# CS 695: Virtualization and Cloud Computing

# Lecture 7: I/O Virtualization Techniques

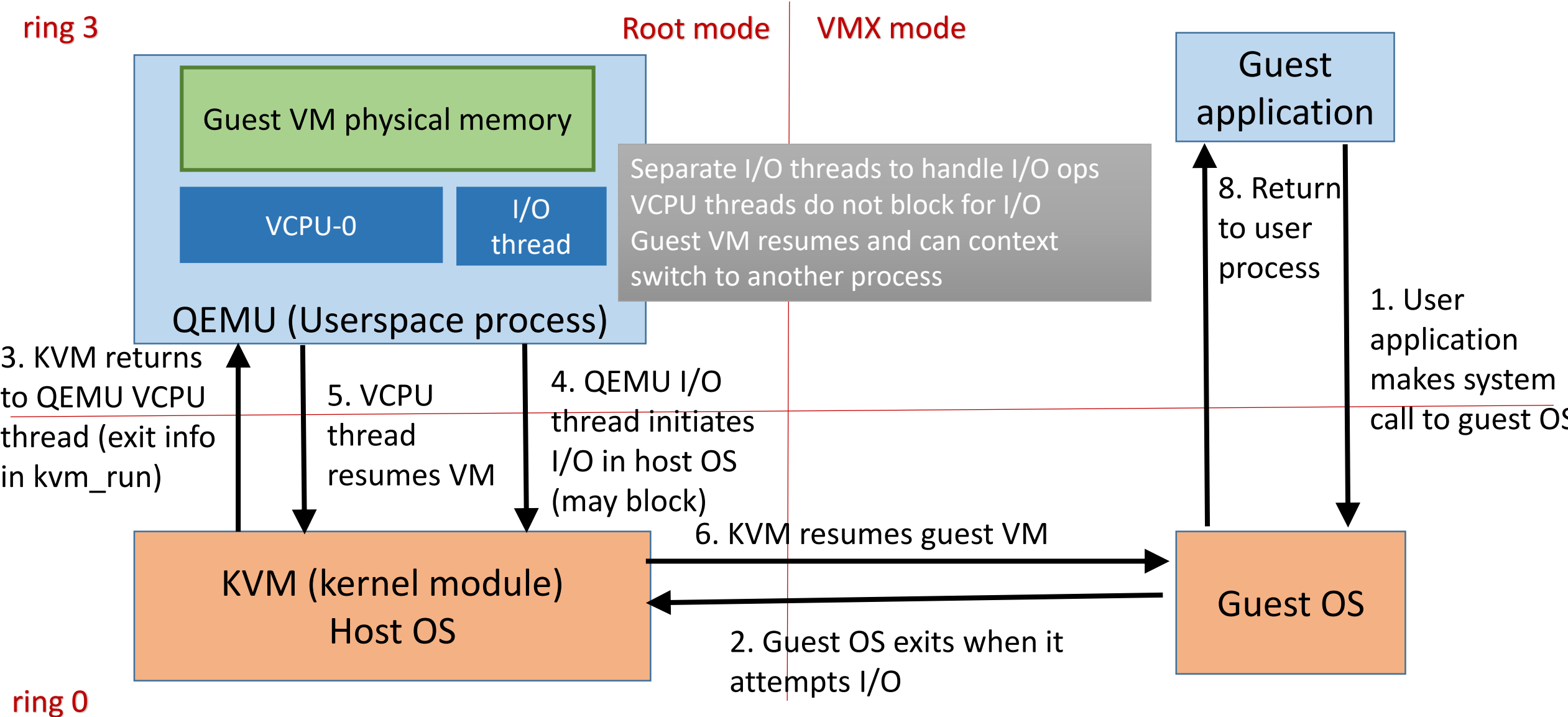Mythili Vutukuru

IIT Bombay

Spring 2021

# Techniques for I/O virtualization

- Guest OS cannot get full access to I/O devices
    - VMM must share I/O device access across guests
- Two ways to virtualize I/O devices:
    - Emulation: I/O access in guest traps to VMM, which performs I/O
    - Direct I/O or device passthrough: a slice of device is assigned directly to guest
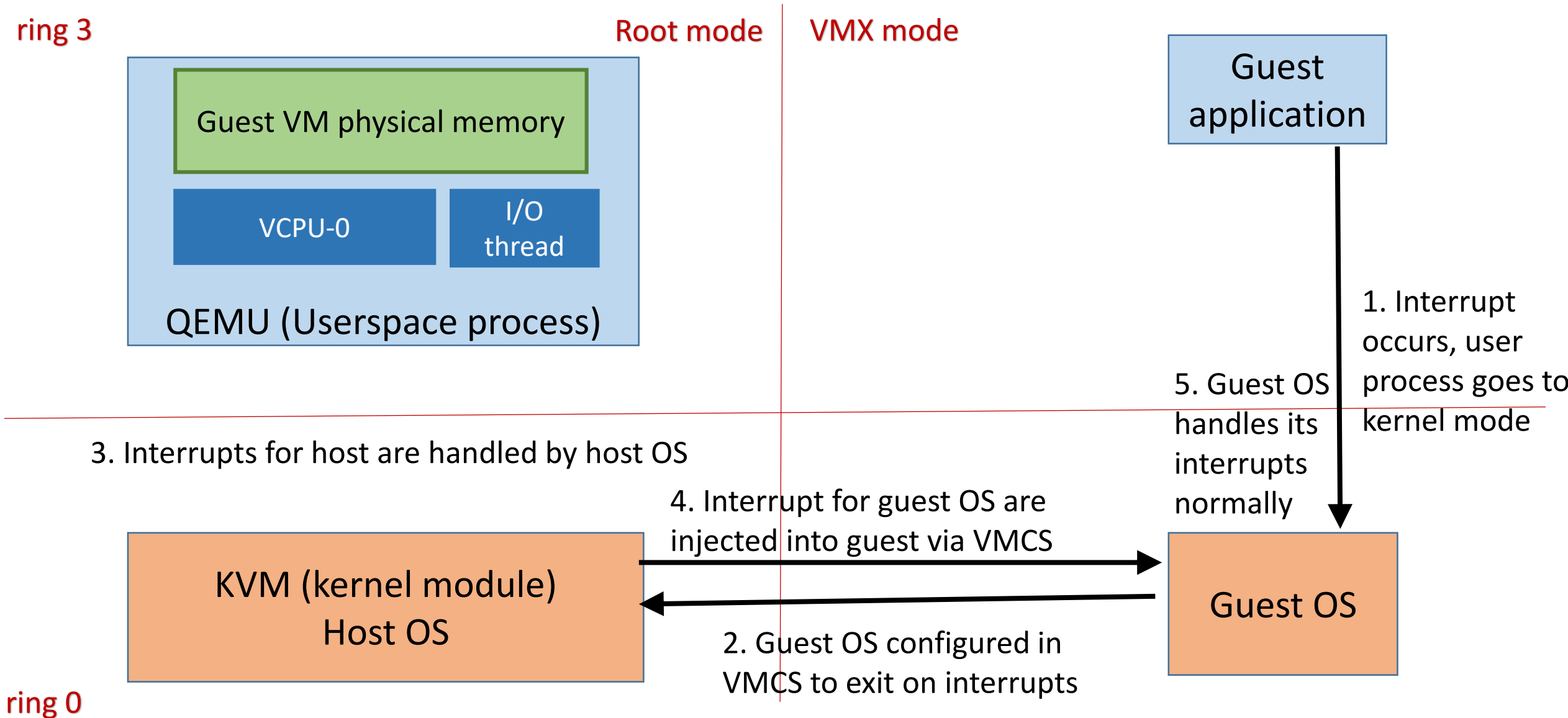- Many optimizations exist, only basics discussed here

# Communication between OS and device

- Device memory exposed as registers (command, status, data etc.)
  - I/O happens by reading/writing this memory
  - E.g., write command into device register to begin I/O
- OS can read/write device registers in two ways:
  - Explicit I/O: in/out instructions in x86 can write to device memory
  - Memory mapped I/O: Some memory addresses are assigned to device memory and are not to RAM. I/O happens by reading/writing this memory.
- Accessing device memory (via explicit I/O or memory mapped I/O) can be configured to trap to VMM
- Device raises interrupt when I/O completes (alternative to polling)
  - Modern I/O devices perform DMA (Direct Memory Access) and copy data from device memory to RAM before raising interrupt
  - Device driver provides physical address of DMA buffers to device

# QEMU/KVM I/O handling

ring 3

Root mode | VMX mode

**Guest VM physical memory**

VCPU-0 | I/O thread

Separate I/O threads to handle I/O ops
VCPU threads do not block for I/O
Guest VM resumes and can context switch to another process

QEMU (Userspace process)

**Guest application**

8. Return to user process

1. User application makes system call to guest OS

3. KVM returns to QEMU VCPU thread (exit info in kvm_run)

5. VCPU thread resumes VM

4. QEMU I/O thread initiates I/O in host OS (may block)

6. KVM resumes guest VM

**KVM (kernel module) Host OS**

**Guest OS**

2. Guest OS exits when it attempts I/O

ring 0

# QEMU/KVM interrupt handling

ring 3

Root mode | VMX mode

Guest application

Guest VM physical memory

VCPU-0 | I/O thread

QEMU (Userspace process)

1. Interrupt occurs, user process goes to kernel mode

5. Guest OS handles its interrupts normally

3. Interrupts for host are handled by host OS

4. Interrupt for guest OS are injected into guest via VMCS

KVM (kernel module) Host OS

Guest OS

2. Guest OS configured in VMCS to exit on interrupts

ring 0

# Full virtualization VMM architecture

Host OS context | VMM context

Guest VM physical memory

VMM userspace process

ring 3

Guest application

1. Guest user app makes system call to perform I/O

ring 1

Guest OS

ring 0

4. VMM kernel driver or userspace process handle I/O requests via emulation

5. Interrupts handled by host and injected into guest

2. Privileged action traps to VMM

VMM kernel driver (Host OS)

3. VMM exits to host OS to handle I/O. Some traps can handled by VMM without world switch, e.g., exit only once per batch of I/O requests

VMM (guest OS traps here)

# QEMU/KVM virtio optimization

ring 3

Root mode    VMX mode

**Guest application**

Guest VM physical memory

VCPU-0    I/O thread

QEMU (Userspace process)

Shared ring

2. Special virtio "front end" device driver places requests in shared memory

1. User application makes system call to guest OS

4. KVM and QEMU "backend" access requests from shared ring. Virtual interrupt raised in guest after batch of responses.

Memory copy avoided
Batching of requests, interrupts
Standardized across devices
High performance

KVM (kernel module) Host OS

Guest OS

3. Guest OS exits after a batch of requests accumulate

ring 0

# Device passthrough or Direct I/O

- More efficient than device emulation

- Example: SR-IOV (Single Root IO Virtualization) in network devices
  - Network card has one physical function (PF) and many virtual functions (VFs)
  - PF managed by host OS, each VF assigned to one guest VM
  - Each VF is like a separate NIC, and is bound to a guest VM
  - Packets destined to the MAC address of VM are switched to corresponding VF

# SR-IOV

- SR-IOV NIC communicates directly with device driver in guest OS
  - Packets do not go to the host OS stack at all
  - Packets switched at Layer-2 using VM virtual device's MAC address
  - Packets DMA'ed directly into guest VM memory, host OS not involved
  - But, interrupts may still cause VM exit (interrupt can be for host too)
- Challenge: when guest device driver provides DMA buffers to VF, it can only provide guest physical addresses (GPA) of the buffer
  - NIC cannot access the DMA buffer memory using GPA alone
- SR-IOV capable NICs have an inbuilt MMU (IOMMU) to translate from GPA to HPA

# Summary

- Techniques for I/O virtualization
  - Device emulation
  - Virtio optimization
  - Device passthrough or direct I/O (SR-IOV)