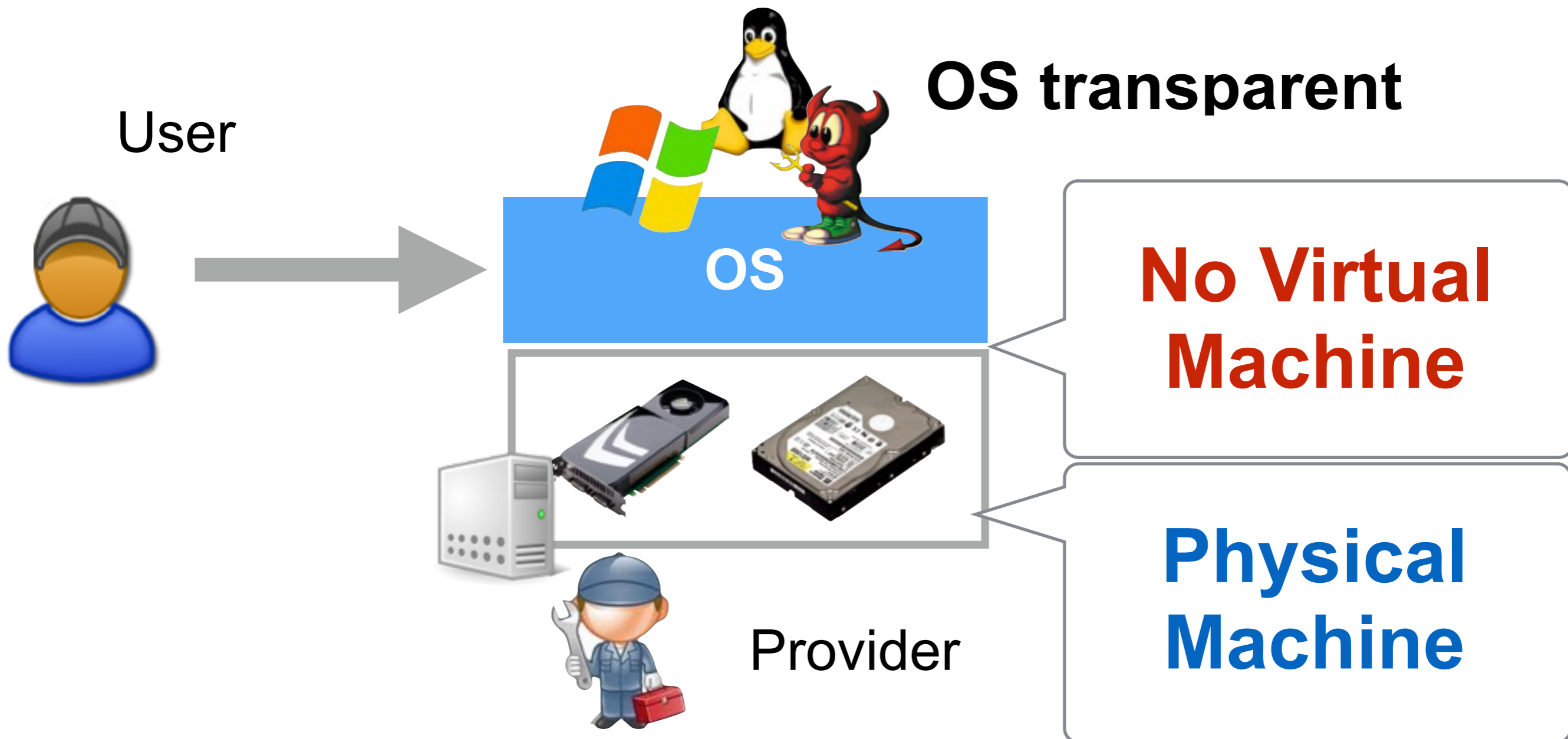# Improving Agility and Elasticity in Bare-metal Clouds

Yushi Omote[†], Takahiro Shinagawa[‡], Kazuhiko Kato[†]

[†]University of Tsukuba, [‡]The University of Tokyo
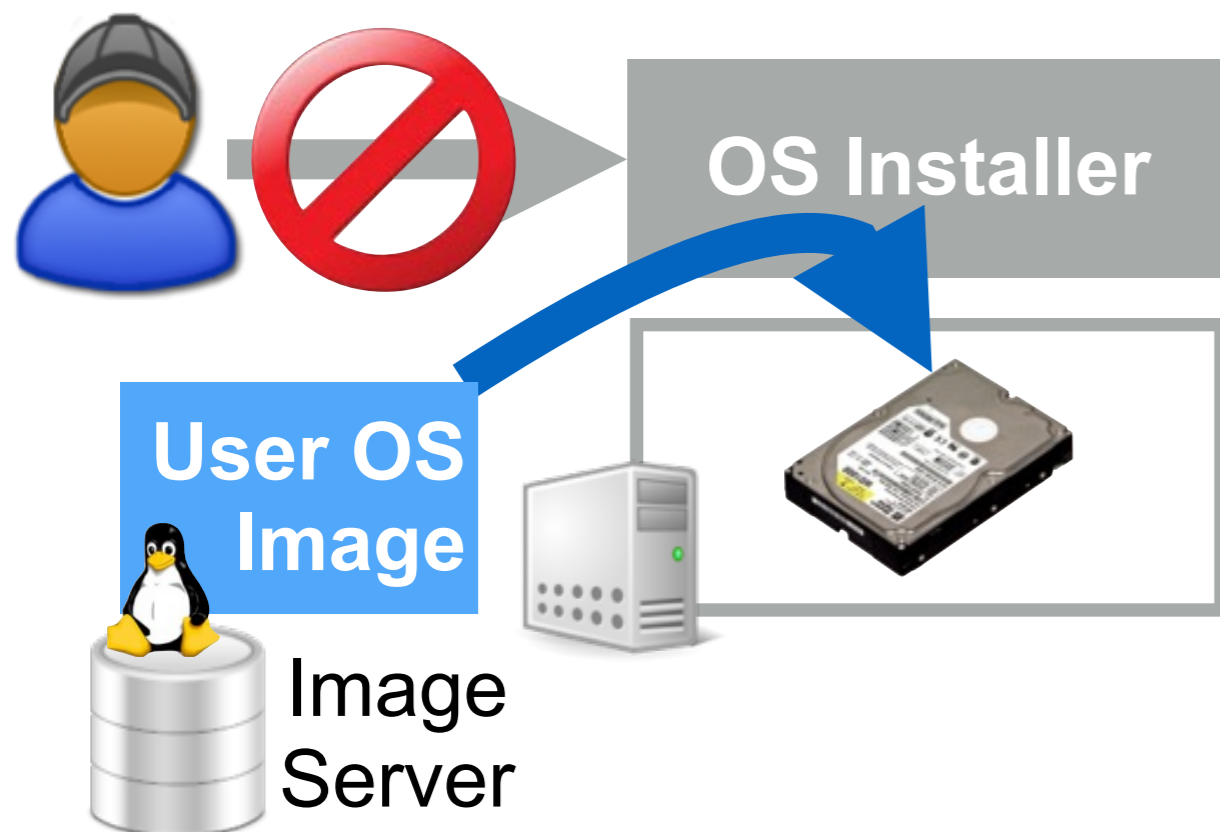
# Bare-metal Clouds
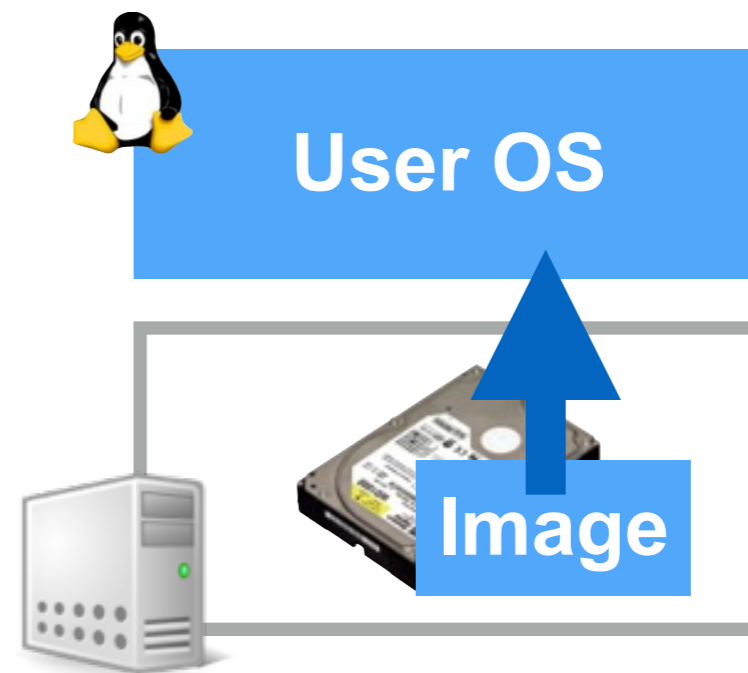
An IaaS for high performance and device functionality

**OS transparent**

User

**OS**

**No Virtual Machine**

Provider

**Physical Machine**

# OS-deployment Problem

Long wait time sacrifices agility and elasticity

## (1) Image Copy
(Tens of minutes)

**OS Installer**

**User OS Image**

Image Server

## (2) Reboot from Local Disk
(A few minutes)

**User OS**

**Image**

# Existing Approach 1
# OS Streaming Deployment

[Clerc et al. IPCCC'10]

**Network Boot + Background Copy**

User OS

Special Driver

Image Server

O Agility and Elasticity

O Performance

**OS-specific drivers are required.**

✗ OS transparency

# Existing Approach 2
# Conventional VMMs

[VMware'01, Xen'03, KVM'07]

**Streaming deployment with VMMs**

User OS

VMM

Image Server

O Agility and Elasticity

O OS transparency

**Continuous virtualization overhead**

✗ Performance
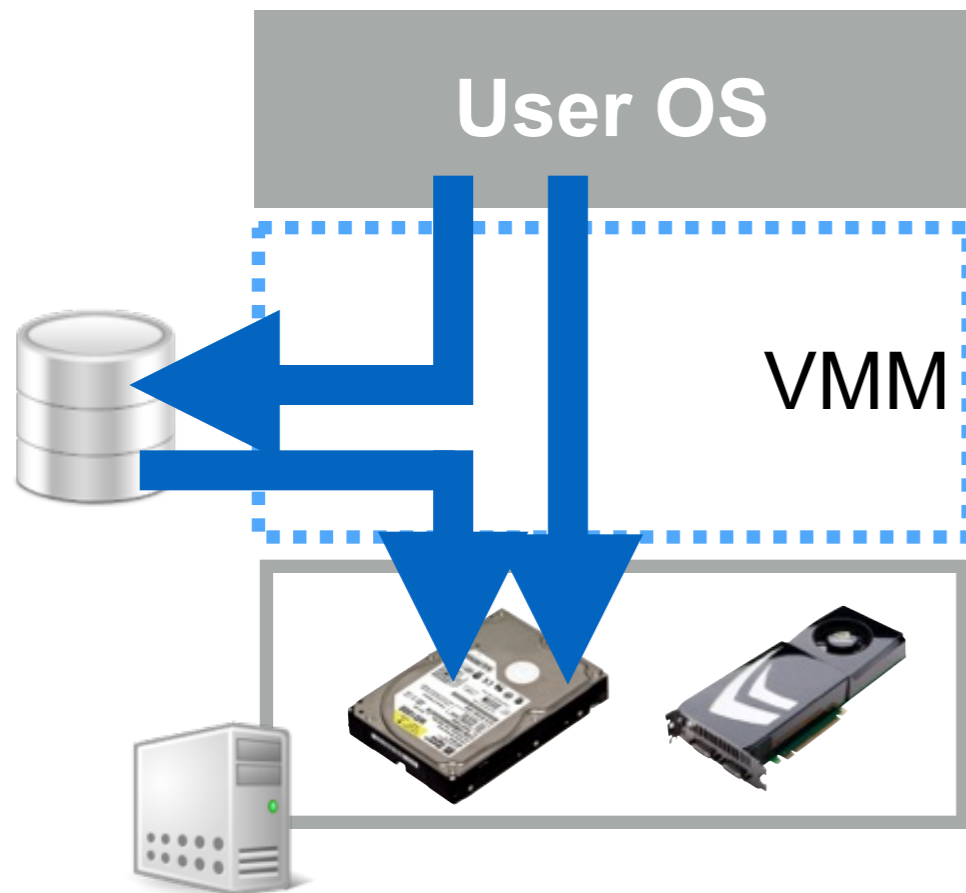
# OS Deployment
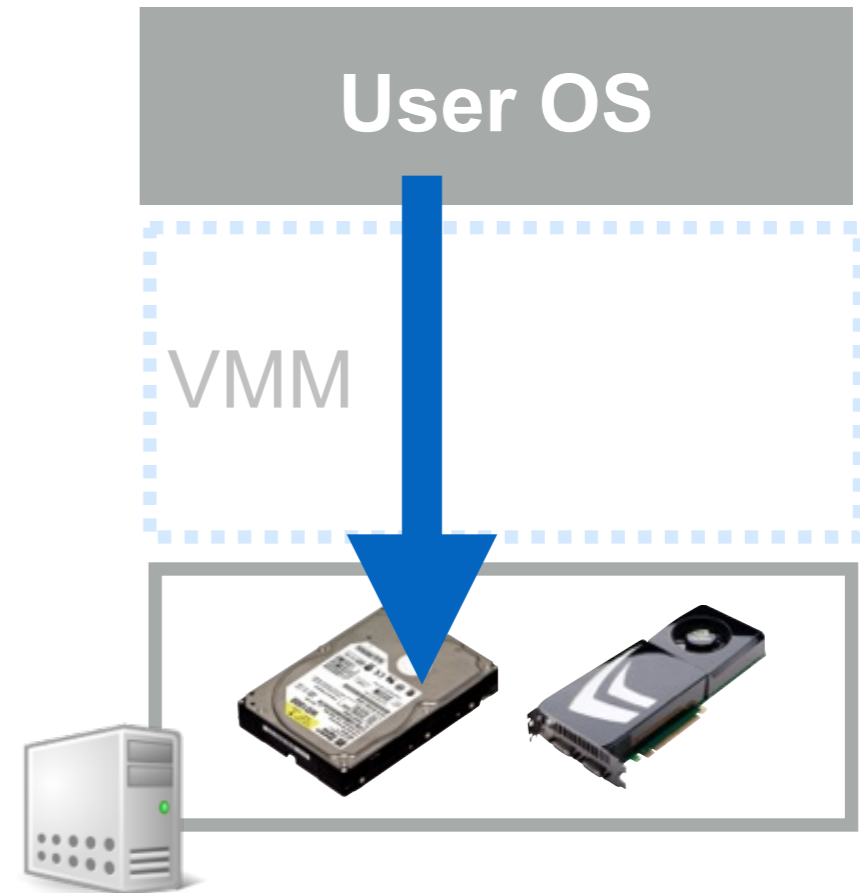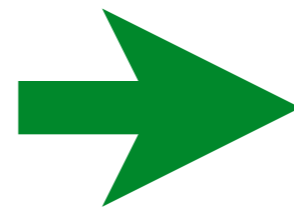# with a Special-purpose VMM

**1) Streaming deployment**

- Agility and Elasticity
- OS transparency

**2) Seamless de-virtualization**
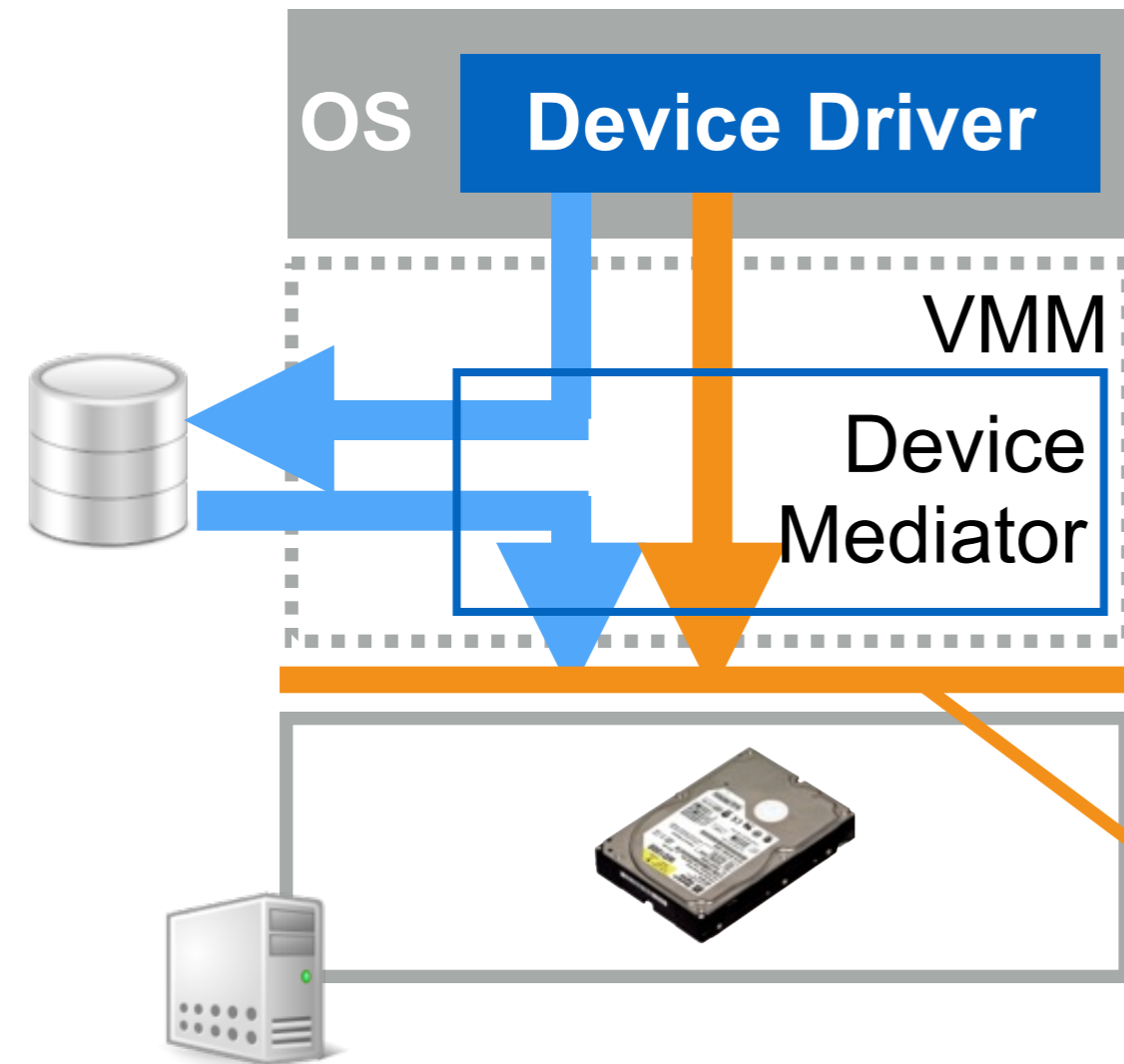
- Performance

# Challenge
# Expose & Control Physical Devices



|  | Virtual Devices? | Direct I/O? |
|---|---|---|
| **Control I/Os** | O | X |
| **Expose physical interface** | X | O |

# Device-interface-level I/O mediation

A device mediator performs:

OS **Device Driver**

VMM

Device Mediator
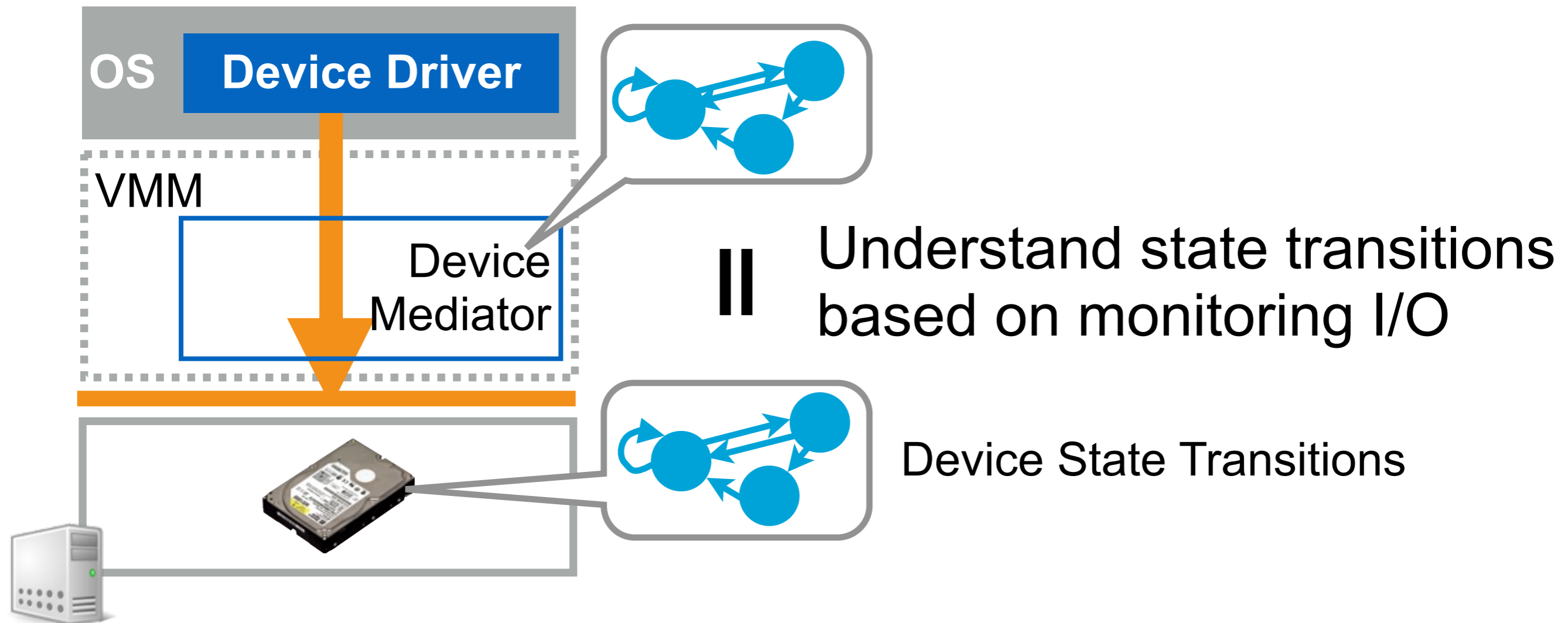
**Physical device interface**

(1) <u>I/O interpretation</u>
   to understand I/O context

(2) <u>I/O redirection</u>
   to perform network booting

(3) <u>I/O multiplexing</u>
   to perform background install

# I/O Interpretation

Determine when/how to mediate I/O requests



Understand state transitions based on monitoring I/O

Device State Transitions

# I/O Redirection

OS

Data

LBA=4
NUM=8

**(1) Interpret**

**(2) Redirect**

LBA=4
NUM=8

Image
Server

**Interrupt**

**Small
Request**

**(3) Restart**

VMM

Disk

# I/O Multiplexing

**Status Check**

**OS Request**

**(1) Request**

**VMM Request**

**Idle State**

**(2) Emulate**

**(3) Queue**

Image Server

VMM

Disk

11

# CPU/Memory Virtualization for De-virtualizable VMM

CPU

Memory

**OS**

VMM

**Guest Physical Address**
**=**
**VMM Physical Address**

## No indirection
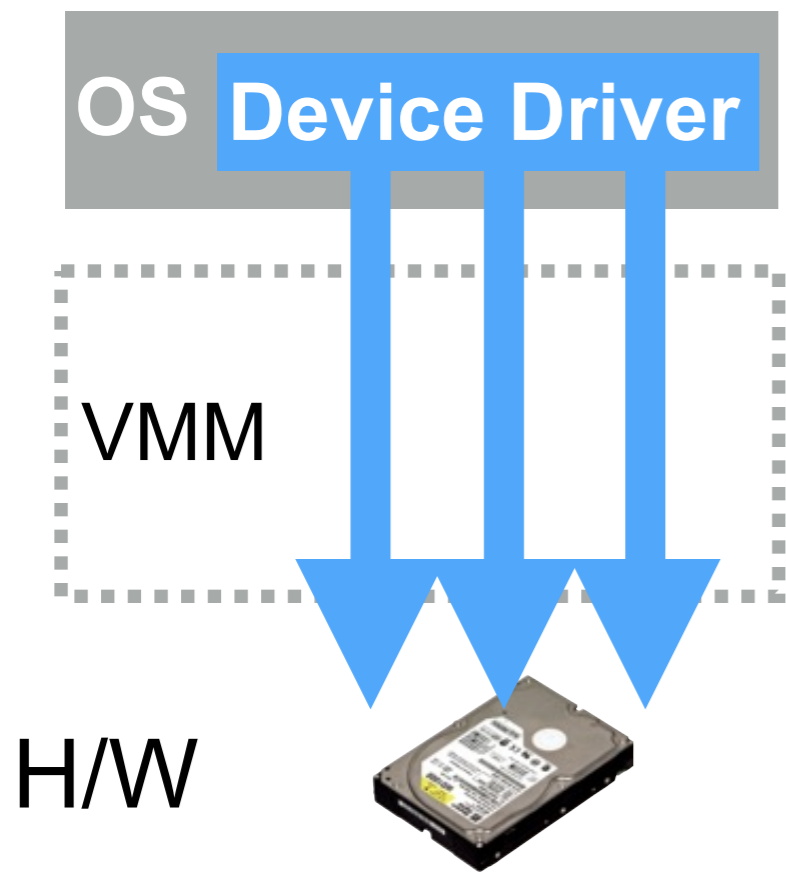
VMM runs passively with VMX

No guest scheduling

## Identity Mapping

VMM exposes physical memory

Mark VMM regions as *reserved*
(via BIOS INT15/e802)

# De-virtualization

**(1) Turns off IO VM exits**

OS **Device Driver**
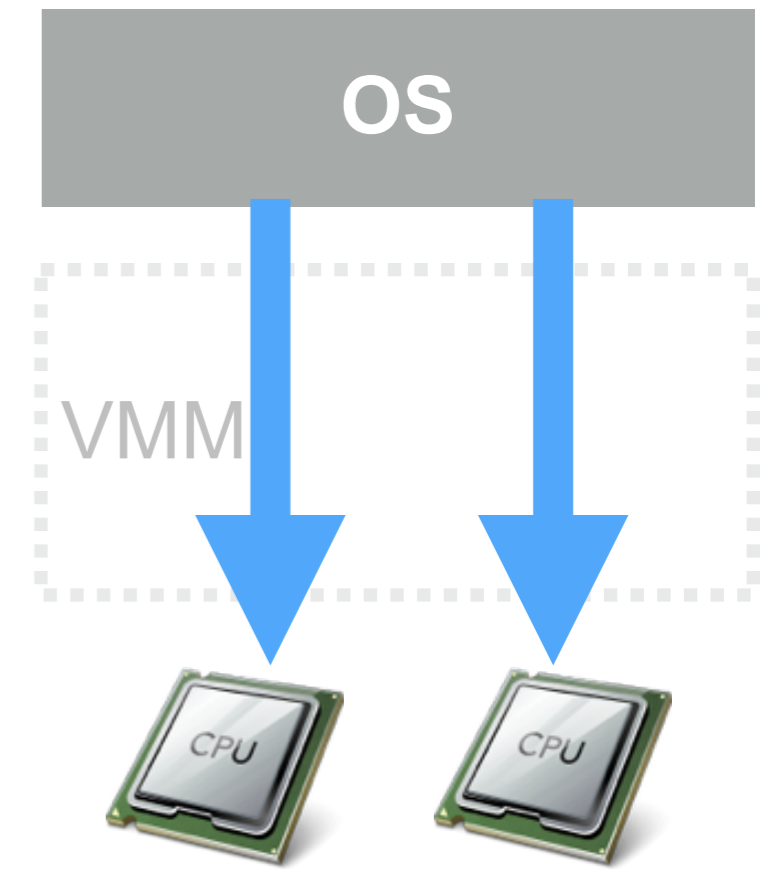
VMM

H/W

Find safe I/O timing

**(2) Turns off nested paging**

OS

VMM

Unsynchronized TLB flush

**(3) Turns off CPU virtualization**

OS

VMM

Ease VM exits condition (VMXOFF Issue)

# Performance Evaluation

- Deployed 32-GB OS Image (Ubuntu 14.04 64-bit)

  - OS-startup Time

  - Cassandra Throughput

  - Storage Throughput

  - InfiniBand Latency

A HPC Cluster

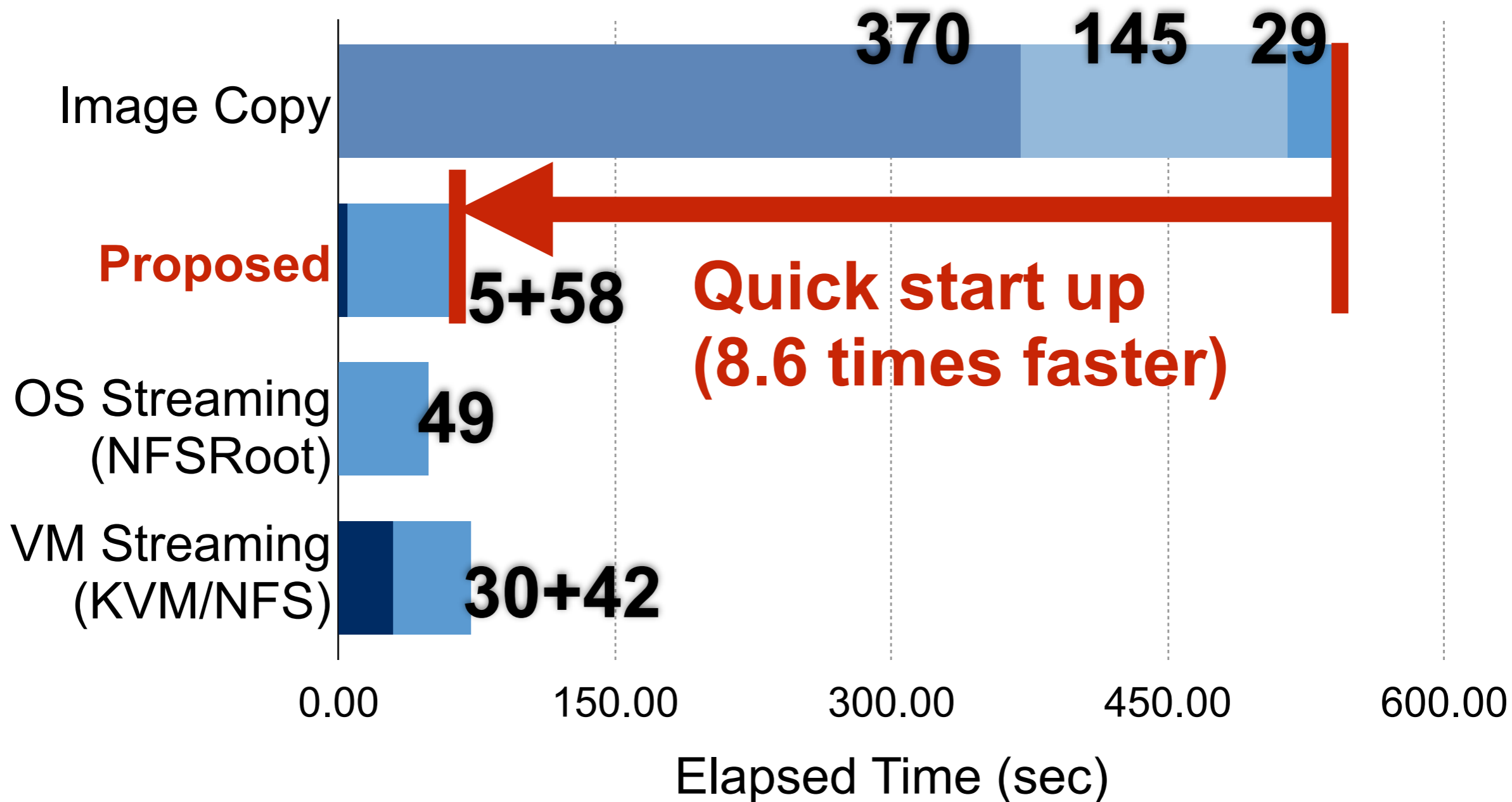| |
|---|
| Intel Xeon X5680 (3.33 GHz) / 96GB RAM |
| HDD 500GB/7200 RPM SATA |
| Mellanox InfiniBand (4X QDR) |
| Intel 82575 EM GbE Network Card |

Interconnected by
A Mellanox Grid Director InfiniBand Switch &
A FUJITSU SR- S348TC1 GbE Switch

# OS-startup Time

**Legend:** Image Copy · Reboot+Firminit. · VMM Boot · OS Boot

| Elapsed Time (sec) | |
|---|---|
| **Image Copy** | 370 · 145 · 29 |
| **Proposed** | 5+58 |
| **OS Streaming (NFSRoot)** | 49 |
| **VM Streaming (KVM/NFS)** | 30+42 |

**Quick start up (8.6 times faster)**

Elapsed Time (sec)

0.00 · 150.00 · 300.00 · 450.00 · 600.00

# Cassandra Throughput
## (Throughout Deployment)

— Proposed     — KVM (No Background Install)

**Eventual bare-metal performance**

**Seamless de-virtualization**

% of Baremetal

Elapsed Time (sec)

# Storage Throughput



Legend: Read (dark blue), Write (light blue)

**Bare-metal performance**

Y-axis: Throughput (MB/sec), values 0.00, 30.00, 60.00, 90.00, 120.00

| Category | Read | Write |
|---|---|---|
| Bare-metal | 117 | 112 |
| Deploy | 112 | 112 |
| Devirt | 112 | 115 |
| KVM/Local | 101 | 100 |

# InfiniBand RDMA latency

# Conclusion

- Improved agility and elasticity in bare-metal clouds

  - De-virtualizable VMM with streaming deployment

    - Device-interface-level I/O mediation

  - Achieved quick startup of an OS

    - 8.6 times faster than image copy

    - Preserved high performance & OS-transparency

# Future work

- Generating device mediators from specification

  - Reduce development cost of device mediators

- More advanced features of IaaS clouds

  - Live migration and checkpointing

# Thank you