

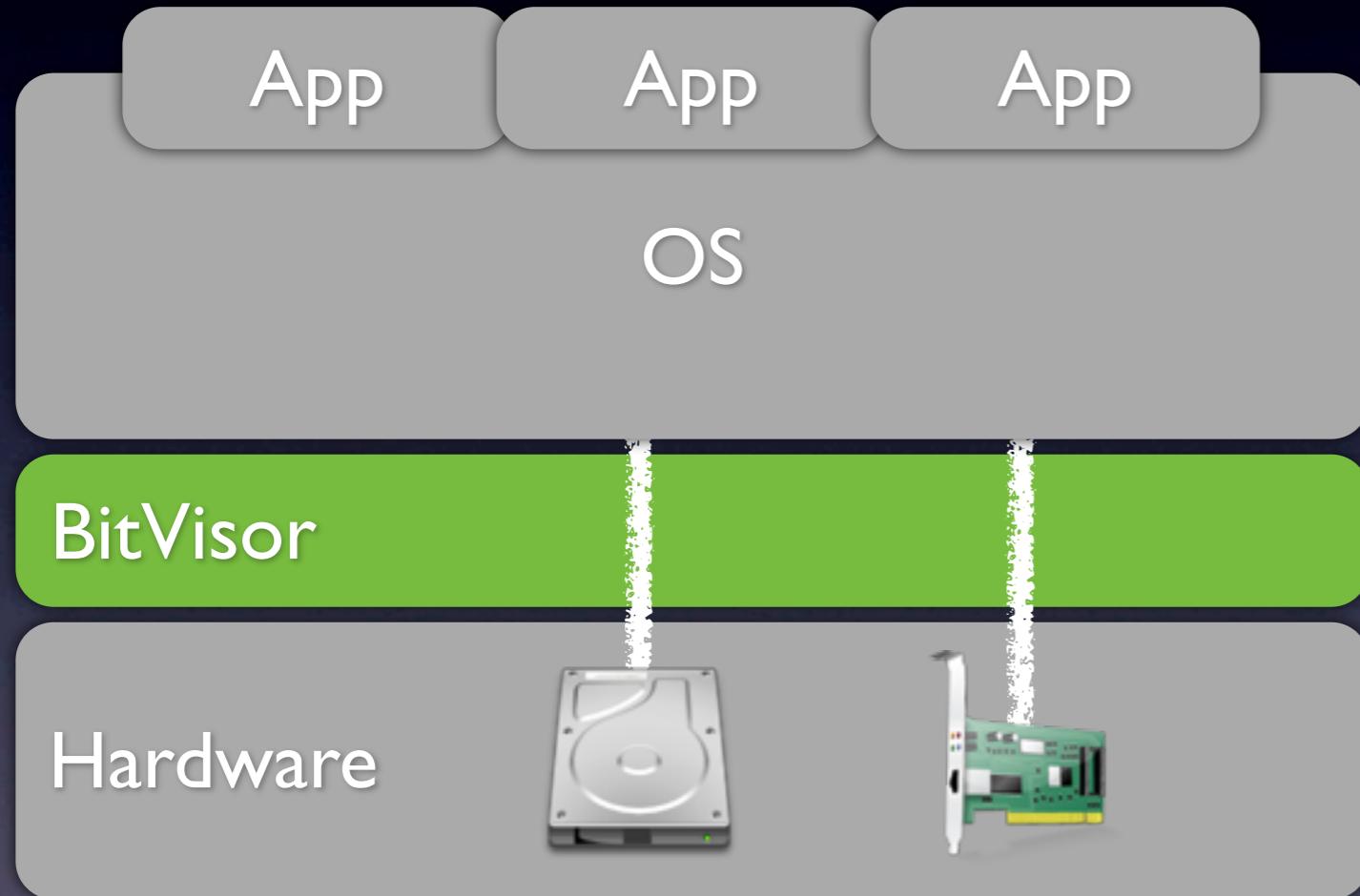
OSb: OSv on BitVisor

(2014/11/21 BitVisor Summit 3)

Yushi Omote
University of Tsukuba

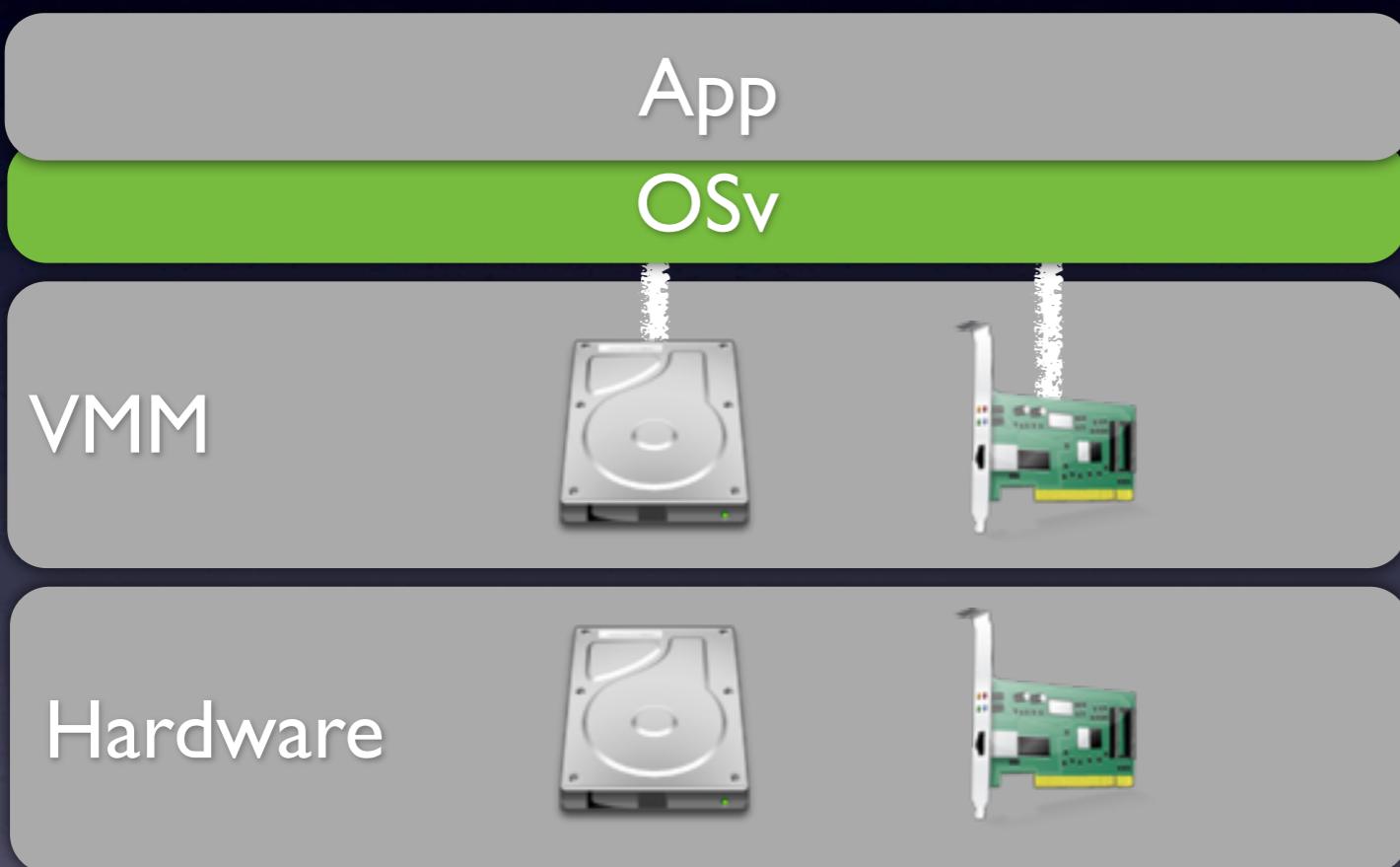
Takahiro Shinagawa
The University of Tokyo

BitVisor



<http://www.justis.as-i.co.jp>

OSv



<http://medical-care.feed.jp>



BitVisor



BitVisor

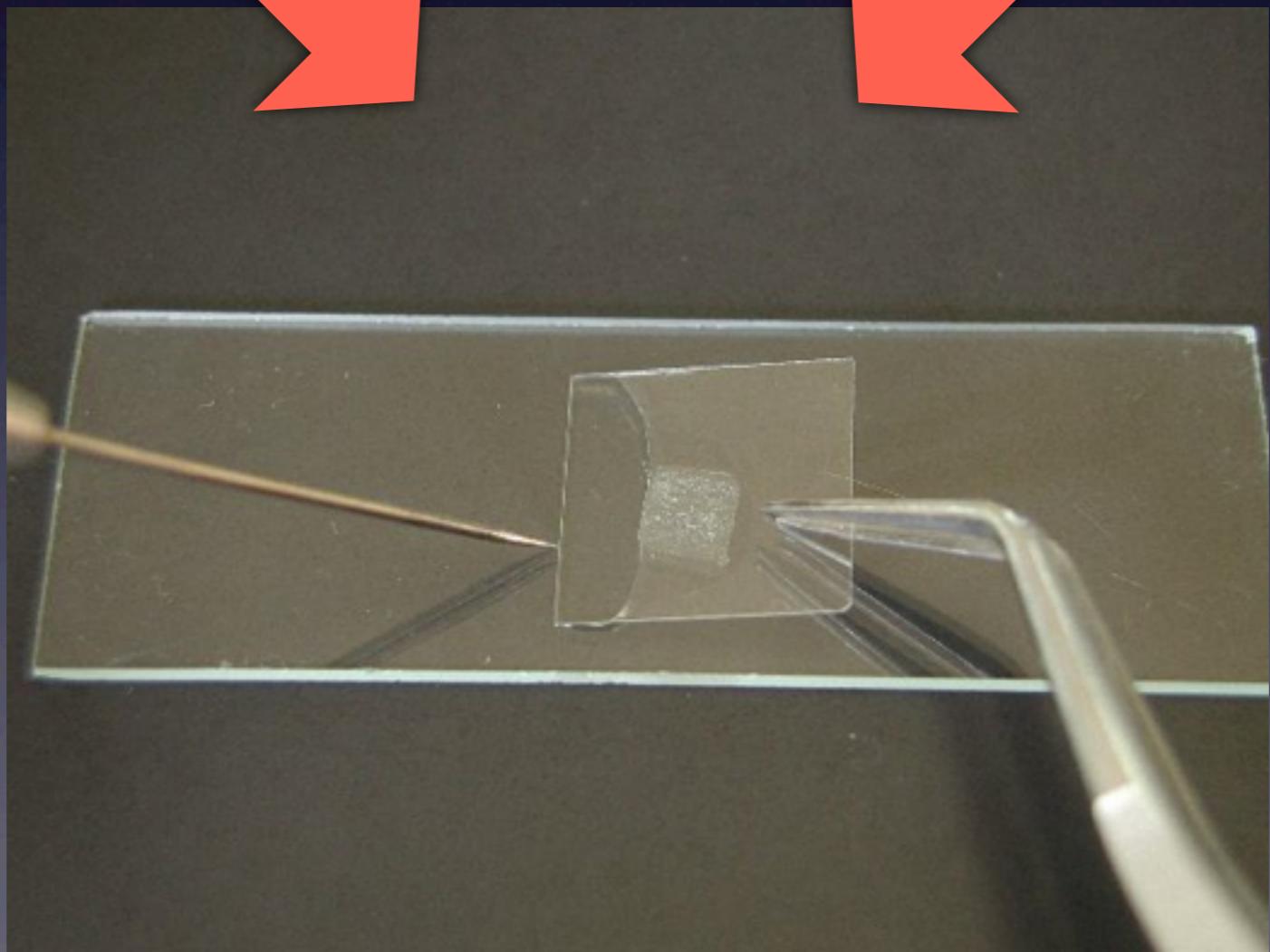


OSv



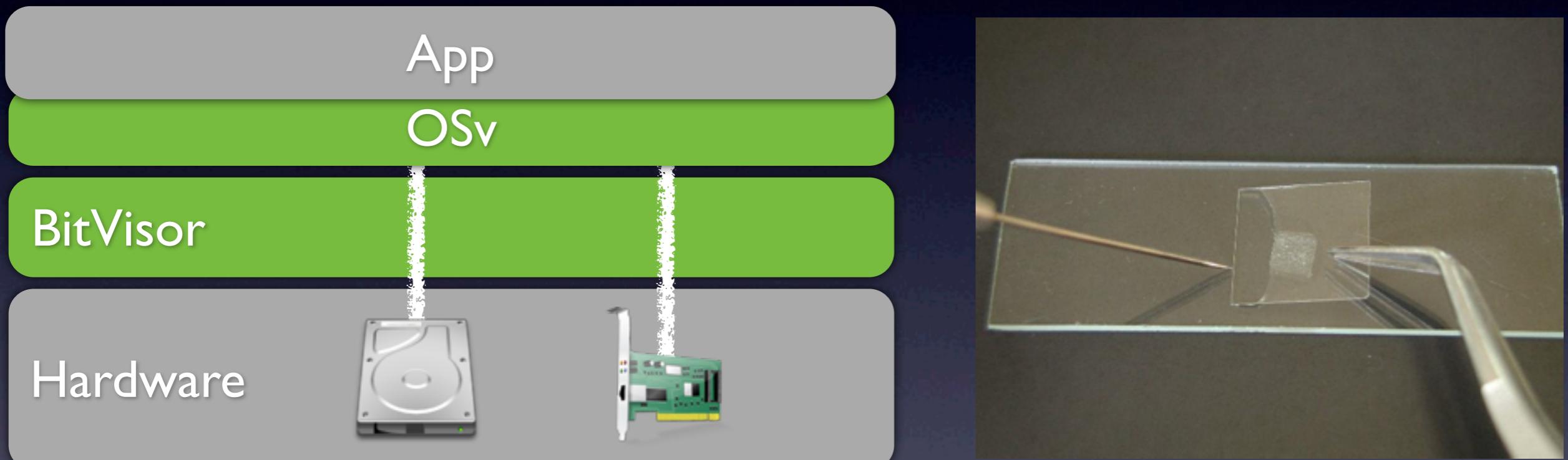
BitVisor

OSv

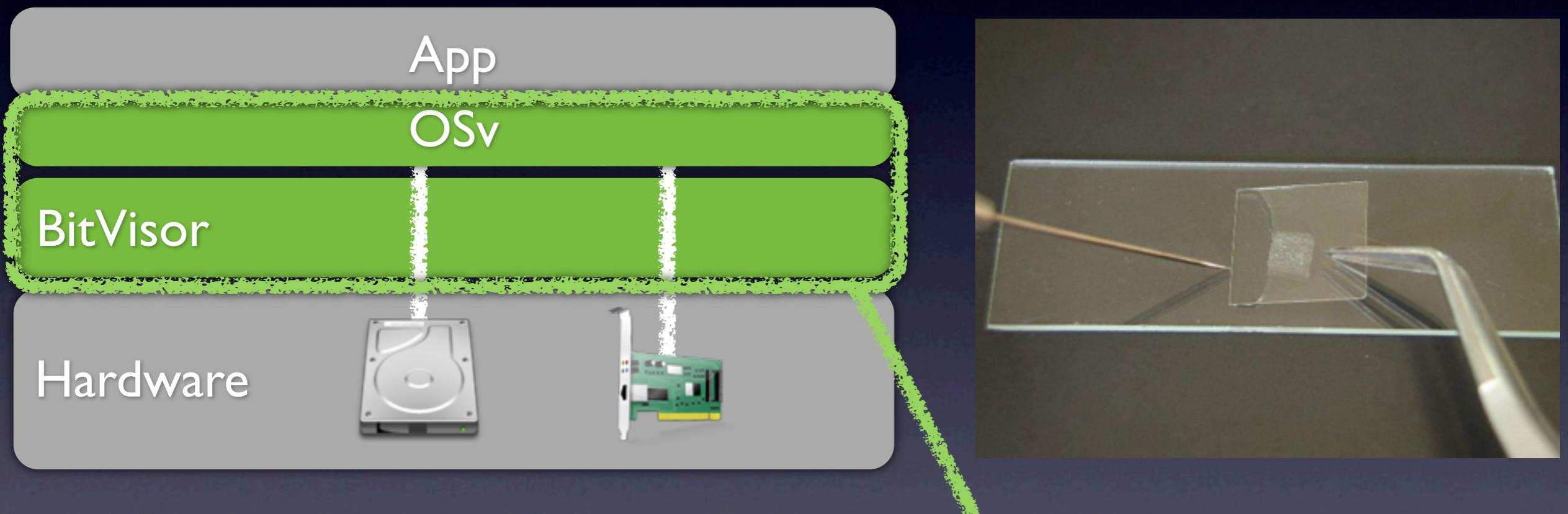


<http://www.root.ne.jp/nishide/shs/>

OSv on BitVisor

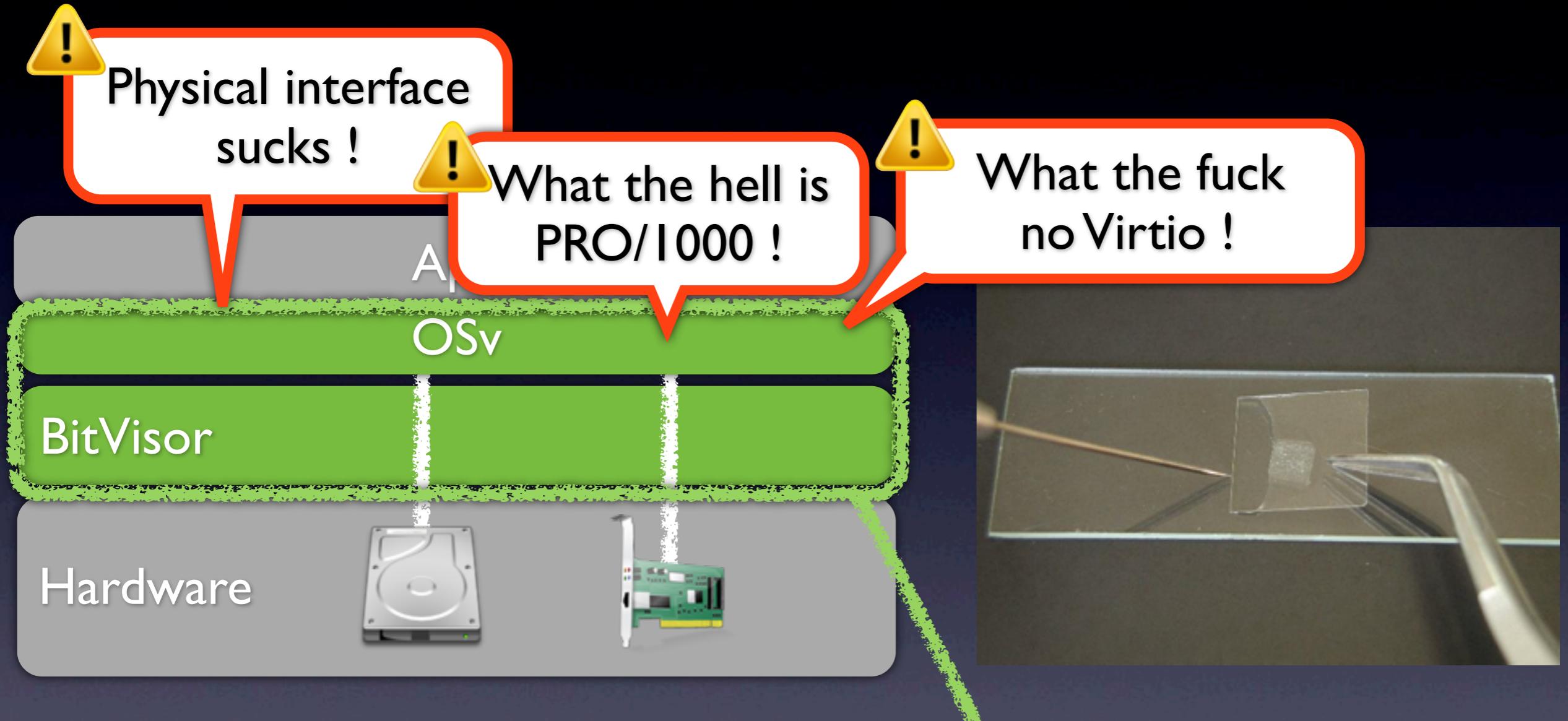


OSv on BitVisor



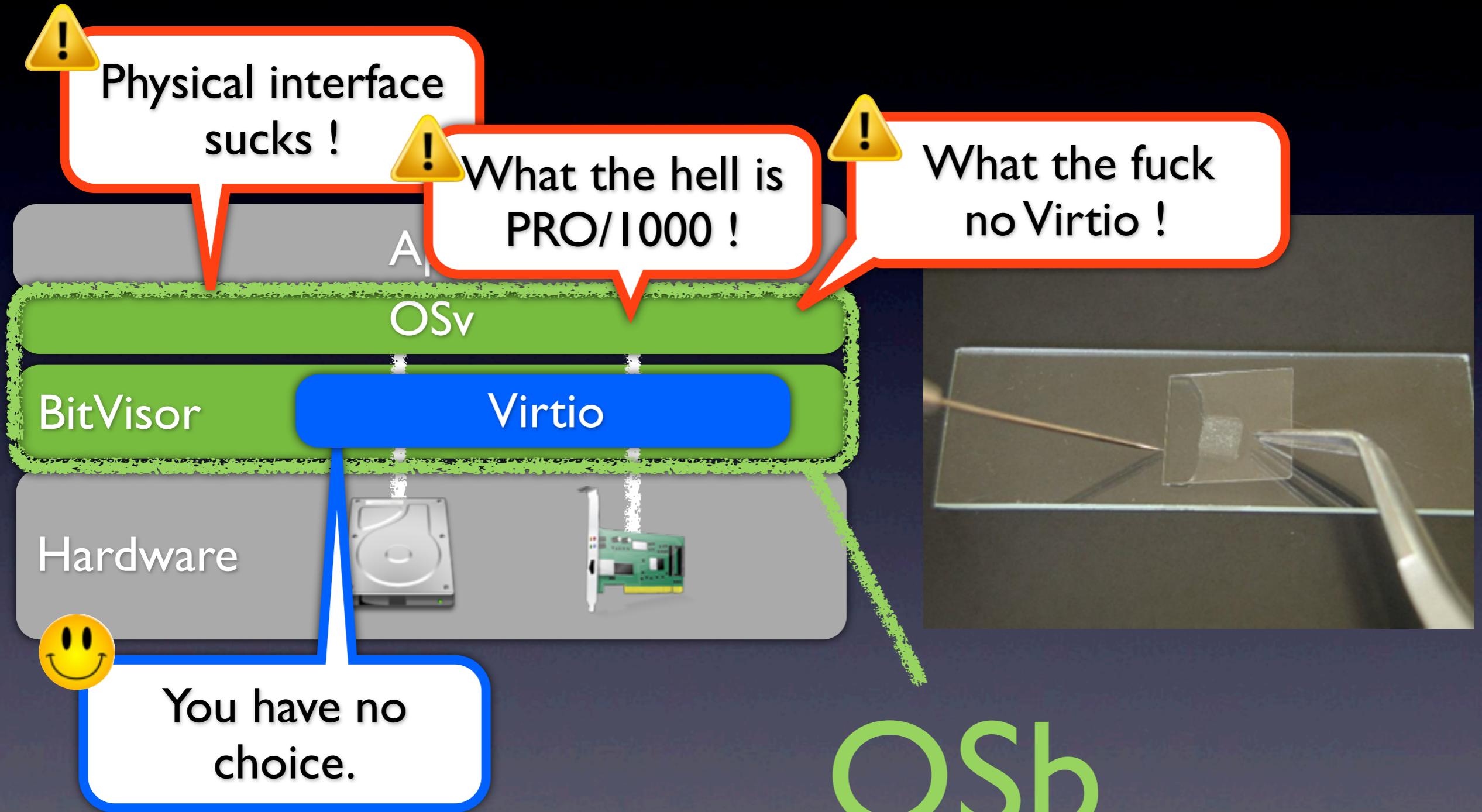
OSb

OSv on BitVisor

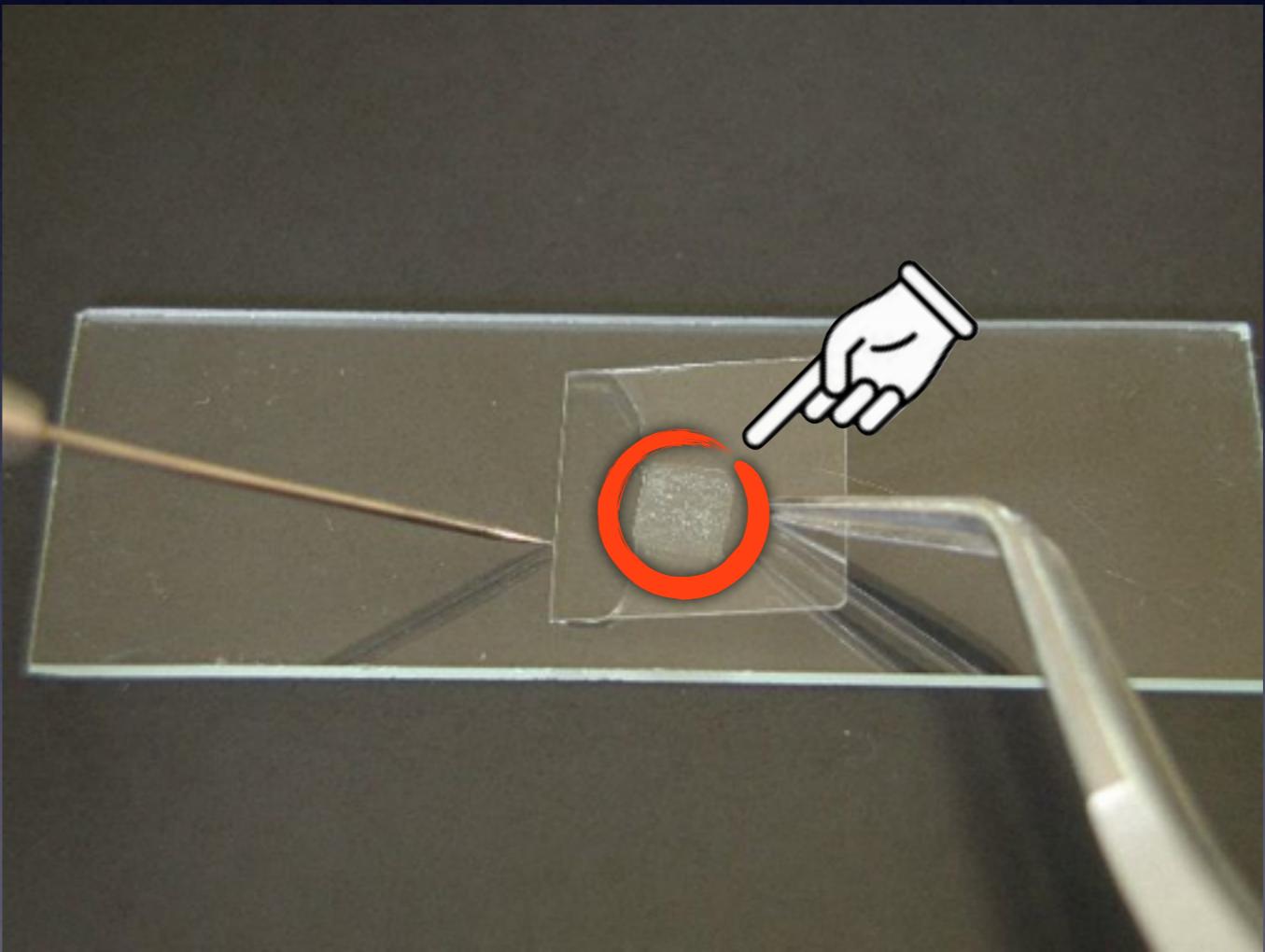


OSb

OSv on BitVisor



~The Road to OSb ~



OSv Code Reading

- Boot Process
- Some Drivers



OSv may run on
physical machine.
(=BitVisor)

- MBR
- Read local disk with INT13/42.
- Load command line, boot loader, OSv kernel.
- Get memory map with INT15/E820.
- Setup segment descriptors/page tables.
- Switch to 64-bit mode.
- premain()
- main()
- ...

- Serial /VGA output.
- SATA.
- Virtio NIC/BLK/RNG/SCSI.
- ACPI.
- APIC.
- ...

Only one thing to fix

Skip on BitVisor

Triple fault 

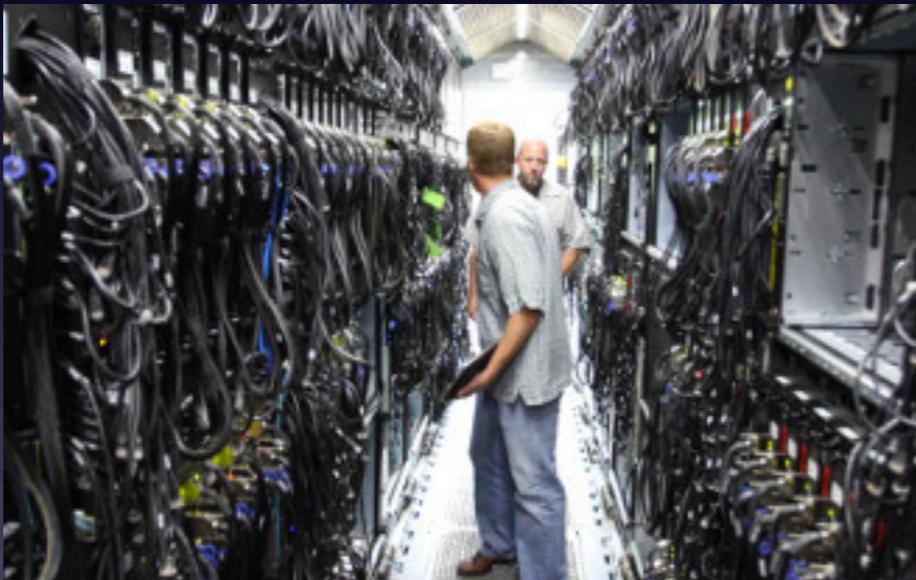
drivers/acpi.c:

```
// Copy the root table list to dynamic memory
if (!is_bitvisor()) {
    status = AcpiReallocateRootTable();
    if (ACPI_FAILURE(status)) {
        acpi_e("AcpiReallocateRootTable failed: %s\n",
        AcpiFormatException(status));
        return;
    }
}
```

Hello world !

On BitVisor,

Whenever updating OSv images,
you update the entire Physical Disk.



<http://a.fsdn.com/sd/articles/14/11/12/1946208-1.jpg>



<http://www.dreamstime.com>

Troublesome...
(especially for development)

Network-boot of OSb

`./scripts/run.py`



Network-boot of OSb

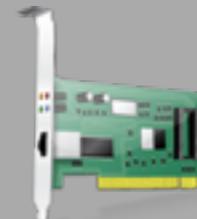
`./scripts/run.py`



Network-boot of OSb

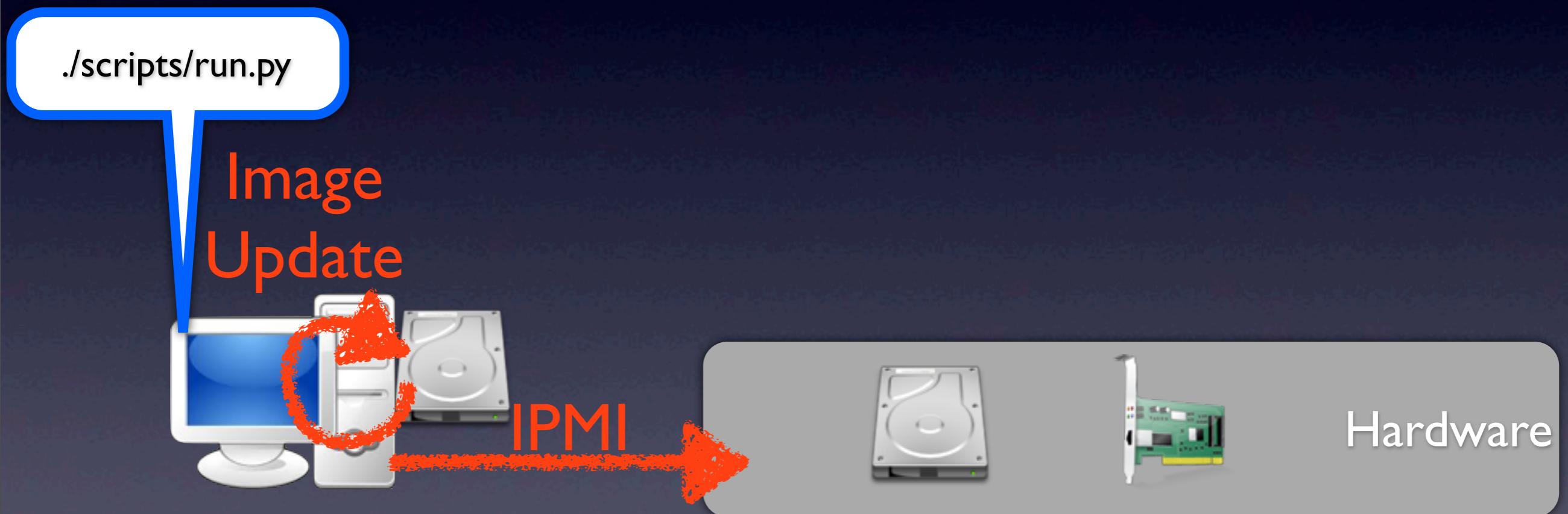
`./scripts/run.py`

Image
Update

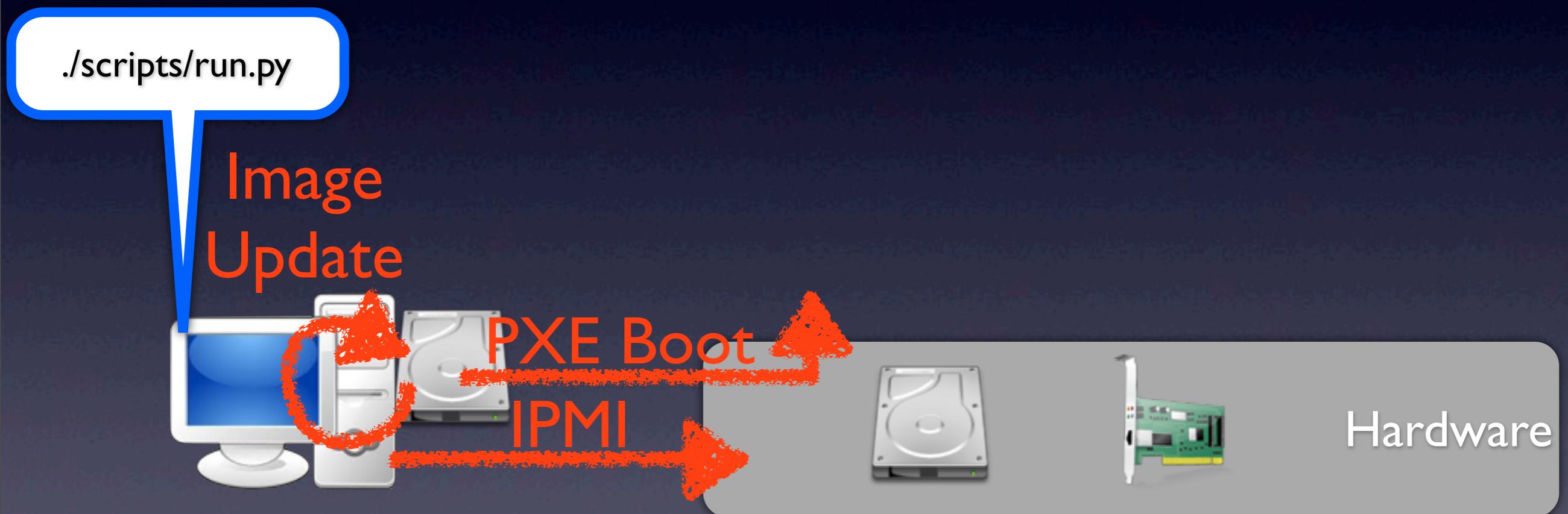


Hardware

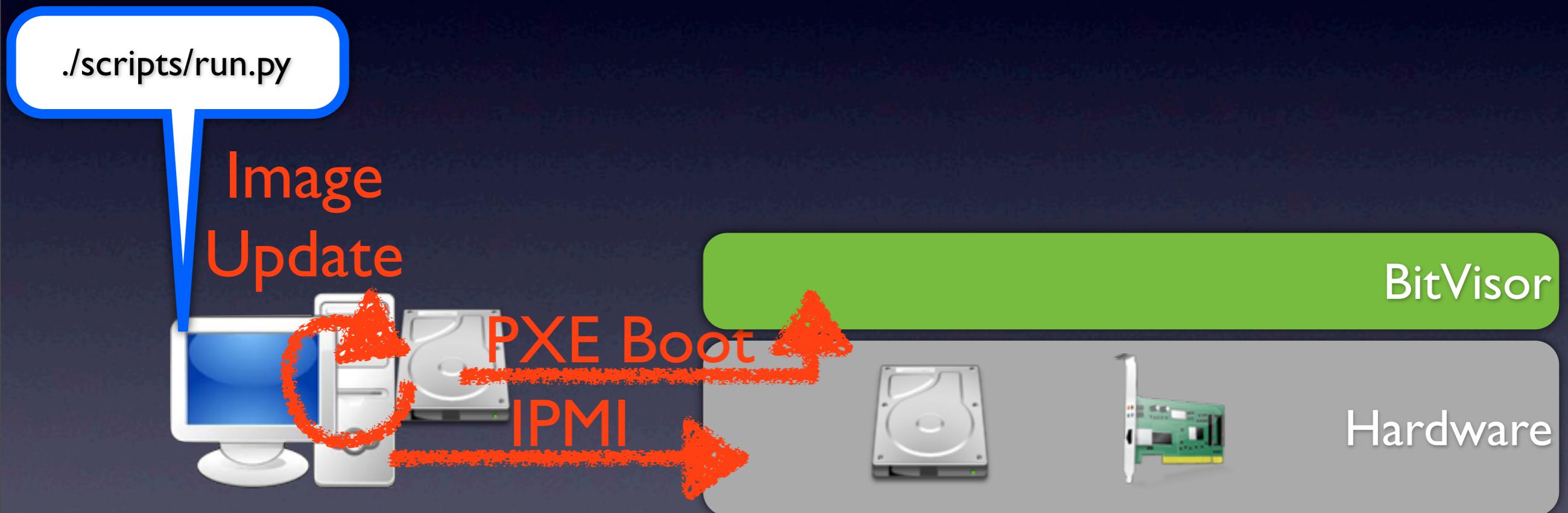
Network-boot of OSb



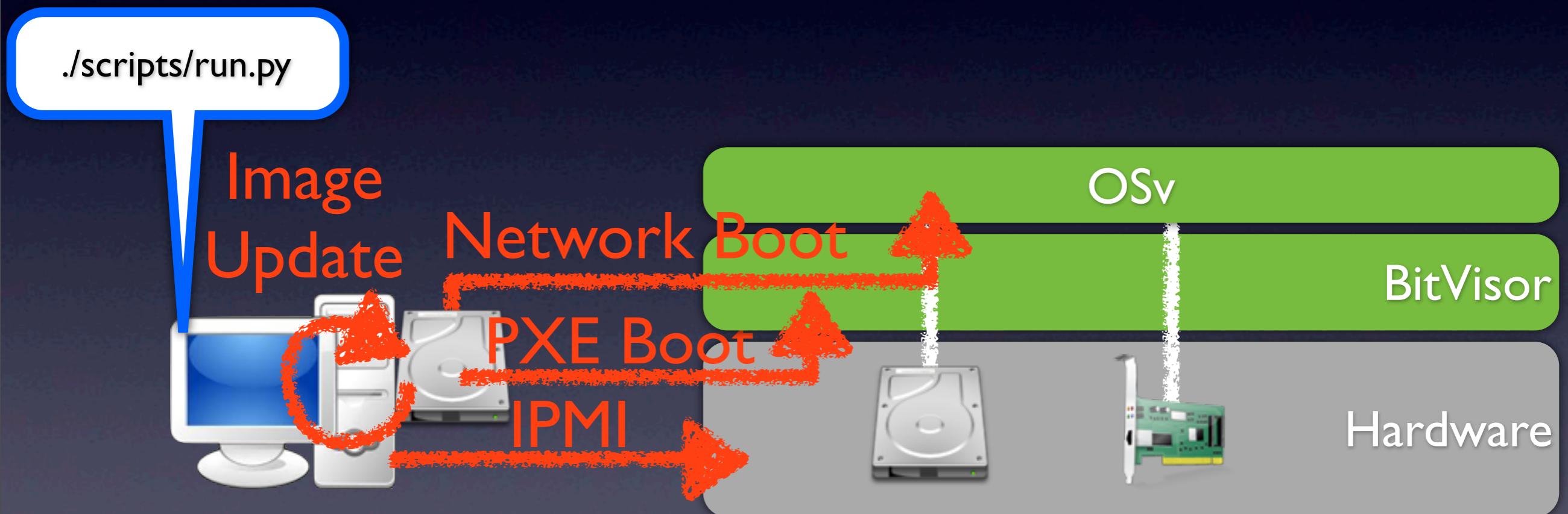
Network-boot of OSb



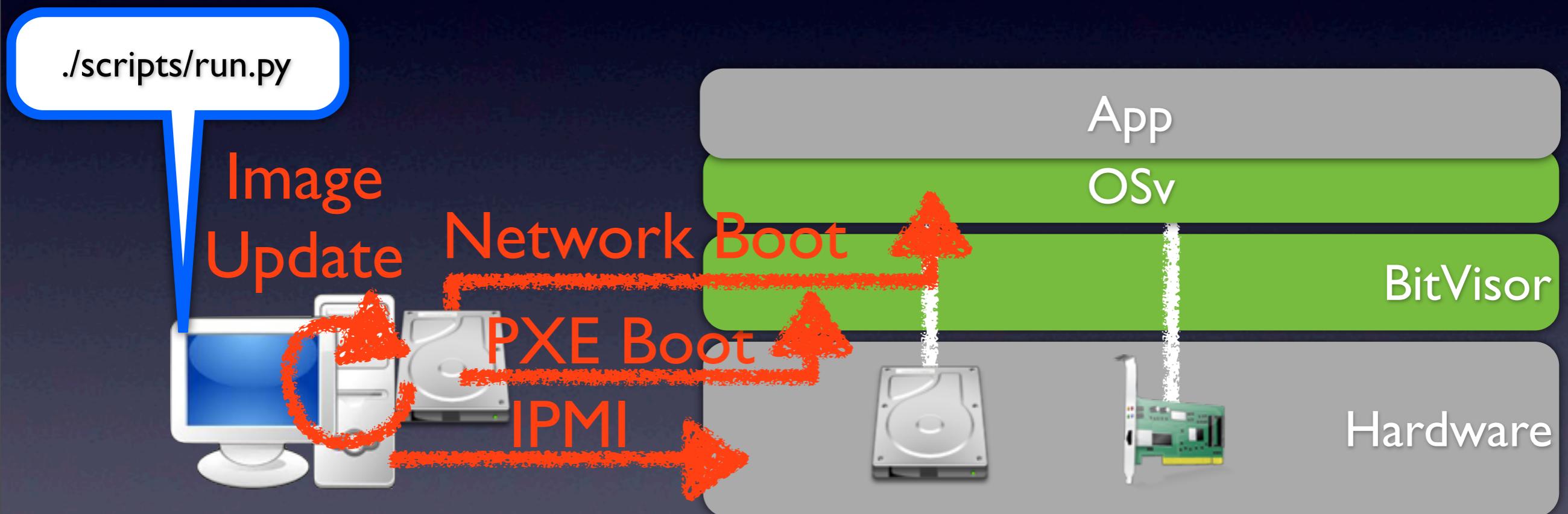
Network-boot of OSb



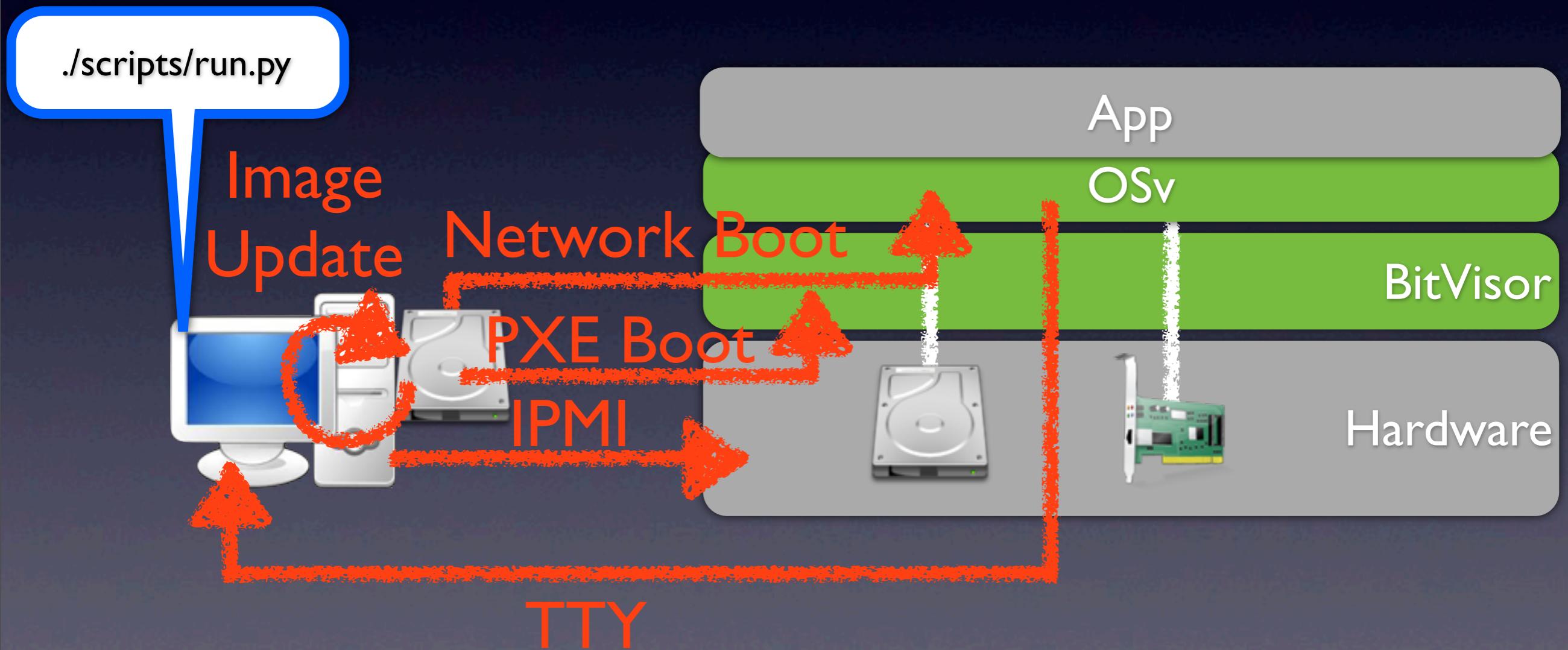
Network-boot of OSb



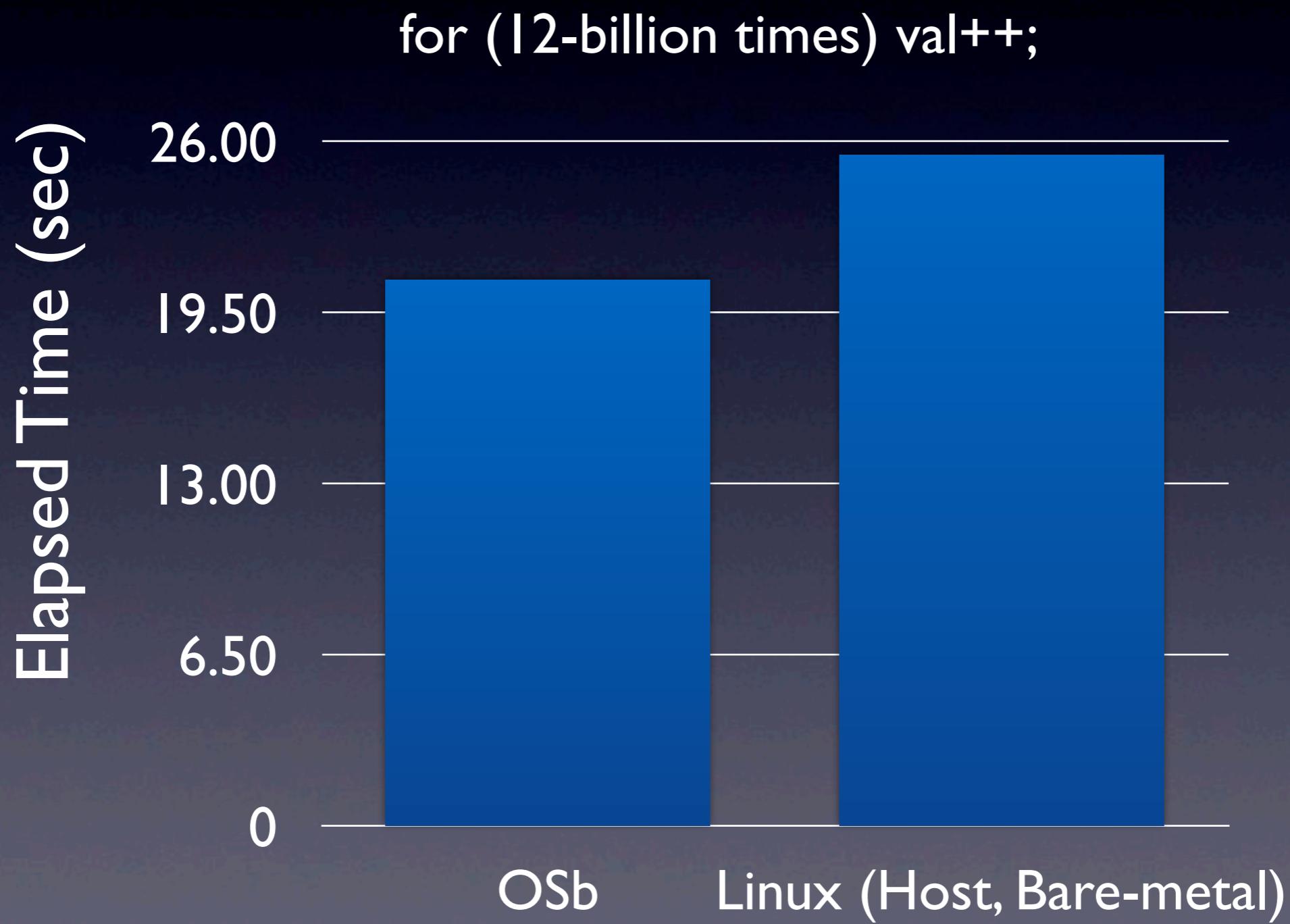
Network-boot of OSb



Network-boot of OSb



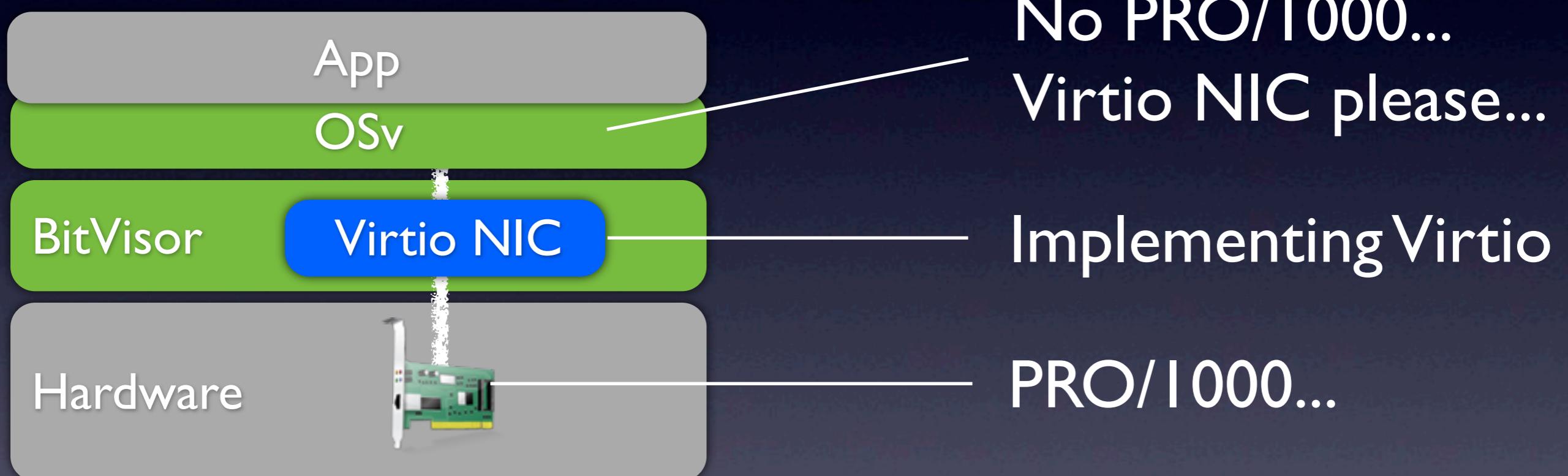
Performance?



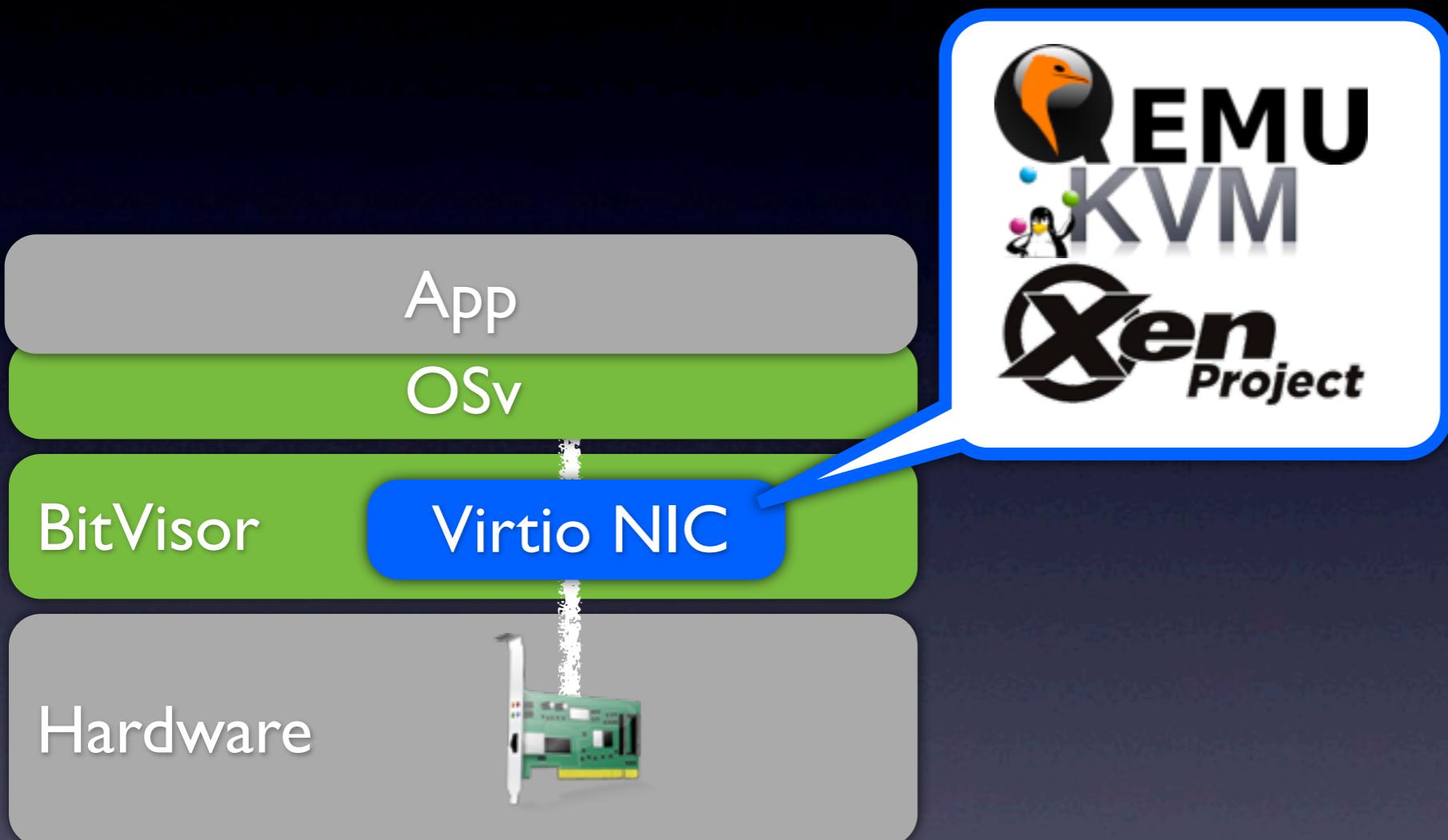
Virtio NIC



Virtio NIC

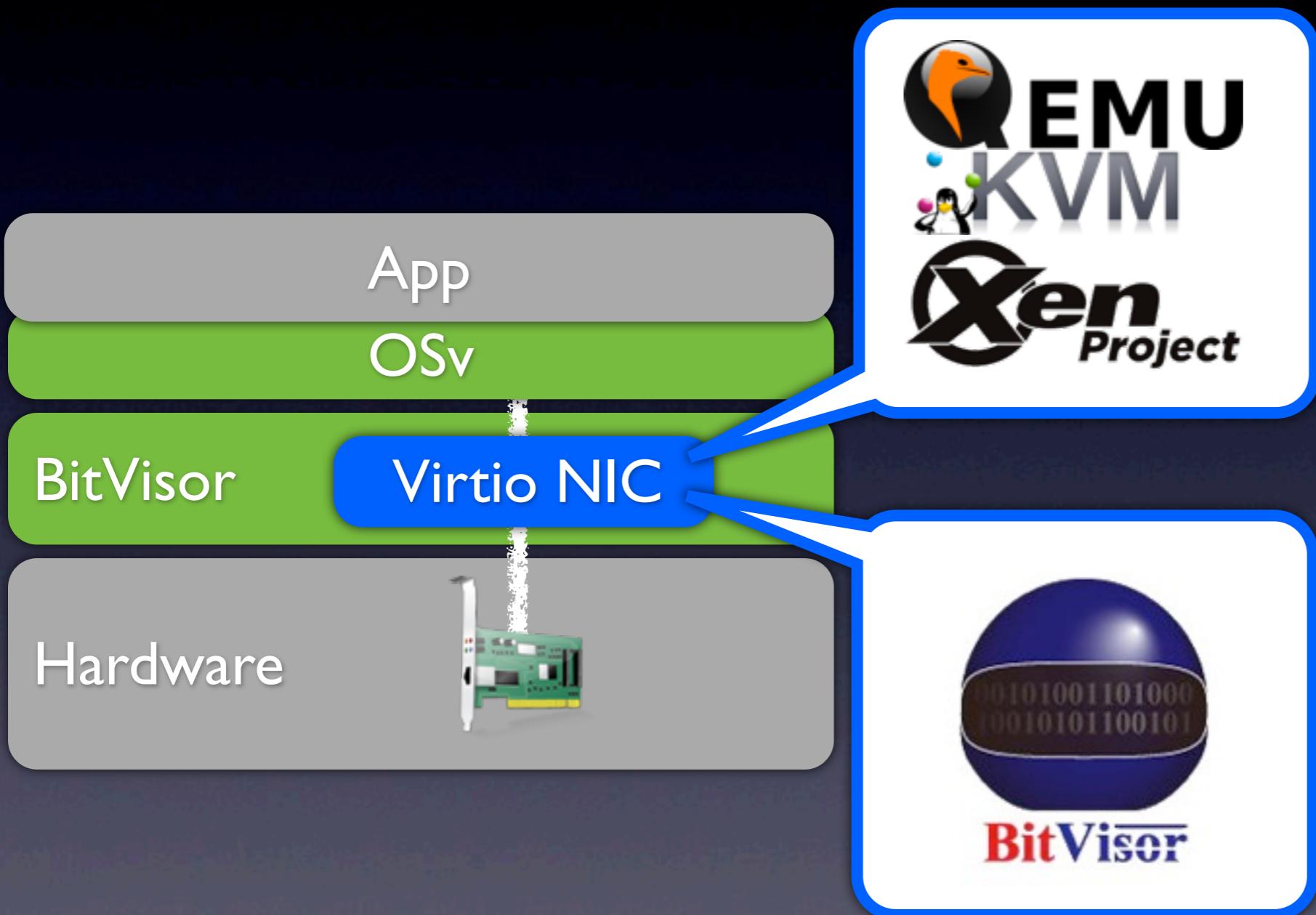


How to implement?



Finally, Our
Virtualization?

How to implement?



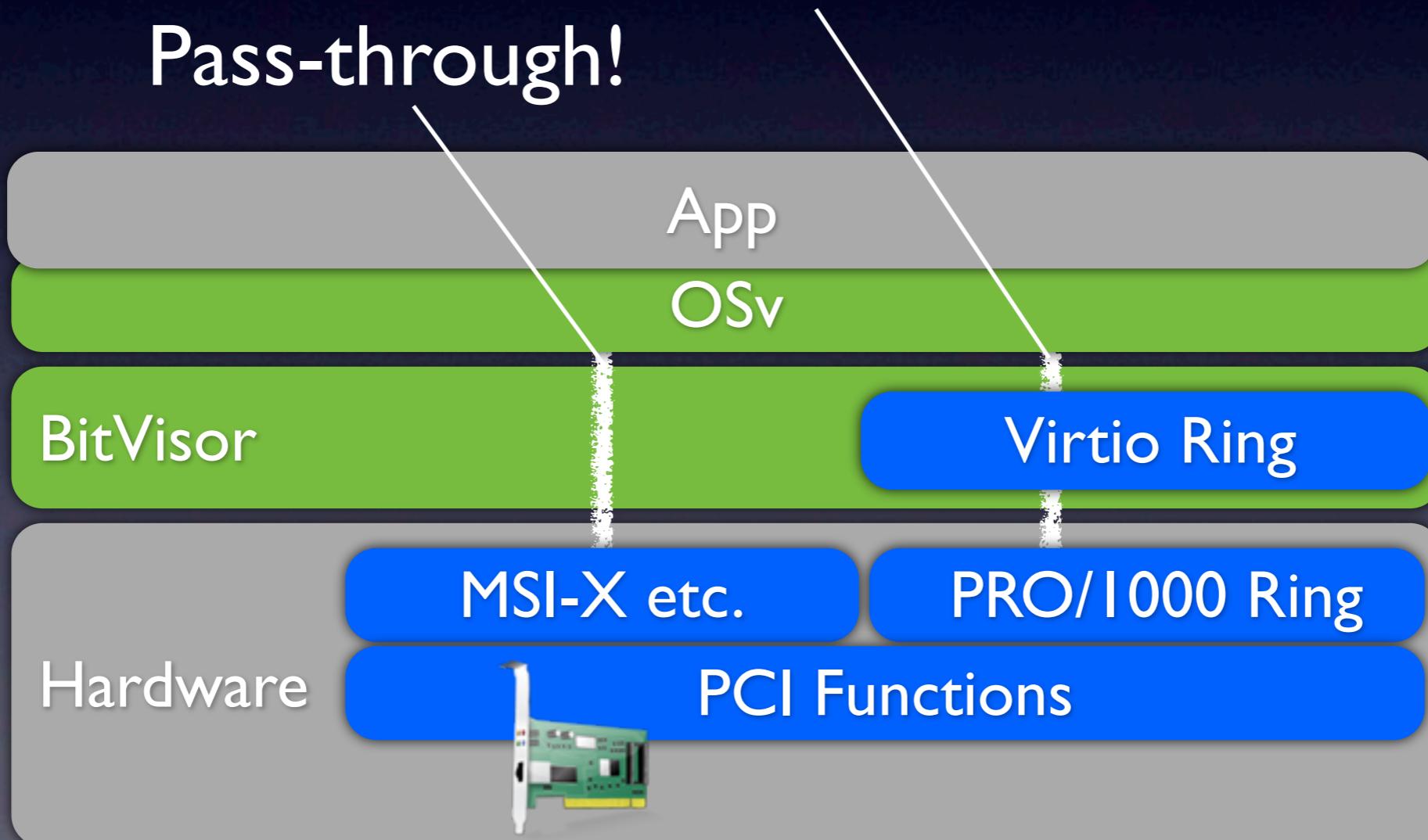
Finally, Our
Virtualization?

No,
BitVisor's Way!

Para Pass-through Virtio

Intercept only interesting I/Os!

Pass-through!



0x1000
(Virtio NIC)

Faking PCI IDs

0x1AF4
(Virtio Device)

31	16 15	00h
	Device ID	Vendor ID
	Status	Command
	Class Code	Revision ID
BIST	Header Type	Lat. Timer
		Cache Line S.
		10h
		1Ch
		20h
		24h
		28h
		2Ch
		30h
		34h
		38h
		3Ch
		18h
		14h
		10h
		0Ch
		08h
		04h
		00h
		1Ch
		20h
		24h
		28h
		2Ch
		30h
		34h
		38h
		3Ch

Base Address Registers

Cardbus CIS Pointer

Subsystem ID **Subsystem Vendor ID**

Expansion ROM Base Address

Reserved **Cap. Pointer**

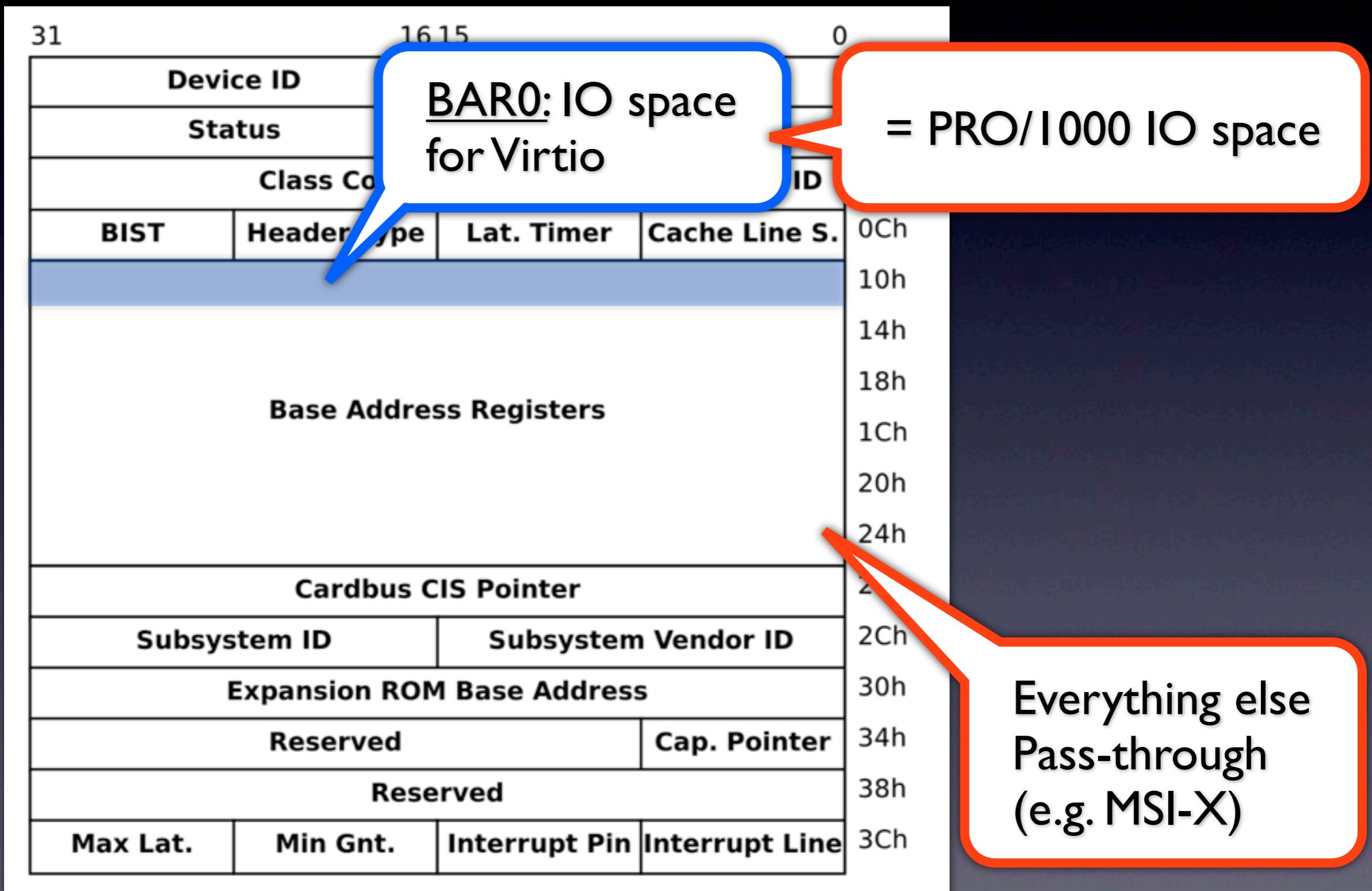
Reserved

Max Lat. **Min Gnt.** **Interrupt Pin** **Interrupt Line**

I (Virtio NIC)

0 (Legacy Virtio)

Faking PCI BARs



Virtio NIC Operations

- Virtio Ring
- Packet Transmission
- Packet Reception

Virtio Ring

BAR0

(IO)



BitVisor emulates
these registers

Guest selects rings:
* Transmission Ring
* Reception Ring

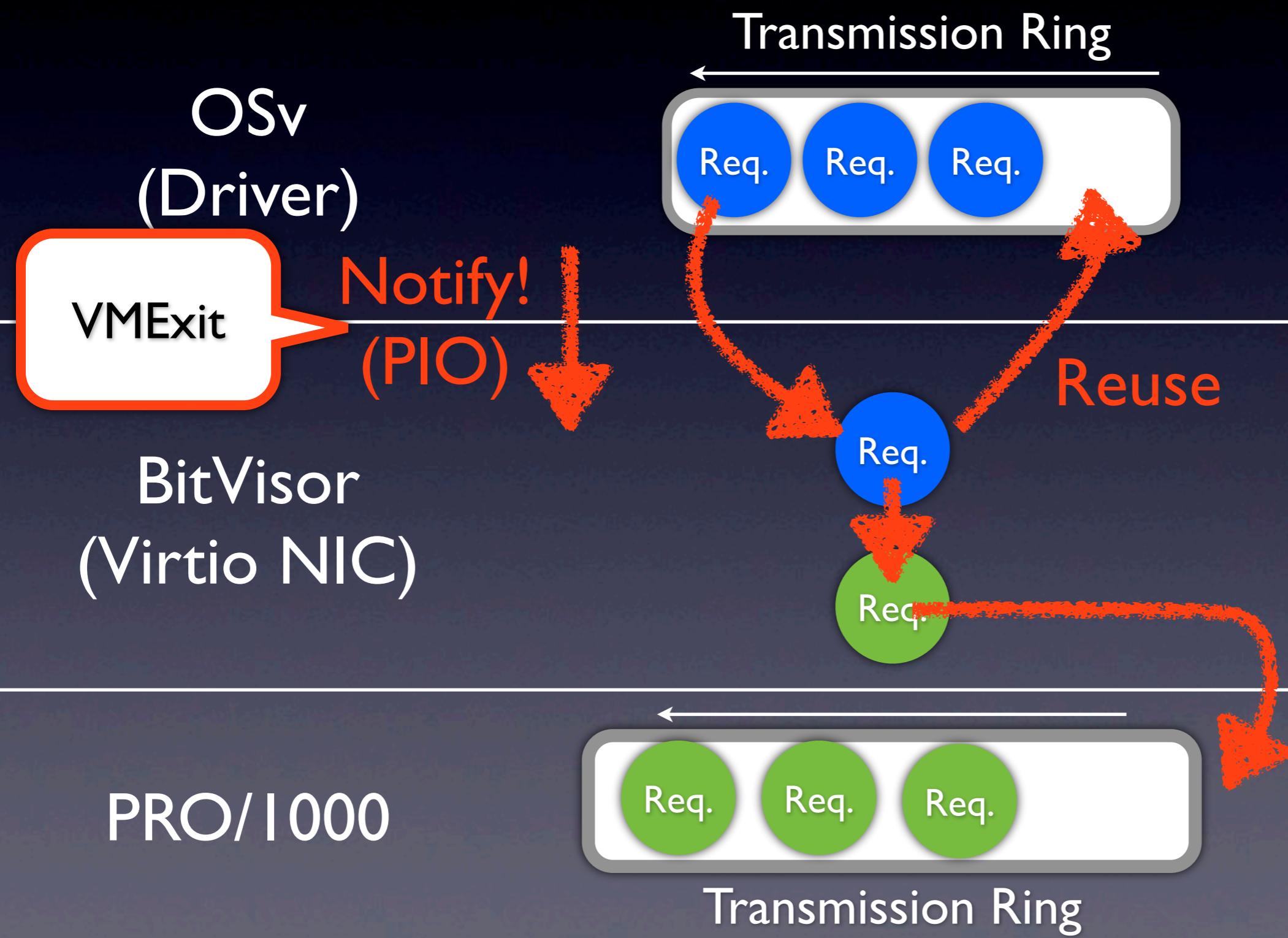
Bits	32	32	32	16		16	8	8
Read / Write	R	R+W	R+W	R	R+W	R+W	R+W	R
Purpose	Device Features bits 0:31	Driver Features bits 0:31	Queue Size	Queue Select	Queue Notify	Queue Address	Device Status	ISR Status

Memory space

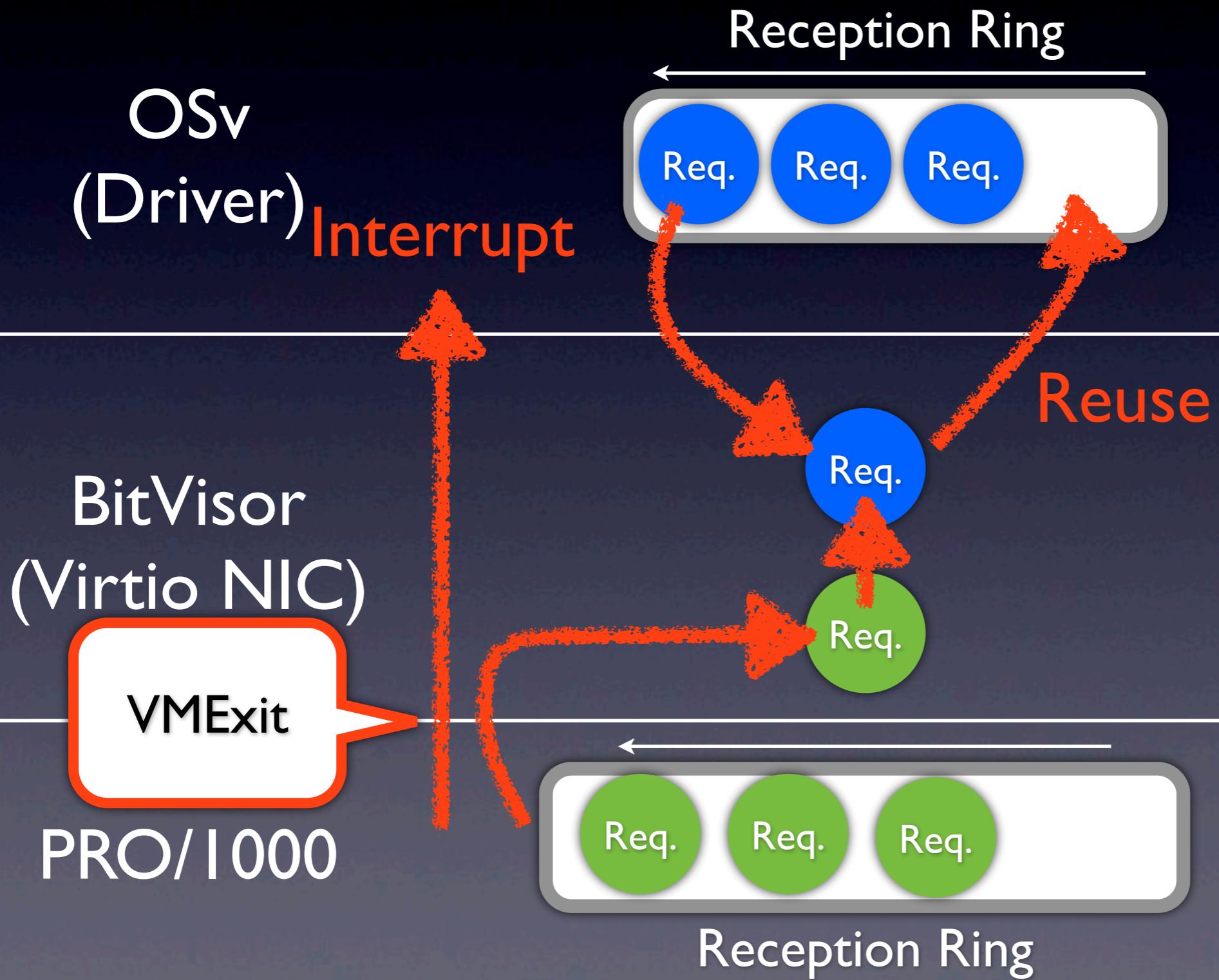
Virtio Ring (Available + Used)



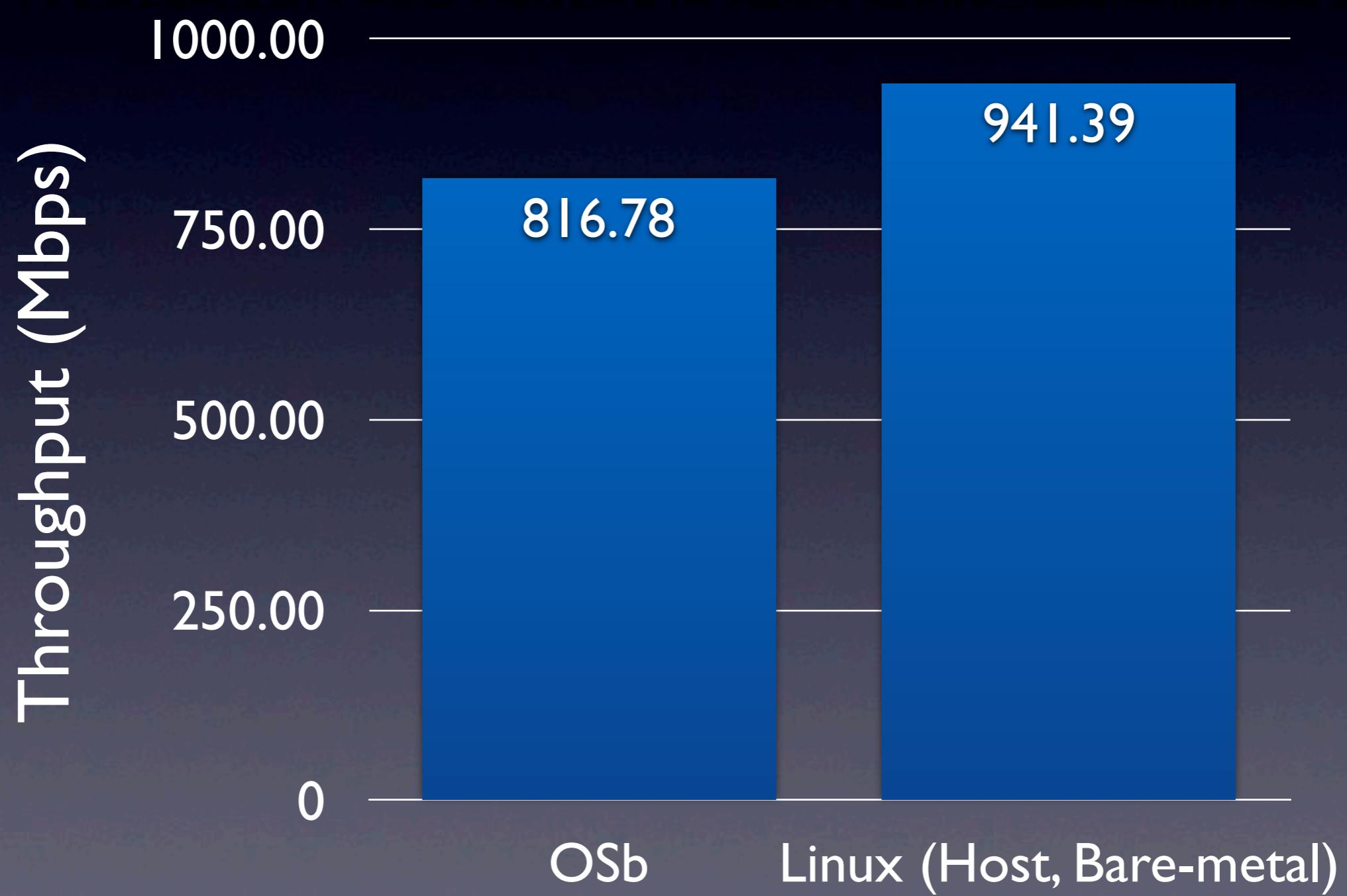
Packet Transmission



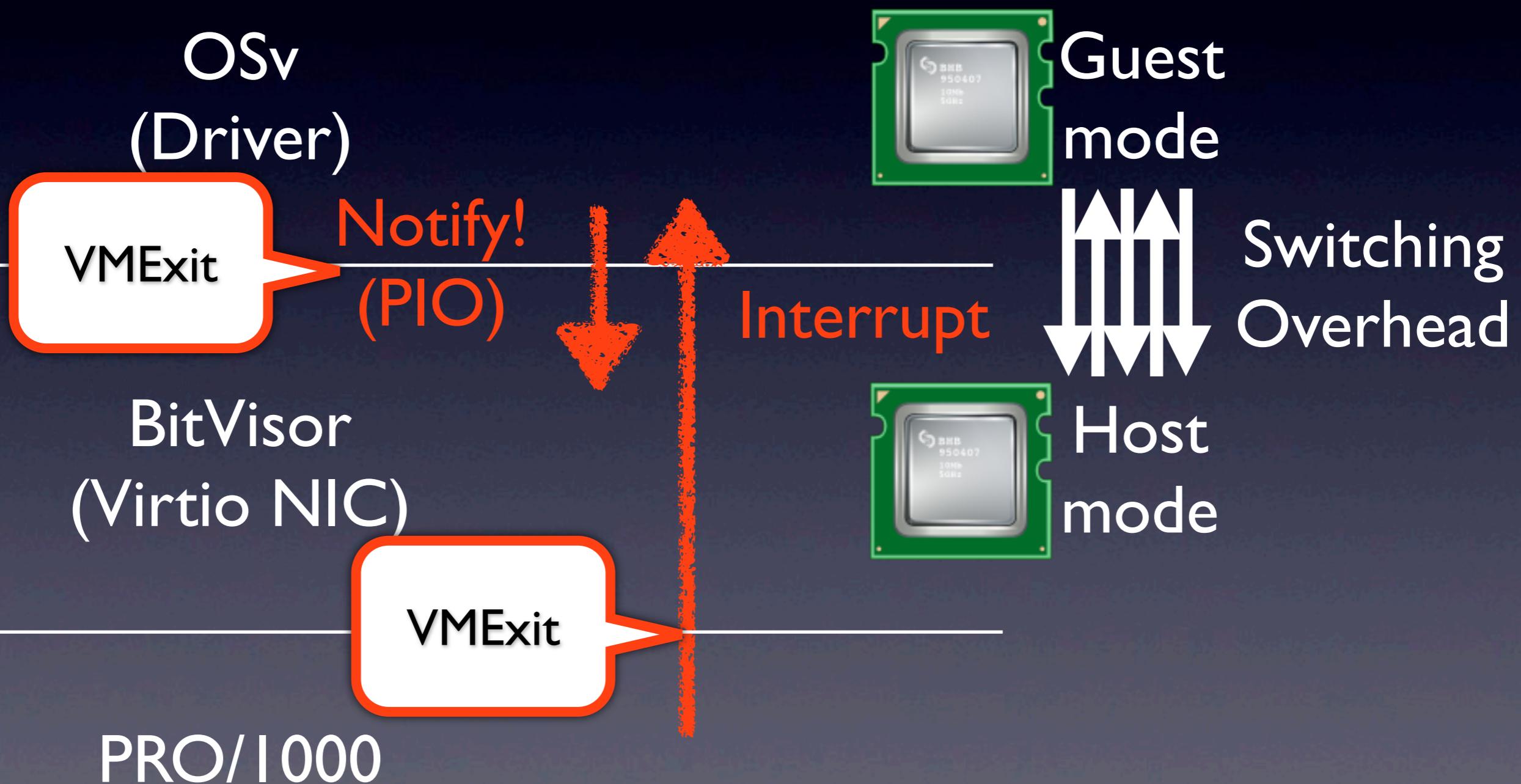
Packet Reception



Netperf TCP_STREAM

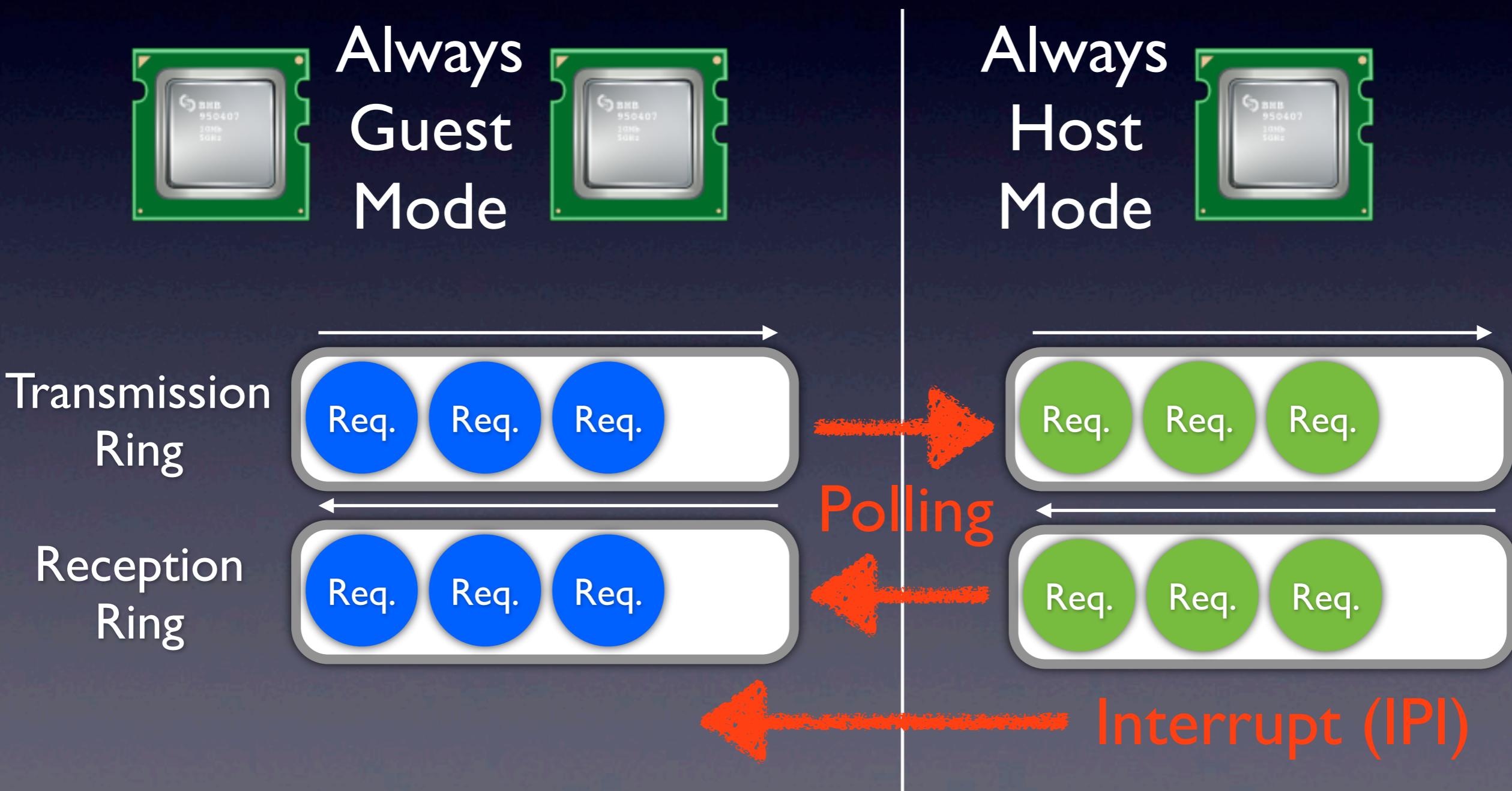


VMExits are Costly...



Exitless Virtio with Dedicated Core

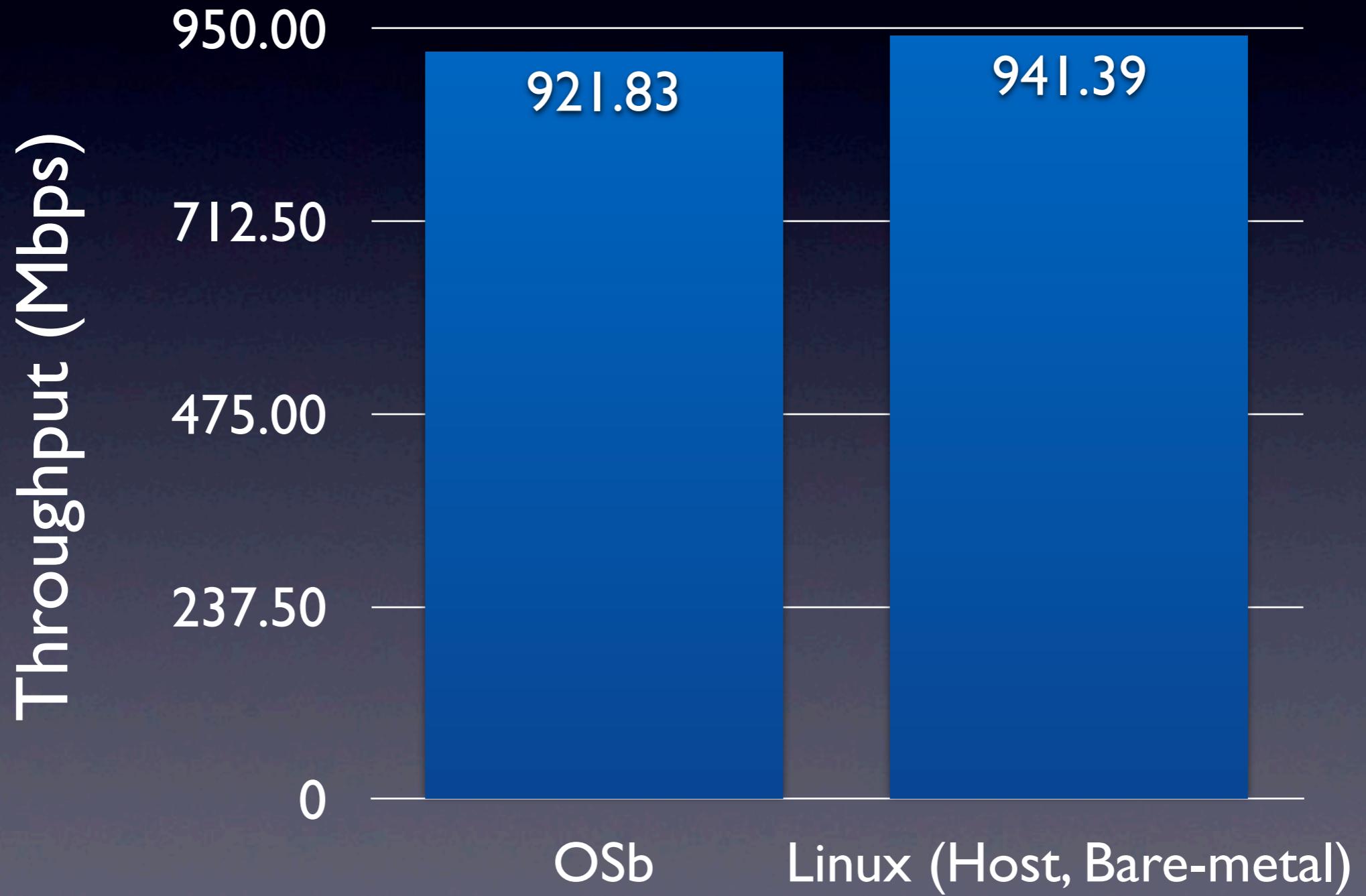
ELVIS [Har'El et al. ATC'13]



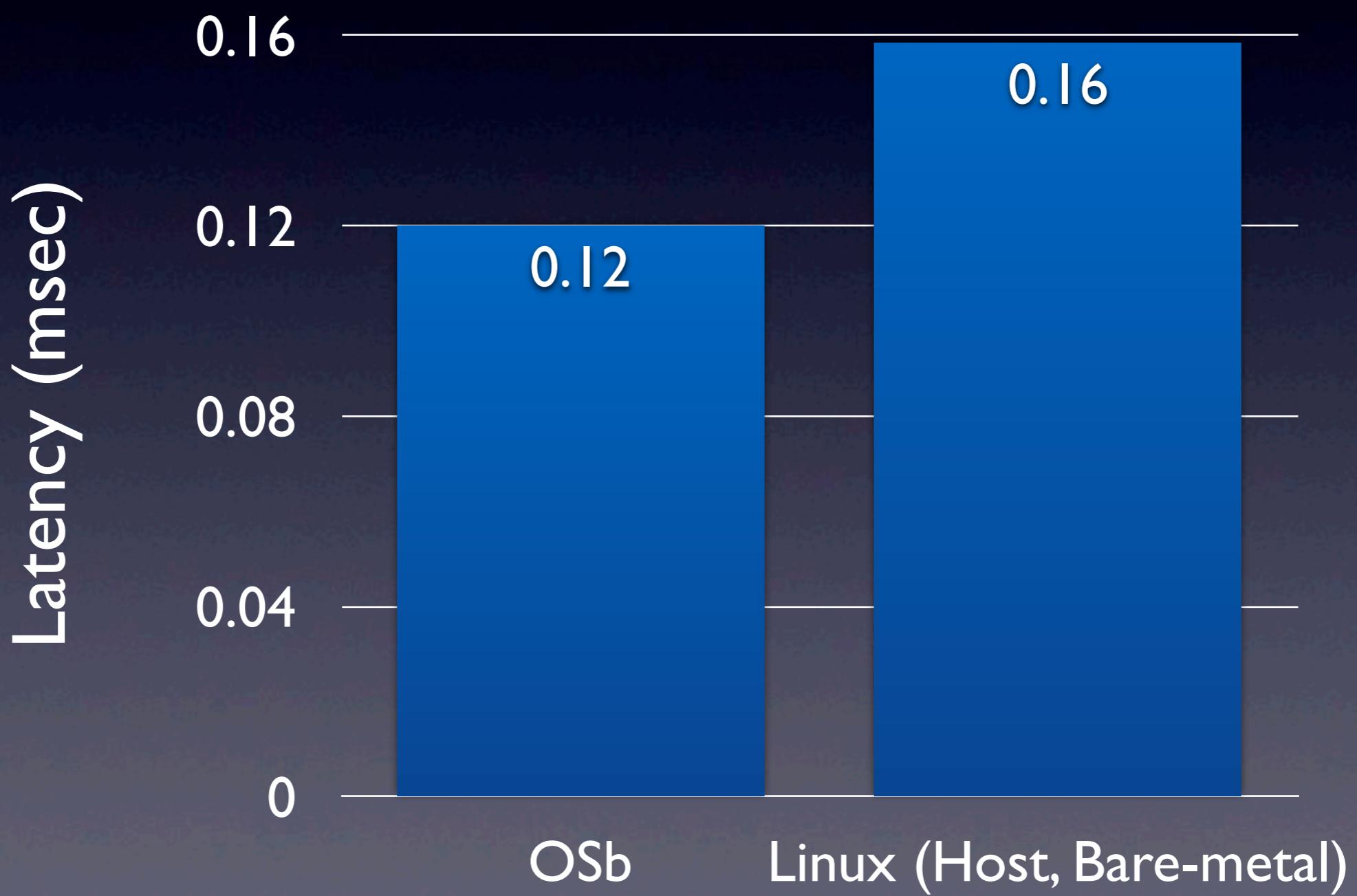
Implementation Summary

- Core Concealing
 - Modifying ACPI MADT Table
 - Notification PIO Pass-through
 - Pass-through to ineffective PIO
 - Interrupt
 - Get vector # from MSI-X table and Send IPI
 - (Or ask PRO/I000 to trigger)

Netperf TCP_STREAM (Exitless Virtio)



Ping (Exitless Virtio)



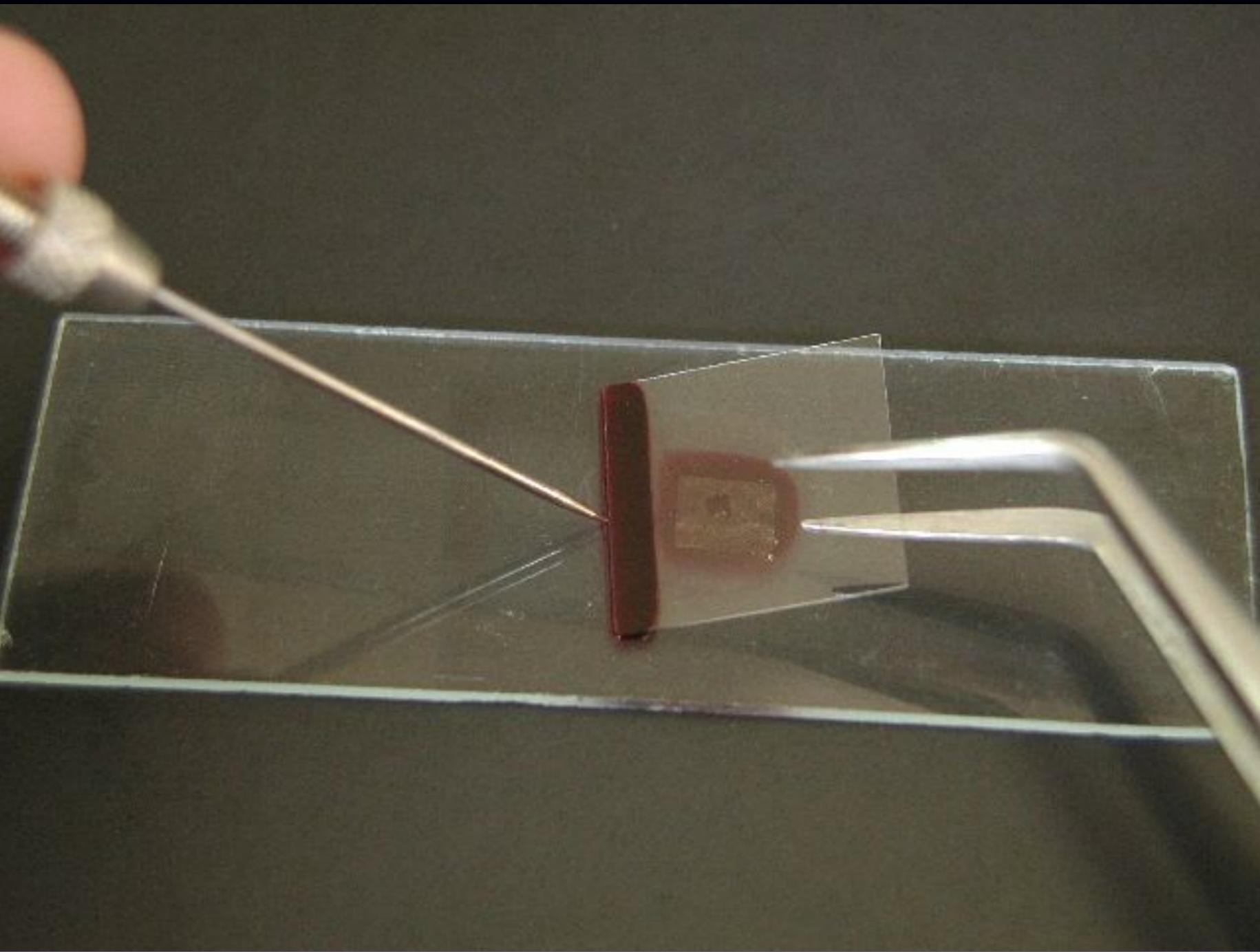
Still Under Optimization

- TODOs
 - Interrupt Moderation
 - Zero Copy/VMDq
 - Advanced PRO/1000 Descriptor
 - TCP Offloading

Summary & Future Work

- Summary
 - Running OSv on BitVisor
 - Virtio NIC + Optimization
- Future work
 - Further Optimization & Evaluation
 - Other Virtio Devices (BLK, RNG...)

Thank you !



<http://www.root.ne.jp/nishide/shs/>