

---

# To PRI or Not To PRI, That's the question

---

Yun Wang<sup>1\*</sup>, Liang Chen<sup>2\*</sup>, Jie Ji<sup>2</sup>, Xianting Tian<sup>2</sup>, Ben Luo<sup>2</sup>, Zhixiang Wei<sup>1</sup>,  
Zhibai Huang<sup>1</sup>, Kailiang Xu<sup>1</sup>, Kaihuan Peng<sup>2</sup>, Kaijie Guo<sup>2</sup>, Ning Luo<sup>2</sup>,  
Guangjian Wang<sup>2</sup>, Shengdong Dai<sup>2</sup>, Yibin Shen<sup>2</sup>, Jiesheng Wu<sup>2</sup>, Zhengwei Qi<sup>1</sup>

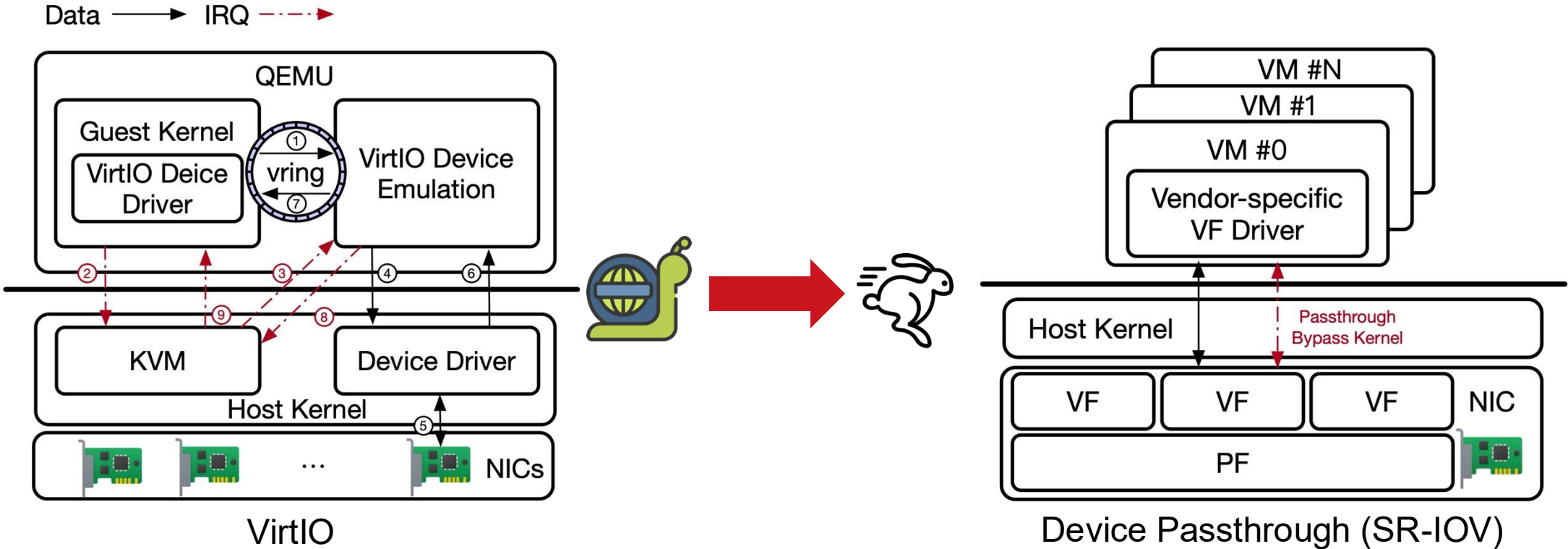


# Motivation

---

- Server memory is **expensive**.
- Cloud providers improve resource usage by taking back **unused/cold** memory pages.
- I/O device passthrough **makes it harder** to dynamically reclaim these resources.

# Evolution in I/O Virtualization



# I/O Devices Passthrough

---



The Goal: **Near-Native** Performance



The Problem: I/O Devices **Can't Handle Page Faults**

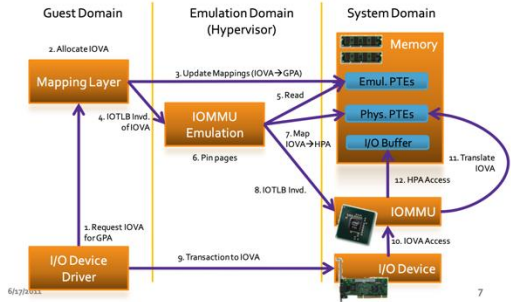


The "Official" Fix: The **Page Request Interface (PRI)**

# Existing Solutions

## Software

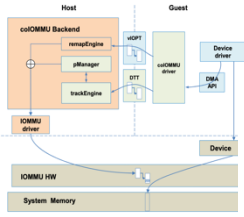
### IOMMU Emulation Flow



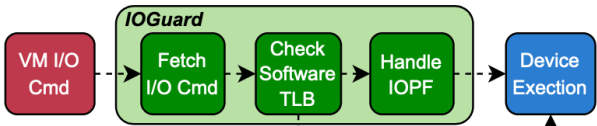
### vIOMMU

### coIOMMU Architecture

- DMA Tracking Table (DTT)
  - > Hold shared DMA buffer info.
- coIOMMU driver
  - > Hook to guest DMA API layer.
- coIOMMU backend
  - > DMA remapping engine (remapEngine)
  - > DMA tracking engine (trackEngine)
  - > Page pinning manager (pManager)



### coIOMMU



### IOGuard

## Hardware



On-Demand Paging

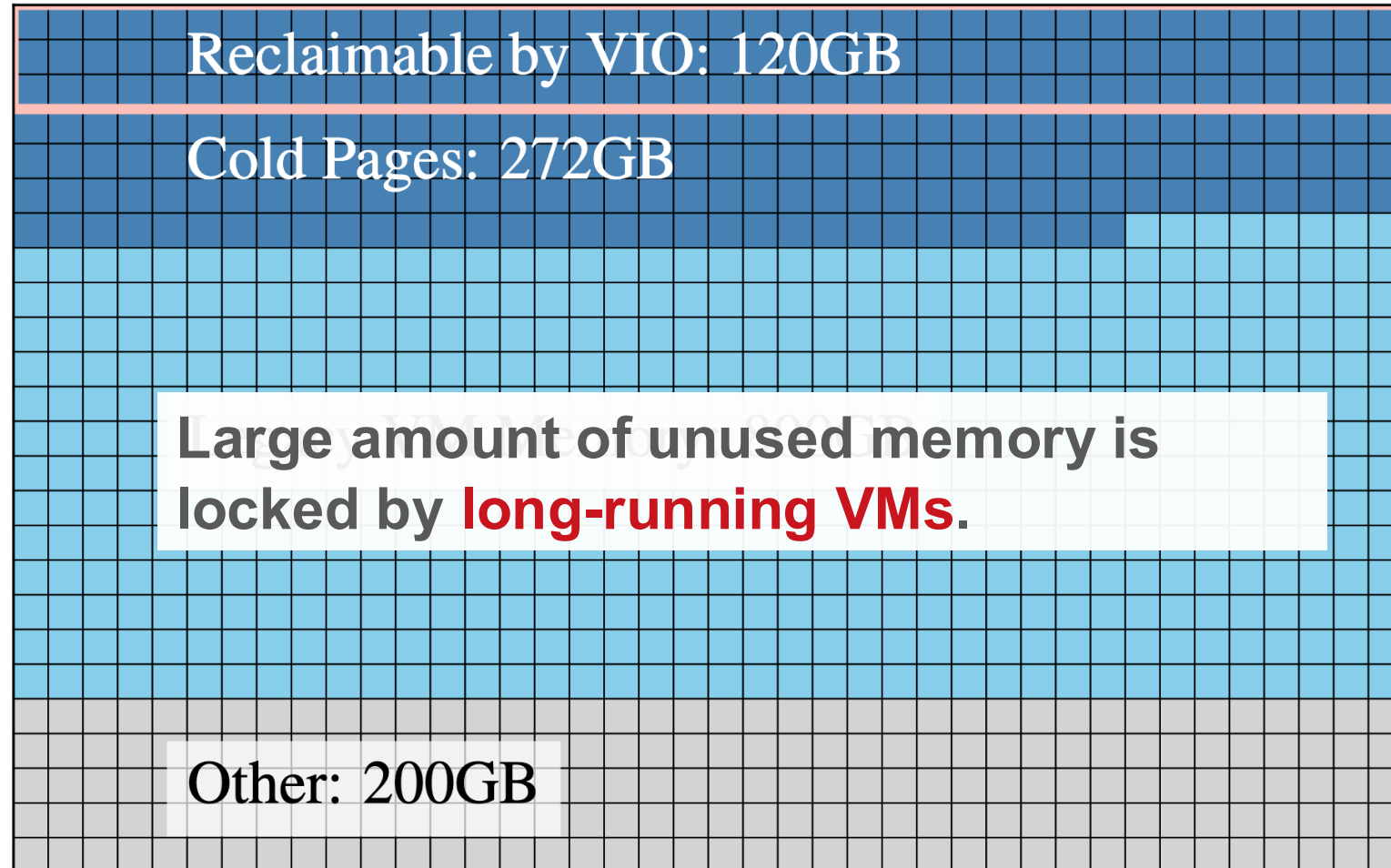


VPRI

# Existing Solutions **Fall Short**

---

Memory Usage Distribution  
(Each Square = 1GB, 1000 in total.)



# Existing Solutions **Fall Shorts**

---

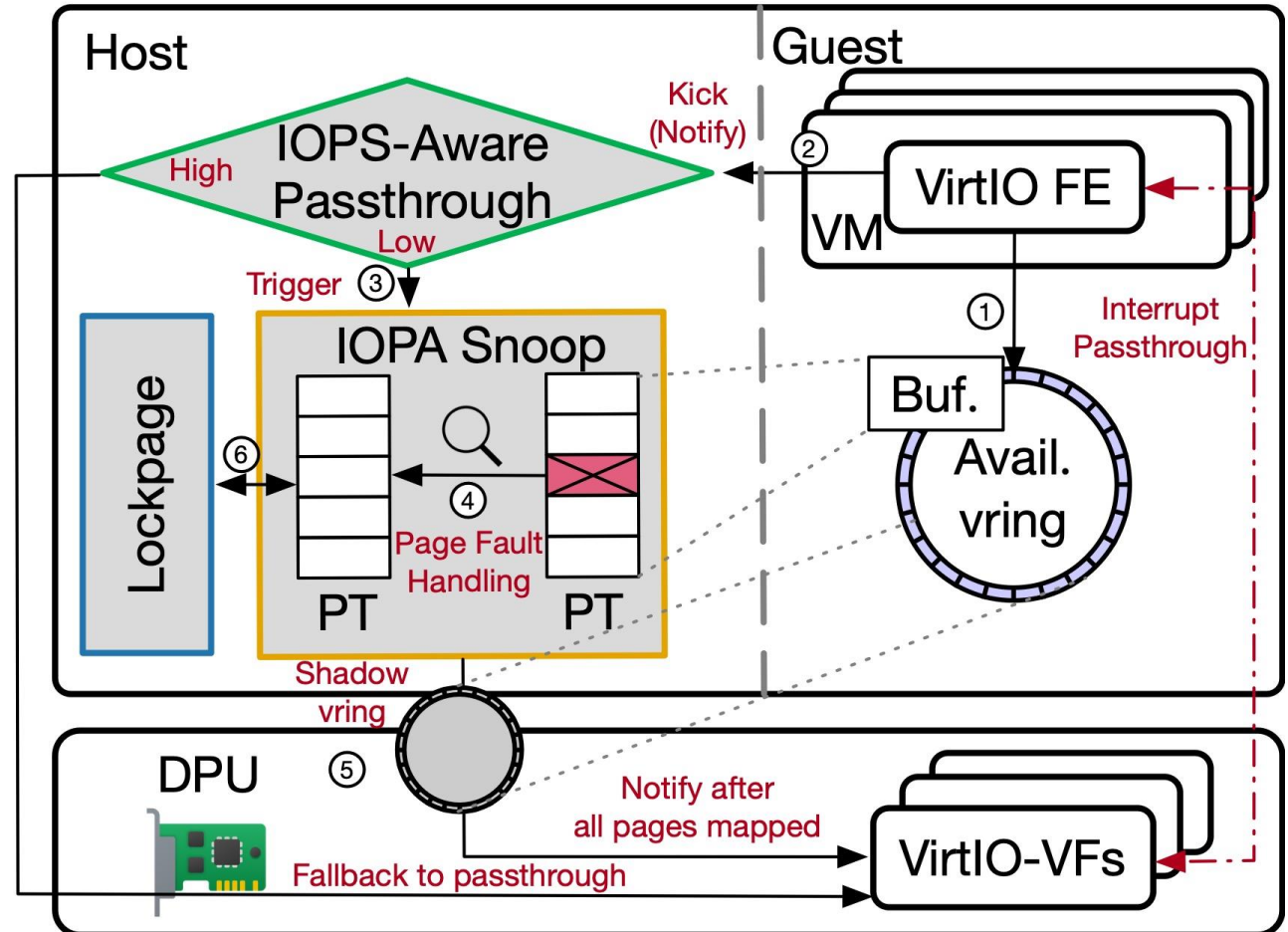
	Software	Hardware	Ours
	vIOMMU[ATC'11] coIOMMU[ATC'20] IOGuard[APSys'24]	ODP[ASPLOS'17] VPRI[SOSP'24]	<b>VIO</b>
Hardware Requirement	✓ General Hardware	✗ Proprietary Hardware	✓ General VirtIO
Guest Modifications	✗ Guest Modification	✓ No Modification	✓ No Modification
Low CPU Overhead	⚠ Non-negligible	✓ Minimal	✓ Minimal

**Limitation:** Adoption for the vast number of existing, **long-running legacy VMs** in production environments.

Can we non-intrusively handle IOPFs from legacy VMs on existing hardware, and with minimal performance overhead?

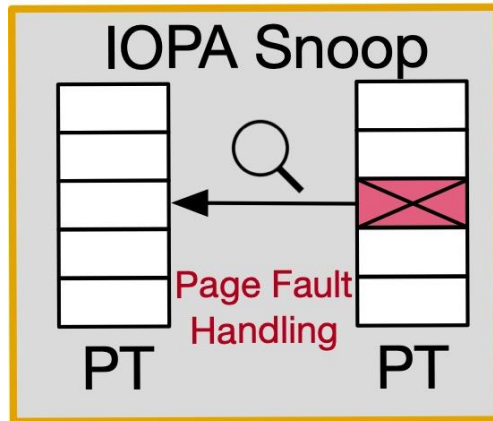
# VIO Overview

- IOPA-Snoop
- Elastic Passthrough
- Lockpage



# Challenges

---



Page Fault  
Detection



Optimal  
Performance

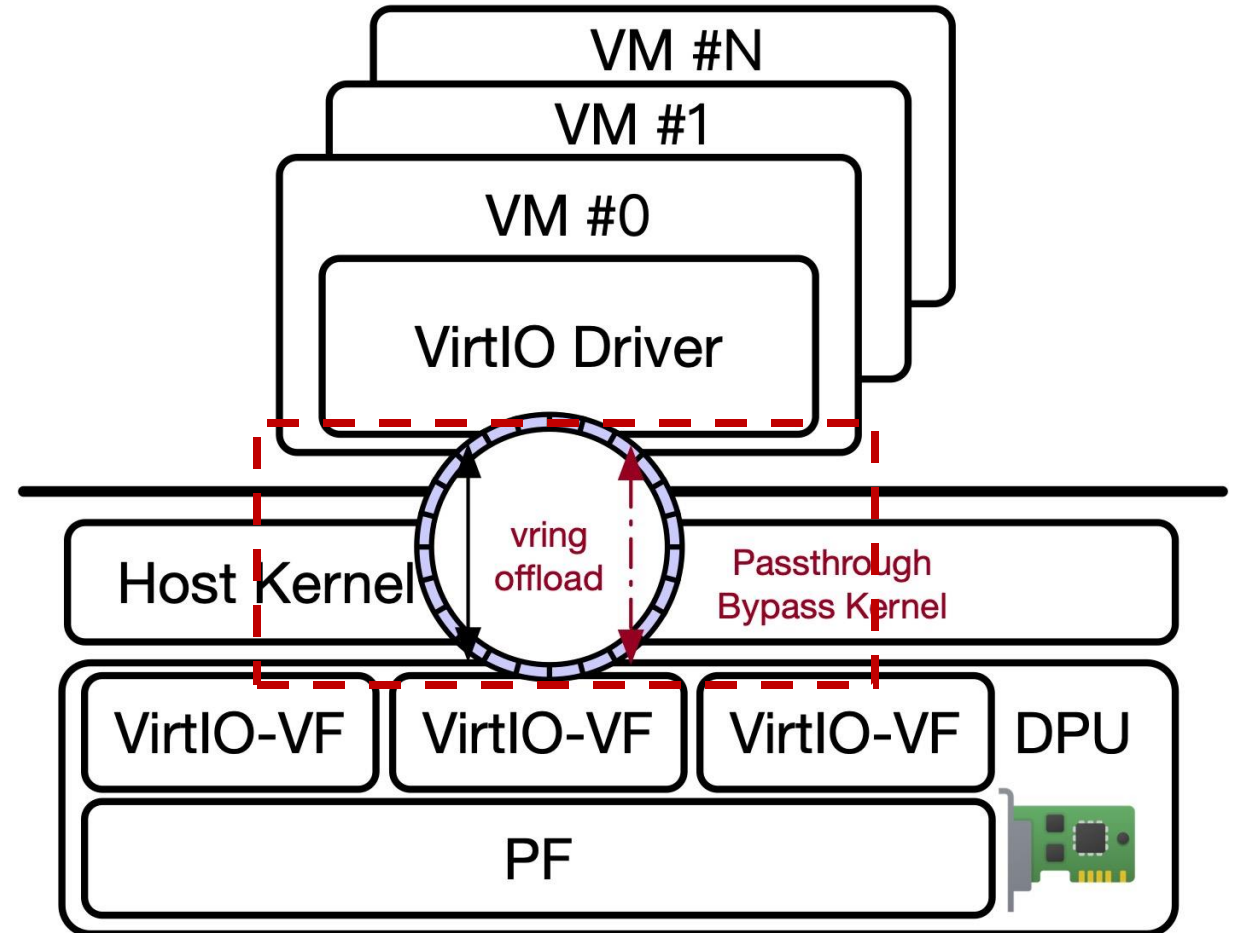


IOPFs  
Reduction

# Loose I/O Interposition

---

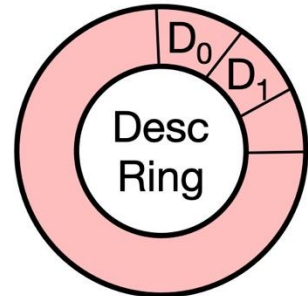
**Root Cause:** I/O operations **bypass** the hypervisor.



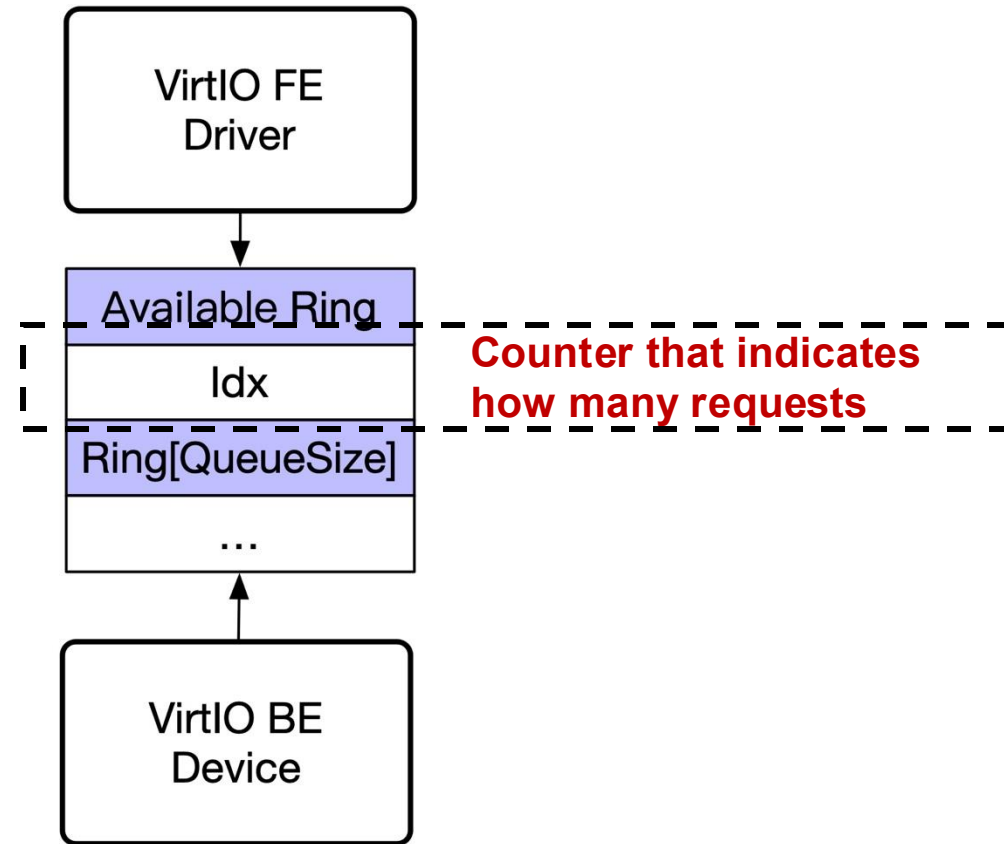
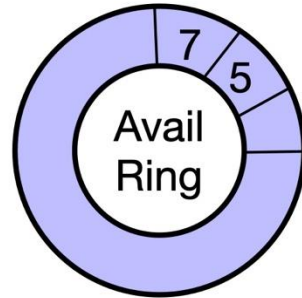
# IOPA-Snoop

---

Contains I/O request descriptors



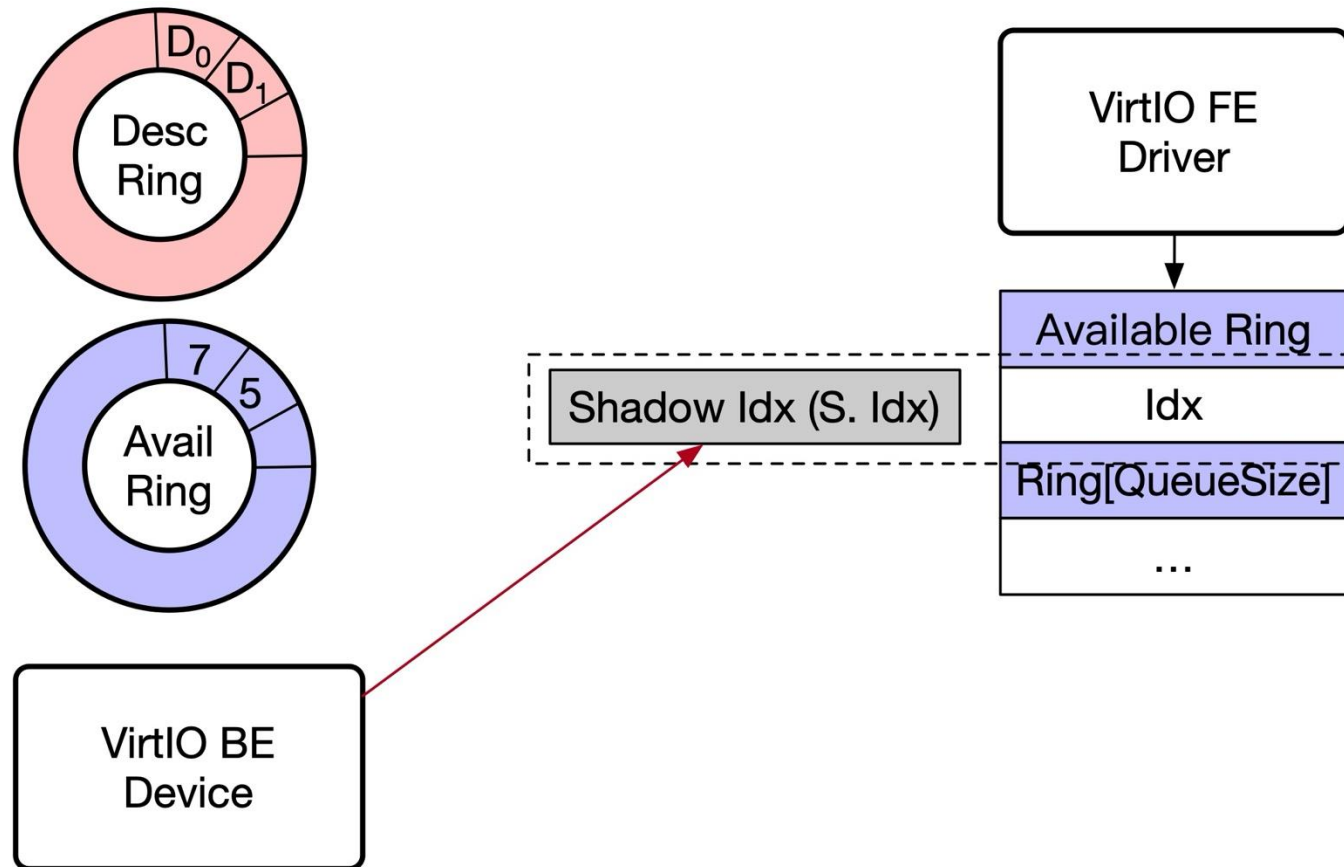
Places descriptor indices to submit requests



VirtIO Model (Simplified)

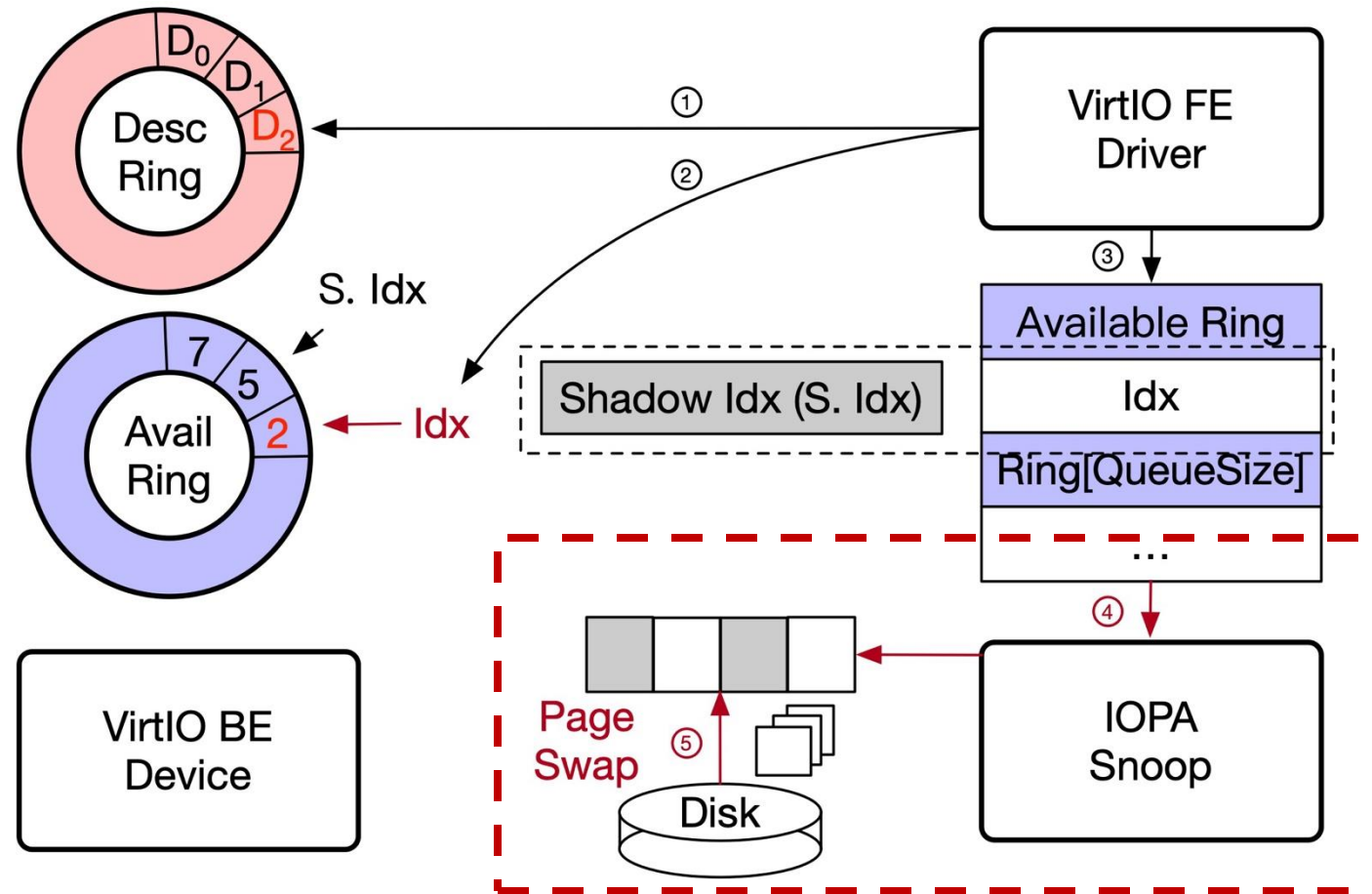
# IOPA-Snoop (cont.)

---



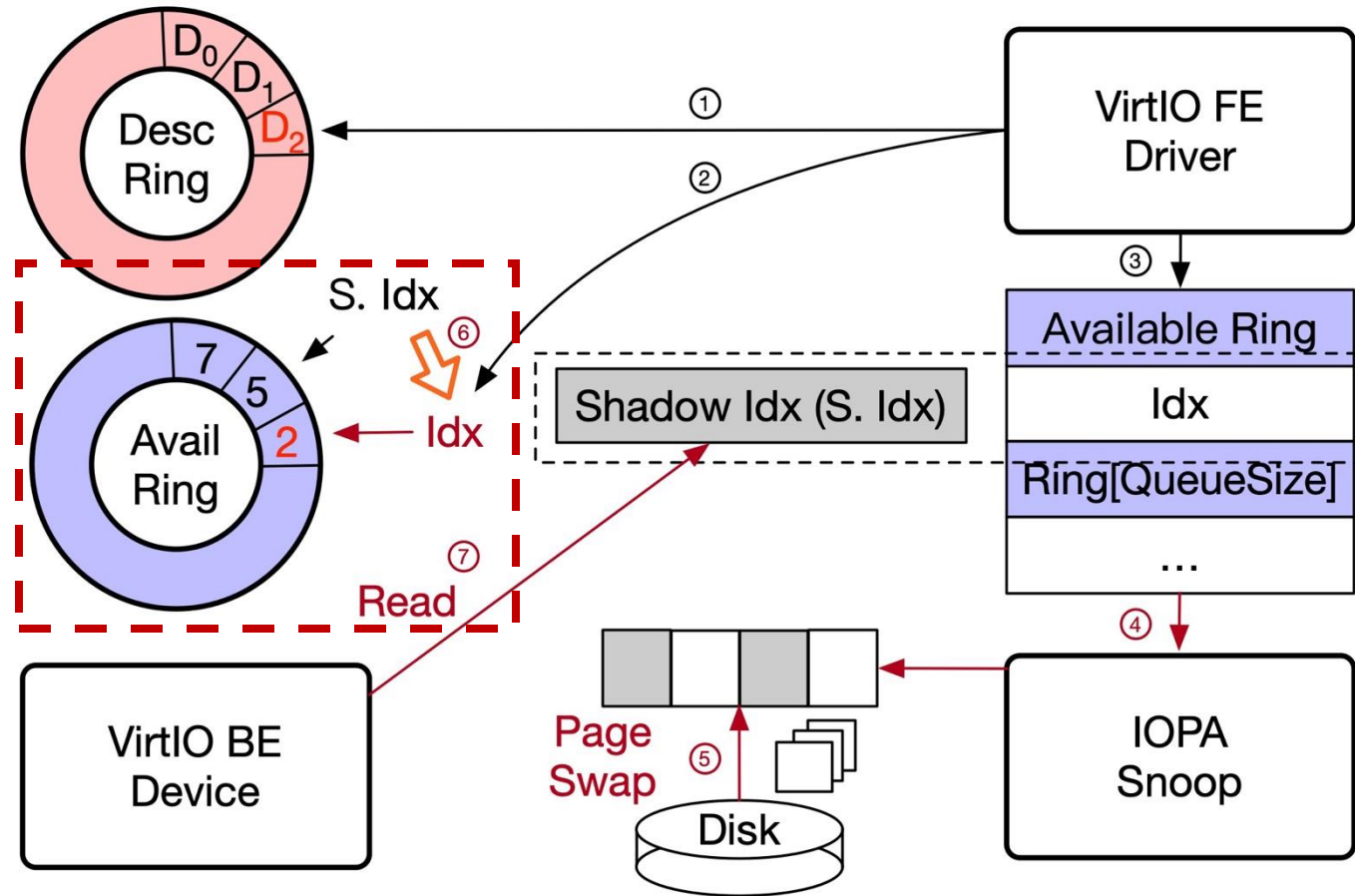
In **VIO**, device reads **Shadow Index** instead of **regular available index**.

# IOPA-Snoop (cont.)



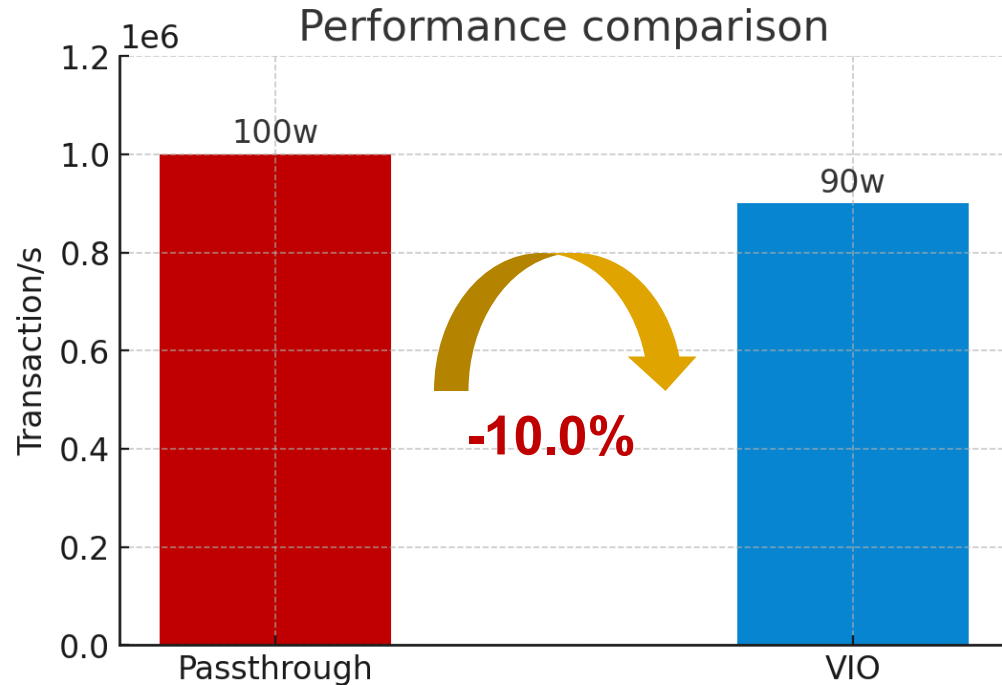
By monitoring DMA addresses in the **vring**, the hypervisor handles page faults on the CPU side.

# IOPA-Snoop (cont.)

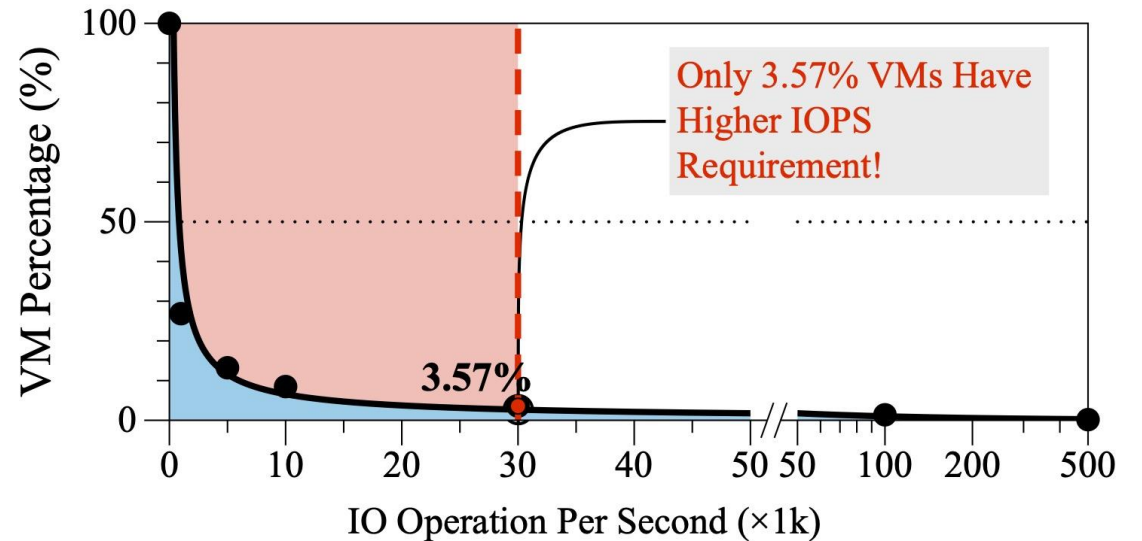


Once the hypervisor updates the **Shadow Index**, the device sees the update and starts the DMA transfer using the descriptor.

# I/O-Intensive Workloads



At high IOPS, cumulative snooping overhead causes a **~10% performance drop**.



High IOPS is rare in production, affecting **less than 4% of VMs**.

# Elastic Passthrough

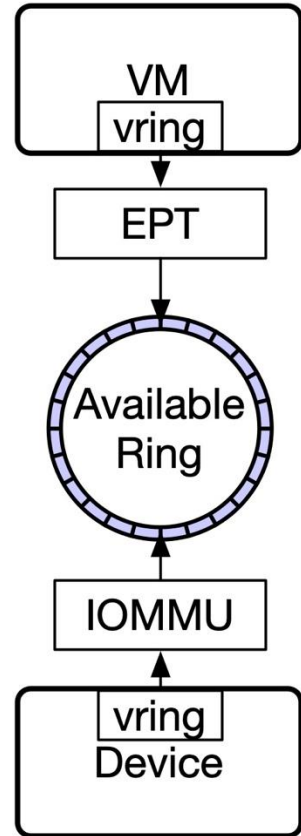
---

- IOPS-Aware Elastic Passthrough, dynamic change mode according to **I/O pressure**.
  - Low IOPS (The Common Case): The VM runs in **VIO snooping mode**.
  - High IOPS (The Rare Case): The VM is automatically switched to **native passthrough mode**.

# Elastic Passthrough (cont.)

---

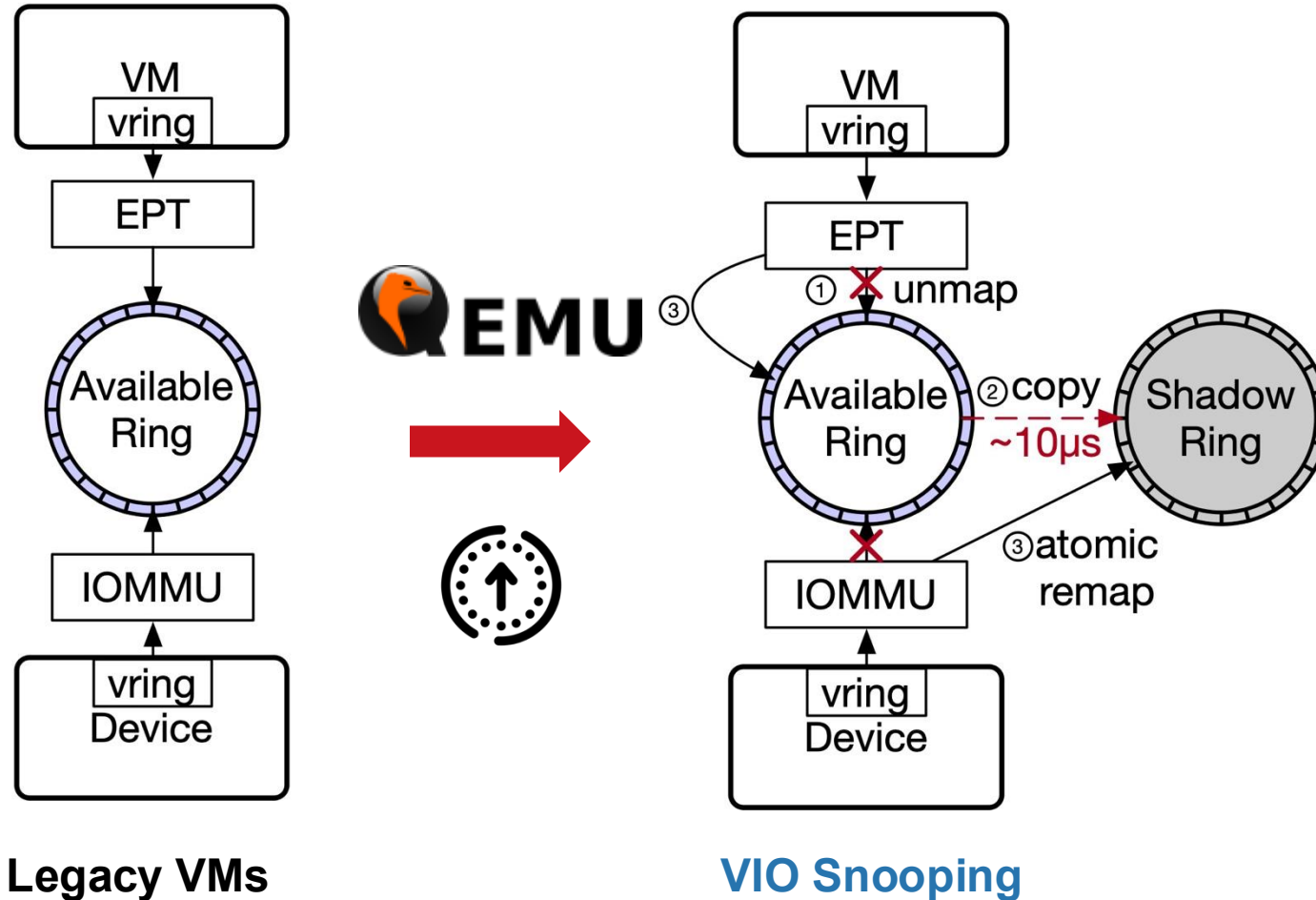
- Passthrough



Legacy VMs

# Elastic Passthrough (cont.)

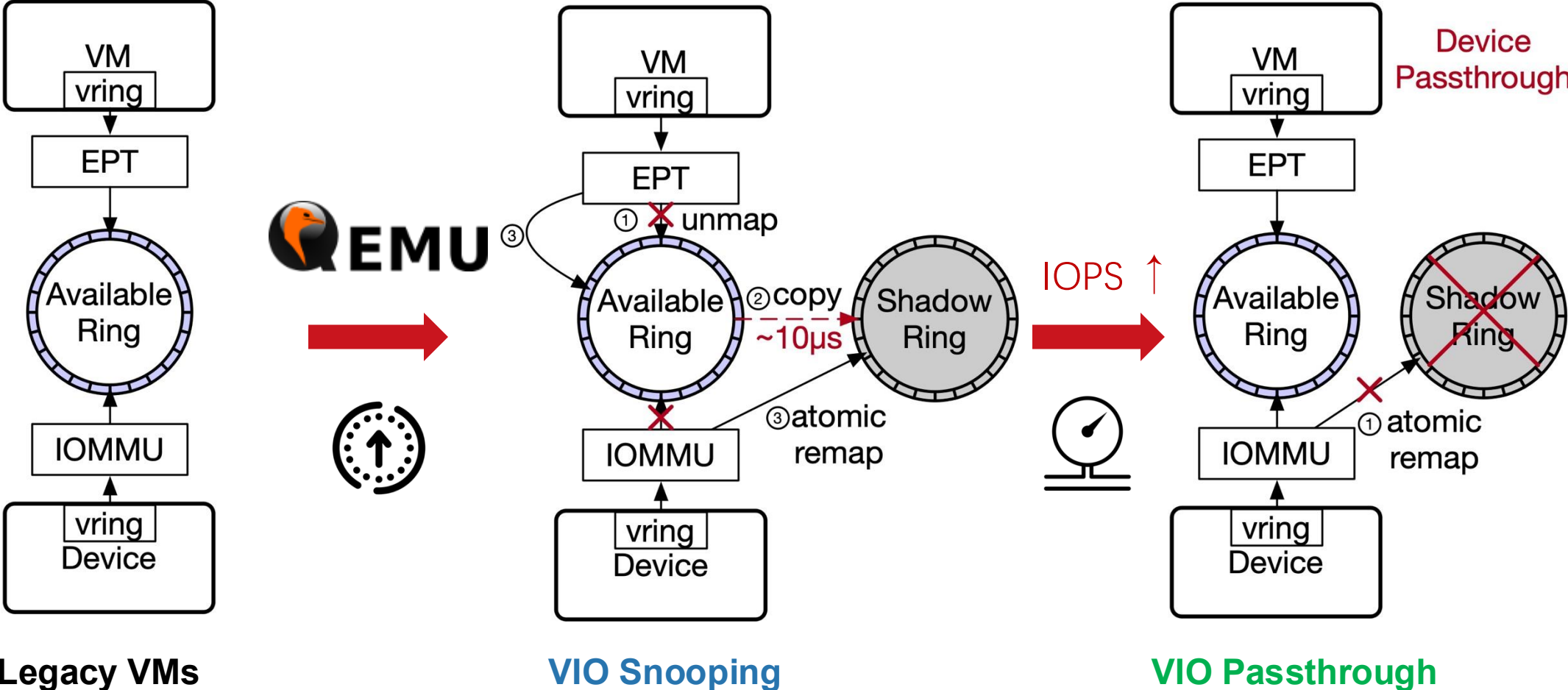
- Passthrough → Snooping (Live Upgrade)\*



[1]Zhang, Xiantao, et al. "Fast and scalable VMM live upgrade in large cloud infrastructure." *Proceedings of the Twenty-Fourth International Conference on Architectural Support for Programming Languages and Operating Systems*. 2019.

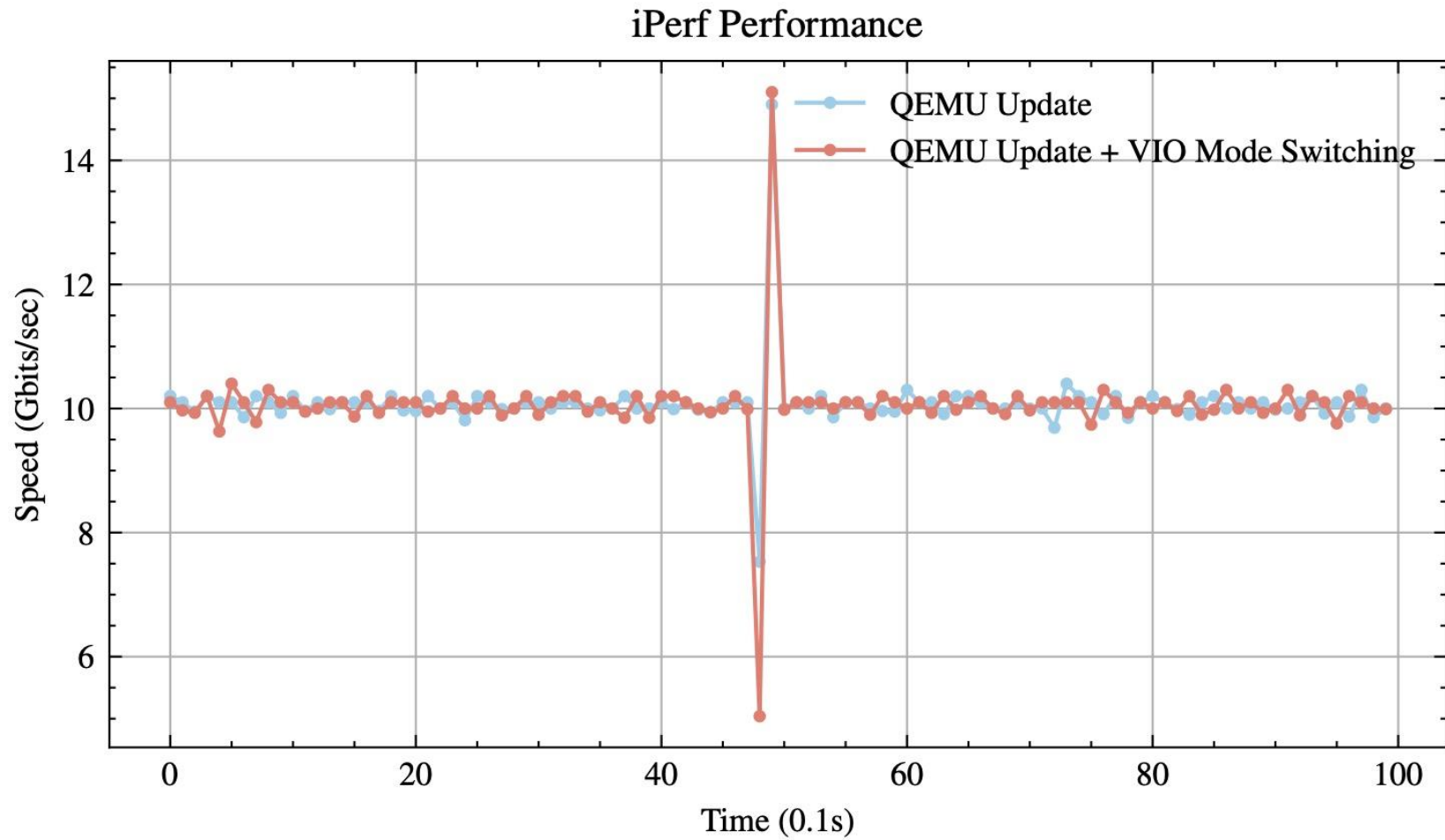
# Elastic Passthrough (cont.)

- Snooping → Passthrough (IOPS increases)



# Elastic Passthrough (cont.)

---

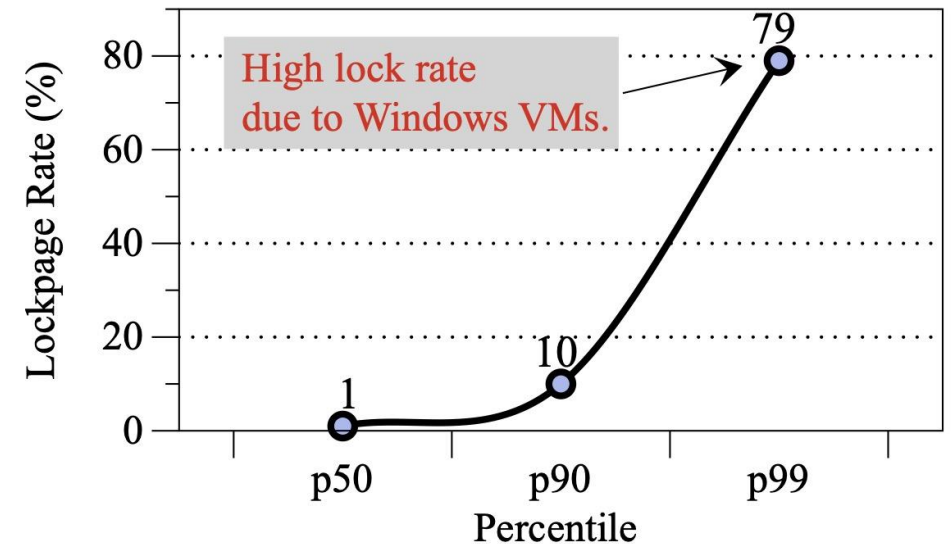


VIO rollout in the Alibaba production

# IOPFs Reduction - Lockpage

---

- Lockpage mechanism
  - Uses **IOPA-snoop** to log page access patterns and identify these **frequently used I/O pages**.
  - Lockpage mechanism **"locks" or "pins" hot pages** in memory.

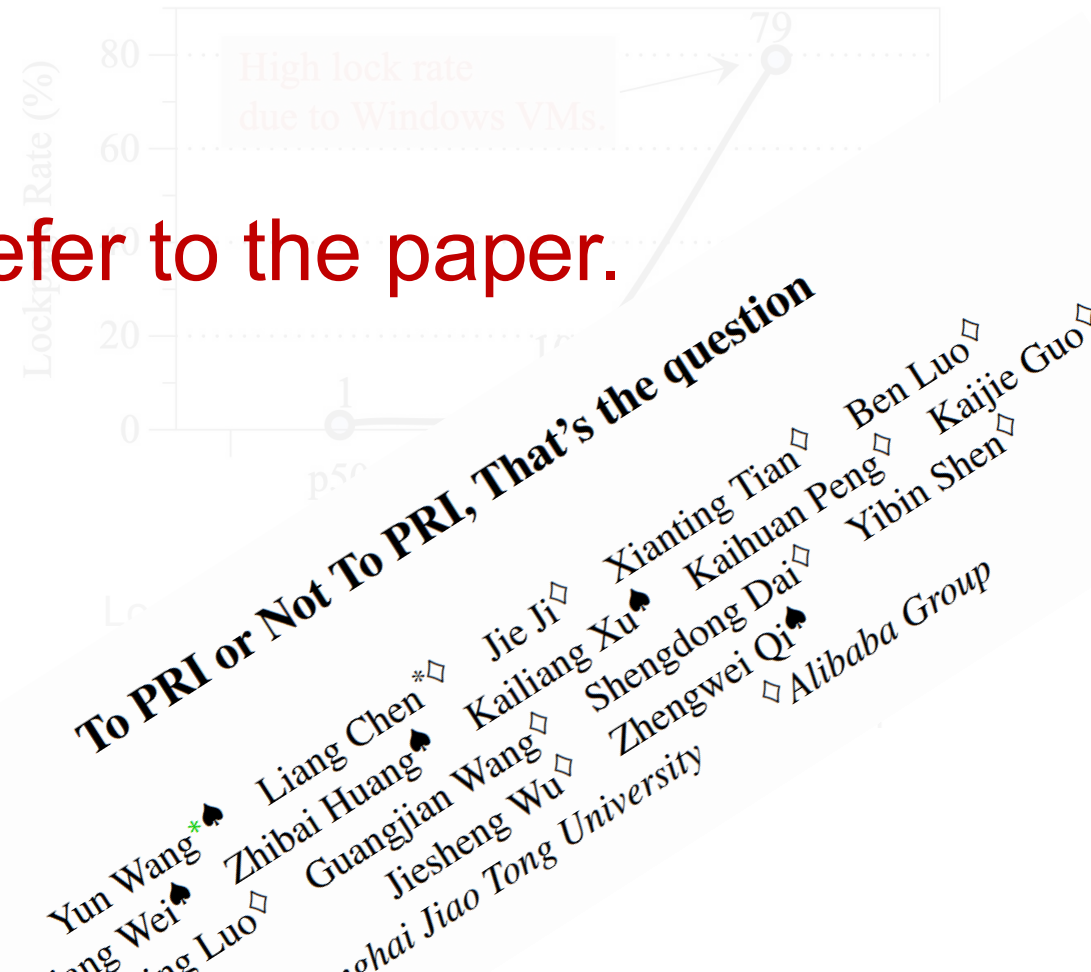


Lockpage rate in a production environment, showing **P50 at 1%**, **P90 at 10%**, and **P99 spiking to 79%** due to poor I/O locality in Windows VirtIO driver.

# IOPFs Reduction - Lockpage

- Lockpage mechanism
  - Uses IOPA-snoop to log page access patterns and identify the most frequently used I/O pages.
  - Lockpage mechanism "locks" or "unlocks" hot pages in memory.

For more details, please refer to the paper.



# Experimental Setup

---

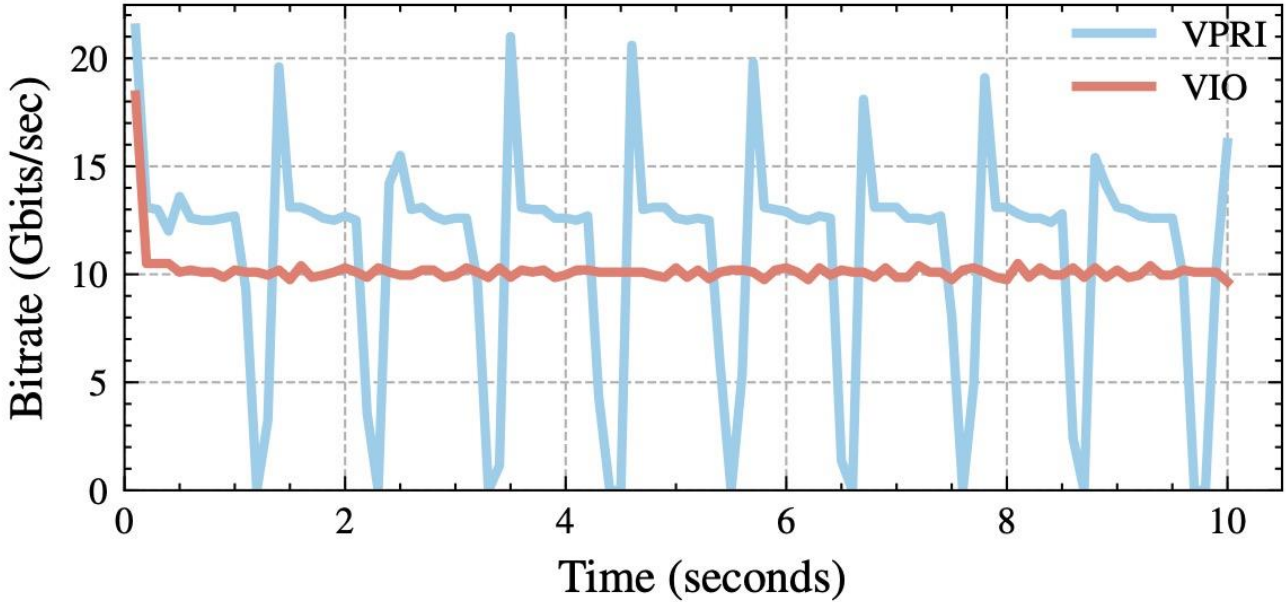
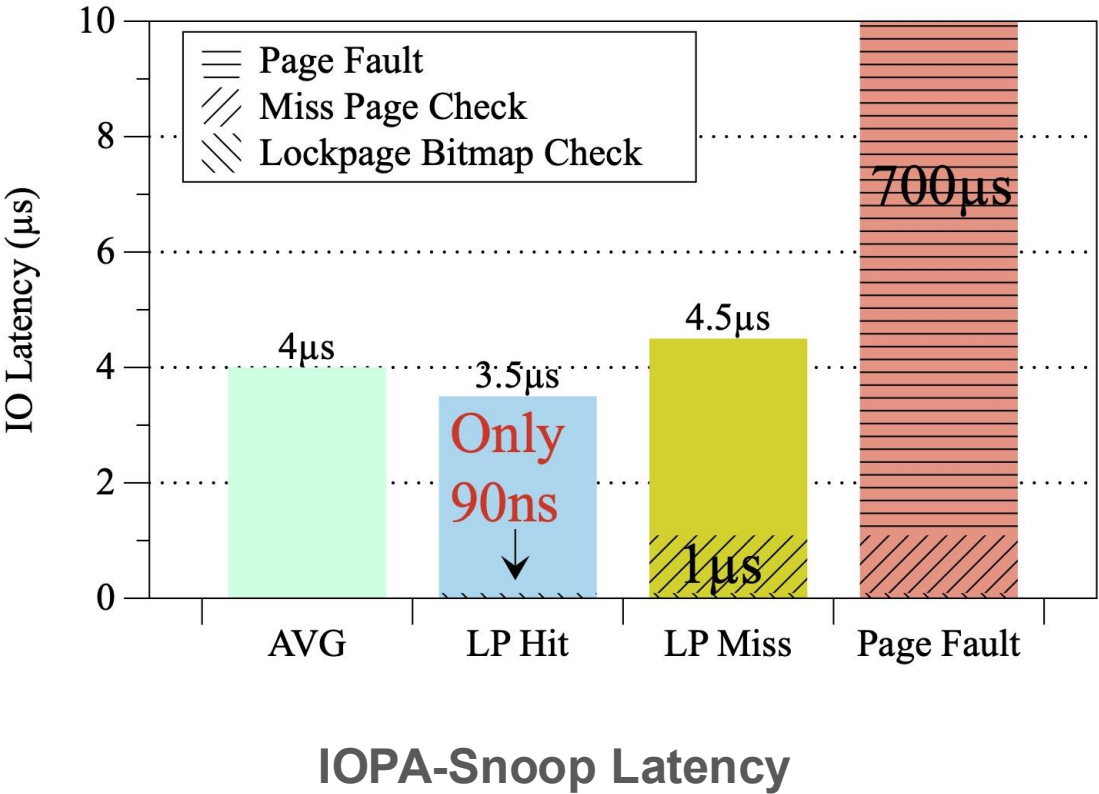
## Software

- Applications: iperf/nginx/Redis/Memcached
- Dataset: YCSB
- Baseline: VPRI[SOSP'24]

## Hardware

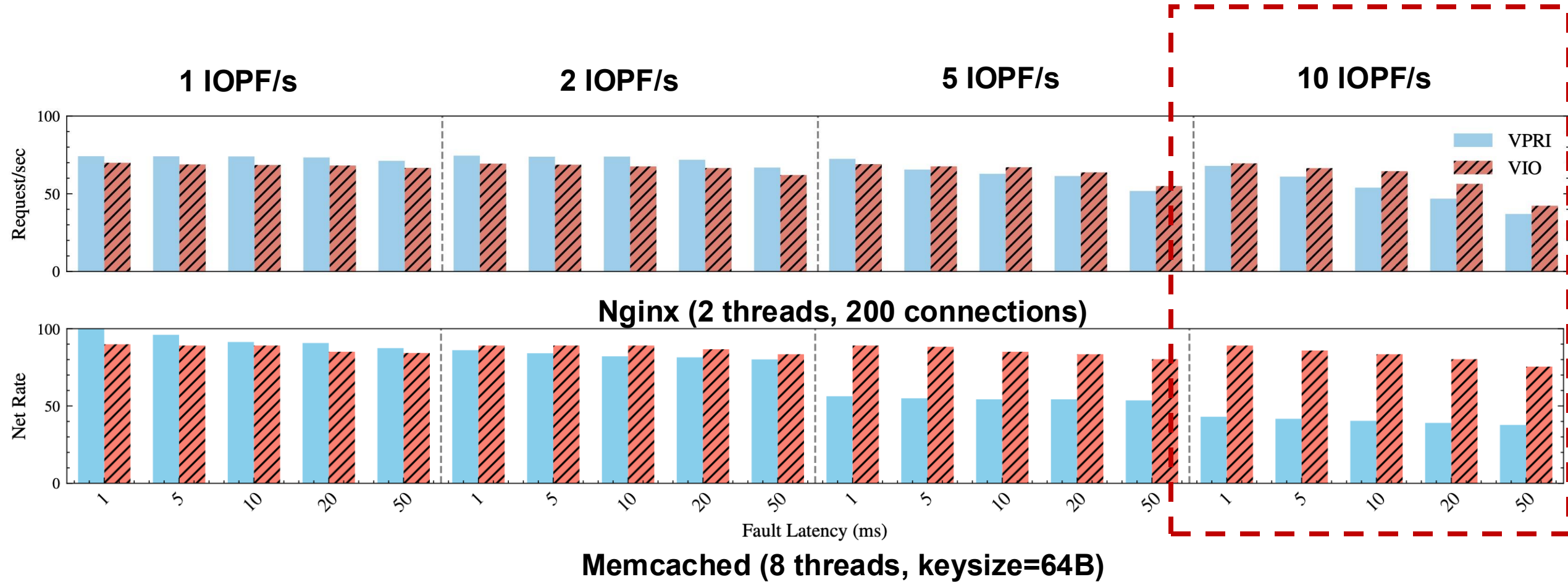
<b>Host</b>	2x Intel Xeon Platinum 8269CY (52C/104T), 192GB RAM, 1TB SSD
<b>DPU</b>	200 Gb/s Bandwidth, PCIe Gen3, Supports 2300 VFs
<b>VM</b>	4 vCPUs, 8GB RAM, CentOS 7.9, 10 Gb/s Network Throughput

# Evaluation



iperf3 performance  
(1 IOPF/s with 5ms latency)

# Evaluation



As the rate of I/O page faults per second (IOPF/s) increases, VIO's performance remains consistent, with **throughput degradation staying below 10%** across all applications.

# Conclusion

---

- VIO is a **practical software alternative** to hardware-dependent approaches like PRI, offering hardware independence and high performance.
- It uses **IOPA-Snooping** and **Elastic Passthrough** to move fault handling off the critical path, achieving both near-native speed and resource efficiency.
- VIO is successfully **deployed on 300K VMs in the production**, reclaiming significant daily memory while maintaining user SLOs.

---

**Thank you!**

---

**Email: [yunwang94@sjtu.edu.cn](mailto:yunwang94@sjtu.edu.cn)**