

本站文章大部分为作者原创，非商业用途转载无需作者授权，但务必在文章标题下面注明作者 刘世民（Sammy Liu）以及可点击的本博客地址超级链接 <http://www.cnblogs.com/sammyliu/>，谢谢合作



世民谈云计算

（声明：本站文章皆基于公开来源信息，仅代表作者个人观点，与作者所在公司无关）

昵称：SammyLiu
园龄：2年6个月
荣誉：推荐博客
粉丝：470
关注：30
+加关注

< 2017年5月 >						
日	一	二	三	四	五	六
30	1	2	3	4	5	6
7	8	9	10	11	12	13
14	15	16	17	18	19	20
21	22	23	24	25	26	27
28	29	30	31	1	2	3
4	5	6	7	8	9	10

常用链接

我的随笔
我的评论
我的参与
最新评论
我的标签

我的标签

GRE(1)
Neutron(1)
Open vSwitch(1)
OpenStack(1)

随笔分类(254)

Ceilometer(3)
Ceph(13)
Cinder(6)
Docker(8)
Glance
Heat(2)
K8S
Keystone(1)
KVM(10)
MessageQueue(4)
MySQL(1)
Neutron(17)
Nova(10)
OpenStack(33)
Sahara
Storage(1)
Swift(3)
Trove
Ubuntu(3)
VMware(3)
安装和配置(1)
版本(4)
备份(1)
大数据(5)
翻译(4)
高可用（HA）(6)
基础知识(19)
监控(1)
容器(4)
容器编排
使用案例(4)
网络(8)
问题定位(3)
行业(14)

博客园 首页 新随笔 订阅 XML 管理

随笔-121 评论-504 文章-45

KVM 介绍（4）：I/O 设备直接分配和 SR-IOV [KVM PCI/PCIe Pass-Through SR-IOV]

学习 KVM 的系列文章：

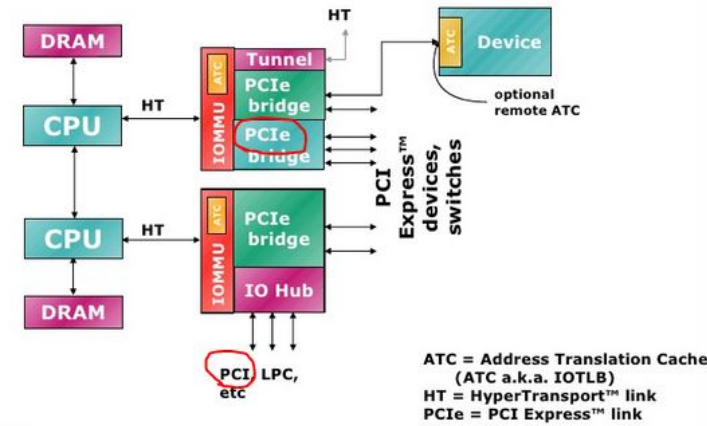
- (1) 介绍和安装
- (2) CPU 和 内存虚拟化
- (3) I/O QEMU 全虚拟化和准虚拟化（Para-virtualization）
- (4) I/O PCI/PCIe设备直接分配和 SR-IOV
- (5) libvirt 介绍
- (6) Nova 通过 libvirt 管理 QEMU/KVM 虚拟机
- (7) 快照（snapshot）
- (8) 迁移（migration）

本文将分析 PCI/PCIe 设备直接分配（Pass-through）和 SR-IOV，以及三种 I/O 虚拟化方式的比较。

1. PCI/PCI-E 设备直接分配给虚拟机（PCI Pass-through）

设备直接分配（Device assignment）也称为 Device Pass-Through。

先简单看看PCI 和 PCI-E 的区别（AMD CPU）：



（简单点看，PCI 卡的性能没有 PCI-E 高，因为 PCI-E 是直接连在 IOMMU 上，而 PCI 卡是连在一个 IO Hub 上。）

主要的 PCI 设备类型：

- Network cards (wired or wireless)
- SCSI adapters
- Bus controllers: USB, PCMCIA, I2C, FireWire, IDE
- Graphics and video cards
- Sound cards

1.1 PCI/PCIe Pass-through 原理

这种方式，允许将主机中的物理 PCI 设备直接分配给客户机使用。较新的x86平台已经支持这种类型，Intel 定义的 I/O 虚拟化技术成为 VT-d，AMD 的称为 AMD-V。KVM 支持客户机以独占方式访问这个宿主机的 PCI/PCI-E 设备。通过硬件支持的 VT-d 技术将设备分给客户机后，在客户机看来，设备是物理上连接在PCI或者PCI-E总线上的，客户机对该设备的I/O交互操作和实际的物理设备操作完全一样，不需要或者很少需要 KVM 的参与。运行在 VT-d 平台上的 QEMU/KVM，可以分配网卡、磁盘控制器、USB控制器、VGA 显卡等设备供客户机直接使用。

几乎所有的 PCI 和 PCI-E 设备都支持直接分配，除了显卡以外（显卡的特殊性在这里）。PCI Pass-through 需要硬件平台 Intel VT-d 或者 AMD IOMMU 的支持。这些特性必须在 BIOS 中被启用。Red Hat Enterprise Linux 6.0 及以上版本支持热插拔的 PCI 设备直接分配到虚拟机。

网卡直接分配：

性能(4)
虚拟化(7)
原理(22)
云Cloud(29)

随笔档案(121)

2017年5月 (1)
2017年3月 (1)
2017年1月 (1)
2016年10月 (7)
2016年9月 (5)
2016年8月 (4)
2016年7月 (1)
2016年6月 (5)
2016年5月 (1)
2016年4月 (1)
2016年3月 (9)
2016年2月 (4)
2016年1月 (2)
2015年12月 (7)
2015年11月 (7)
2015年10月 (4)
2015年9月 (4)
2015年8月 (5)
2015年7月 (9)
2015年6月 (10)
2015年5月 (3)
2015年4月 (11)
2015年3月 (2)
2015年2月 (6)
2015年1月 (5)
2014年12月 (6)

文章分类(21)

Ceph(1)
GlusterFS
Web 服务器(2)
操作系统(1)
存储
大数据(2)
分布式系统
服务器(1)
网络(11)
虚拟化(3)
云

文章档案(42)

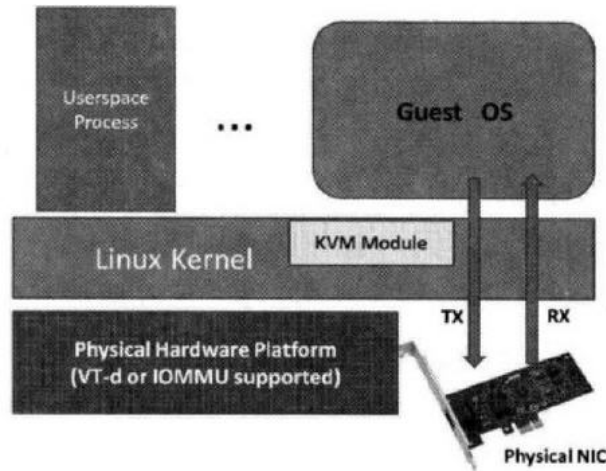
2016年10月 (2)
2016年9月 (1)
2016年6月 (1)
2016年5月 (3)
2015年12月 (4)
2015年10月 (5)
2015年9月 (2)
2015年6月 (1)
2015年4月 (23)

积分与排名

积分 - 286831
排名 - 535

最新评论

1. Re: Neutron 理解 (1): Neutron 所实现的虚拟化网络 [How Neutron Virtualizes Network]
eth1 - 公共网络 (untagged), 管理网络 (tag=102), 存储网络 (tag=103) 不好意思, 大家共用同一个eth1端口的话, 请问这里交



硬盘直接分配:

- 一般 SATA 或者 SAS 等类型的硬盘的控制器都是直接接入到 PCI 或者 PCI-E 总线的, 所以也可以将硬盘作为普通的PCI设备直接分配个客户机。需要注意的是, 当分配硬盘时, 实际上将其控制器作为一个整体分配到客户机中, 因此需要在硬件平台上至少有另两个或者多个SATA或者 SAS控制器。

1.2 在 RedHat Linux 6 上使用 virt-manger 分配一个光纤卡给虚拟机

准备工作:

- (1) 在 BIOS 中打开 Intel VT-d
- (2) 在 Linux 内核中启用 PCI Pass-through

添加 intel_iommu=on 到 /boot/grub/grub.conf 文件中。(在我的 RedHat Linux 6上, 该文件是 /boot/grub.conf)

- (3) 重启系统, 使得配置生效

实际分配:

- (1) 使用 lspci -nn 命令找到待分配的 PCI 设备。这里以一个 FC 卡为例:

```
[root@rh65 ~]# lspci -nn | grep Fibre
1f:00.0 Fibre Channel [0c04]: QLogic Corp. ISP2532-based 8Gb Fibre Channel to PC
I Express HBA [1077:2532] (rev 02)
1f:00.1 Fibre Channel [0c04]: QLogic Corp. ISP2532-based 8Gb Fibre Channel to PC
I Express HBA [1077:2532] (rev 02)
```

使用 Lspci 命令得到的 PCI 数字的含义, 以后使用 libvirt API 分配设备时会用到:

```
02:01.0 CardBus bridge: Ricoh Co Ltd RL5c476 II (rev b4)
02:01.1 FireWire (IEEE 1394): Ricoh Co Ltd R5C552 IEEE 1394 Controller
02:01.2 SD Host controller: Ricoh Co Ltd R5C822 SD/SDIO/MMC/MS/MSPro
```

Function number
PCI device number
PCI bus number

- (2) 使用 virsh nodev-list 命令找到该设备的 PCI 编号

```
[root@rh65 ~]# virsh nodev-list --tree | grep 1f
| +- pci_0000_1f_00_0
|   +- pci_0000_1f_00_1
|     +- pci_0000_00_1f_0
|       +- pci_0000_00_1f_2
|         +- pci_0000_00_1f_3
|           +- pci_0000_00_1f_5
```

- (3) 将设备从主机上解除

```
[root@rh65 ~]# virsh nodev-detach pci_0000_1f_00_0
Device pci_0000_1f_00_0 detached
```

- (4) 使用 virt-manager 将设备直接分配给一个启动了的虚拟机

交换机端口是配置为tagged还是untagged.....

--xianke9

2. Re:理解Docker (5) : Docker 网络

1.12版本上网络的表现如何?

--幽灵狼

3. Re:理解Docker (5) : Docker 网络

我想请问一下运行docker quickstart terminal时一直卡在"waiting for an IP"应该如何解决呢?希望楼主能解答一下。

--silentbell

4. Re:理解Docker (6) : 若干企业生产环境中的容器网络方案写得好好! 加油。

--itbj00

5. Re:理解Docker (5) : Docker 网络

非常好,写得很详细。加油!

--itbj00

阅读排行榜

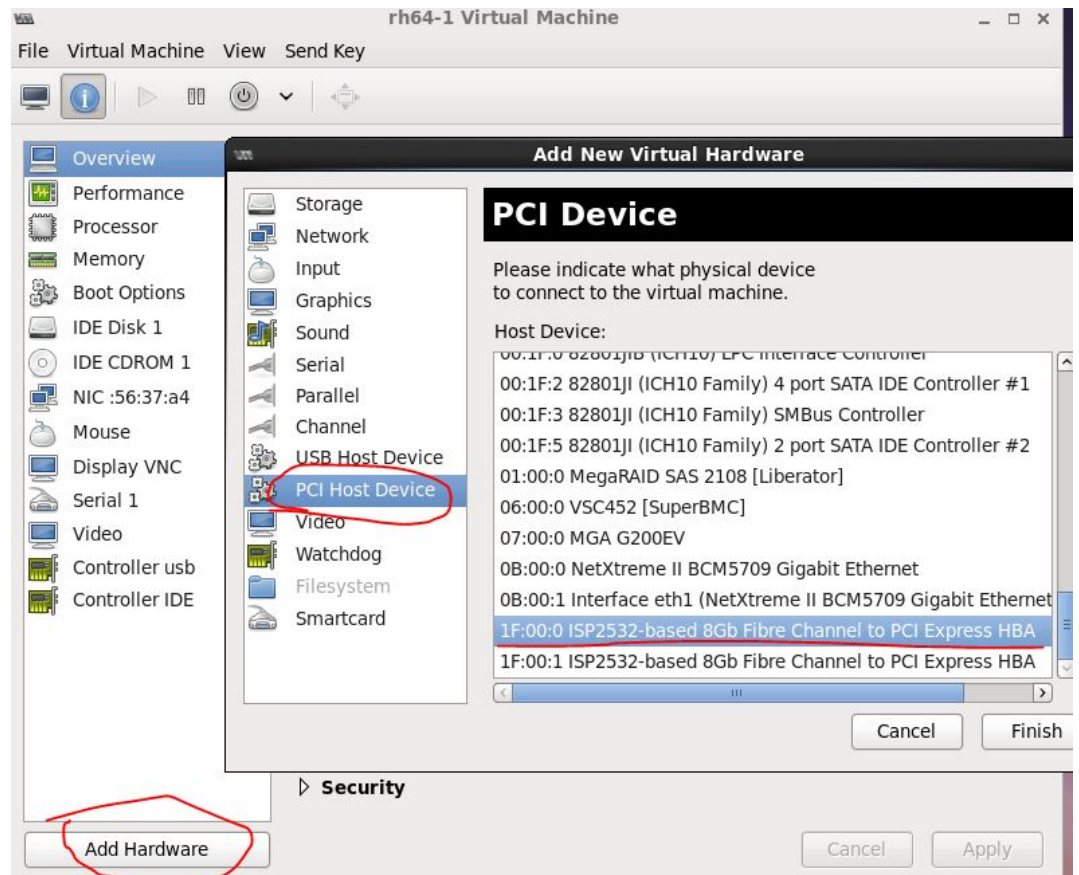
1. Neutron 理解 (1): Neutron 所实现的虚拟化网络 [How Neutron Virtualizes Network](22087)
2. 理解 OpenStack 高可用 (HA) (1) : OpenStack 高可用和灾备方案 [OpenStack HA and DR](13707)
3. Neutron 理解 (3): Open vSwitch + GRE/VxLAN 组网 [Neutron Open vSwitch + GRE/VxLAN Virtual Network](13434)
4. 探索 OpenStack 之 (9) : 深入块存储服务Cinder (功能篇) (12921)
5. 理解 OpenStack + Ceph (1) : Ceph + OpenStack 集群部署和配置 (12444)

评论排行榜

1. Neutron 理解 (1): Neutron 所实现的虚拟化网络 [How Neutron Virtualizes Network](63)
2. Neutron 理解 (14) : Neutron ML2 + Linux bridge + VxLAN 组网 (54)
3. Neutron 理解 (8): Neutron 是如何实现虚拟机防火墙的 [How Neutron Implements Security Group](34)
4. Neutron 理解 (3): Open vSwitch + GRE/VxLAN 组网 [Neutron Open vSwitch + GRE/VxLAN Virtual Network](25)
5. Neutron 理解 (5) : Neutron 是如何向 Nova 虚拟机分配固定IP地址的 (How Neutron Allocates Fixed IPs to Nova Instance) (21)

推荐排行榜

1. Neutron 理解 (1): Neutron 所实现的虚拟化网络 [How Neutron Virtualizes Network](9)
2. 我所了解的 京东、携程、eBay、小米的 OpenStack 云(6)
3. 理解 OpenStack 高可用 (HA) (1) : OpenStack 高可用和灾备方案 [OpenStack HA and DR](6)
4. Neutron 理解 (2): 使用 Open vSwitch + VLAN 组网 [Neutron Open vSwitch + VLAN Virtual Network](6)
5. 理解 OpenStack 高可用 (HA) (2) : Neutron L3 Agent HA 之虚拟路由冗余协议 (VRRP) (5)



(5) 添加好了后的效果



(6) 在虚拟机中查看该PCI设备

```
[root@rh64-1 ~]# lspci
00:00.0 Host bridge: Intel Corporation 440FX - 82441FX PMC [Natoma] (rev 02)
00:01.0 ISA bridge: Intel Corporation 82371SB PIIX3 ISA [Natoma/Triton II]
00:01.1 IDE interface: Intel Corporation 82371SB PIIX3 IDE [Natoma/Triton II]
00:01.2 USB controller: Intel Corporation 82371SB PIIX3 USB [Natoma/Triton II] (rev 01)
00:01.3 Bridge: Intel Corporation 82371AB/EB/MB PIIX4 ACPI (rev 03)
00:02.0 VGA compatible controller: Cirrus Logic GD 5446
00:03.0 Ethernet controller: Realtek Semiconductor Co., Ltd. RTL-8139/8139C/8139C+ (rev 20)
00:04.0 RAM memory: Red Hat, Inc Virtio memory balloon
00:05.0 Fibre Channel: QLogic Corp. ISP2532-based 8Gb Fibre Channel to PCI Express HBA (rev 02)
```

(7) 不再使用的话,需要在 virt-manager 中首先将该设备移除,然后在主机上重新挂载该设备

```
[root@rh65 s1]# virsh nodedev-reattach pci_0000_0b_00_0
Device pci_0000_0b_00_0 re-attached
```

1.3 在 RedHat Linux 6 上使用 qemu-kvm 分配一个光纤卡给虚拟机

除了步骤 (4), 其他步骤同上面。


```
[root@rh65 domains]# kvm -smp 2 -m 2048 -drive file=../domains/rh64-9.qcow2,if=i
de,media=disk,format=qcow2 -boot c -name rh64-7 -device pci-assign,host=1f:00.0
VNC server running on ':1:5900'
```

```
00:03.0 Ethernet controller: Realtek Semiconductor Co., Ltd. RTL-81
C+ (rev 20)
00:04.0 Fibre Channel: QLogic Corp. ISP2532-based 8Gb Fibre Channel
ss HBA (rev 02)
[root@rh64-2 ~]# _
```

1.4 设备直接分配让客户机的优势和不足

- 好处：在执行 I/O 操作时大量减少甚至避免 VM-Exit 陷入到 Hypervisor 中，极大地提高了性能，可以达到几乎和原生系统一样的性能。VT-d 克服了 virtio 兼容性不好和 CPU 使用频率较高的问题。
- 不足：（1）一台服务器主板上的空间比较有限，因此允许添加的 PCI 和 PCI-E 设备是有限的。大量使用 VT-d 独立分配设备给客户机，让硬件设备数量增加，这会增加硬件投资成本。（2）对于使用 VT-d 直接分配了设备的客户机，其动态迁移功能将受限，不过也可以使用热插拔或者libvirt 工具等方式来缓解这个问题。
- 不足的解决方案：（1）在一台物理宿主机上，仅少数 I/O 如网络性能要求较高的客户机使用 VT-d 直接分配设备，其他的使用纯模拟或者 virtio 已达到多个客户机共享同一个设备的目的（2）对于网络 I/O 的解决办法，可以选择 SR-IOV 是一个网卡产生多个独立的虚拟网卡，将每个虚拟网卡分配个一个客户机使用。

2. SR-IOV 设备分配

2.1 原理

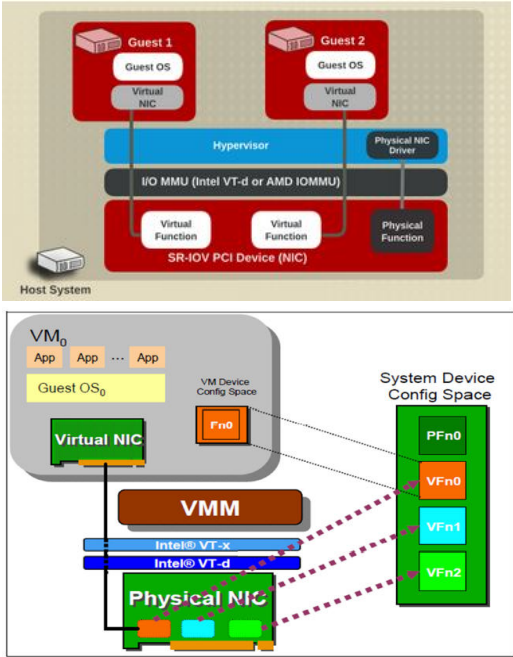
VT-d 的性能非常好，但是它的物理设备只能分配给一个客户机使用。为了实现多个虚拟机共享一个物理设备，并且达到直接分配的目的，PCI-SIG 组织发布了 SR-IOV（Single Root I/O Virtualization and sharing）规范，它定义了一个标准化的机制用以原生地支持实现多个客户机共享一个设备。不过，目前 SR-IOV（单根 I/O 虚拟化）最广泛地应用还是网卡上。

SR-IOV 使得一个单一的功能单元（比如，一个以太网端口）能看起来像多个独立的物理设备。一个带有 SR-IOV 功能的物理设备能被配置为多个功能单元。SR-IOV 使用两种功能（function）：

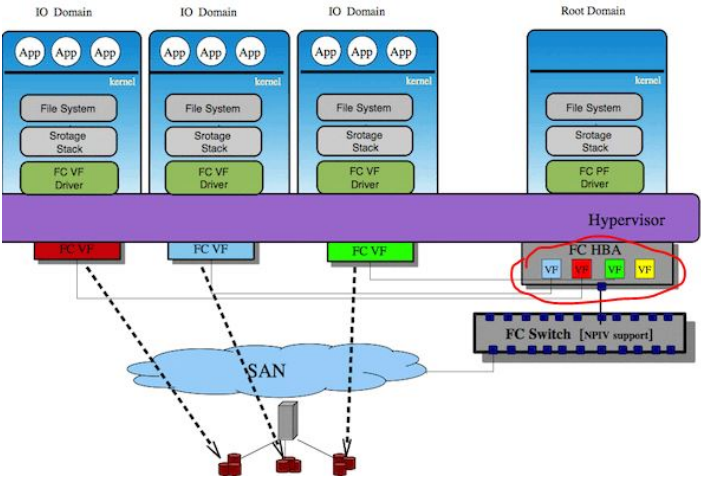
- 物理功能（Physical Functions，PF）：这是完整的带有 SR-IOV 能力的 PCIe 设备。PF 能像普通 PCI 设备那样被发现、管理和配置。
- 虚拟功能（Virtual Functions，VF）：简单的 PCIe 功能，它只能处理 I/O。每个 VF 都是从 PF 中分离出来的。每个物理硬件都有一个 VF 数目的限制。一个 PF，能被虚拟成多个 VF 用于分配给多个虚拟机。

Hypervisor 能将一个或者多个 VF 分配给一个虚拟机。在某一时刻，一个 VF 只能被分配给一个虚拟机。一个虚拟机可以拥有多个 VF。在虚拟机的操作系统看来，一个 VF 网卡看起来和一个普通网卡没有区别。SR-IOV 驱动是在内核中实现的。

网卡 SR-IOV 的例子：



光纤卡 SR-IOV 的例子：



2.2 SR-IOV 的条件

- 1. 需要 CPU 支持 Intel VT-x 和 VT-D（或者 AMD 的 SVM 和 IOMMU）
- 2. 需要有支持 SR-IOV 规范的设备：目前这种设备较多，比如Intel的很多中高端网卡等。
- 3. 需要 QEMU/KAM 的支持。

RedHat Linux 6.0 官方只完整测试了下面的几款 SR-IOV 网卡：

- Intel® 82576NS Gigabit Ethernet Controller (igb 驱动)
- Intel® 82576EB Gigabit Ethernet Controller (igb 驱动)
- Intel® 82599ES 10 Gigabit Ethernet Controller (ixgbe 驱动)
- Intel® 82599EB 10 Gigabit Ethernet Controller (ixgbe 驱动)

2.3 分配 SR-IOV 设备的步骤

手头没有支持SR-IOV的设备。这是 RedHat 上 SR-IOV 的配置步骤：[Using SR-IOV](#)。

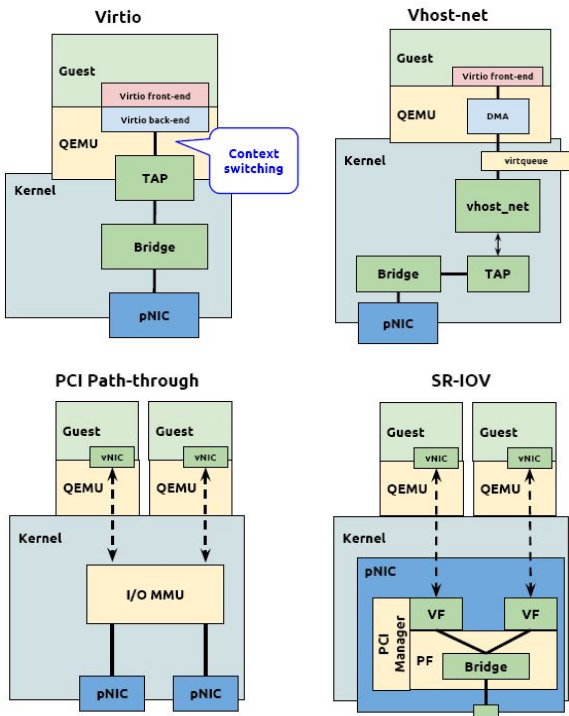
简单来说，SR-IOV 分配步骤和设备直接分配相比基本类似，除了要使 PF 虚拟化成多个 VF 以外。

2.4 优势和不足

优势	不足
<div>1. 真正实现设备共享（多个客户机共享一个 SR-IOV 设备的物理端口）</div> <div>2. 接近原生性能</div> <div>3. 相比 VT-d，SR-IOV 可以使用更少的设备来支持更多的客户机，可以提高数据中心的空</div> <div>间利用率。</div>	<div>1. 对设备有依赖，目前只有部分设备支持 SR-IOV。RedHat Linux 只是测试了 Intel 的几款高端网卡。</div> <div>2. 使用 SR-IOV 时不方便动态迁移客户机。这是因为这时候虚拟机直接使用主机上的物理设备，因此虚拟机的迁移（migration）和保存（save）目前都不支持。这个在将来有可能被改变。</div>

3. 各种设备虚拟化方式的比较

3.1 架构上的比较（以网卡为例）



3.2 性能上的比较（以网卡为例）

纯模拟网卡和物理网卡的比较：

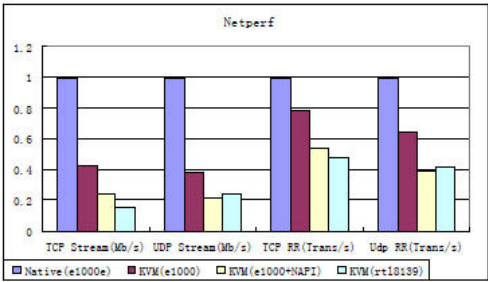


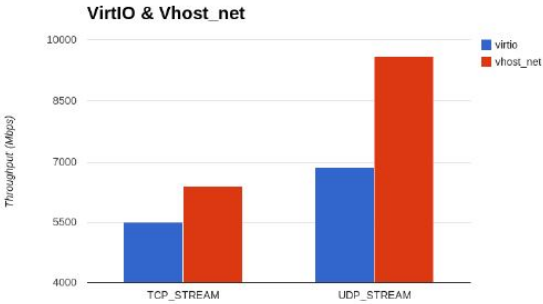
Fig. 4. Netperf benchmark test, Guest OS used different NIC configurations - virtual e1000 NIC (e1000 driver, e1000 driver with NAPI support), virtual rtl8139 NIC

（来源：Evaluating and Optimizing I/O Virtualization in Kernel-based Virtual Machine (KVM), Binbin Zhang, Xiaolin Wang, Rongfeng Lai, Liang Yang, Zhenlin Wang, Yingwei Luo, Xiaoming Li）

（测试环境：两台物理服务器 HostA 和 HostB，都使用GB以太网。HostA 使用 82566DC 网卡，HostB 使用 82567LM-2 网卡，一台虚拟机运行在 HostB 上，使用 KVM-76。）

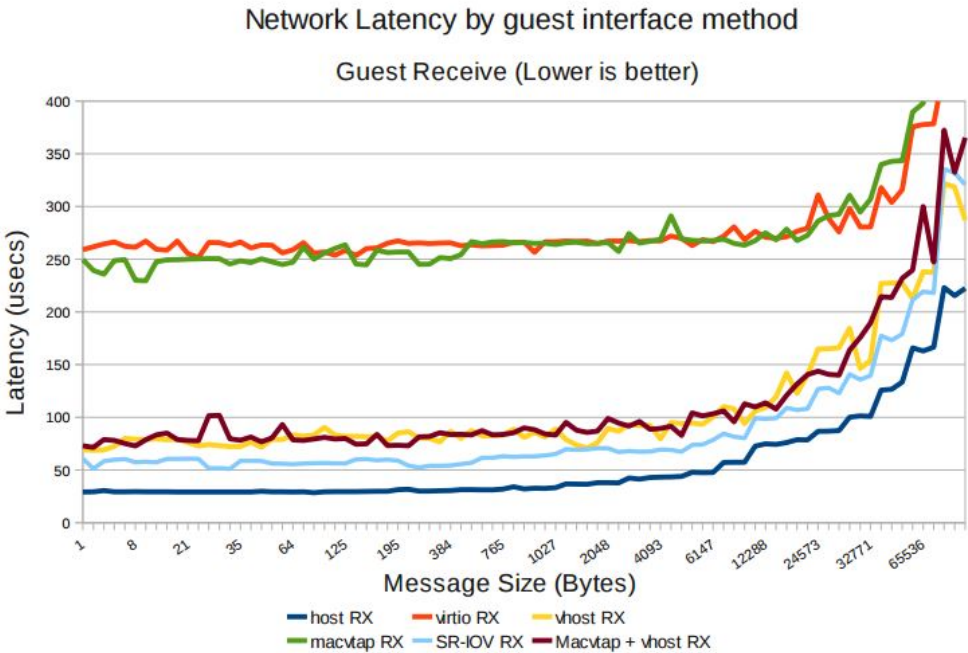
- 结论：
- 纯模拟网卡的性能只有物理网卡的四成到六成
 - 纯模拟网卡的 UDP 性能比 TCP 性能高 50% 到 100%
 - 在虚拟网卡上使用 NAPI，不但不会提高性能，反而会是性能下降
 - e1000 的性能比 rtl8139 的性能高不少（为什么 RedHat Linux KVM 上默认的网卡是 rt18139 呢？）

Virtio 和 vhost_net 的吞吐量比较：



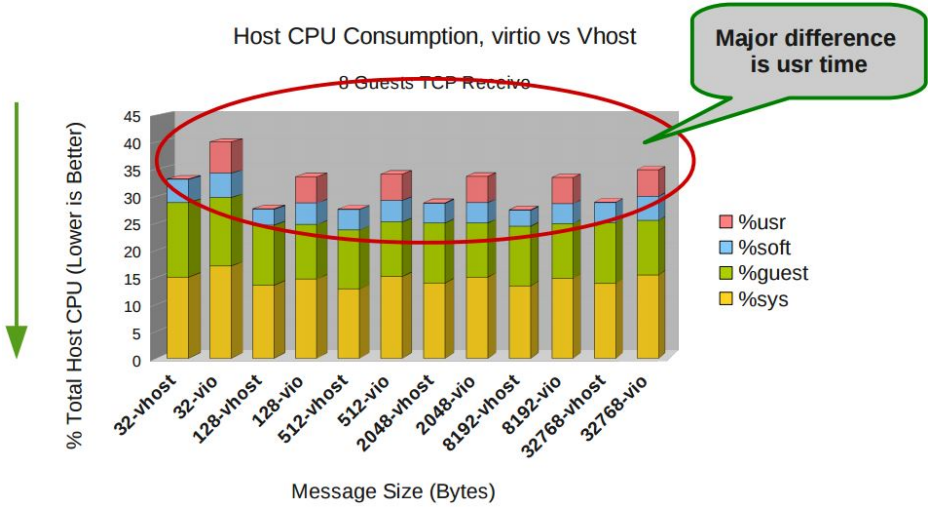
- 来源：CANONICAL, KVM Performance Optimization, Paul Sim,Cloud Consultant, paul.sim@canonical.com
- 结论：vhost_net 比 virtio 的 UDP 和 TCP 性能高 20% 左右。

RedHat Linux 6 上 virtio，vhost_net，SR-IOV 和物理设备网络延迟的比较：



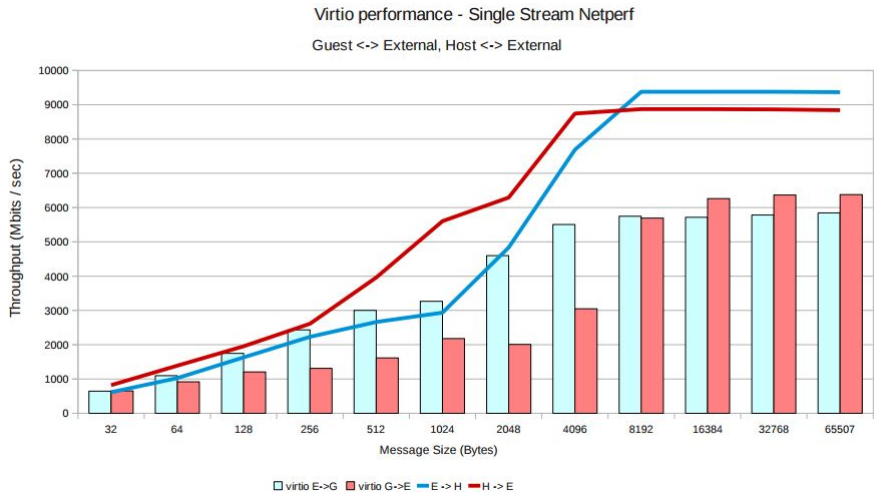
（来源：RedHat 官网）

RedHat Linux 6 上 virtio 和 vhost_net 所消耗的主机CPU资源的比较：



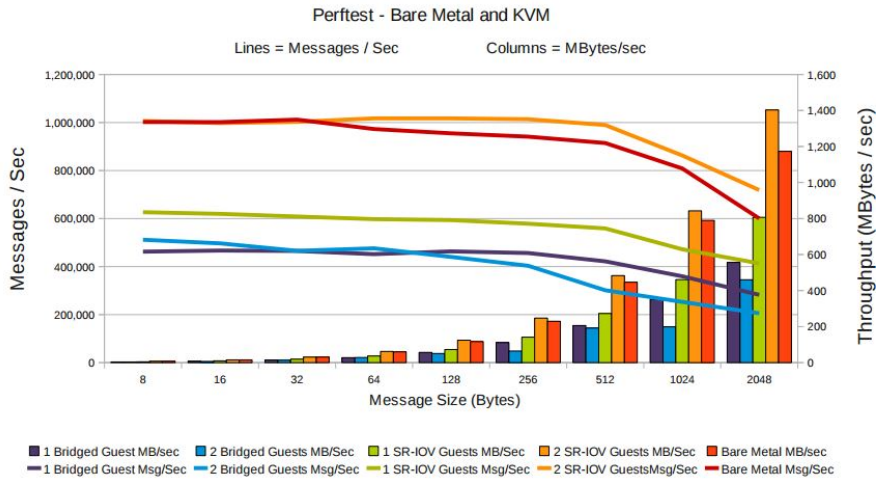
(来源同上)

使用 virtio 的 KVM 与物理机的 TCP 吞吐量对比：



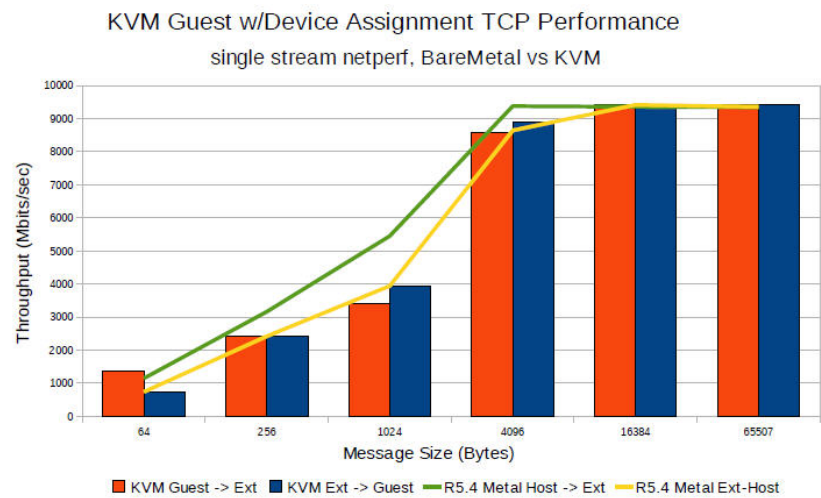
(数据来源：[RedHat 官网](#))

物理机与使用 SR-IOV 的 KVM 的网络性能对比：



(来源：同上)

物理机与使用 Pass-through 的KVM 的 TCP 性能对比：



(资料来源：Open Source Virtualization: KVM and Linux, Chris Wright, Principal Software Engineer, Red Hat, September 4, 2009)

3.3 Virtio 和 Pass-Through 的详细比较

Basic Operation	<ul style="list-style-type: none">- Backend/Guest direct access to shared Vring buffers - PIO- Switching at software level- Management Flexibility - internal SDN support ovs-vsctl add-port br0 <phys-intfc> - vswitch ovs-ofctl control flows- IRQ bottleneck - QEMU - call into kvm inject Kernel - inject directly	<ul style="list-style-type: none">- Direct access to hw memory regions- DMA Support- Switching at hw level - SR-IOV depends on #of Queues- Management Flexibility - external SDN capable- IRQ bottleneck - hw enhancements, posted interrupts, exitless EOI improve things - closer to native
Migration	<ul style="list-style-type: none">- Virtio lockless- Saves device state, tracks dirty pages	<ul style="list-style-type: none">- QEMU sets 'unmigratable', or installs migration blocker- Guest can be holding a lock - deadlock, hw state,
Scalability	<ul style="list-style-type: none">- Practical limitations - primarily per formance	<ul style="list-style-type: none">- Number of Devices limited, limits #VMs- SR-IOV - #of VF - # of queues
Network Performance	<ul style="list-style-type: none">- Soft switching - bridge, vSwitch- Several IO HOPS- Can approach near native - 10Ge for few bridged Guest	<ul style="list-style-type: none">- Switching done at HW level - hw queues- Performance scales with # of Guests- DMA support- IRQ Passthrough still a problem
Host Performance	<ul style="list-style-type: none">- PIO - takes cpu cycles- Exits - few but still- Guest pages swappable	<ul style="list-style-type: none">- Guest pinned - can't swap- Fewer exits- Less PIO
Cloud Environment	<ul style="list-style-type: none">- Cloud friendly - migration, SDN, paging	<ul style="list-style-type: none">- Not Cloud friendly, great for NFV/RT DDPK, run to completion

(来源：Reconnaissance of Virtio: What's new and how it's all connected? by Mario Smarduch)

4. 综合结论

KVM 依赖的Intel/AMD 处理器的各种虚拟化扩展：

处理器	CPU 虚拟化	内存虚拟化	PCI Pass-through
Intel	VT-x	VPID , EPT	VT-d
AMD	AMD-V	ASID , NPT	IOMMU

I/O 虚拟化方案的选择：

- I/O设备尽量使用准虚拟化（virtio 和 vhost_net）
- 如果需要实时迁移，不能使用 SR-IOV
- 对更高I/O要求又不需要实时迁移的，可以使用 SR-IOV
- 每种方案都有优势和不足，在特定环境下其性能有可能反而下降，因此在生产环境中使用各种虚拟化方式前需要经过完整测试

其它参考资料：

- RedHat Linux 6 官方文档
- KVM 官方文档
- KVM 虚拟化技术实战与解析 任永杰、单海涛 著
- KVM 虚拟化技术在 AMD 平台上的实现

分类: KVM,虚拟化

好文要顶

关注我

收藏该文

SammyLiu

关注 - 30

粉丝 - 470

荣誉：推荐博客

+加关注

2

推荐

0

反对

« 上一篇：KVM 介绍（3）：I/O 全虚拟化和准虚拟化 [KVM I/O QEMU Full-Virtualizaiton Para-virtualization]
» 下一篇：KVM 介绍（5）：libvirt 介绍 [Libvirt for KVM/QEMU]

posted on 2015-06-05 08:06 SammyLiu 阅读(7732) 评论(0) 编辑 收藏

[刷新评论](#) [刷新页面](#) [返回顶部](#)

注册用户登录后才能发表评论，请 [登录](#) 或 [注册](#)，[访问网站首页](#)。

- 【推荐】50万行VC++源码: 大型组态工控、电力仿真CAD与GIS源码库
- 【报表】Excel 报表开发18 招式，人人都能做报表
- 【活动】阿里云海外云服务全面降价助力企业全球布局
- 【实用】40+篇云服务器操作及运维基础知识！

每周末去

“硅谷大学”

Google 机器学习认证项目

首推价末班车

立即加入

Google

认证

14

优达学城

- 最新IT新闻:
- 知乎上线视频功能，以后看教程更方便了
 - 一年只赚2万元：乐视游戏或被出售
 - Unity获得4亿美元投资，现估值为26亿美元
 - 直播对陌陌的意义，就像王者荣耀之于腾讯游戏
 - 死磕支付宝？苏宁金融发布“星辰计划”：扫码支付返888元
- » 更多新闻...

阿里云

云服务器

降破底价

30元/月

轻松搭建网站/应用

- 最新知识库文章:
- 程序员的工作、学习与绩效
 - 软件开发为什么很难
 - 唱吧DevOps的落地，微服务CI/CD的范本技术解读
 - 程序员，如何从平庸走向理想？
 - 我为什么鼓励工程师写blog
- » 更多知识库文章...