

I/O in Linux Hypervisors and Virtual Machines Lecture for the Embedded Systems Course CSD, University of Crete (May 7 & 9, 2025)

Manolis Marazakis (maraz@ics.forth.gr)

FORTH-ICS Institute of Computer Science (ICS)

Foundation for Research and Technology – Hellas (FORTH)

Outline

- Elements of I/O virtualization
 - Machine emulator, Hypervisor, Transport
- Alternative designs for the device I/O path
 - Device emulation (fully virtualized)
 - Para-virtualized devices
 - Direct device assignment (pass-through access)
- Follow the I/O path of hypervisor technologies most commonly used on Linux servers (x86 platform)
 - xen, kvm
- Provide a glimpse of hypervisor internals *
 - xen: device channels, grant tables
 - kvm: virtio
 - * VMware: market leader ... but out-of-scope for this lecture



I/O device types

- Dedicated
 - E.g. Display
- Partitionable
 - E.g. Disk
- Shared
 - E.g. NIC
- Devices that can be enumerated (PCI, PCI-Express)
 - VMM needs to emulate 'discovery' over a system bus/interconnect
- Devices with hard-wired addresses (e.g. PS/2)
 - VMM should maintain status information on virtual device ports
- Emulated (e.g. experimental hardware)
 - VMM must define & emulate <u>all</u> H/W functionality
 - GuestOS needs to load corresponding device drivers



Why is I/O hard to virtualize ?

- Multiplexing/de-multiplexing for guests
 - Programmed I/O (PIO): <u>privileged</u> CPU instructions specifically for I/O
 - I/O devices have a separate address space from general memory
 - Memory-mapped I/O (MMIO): CPU instructions for memory access are also used for accessing devices
 - ► The memory-mapped I/O region is protected
- Direct Memory Access (DMA)
 - Allow hardware subsystems within the computer to access system memory for reading and/or writing <u>independently of the CPU</u>
 - Synchronous: triggered by software
 - Asynchronous: triggered by devices (e.g. NIC)
- Implementation layers:
 - ▶ system calls (application to GuestOS) \rightarrow trap to VMM
 - driver calls (GuestOS) \rightarrow paravirtualization
 - Hypercall between modified driver in GuestOS and VMM
 - I/O operations (GuestOS driver to VMM)



Processing an I/O Request from a VM



I/O request issued by an application is processed first by the guest operating system I/O stack running within the VM, and then by the hypervisor I/O stack managing physical hardware.



Source: Mendel Rosenblum, Carl Waldspurger: "I/O Virtualization" ACM Queue, Volume 9, issue 11, November 22, 2011



Device Front-Ends & Back-Ends



Source: Mendel Rosenblum, Carl Waldspurger: "I/O Virtualization" ACM Queue, Volume 9, issue 11, November 22, 2011



I/O stack: Host side



Storage Controller





I/O stack: Host + Guest





Device Emulation



- Hypervisor emulates devices → traps I/O accesses (PIO, MMIO)
- Emulation of DMA and interrupts



Para-virtualized drivers



- Hypervisor-specific virtual device drivers (front-end) in Guest OS
- Involvement of Hypervisor (for Back-end)



Direct device assignment



- <u>Bypass</u> Hypervisor to directly access
 I/O devices
- Security & safety concerns
 - IO-MMU for address translation
 & isolation (DMA restrictions)
 - SR-IOV for shared access



xen history

- Developed at Systems Research Group, Cambridge University (UK)
- Creators: Keir Fraser, Steven Hand, Ian Pratt, et al (2003)
- Broad scope, for both Host and Guest
- Merged in Linux mainline: 2.4.22
- Company: Xensource.com
 - Acquired by Citrix Systems (2007)

Xen VMM: Paravirtualization (PV) 1/2



source: https://wiki.xenproject.org/wiki/Paravirtualization_(PV)

I/O in Linux Hypervisors and Virtual Machines



Xen VMM: Paravirtualization (PV) 2/2





kvm history

- Prime creator: Avi Kivity, Qumranet, circa. 2005 (IL)
 - Company acquired by RedHat (2008)
- "Narrow" focus: x86 platform, Linux host
 - Assumes Intel VT-x or AMD svm
- Merged in Linux kernel mainline: 2.6.20
 - ... < 4 months after 1st announcement !



kvm VMM: tightly integrated in the Linux kernel



- Hypervisor:
 Kernel module
- Guest OS: User-space
 process
 (QEMU)
- Requires H/W virtualization extensions

xen vs kvm

Xen

- Strong support for paravirtualization with modified host-OS
 - → Near-native performance for I/O devices
- Separate code based for DOM0 and device drivers
- Security model: Rely on DOM0
- Maintainability: Hard to keep up with all versions of possible guests due to PV

KVM

- Requires H/W virtualization extension – Intel VT, AMD Pacifica (AMD-V)
- Limited support for paravirtualization (via virtio)
- Code-base integrated into Linux kernel source tree
- Security model: Rely on Commodity/Casual Linux systems
- Maintainability: Easy – Integrated well into infrastructure, code-base

I/O virtualization in xen



- <u>Bridge</u> in driver domain: multiplex/de-multiplex network I/Os from guests
- I/O Channel
 - Zero-copy transfer with Grant-copy
 - Enable driver domain to access I/O buffers in guest memory

- Xen network I/O extension schemes
 - Multiple RX queues, SR-IOV ...

Source:"Bridging the gap between software and hardware techniques for i/o virtualization" Usenix Annual Technical conference, 2008



Xen device channels

- Asynchronous shared-memory transport
- Event ring (for interrupts)
- Xen "peer domains"
 - Inter-guest communication
 - Mapping one guest's buffers to another
 - Grant tables for "DMA" (bulk transfers)
- Xen dom0 (privileged domain) can access all devices
 - Exports subset to other domains
 - Runs back-end of device drivers (e.g. net, block)



Xen grant tables

- Share & Transfer pages between domains
 - a software implementation of certain IOMMU functionality

Transferred pages:

- ▶ Driver in local domain "advertises" buffer \rightarrow notify hypervisor
- Driver then transfers page to remote domain _and_ takes a free page from a producer/consumer ring ("page-flip")
- ► Use case: network drivers → receive their data asynchronously, i.e. may not know origin domain (need to inspect network packet before actual transfer between domains)
- With RDMA NICs, we can transfer (DMA) directly into domains ...

Shared pages:

- ▶ Driver in local domain "advertises" buffer → notify hypervisor that this page can be access by other domains
- ► Use case: block drivers → receive their data synchronously, i.e. know which domain requested data to be transferred via DMA



kvm run-time environment







Run-time view of a kvm guest



- Guest Host switch via scheduler
- Use of Linux
 subsystems:
 scheduler, memory
 management, ...
- Re-use of user-space tools
 - VM images
 - Network configuration



- Processes create virtual machines
 - A process together with /dev/kvm is in essence the Hypervisor
- VMs contain memory, virtual CPUs, and (in-kernel) devices
- Guest ("physical") memory is part of the (virtual) address space of the creating process
 - Virtual MMU+TLB, and APIC/IO-APIC (in-kernel)
 - Machine instruction interpreter (in-kernel)
- vCPUs run in process context (i.e. as threads)
 - To the host, the process that started the guest _is_ the guest!

kvm API demonstration:

- https://lwn.net/Articles/658511/
- <u>https://github.com/soulxu/kvmsample</u>



kvm API sample

```
struct kvm *kvm = malloc(sizeof(struct kvm));
#include <linux/kvm.h>
                                                 kvm->dev fd = open(KVM DEVICE, O RDWR);
struct kvm {
                                                 kvm->vm fd = ioctl(kvm->dev fd, KVM CREATE VM, 0);
 int dev fd;
                                                 kvm->ram size = ram size;
 int vm fd;
                                                 kvm->ram start = ( u64)mmap(NULL, kvm->ram size,
  u64 ram size;
                                                 PROT READ | PROT WRITE, MAP PRIVATE | MAP ANONYMOUS | MAP NORESERVE, -
  u64 ram start;
                                                 1, 0);
 int kvm version;
 struct kvm userspace memory region mem;
                                                 kvm->mem.slot = 0;
 struct vcpu *vcpus;
                                                 kvm->mem.guest phys addr = 0;
 int vcpu number;
                                                 kvm->mem.memory size = kvm->ram size;
};
                                                 kvm->mem.userspace addr = kvm->ram start;
struct vcpu {
                                                 ret = ioctl(kvm->vm fd, KVM SET USER MEMORY REGION, &(kvm->mem));
  int vcpu id;
  int vcpu fd;
                                                 struct vcpu *vcpu = malloc(sizeof(struct vcpu));
  pthread tvcpu thread;
                                                 vcpu->vcpu id = 0;
  struct kvm run *kvm run;
                                                 vcpu->vcpu fd = ioctl(kvm->vm fd, KVM CREATE VCPU, vcpu->vcpu id);
  int kvm run mmap size;
                                                 vcpu->kvm run mmap size = ioctl(kvm->dev fd, KVM GET VCPU MMAP SIZE, 0);
  struct kvm regs regs;
                                                 vcpu->kvm run = mmap(NULL, vcpu->kvm run mmap size, PROT READ |
  struct kvm sregs sregs;
                                                 PROT WRITE, MAP SHARED, vcpu->vcpu fd, 0);
  void *(*vcpu thread func)(void *);
                                                 pthread create(&(kvm->vcpus->vcpu thread), (const pthread attr t *)NULL, kvm-
                                                 >vcpus[i].vcpu thread func, kvm);
};
                                                 pthread join(kvm->vcpus->vcpu thread, NULL);
```



I/O virtualization in kvm



- Native KVM I/O model
- PIO: Trap
- MMIO: The machine emulator executes the faulting instruction
- Slow due to mode-switch!
- Extensions to support PV
 - VirtIO: An API for Virtual I/O aims to support many hypervisors (of all types)



virtio

- A family of drivers which can be adapted for various hypervisors, by porting a shim layer
- Related: VMware tools, Xen para-virtualized drivers
- Explicit separation of Drivers, Transport, Configuration





virtio architecture



Front-end

A kernel module in guest OS. Accepts I/O requests from user process. Transfer I/O requests to back-end.

Back-end

A device in QEMU. Accepts I/O requests from front-end. Perform I/O operation via physical device.

- Virtqueues (per device)
- Vring (per virtqueue)
- Queue requests

vring & virtqueue

- vring: transport implementation (ring-buffer)
 - shared (memory-mapped) between Guest and QEMU
 - Reduce the number of MMIOs
 - published & used buffers
 - descriptors
- virtqueue API:
 - add_buf: expose buffer to other end
 - ø get_buf: get next used buffer
 - kick: (after add_buf) notify QEMU to handle buffer
 - disable_cb, enable_cb: disable/enable callbacks
- "buffer" := scatter/gather list \rightarrow (address, length) pairs
- QEMU: virtqueue_pop, virtqueue_push
- virtio-blk: 1 queue
- virtio-net: 2 queues



kvm with virtio





virtio processing flow







Virtual interrupts





Sources

- Mendel Rosenblum, Carl Waldspurger: "I/O Virtualization" ACM Queue, Volume 9, issue 11, November 22, 2011 URL: <u>http://queue.acm.org/detail.cfm?id=2071256</u>
- Avi Kivity, et al: kvm: The Linux Virtual Machine Monitor, Proceedings of the Linux Symposium, 2007
 - http://www.linux-kvm.com/sites/default/files/kivity-Reprint.pdf
 - http://kerneltrap.org/node/8088
- Paul Barham, Boris Dragovic, Keir Fraser, Steven Hand, Tim Harris, Alex Ho, Rolf Neugebauer, Ian Pratt, Andrew Warfield, Xen and the Art of Virtualization, SOSP'03
- Xen (v.3.0. for x86) Interface Manual
 - http://pdub.net/proj/usenix08boston/xen_drive/resources/developer_ manuals/interface.pdf
- Jun Nakajima, Asit Mallick, Ian Pratt, Keir Fraser, X86-64 Xen-Linux: Architecture, Implementation, and Optimizations, OLS 2006
- > Xen: finishing the job, lwn.net 2009



NVIDIA Virtual GPU (vGPU) + VMware Horizon



source: https://www.nvidia.com/en-eu/data-center/virtual-gpu-technology/



VMware Virtual Desktop Infrastructure (VDI)









QEMU machine emulator

- Creator: Fabrice Bellard (circa. 2006)
- Machine emulator using a dynamic binary translator
 - Run-time conversion of target CPU instructions to host ISA
 - Translation cache
- Emulated Machine := { CPU emulator, Emulated Devices, Generic Devices } + "machine description"
 - Link between emulated devices & underlying host devices
 - ▶ Alternative storage back-ends e.g. POSIX/AIO (\rightarrow thread pool)
 - Caching modes for devices:
 - ▶ cache=none \rightarrow O_DIRECT
 - cache=writeback \rightarrow buffered I/O
 - ► Cache=writethrough → buffered I/O, O_SYNC

Emulated Platform







Virtualization components

Machine emulation

- CPU, Memory, I/O
- Hardware-assisted

Hypervisor

- Hyper-call interface
- Page Mapper
- ► I/O
- Interrupts
- Scheduler
- Transport
 - Messaging, and bulk-mode

Virtual disks

- Exported by host
- Physical device/partition
- ... or logical device
- ... or file (image) (e.g. : .qcow2, .vmdk, "raw" .img)
 - file format that describes containers for virtual hard disk drives
 - features: compression, encryption, copy-on-write snapshots



Device model (with full-system virtualization)

- VMM intercepts I/O operations from GuestOS and passes them to device model at the host
- Device model emulates I/O operation interfaces:
 - PIO, MMIO, DMA, ...
- Two different implementations:
 - Part of the VMM
 - User-space standalone service



Virtualized I/O flow





Device emulation: full virtualization vs para-virtualization

GuestOS		
VMM (full virtualization)	Traps	
	Device Emulation	
Hardware		

GuestOS	PV drivers	
VMM (full virtualization)	Traps	
	•	
	Device Emulation	
Hardware		

kvm guest initialization (command-line)

- > qemu-system-x86_64 -L /usr/local/kvm/share/qemu
 - -hda /mnt/scalusMar2013/01/scalusvm.img

-drive

file=/mnt/scalusMar2013/01/datavol.img,if=virtio,index=0
,cache=writethrough

```
-net nic -net user, hostfwd=tcp::12301-:22
```

-nographic

-m 4096 -smp 2



Device I/O in kvm

- Configuration via MMIO/PIO
- eventfd for events between host/guest
 - irqfd : host \rightarrow guest
 - ▶ ioeventfd : guest \rightarrow host
- virtio: abstraction for virtualized devices
 - Device types: PCI, MMIO
 - Configuration
 - Queues



IO-MMU Emulation

- Shortcomings of device assignment for <u>unmodified</u> guests:
 - Requires pinning all of the guest's pages, thereby disallowing memory over-commitment
 - Exposes the guest's memory to buggy device drivers
- A single physical IO-MMU can emulate multiple IO-MMU's (for multiple guests)
- Why?
 - Allow memory over-commitment (pin/unpin during map/unmap of I/O buffers)
 - Intra-guest protection, redirection of DMA transactions
 - ... without compromising inter-guest protection





I/O in Linux Hypervisors and Virtual Machines

