# *Hardware-Assisted Virtualization*
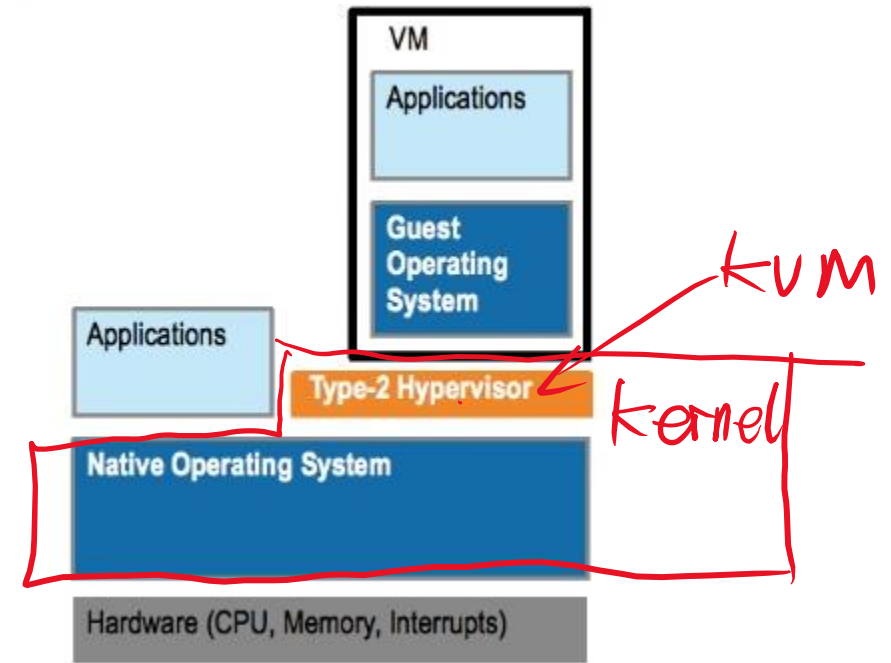
温瑞林

# CONTENTS

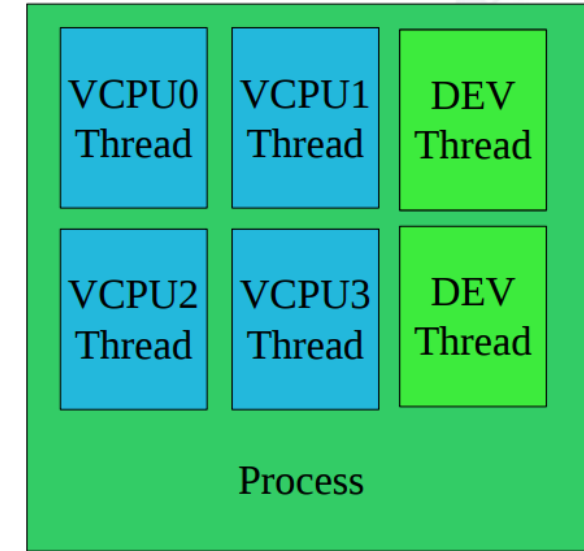What's KVM? **/01**

# What's KVM?

## Type 2 Hypervisor

# Abstraction Model

Design
1. One process per VM
2. One thread per-VCPU
3. Device models run concurrent in VCPU thread
4. Long running operations run in additional device thread



| | | |
|---|---|---|
| VCPU0 Thread | VCPU1 Thread | DEV Thread |
| VCPU2 Thread | VCPU3 Thread | DEV Thread |

Process

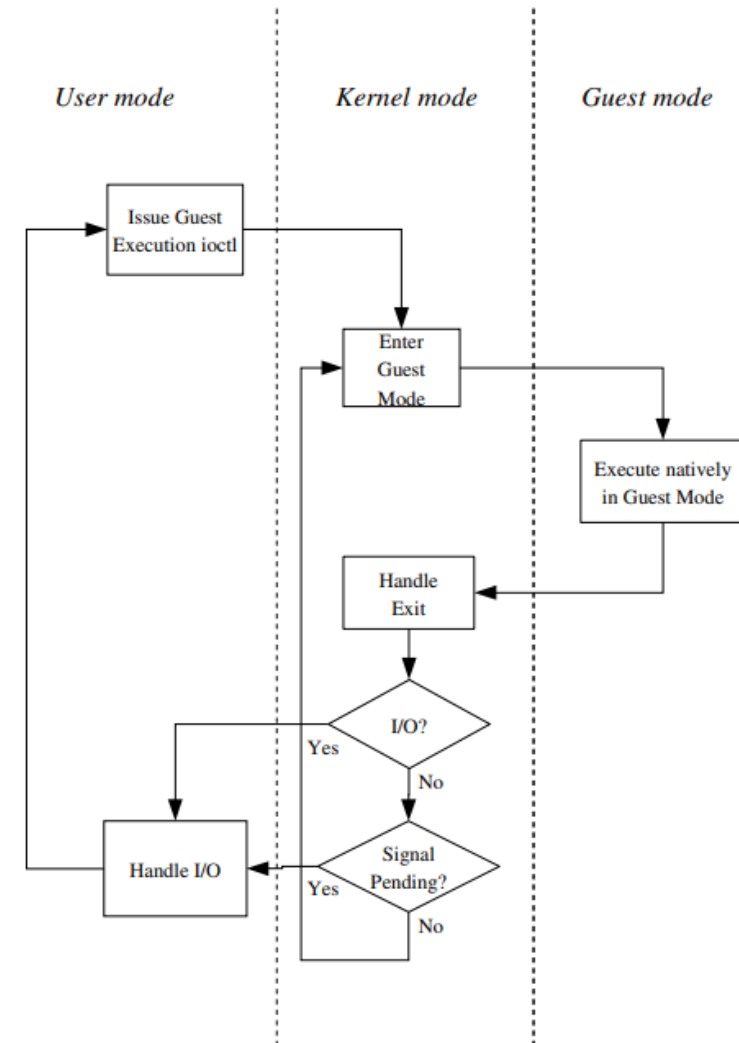| | | |
|---|---|---|
| Memory | 1.00 GiB/3.00 GiB | |
| Processors | 2 (1 sockets, 2 cores) | |
| BIOS | Default (SeaBIOS) | |
| Display | Standard VGA (std) | |
| Machine | q35 | |
| SCSI Controller | VirtIO SCSI | |
| Hard Disk (scsi0) | local-lvm:vm-101-disk-0,backup=0,size=25G,ssd=1 | |
| Hard Disk (scsi1) | BackupSsd:vm-101-disk-1,size=55064M | |
| Network Device (net0) | virtio=A2:43:2D:23:3C:3E,bridge=vmbr0,firewall=1 | |
| PCI Device (hostpci0) | 00:02.0,pcie=1,x-vga=1 | |
| Serial Port (serial0) | socket | |

| NI | VIRT | RES | SHR | S | CPU% | MEM% | TIME+ | Command |
|---|---|---|---|---|---|---|---|---|
| 3 | 167M | 11088 | 7916 | S | 0.0 | 0.1 | 1:24.84 | /sbin/init |
| 0 | 86172 | 2412 | 2224 | S | 0.0 | 0.0 | 0:11.71 | ─ /usr/sbin/pvefw-logger |
| 0 | 86172 | 2412 | 2224 | S | 0.0 | 0.0 | 0:11.71 | └ /usr/sbin/pvefw-logger |
| 0 | 21404 | 9372 | 7772 | S | 0.0 | 0.1 | 0:00.05 | ─ /lib/systemd/systemd --user |
| 0 | 168M | | 52 | S | 0.0 | 0.0 | 0:00.00 | └ (sd-pam) |
| 0 | 3910M | 3120M | 10932 | S | 4.0 | 39.7 | 9h47:08 | ─ /usr/bin/kvm -id 101 -name ubuntu-20.04 -no- |
| 0 | 3910M | 3120M | 10932 | S | 0.0 | 39.7 | 0:00.00 | ├ /usr/bin/kvm -id 101 -name ubuntu-20.04 - |
| 0 | 3910M | 3120M | 10932 | S | 1.3 | 39.7 | 4h51:31 | ├ /usr/bin/kvm -id 101 -name ubuntu-20.04 - |
| 0 | 3910M | 3120M | 10932 | S | 2.0 | 39.7 | 4h35:28 | ├ /usr/bin/kvm -id 101 -name ubuntu-20.04 - |
| 0 | 3910M | 3120M | 10932 | S | 0.0 | 39.7 | 0:00.05 | ├ /usr/bin/kvm -id 101 -name ubuntu-20.04 - |

# Programming Model

1. Event-Loop like model
2. Hide the detail of cross-platform hardware features

```
int runsz = ioctl(kvm_fd, KVM_GET_VCPU_MMAP_SIZE, 0);
struct kvm_run *run = (struct kvm_run *)
        mmap(       NULL, runsz,
                PROT_READ | PROT_WRITE,
                MAP_SHARED, vcpu_fd, 0);

for (;;) {

        ioctl(vcpu_fd, KVM_RUN, 0);
        switch (run->exit_reason) {
        case KVM_EXIT_IO:
                printf(……);
                break;
        case KVM_EXIT_SHUTDOWN:
                return;
        }

}
```
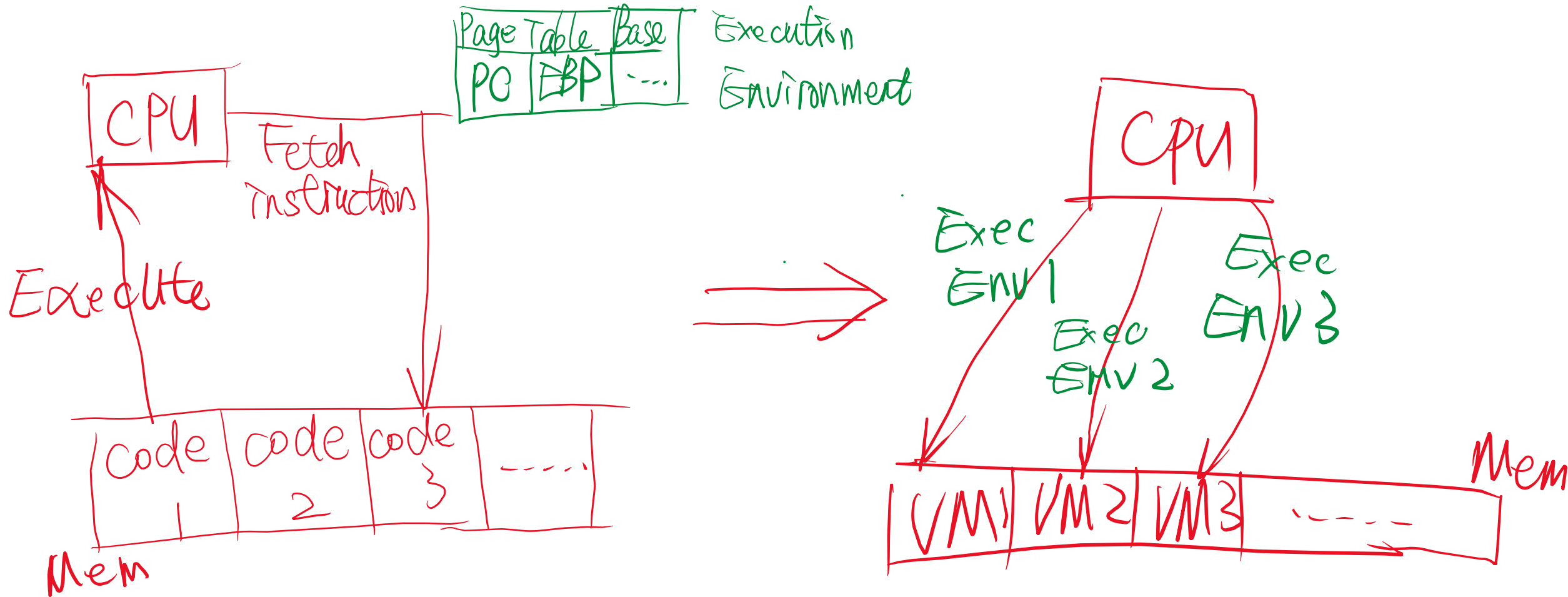
Hardware Virtualization and Implementation of KVM /02

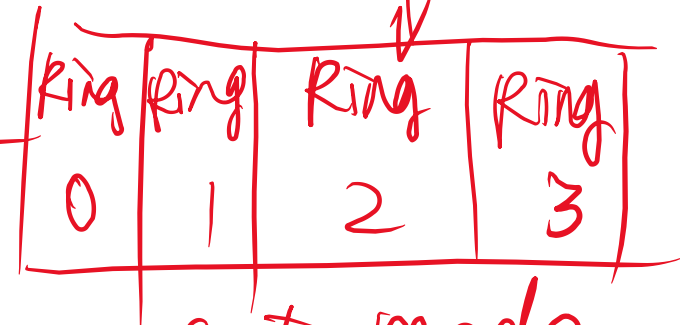# Hardware Virtualization and Implementation of KVM



| Page Table | Base |
|---|---|
| PC | EBP | .... |

Execution Environment

CPU

Fetch Instruction

Execute

Mem

| Code 1 | Code 2 | Code 3 | .... |

$\Rightarrow$

CPU

Exec Env 1

Exec Env 2

Exec Env 3

Mem

| VM1 | VM2 | VM3 | .... |

# Intel VT-x

Root Mode

Ring 3
Ring 2
Ring 1
Ring 0

(Execution Environment)

VMCS (Virtual Machine Control Structure)

VM Exit

Ring 3
Ring 2
Ring 1
Ring 0

Without Virtualization

Caused by sensitive operations or external interrupts
e.g. page fault, timer interrupt

Ring 0 | Ring 1 | Ring 2 | Ring 3
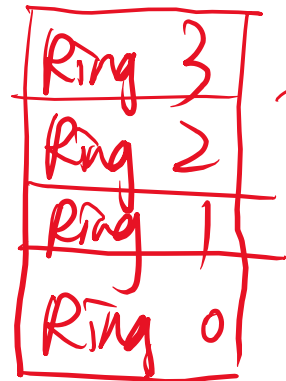
non-root mode

# Intel VT-x

Root Mode

| Ring 3 |
|--------|
| Ring 2 |
| Ring 1 |
| Ring 0 |

(Execution Environment)

VMCS (Virtual Machine Control Structure)

VM Exit

| Ring 0 | Ring 1 | Ring 2 | Ring 3 |
|--------|--------|--------|--------|

non-root mode

Caused by sensitive operations or external interrupts
e.g. page fault, timer interrupt

← root mode → non-root mode

| User mode | Kernel mode | Guest mode |
|-----------|-------------|------------|

Issue Guest Execution ioctl

Enter Guest Mode

Execute natively in Guest Mode

Handle Exit

I/O?
Yes / No

Signal Pending?
Yes / No

Handle I/O

libvirt | kernel module /dev/kvm

# ARM

Non-Secure state

| PL0 User |
| PL1 Kernel |
| PL2 Hyp |

Secure state

| PL0 User |
| PL1 Kernel |

*(handwritten)* PLO PLI

*(handwritten)* PLO PL2

*(handwritten)* PLO PLI

## ARM without virtualization

|  |  |  |
|---|---|---|
| PL 0 (User) | Host User | QEMU | VM User |
| PL 1 (Kernel) | Host Kernel | KVM **Highvisor** | VM Kernel |
|  |  | Trap | Trap |
| PL 2 (Hyp) |  | **Lowvisor** |  |

Figure 2: KVM/ARM System Architecture

*(handwritten)* Arm v7 & Arm v8

*(handwritten)* PLO Host User       VM User

*(handwritten)* PL1       VM Kernel

*(handwritten)* PL2 Host Kernel  KVM

*(handwritten)* Arm v8.1 VHE

Dune: Safe User-level Access
to Privileged CPU Features **/03**

# Idea



non-root mode

Guest OSes run at intended rings

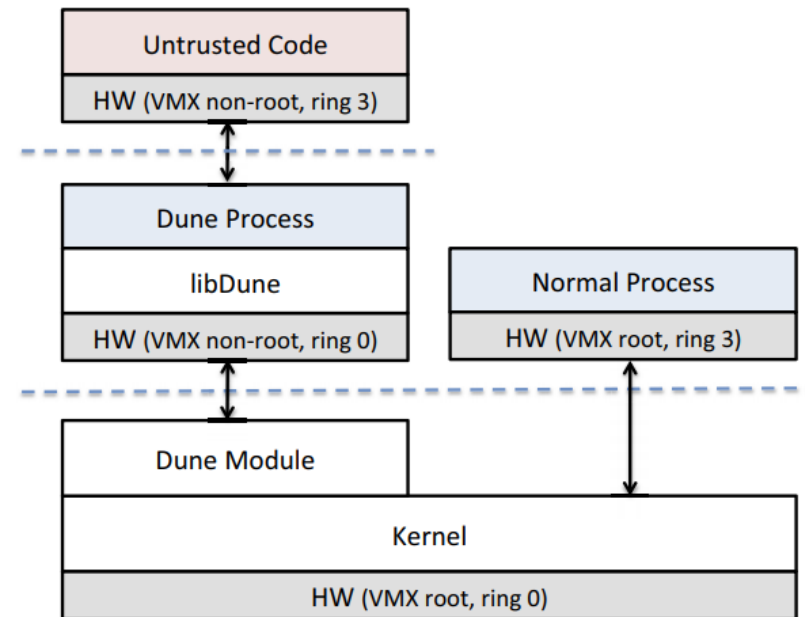Can we replace it with a host process?

Figure 1: The Dune system architecture.

# Benefits

Hardware features exposed by Dune and their corresponding privileged x86 instructions.

| Mechanism | Privileged Instructions |
|-----------|------------------------|
| Exceptions | LIDT, LTR, IRET, STI, CLI |
| Virtual Memory | MOV CRn, INVLPG, INVPCID |
| Privilege Modes | SYSRET, SYSEXIT, IRET |
| Segmentation | LGDT, LLDT |

Allow some applications to improve their performance by exploiting the exposed hardware features
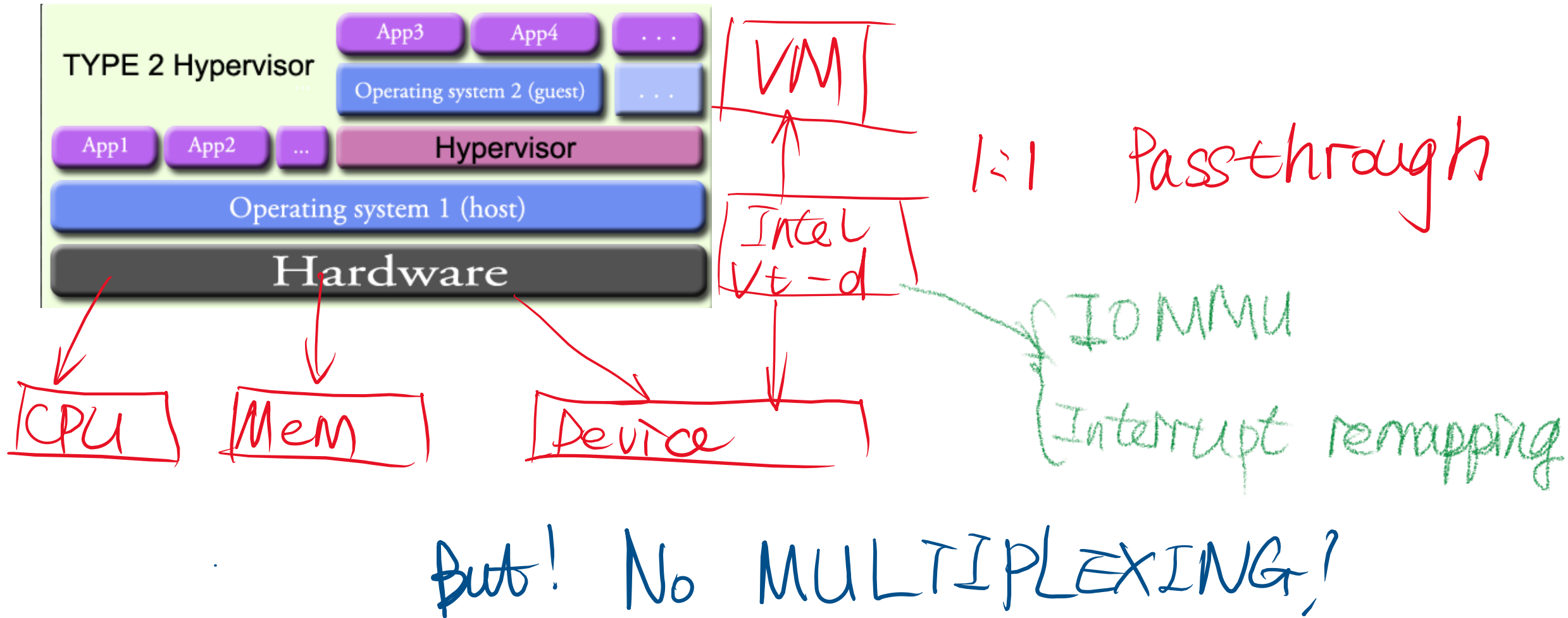
# Application?

1. Faster GC
   1. No syscall or IRQ overhead
   2. More exploitable information(Dirty bits…)
   3. TLB flush batching
2. Sandbox
   1. Faster context switching
   2. More flexible filtering policies on sensitive operations without kernel modification
3. Trace
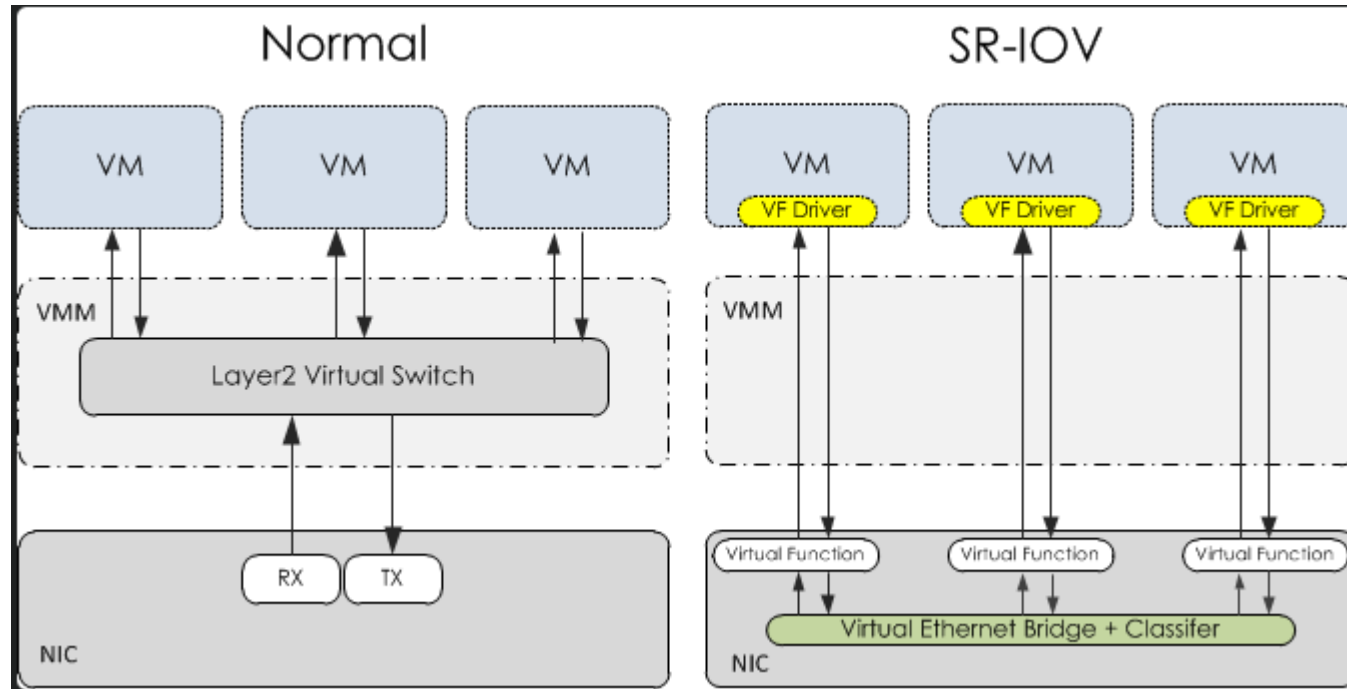   1. A better way to bypass anti-debugging mechanism?

# Arrakis: The Operating System is the Control Plane /04

# Idea



TYPE 2 Hypervisor

App3   App4   . . .

Operating system 2 (guest)   . . .

App1   App2   . . .   Hypervisor

Operating system 1 (host)

Hardware

VM

1:1  Passthrough

Intel
Vt-d

CPU   Mem   Device

IOMMU

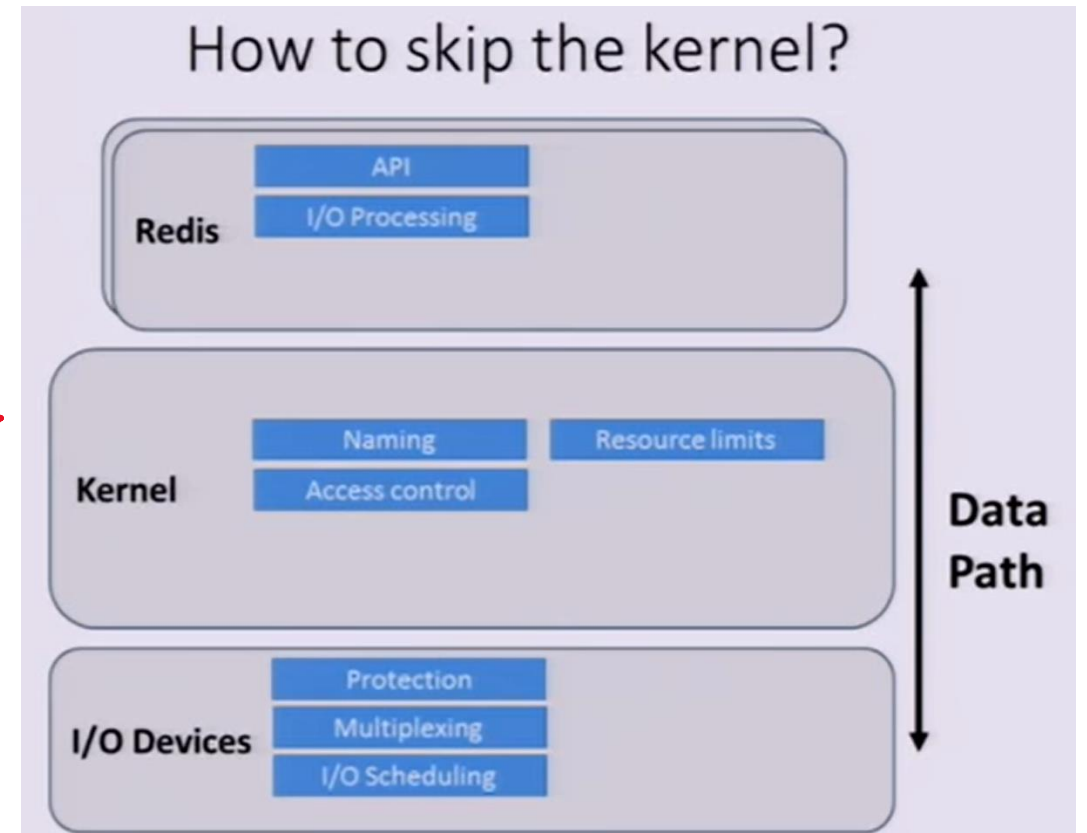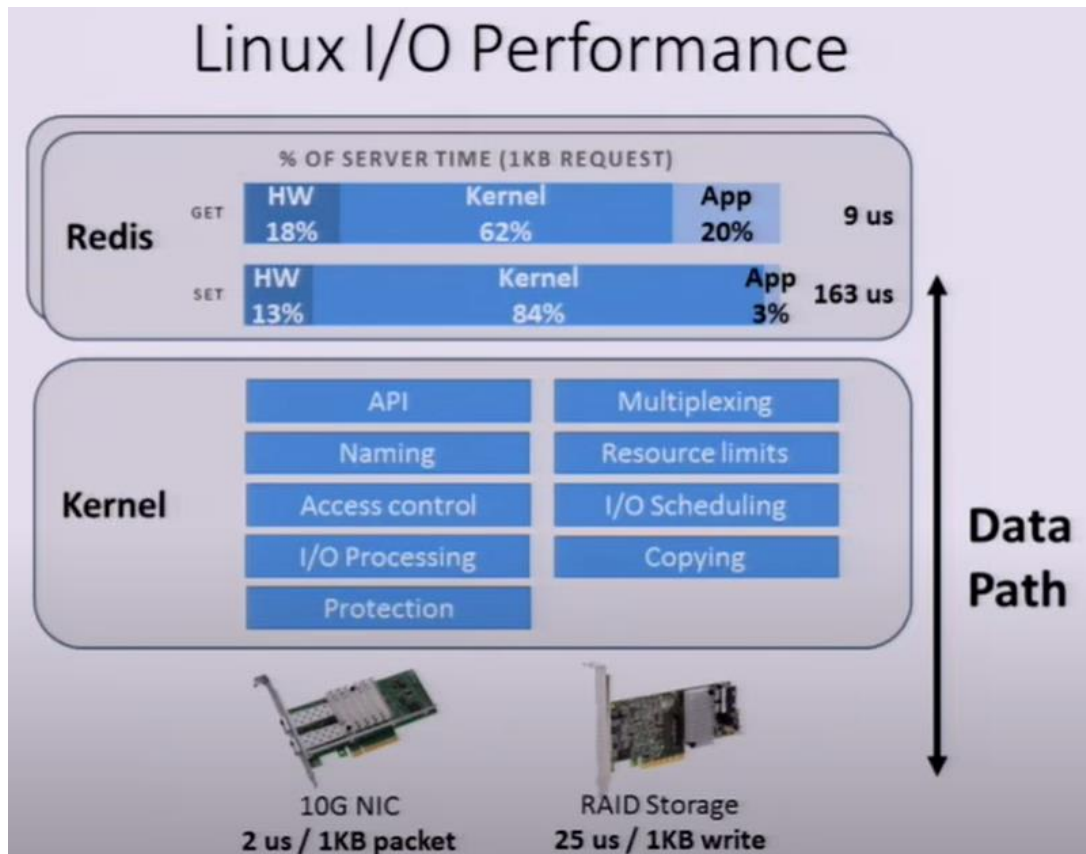Interrupt remapping

But! No MULTIPLEXING?

# Idea



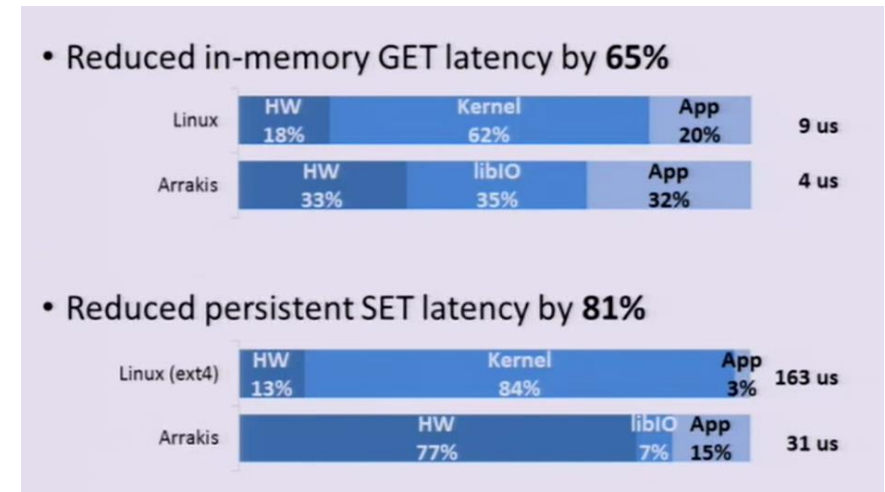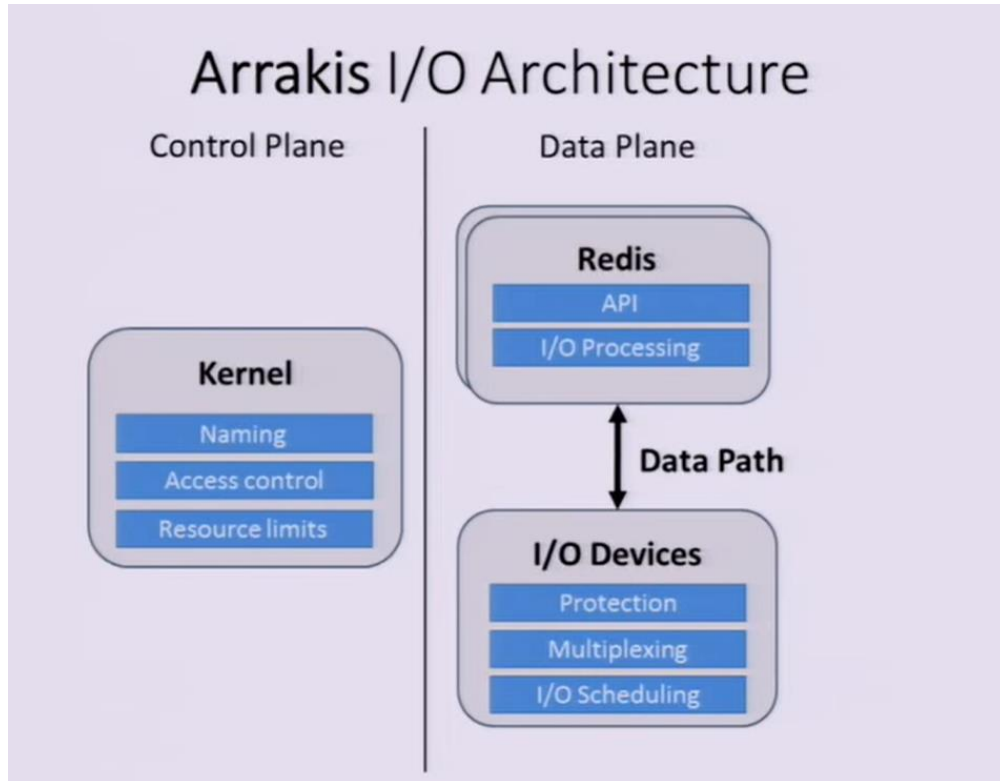With SR-IOV, all hardware resources are able to be multiplexed at the hardware level.

So can we reorganize the OS now?

# Idea

# Design and Result



Arrakis I/O Architecture

Control Plane | Data Plane

Kernel: Naming, Access control, Resource limits

Redis: API, I/O Processing

Data Path

I/O Devices: Protection, Multiplexing, I/O Scheduling

- Reduced in-memory GET latency by **65%**

| | HW | Kernel/libIO | App | |
|---|---|---|---|---|
| Linux | HW 18% | Kernel 62% | App 20% | 9 us |
| Arrakis | HW 33% | libIO 35% | App 32% | 4 us |

- Reduced persistent SET latency by **81%**

| | HW | Kernel/libIO | App | |
|---|---|---|---|---|
| Linux (ext4) | HW 13% | Kernel 84% | App 3% | 163 us |
| Arrakis | HW 77% | libIO 7% | App 15% | 31 us |

*DPDK? Openvurt? ⟹ Arrakis provides a generic architechture*

20

Summary **/05**

# Summary

1. KVM and Dune - How to integrate with the existing system gracefully
2. Arrakis – How to reorganize the existing system gracefully

Thanks