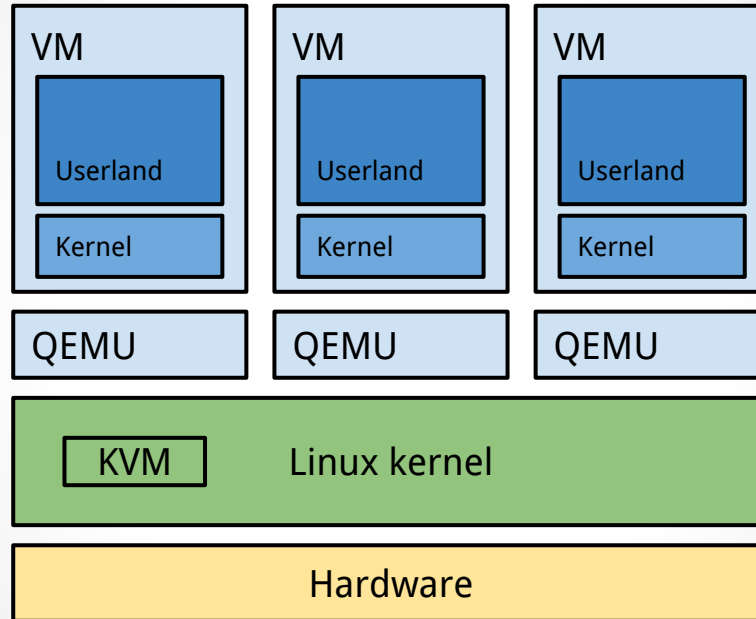


x86: Device Virtualization

Gabriel Laskar <gabriel@lse.epita.fr>

Hardware Virtualization



HW Virtualization : QEMU/KVM

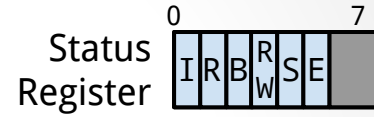
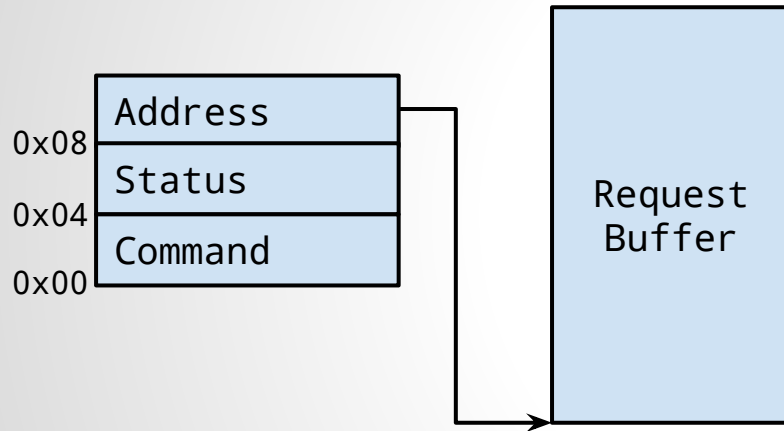
I/O Virtualization

- HW emulation
- Paravirtualized Devices
 - Xen
 - Virtio
 - Other (vmxnet, synthetic devices, ...)
- Hardware Pass through
 - Full Device (pci, vga)
 - Protocol (usb, serial)
 - Other way

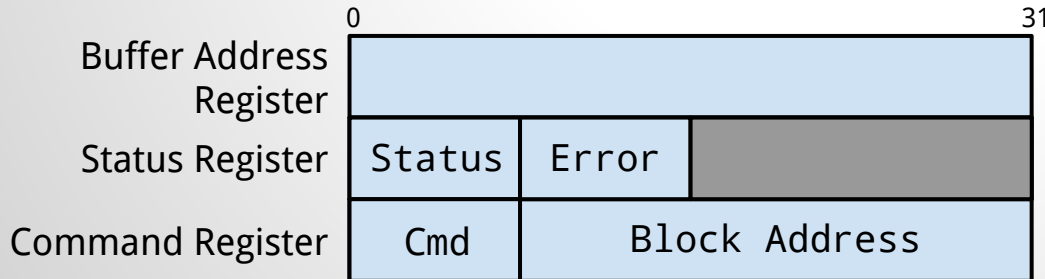
Devices

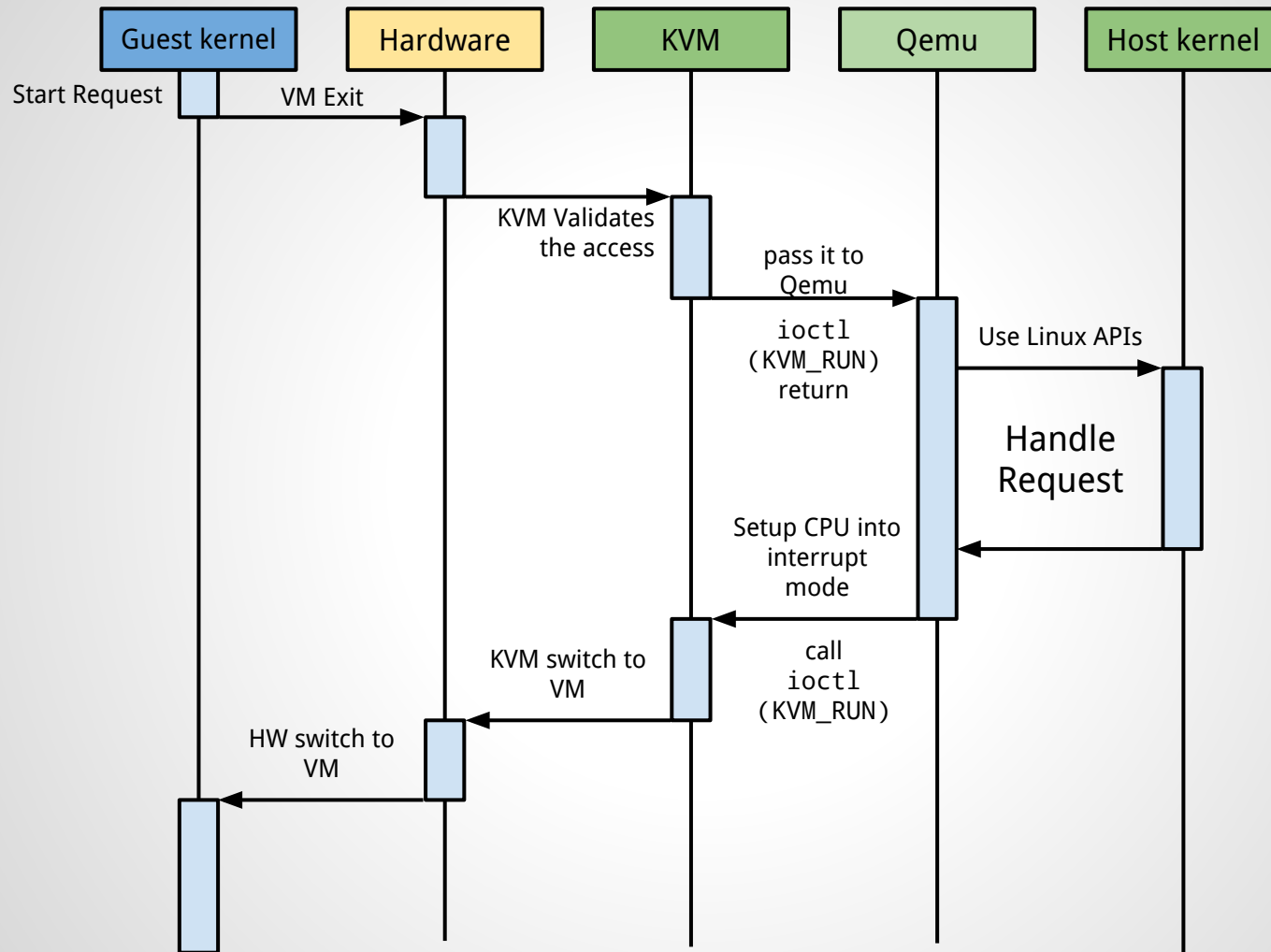
- Registers accessible to CPU :
 - MMIO
 - PIO (in, out)
- Access to Memory (DMA)
- Interrupts (irq, msi)

Example of a simple device: a block device



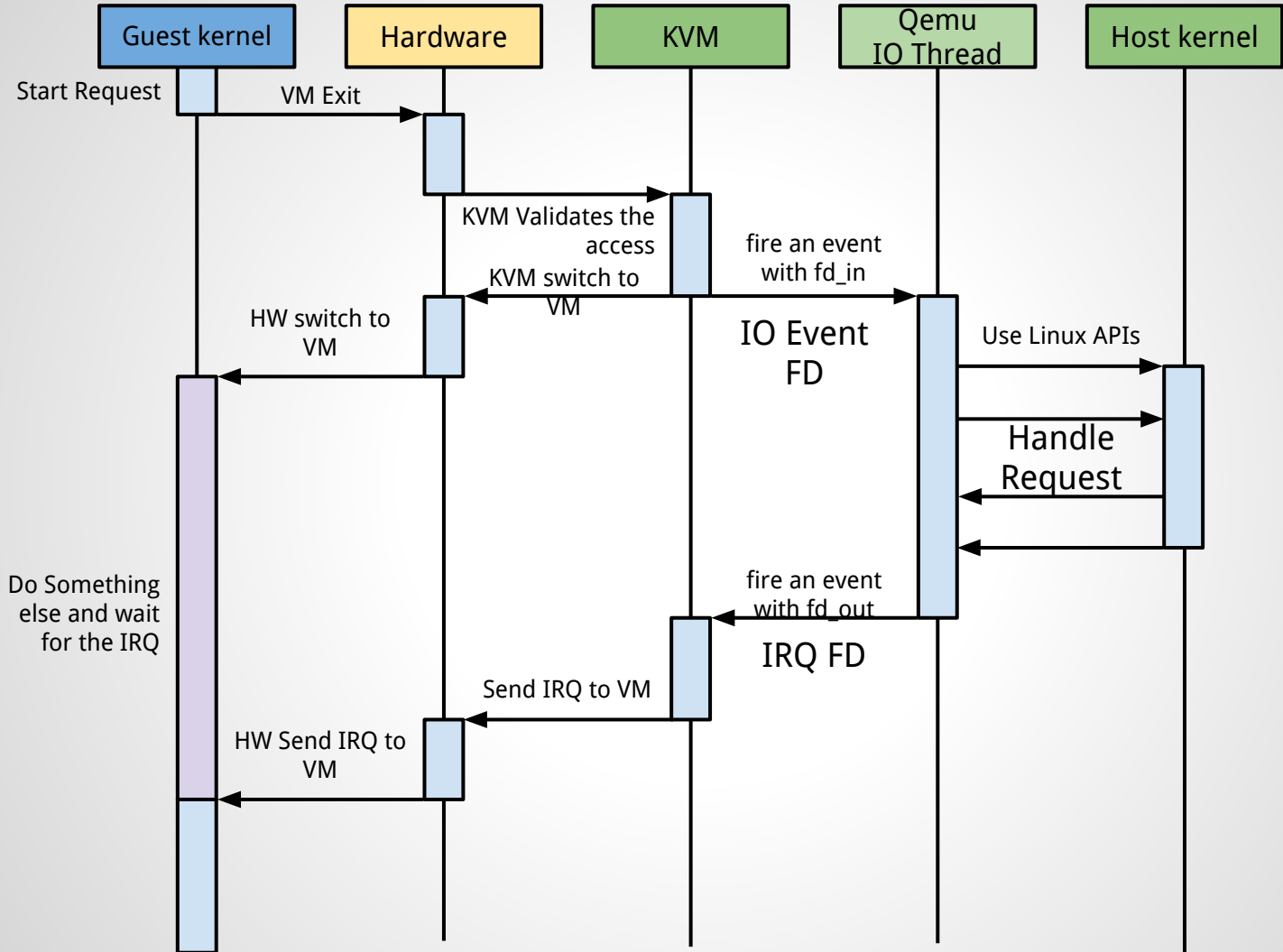
- I : Initialize device
- R : Ready to receive command
- B : Busy
- RW : Read/Write (clear if read, set if write)
- S : Start Command
- E : Error





MMIO, PIO : How fast ?

- For each mmio access, there is an exit
- We have to assert the read/write, and process the command
- Can't be asynchronous, ie : we can't do that with the vcpu guest running
- What solutions do we have ?



How can we solve the IO Problem ?

- **KVM_IOEVENTFD**
 - Attach an ioeventfd to a pio/mmio guest address
 - When guest write into this address, it fire an event instead of an exit
- **KVM_IRQFD**
 - Allow setting an eventfd that will trigger a guest interrupt
- With eventfd and irqfd, we can offload io traffic into another thread, and just listen/fire event through fds.

Example : handling device

```
void handle_device(void *device, int eventfd, int irqfd)
{
    struct pollfd input_queue = {
        .fd = eventfd,
        .events = POLLIN;
    };

    for (;;) {
        int ret = poll(input_queue, 1, timeout);

        if (ret > 0) {
            uint64_t event_value;
            read(eventfd, &event_value, sizeof(event_value));
            uint64_t res = do_something(device, event_value);
            write(irqfd, &res, sizeof(res));
        }
    }
}
```

Virtio

- Abstraction layer for virtualized devices
 - multiple device types: PCI, MMIO
 - Split in 2 parts: Configuration Space, Queues
- Multiple Devices supported:
 - block
 - network
 - channels (for serial line, 9p, or custom data channel)
 - scsi devices
 - RNG
 - Ballooning
 - GPU and Input devices (not completely merged yet)

Direct Pass-Through. VFIO, VT-D and IOMMU

- In order to pass-through a device, we must have some kind of IOMMU support
- VFIO is the linux API allowing to use it
- VFIO allows to use PCI (or other) devices in userland

Notes for the project

- Create a VM that can boot stos without stos . grub
 - Boot directly into PM mode
 - Load the kernel/stos:
 - kernel/stos
 - initramfs
 - commandline
 - modules
 - memory map
 - pass the boot param structure
 - Devices:
 - 8250
 - 8259a
 - 8254
- All your code should be working in your VM