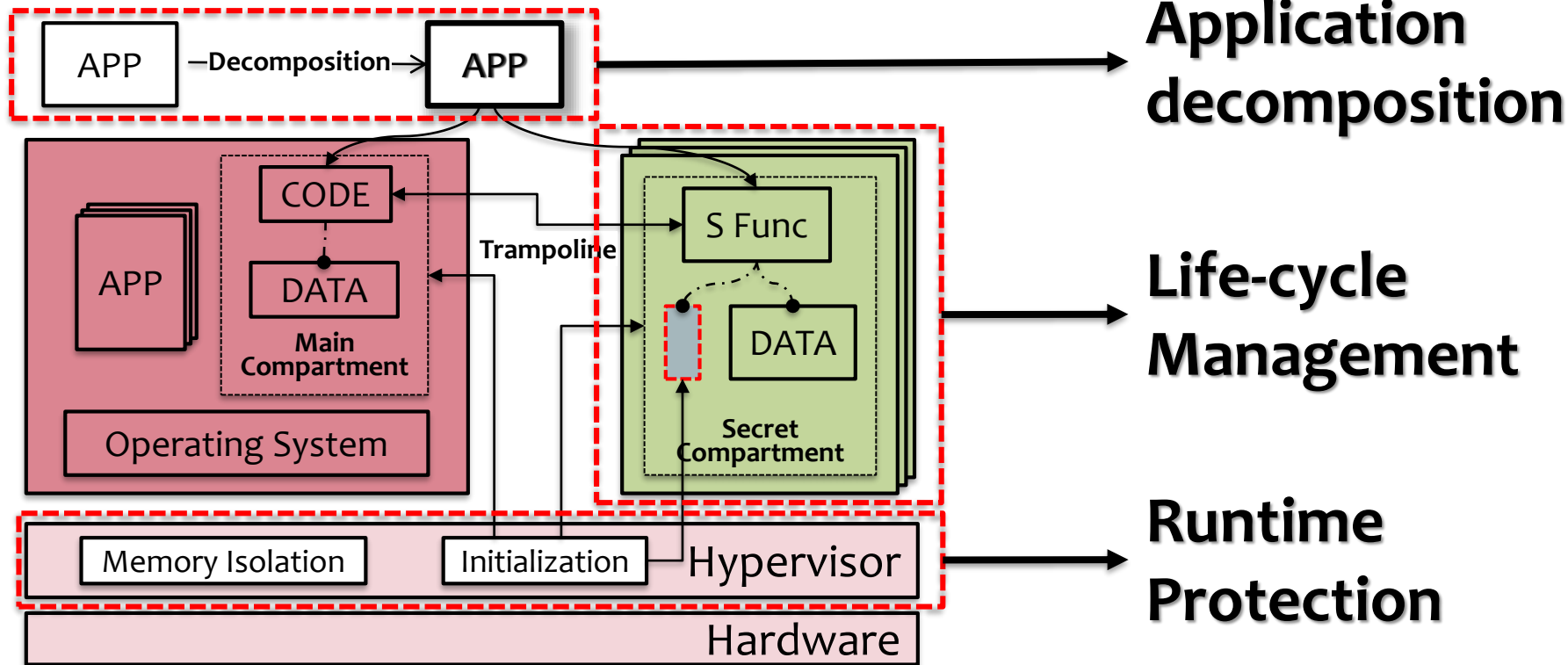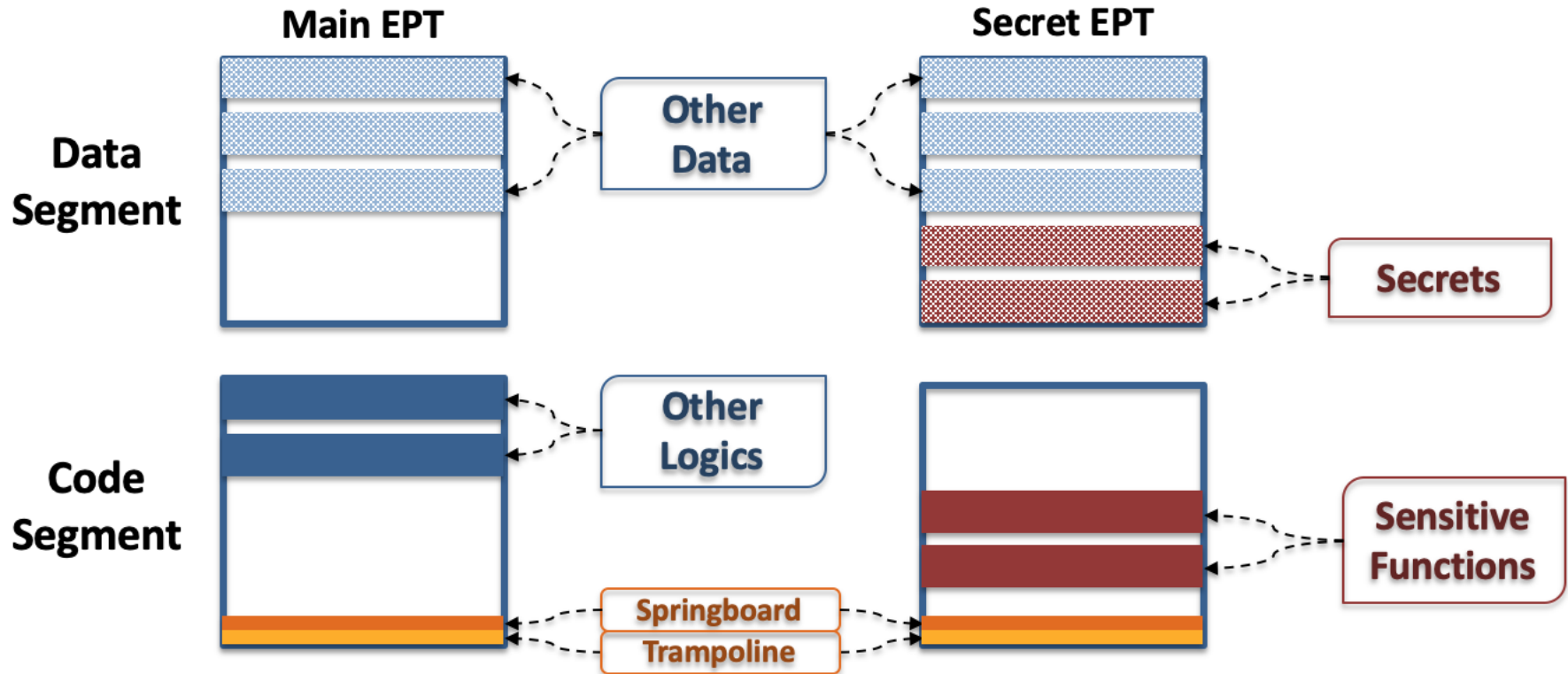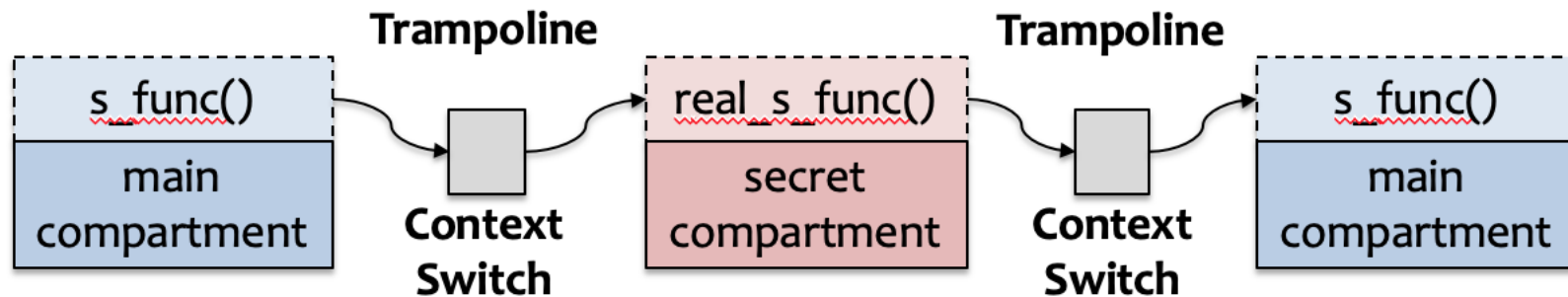  - It can access kernel data or code!

# Segmentation to the Rescue



**Execution Privilege** | **Memory Access Map in LOTR-x86**

User-mode (Ring3): User | U U U U U S S U U U | Kernel | S S S S S

PrivUser-mode (Ring2-x32): User | U U U U U | PrivUser | S S | U U U U | Kernel | S S S S S

Gate/Kernel mode (Ring1/Ring0): User | U U U U U S S U U U | Kernel | S S S S S

Legend:
- U — User Page
- S — Supervisor Page
- M — Access Denied By Paging
- M — Access Denied By Segmentation

# Lord of the x86 Rings

# SeCage



**Application decomposition**

**Life-cycle Management**

**Runtime Protection**

# SeCage: Different EPTs for two Parts

# SeCage: Using VMFUNC to Switch EPTs

- **Trampoline control flow**

# Hodor

# Hodor

- **Leverage MPK to provide different memory views**

Table 1: Latency of Basic Operations

| Instruction or Operation | Cycles* |
|---|---|
| write to `CR3` with `CR3_NOFLUSH` | $186 \pm 9$ |
| `vmfunc` | $109 \pm 1$ |
| `wrpkru` | $26 \pm 2$ |
| no-op system call w/ KPTI | $433 \pm 12$ |
| no-op system call w/o KPTI | $96 \pm 2$ |
| no-op VM call | $1694 \pm 131$ |
| user-space context switch | $748 \pm 8$ |
| process context switch using semaphore | $4426 \pm 41$ |

- **Have to replace illegal wrpkru instructions**

# Principles for System Isolation

- Fine-grained isolation

- **Reduced attack surface**

- Small TCB and Reference monitor

- Defense in depth

# **Reduced Attack Surface**

- **Attack surface**
  - Points where an unauthorized user (the "attacker") can try to enter data to or extract data
  - The interface between two entities
    - API level
    - System call level
    - ...

- **Keep attack surface as small as possible**
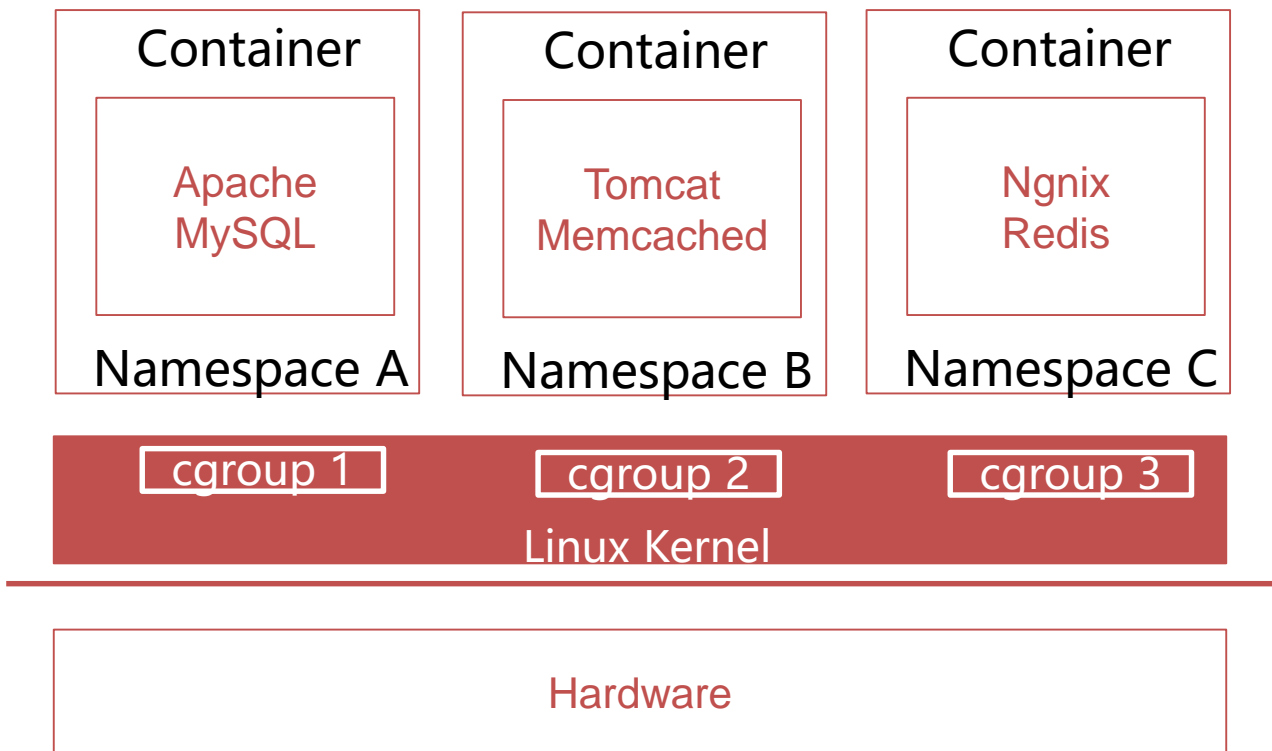  - A fundamental isolation principle

# Attack Surface

- **Interface number**
  - More than 350 system calls in Linux
  - Only 10 hypercalls in KVM
  - KVM has a much smaller attack surface if the call number matters

- **Invoked Code**
  - Small interface may cover a large number of code

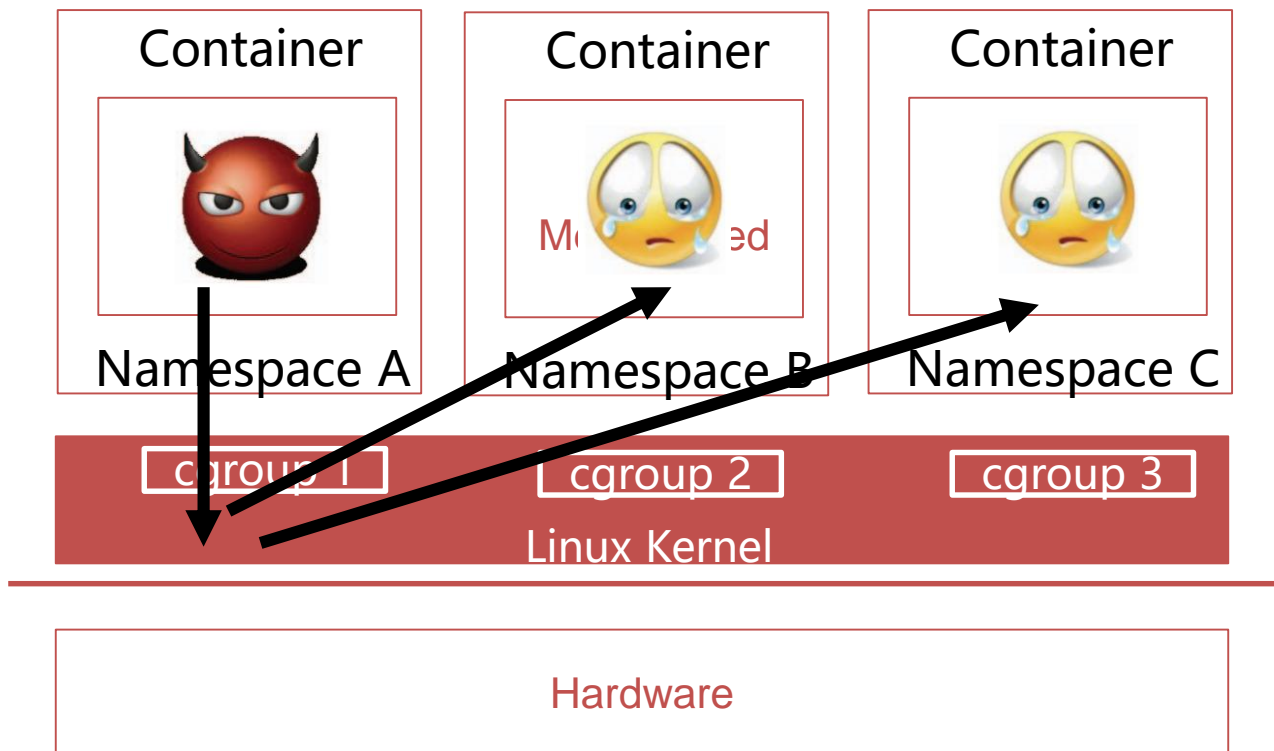| Proc 1 | Proc 2 | Proc 3 | Proc 4 |

| VM-1 | VM-2 |

| VMM |

# Container

- **Lightweight virtualization: container**
  - Process isolation
  - Namespace
  - Cgroups
- **All containers share the same Linux kernel**
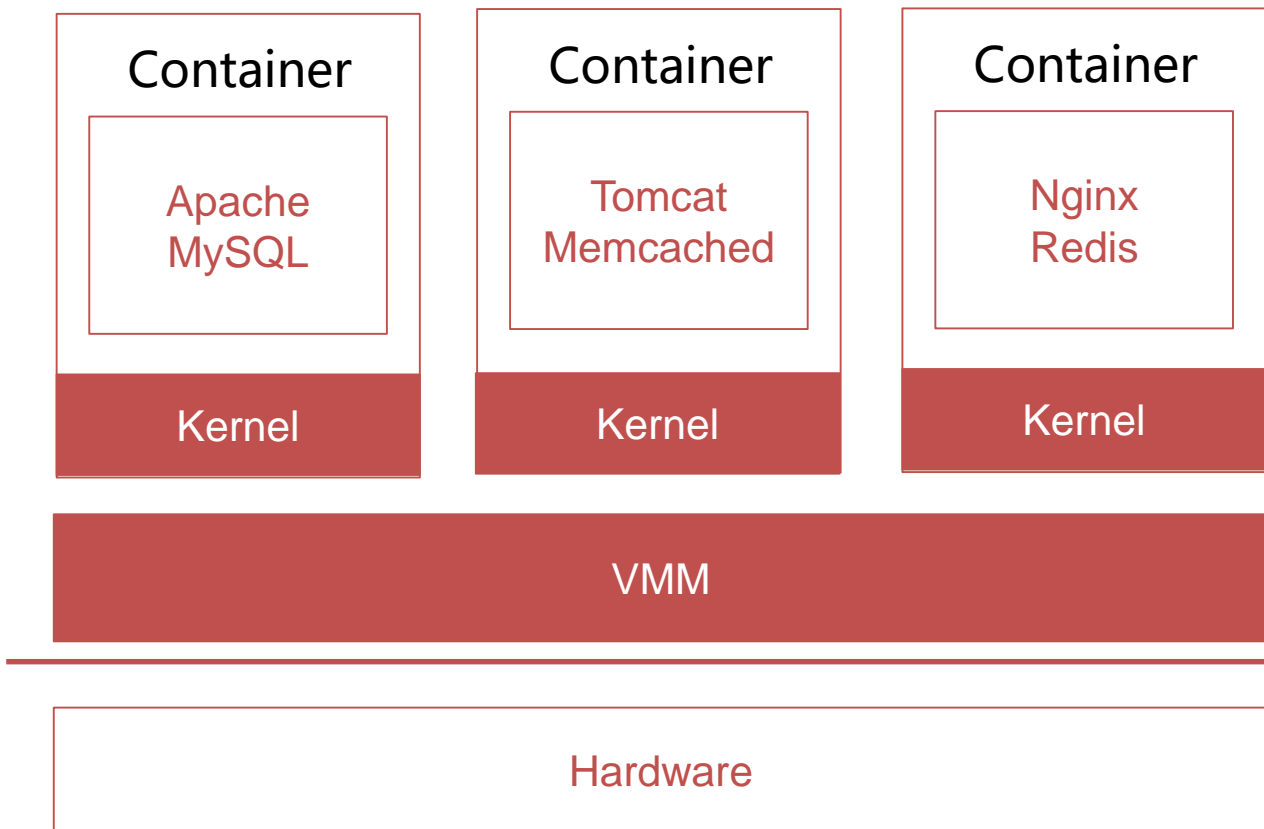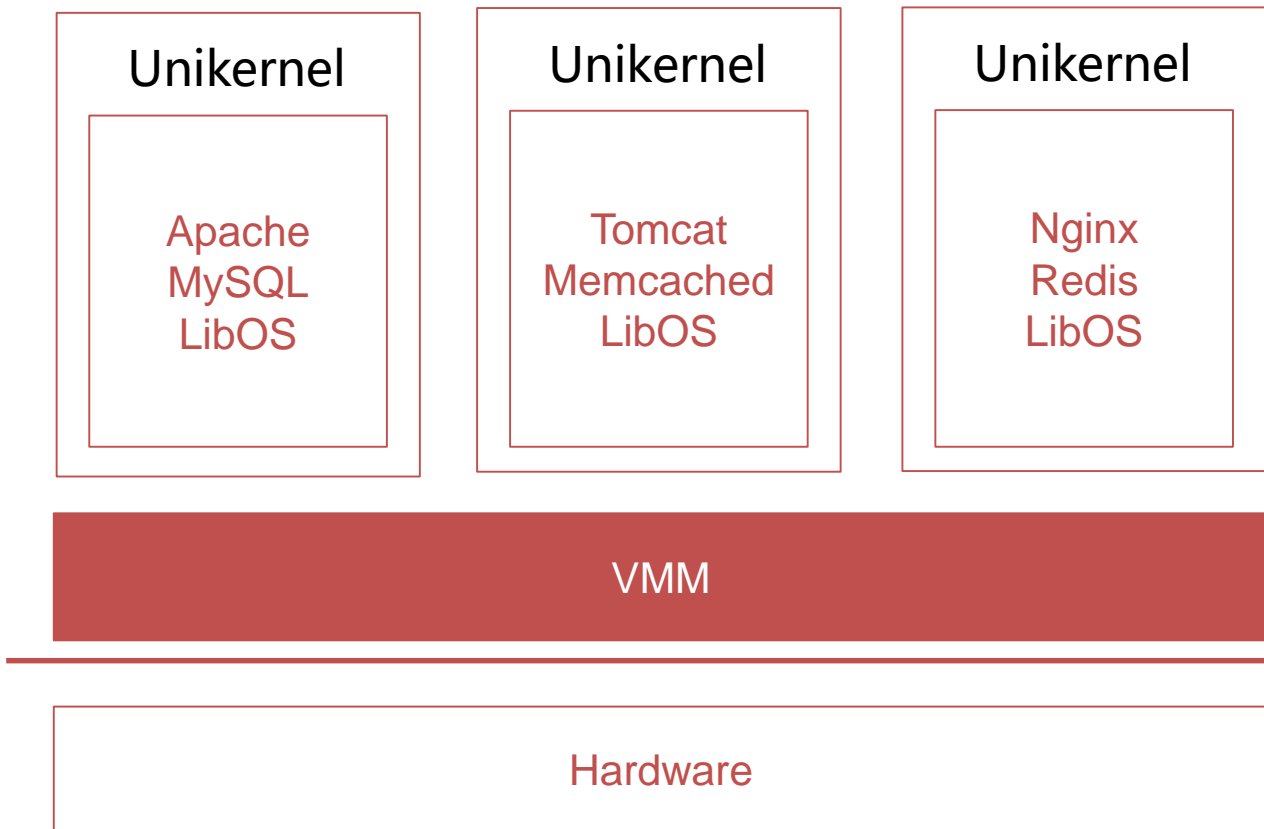- **Large attack surface**
  - System call

# Container

# Container: Weak Isolation
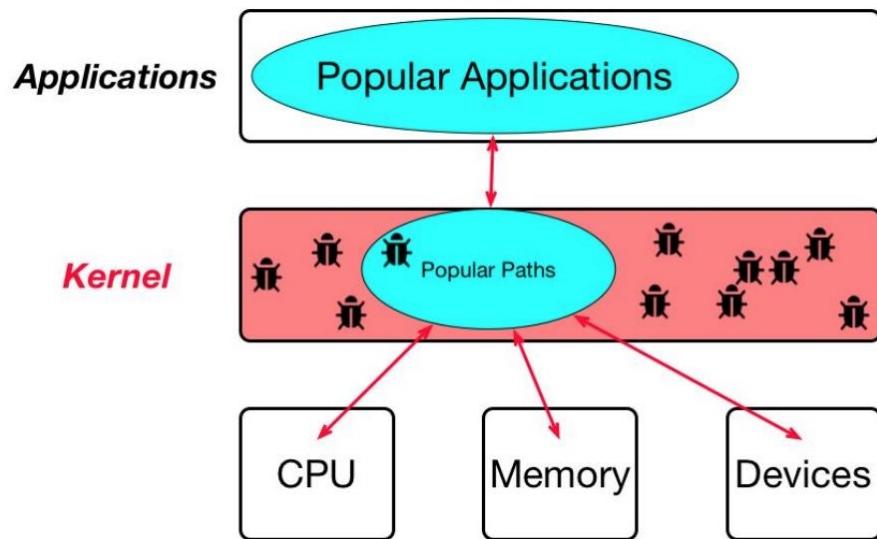
# Hyper Container

# Unikernel

| Unikernel | Unikernel | Unikernel |
|---|---|---|
| Apache MySQL LibOS | Tomcat Memcached LibOS | Nginx Redis LibOS |

**VMM**

**Hardware**

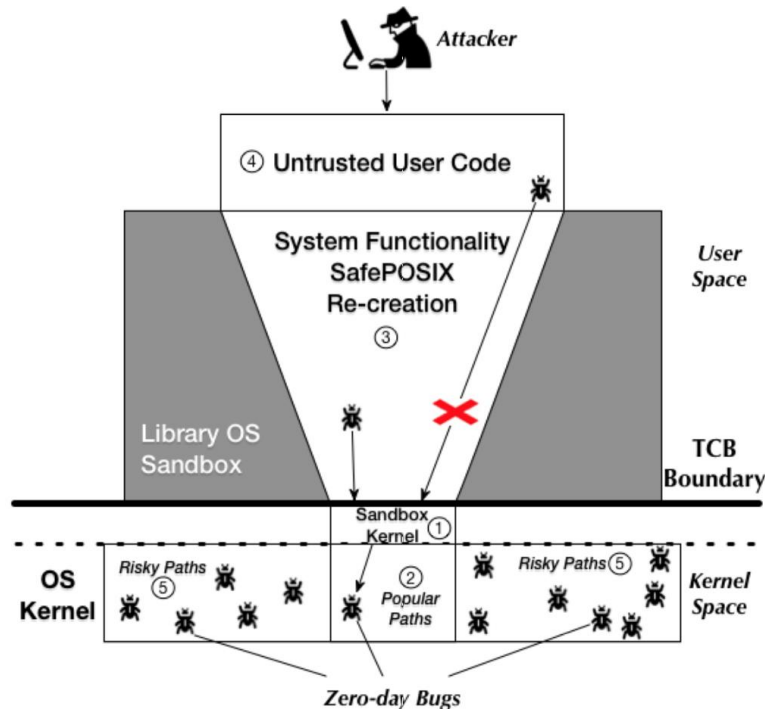# Bug Density in Linux Kernel

- **Bug density in Linux kernel**

# Popular Path

- **Definition**
  - lines of code in the kernel source files, which are commonly executed in the system's normal workload

- **Key insight**
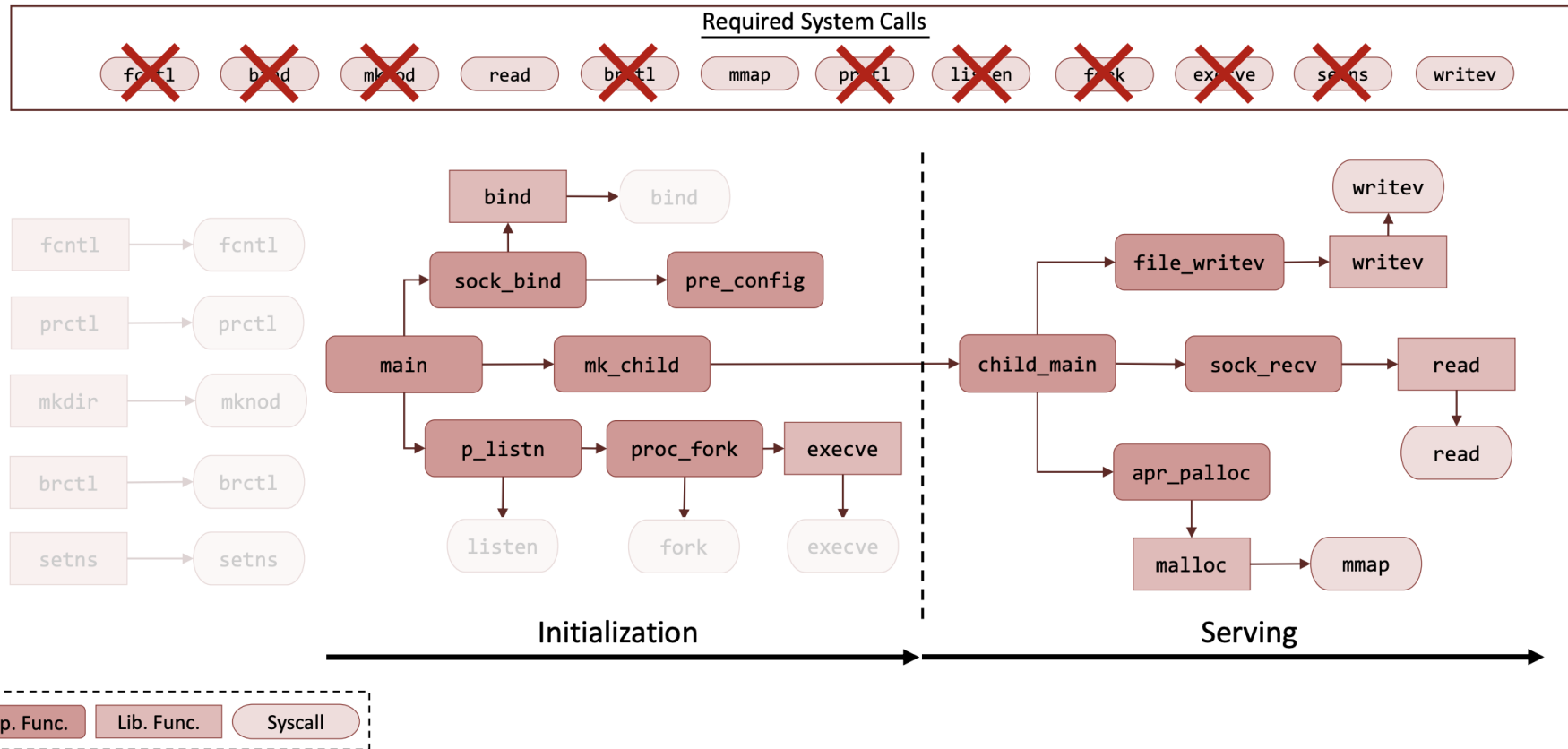  - the <span style="color:red">popular paths</span> contain many fewer bugs!

# Lock-in-Pop

- **Lock applications into using only popular paths**
- **Safely re-create file directories with basic calls like open(), read(), write(), close() to avoid using unpopular paths**

# Considering Timeline of a Process

- **Initialization phase**
  - Read configuration files
  - Fork worker processes
  - Execute other programs
  - Create files and set their permissions

- **Serving phase**
  - Handle client requests
  - Establish connections
  - ...

# Example: Apache Web Server

# Temporal System Call Specialization

- **Disable additional system calls that are needed only during the initialization phase, after entering the serving phase**

- Disables **51%** more security-critical system calls, breaking **218** more shellcodes and ROP payloads
- Mitigates **13** more Linux kernel CVEs

# Principles for System Isolation

- **Fine-grained isolation**

- **Reduced attack surface**

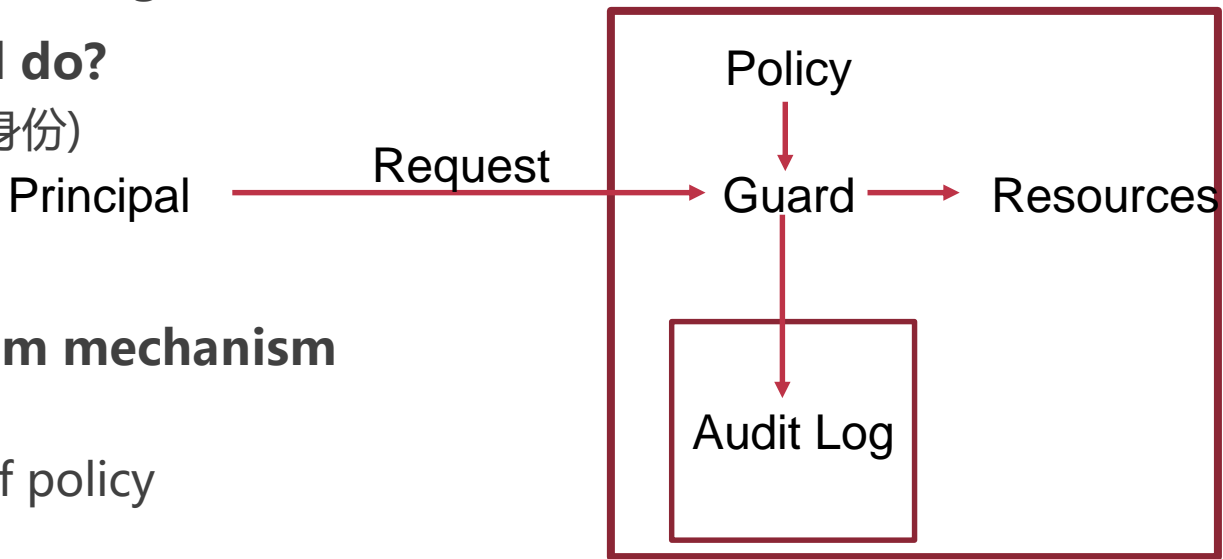- **Small TCB and Reference monitor**

- **Defense in depth**

# TCB

- **The Trusted Computing Base (TCB)**
  - A trusted component is a part of the system that we rely upon to operate correctly

- **TCB lets us separate the system into two parts**
  - The part that is security-critical (the TCB)
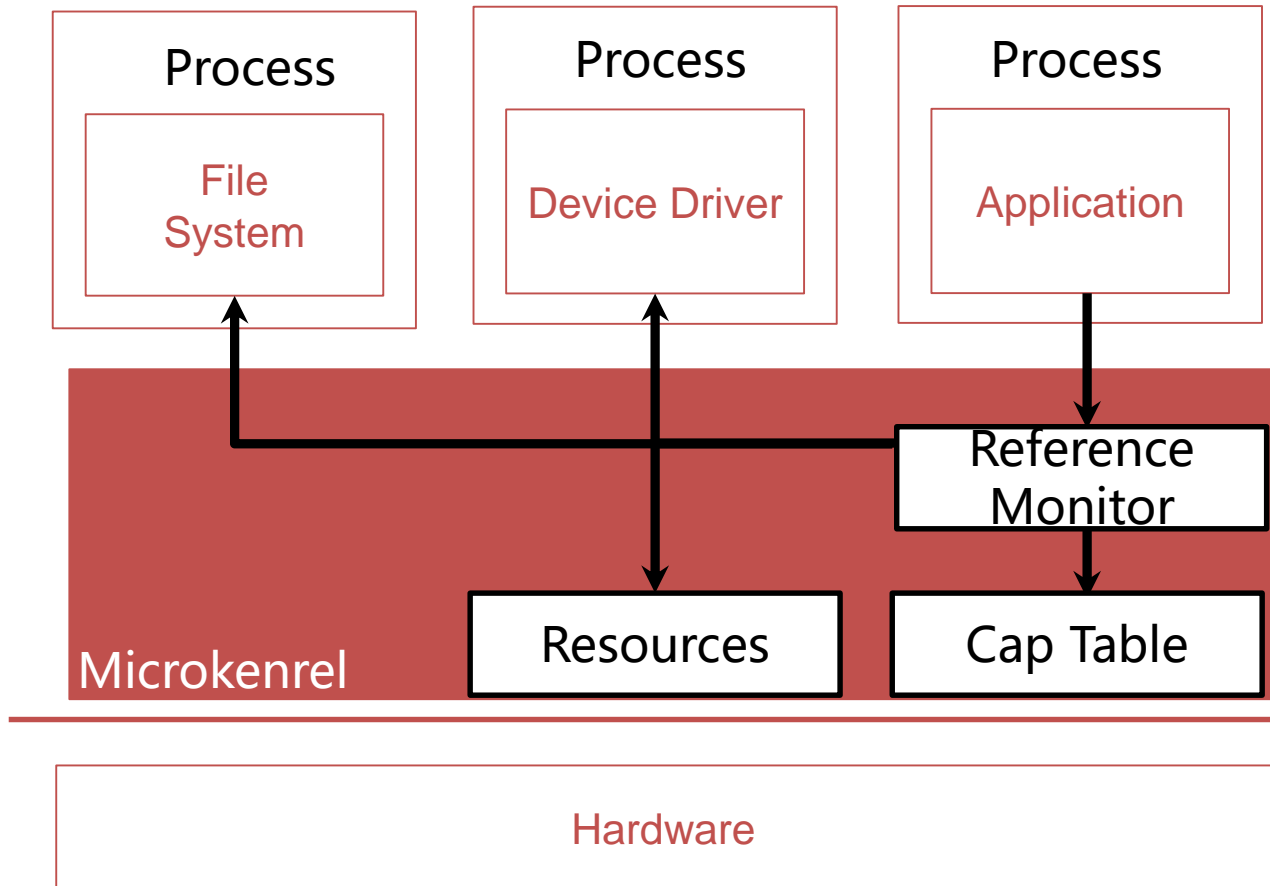  - Everything else.

# Reference Monitor

- **100% isolation is usually not what we want**

- **We need controlled sharing/interaction as well**

- **What does the guard do?**
  - Authenticate (验证身份)
  - Authorize (授权)
  - Audit (审计)

- **Separation policy from mechanism**
  - To ease reasoning
  - To ease evolution of policy

Policy

Principal → Request → Guard → Resources

Audit Log

# Reference Monitor

- **Reference monitor concept was defined in 1972 by James Anderson to describe design requirements on a "reference validation mechanism"**
    - The reference validation mechanism must always be invoked (complete mediation).
    - The reference validation mechanism must be tamperproof (tamperproof).
    - The reference validation mechanism must be small enough to be subject to analysis and tests, the completeness of which can be assured (verifiable).
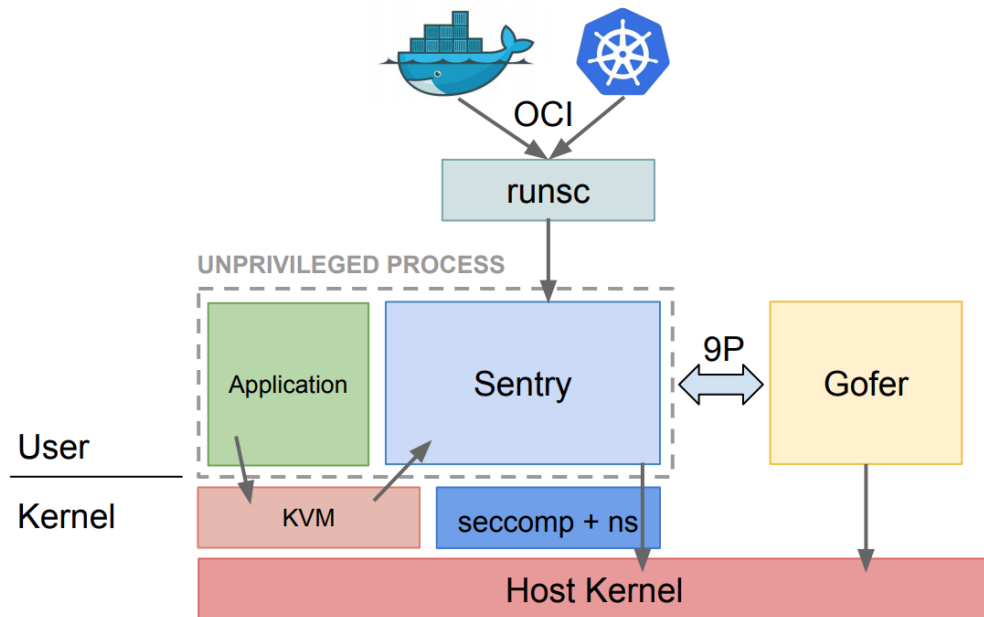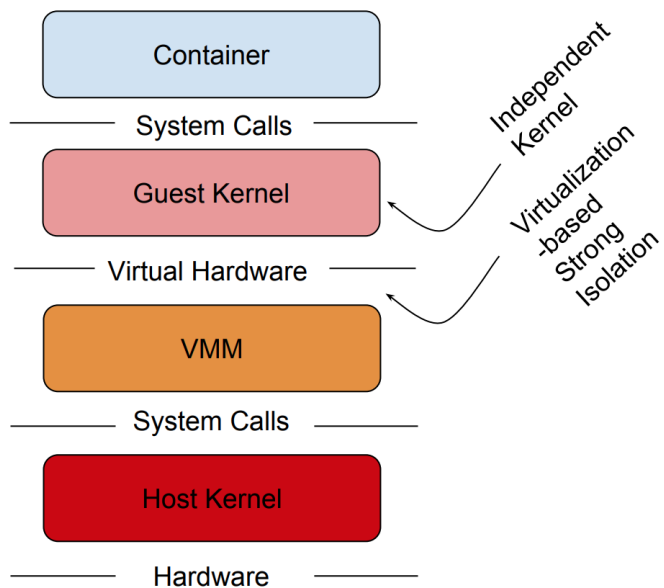
# Microkernel



Process

File System

Process

Device Driver

Process

Application

Microkenrel

Reference Monitor

Resources

Cap Table

Hardware

# Principles for System Isolation

- **Fine-grained isolation**

- **Reduced attack surface**

- **Small TCB and Reference monitor**
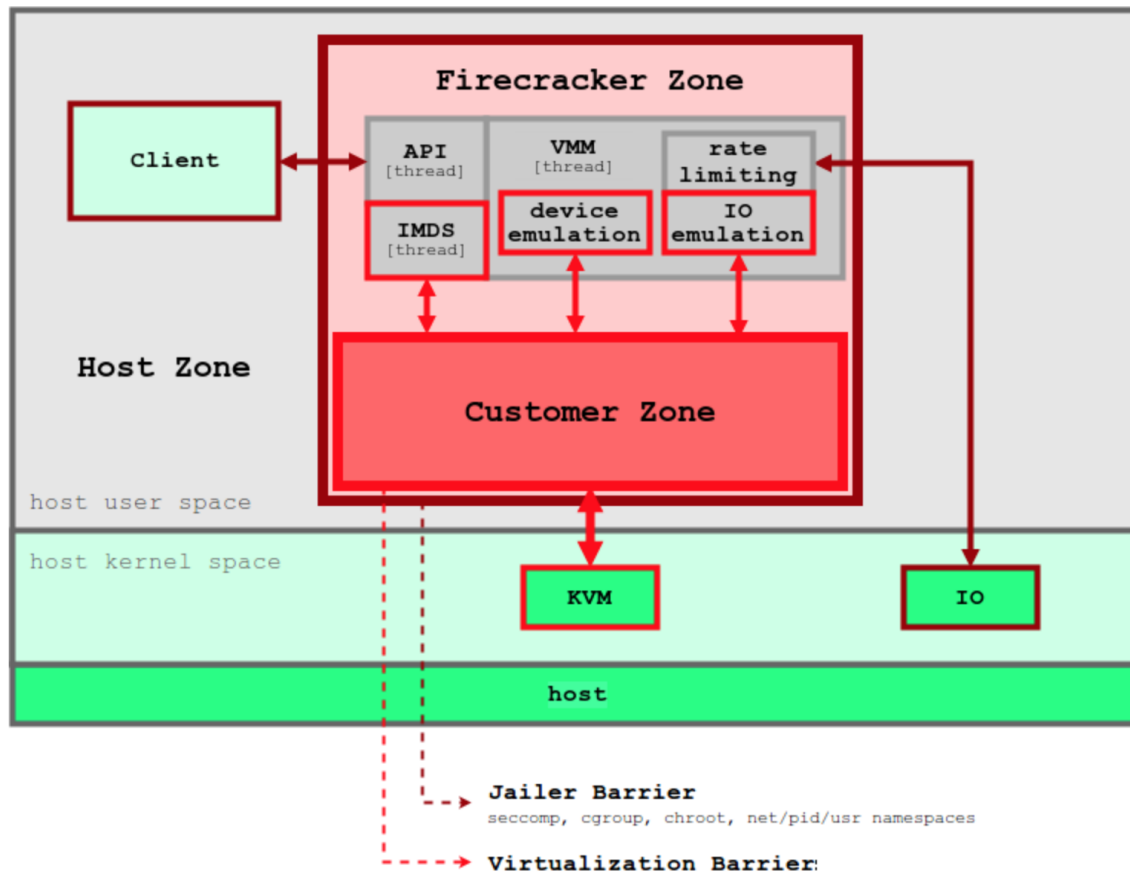
- **Defense in depth**

# Defense in Depth

- **The notion of layering multiple types of protection together**

- **Hypothesis is that attacker needs to breech all the defenses**

- **But defense in depth isn't free:**
  - You are throwing more resources at the problem
  - And although it can be better, it is less than the sum of the parts…

# Google gVisor



Container

System Calls

Guest Kernel

Virtual Hardware

VMM

System Calls

Host Kernel

Hardware

Independent Kernel

Virtualization -based Strong Isolation

OCI

runsc

UNPRIVILEGED PROCESS

Application

Sentry

9P

Gofer
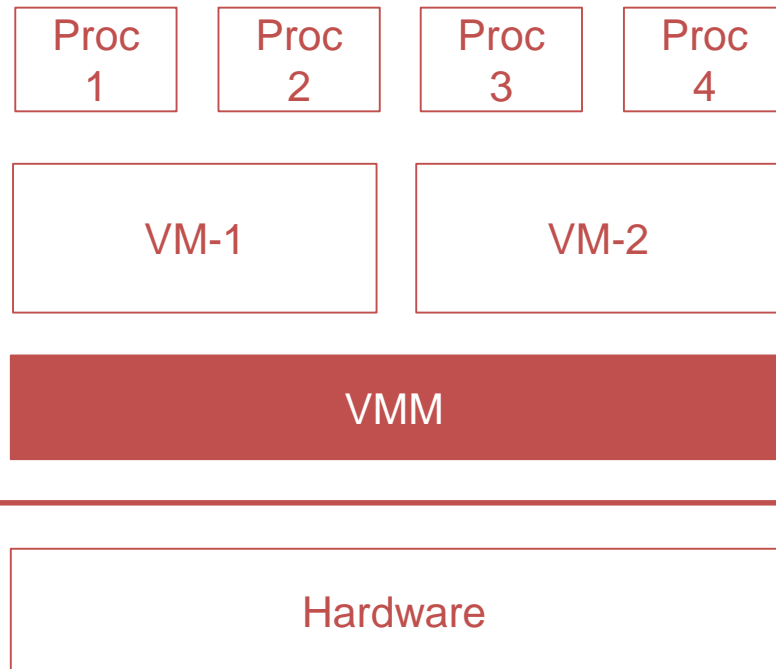
User

Kernel

KVM

seccomp + ns

Host Kernel

# Amazon Firecracker

# Principles for System Isolation

- **Fine-grained isolation**
  - Light-Weight Context
  - Lord of the x86 Ring
  - SeCage
  - Hodor

- **Reduced attack surface**
  - Hyper Container
  - Unikernel
  - Lock-in-Pop

- **Small TCB and Reference monitor**
  - Microkernel
  - Nested kernel

- **Defense in depth**
  - gVisor
  - Firecracker

| Proc 1 | Proc 2 | Proc 3 | Proc 4 |
|--------|--------|--------|--------|

| VM-1 | VM-2 |
|------|------|

**VMM**

Hardware

# Conclusion

- What is system isolation?

- Operating system and hypervisor

- Isolation tools

- Principles for system isolation

**Thanks!**