# Is Bare-metal I/O Performance with User-defined Storage Drives Inside VMs Possible?

## Benchmarking libvfio-user vs. Common Storage Virtualization Configurations

Sebastian Rolon & Oana Balmau

**CHEOPS '23**
May 8th, 2023
Rome, Italy

# Virtualization is everywhere

- Datasets keep growing
- We want storage to be efficient



Source: Intel [1]. An Intel Optane PCIe NVMe SSD.



Source: Google [2]. Inside a Google Datacenter.

# What is missing from storage virtualization?

1. Good performance with loose coupling
2. Rapid device prototyping
3. Live migration
4. Userspace drivers

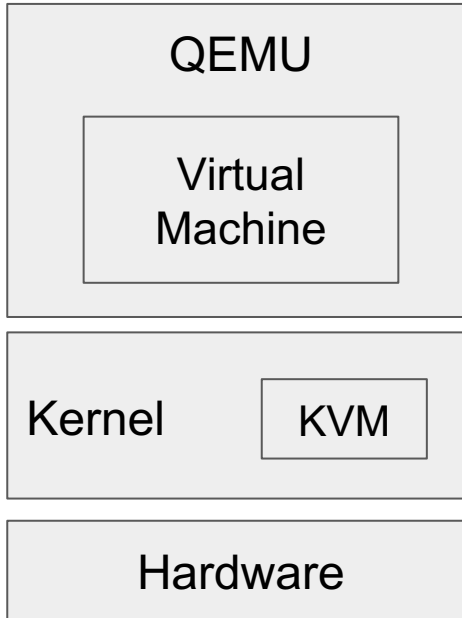# What is missing from storage virtualization?

1. Good performance with loose coupling
2. Rapid device prototyping
3. Live migration
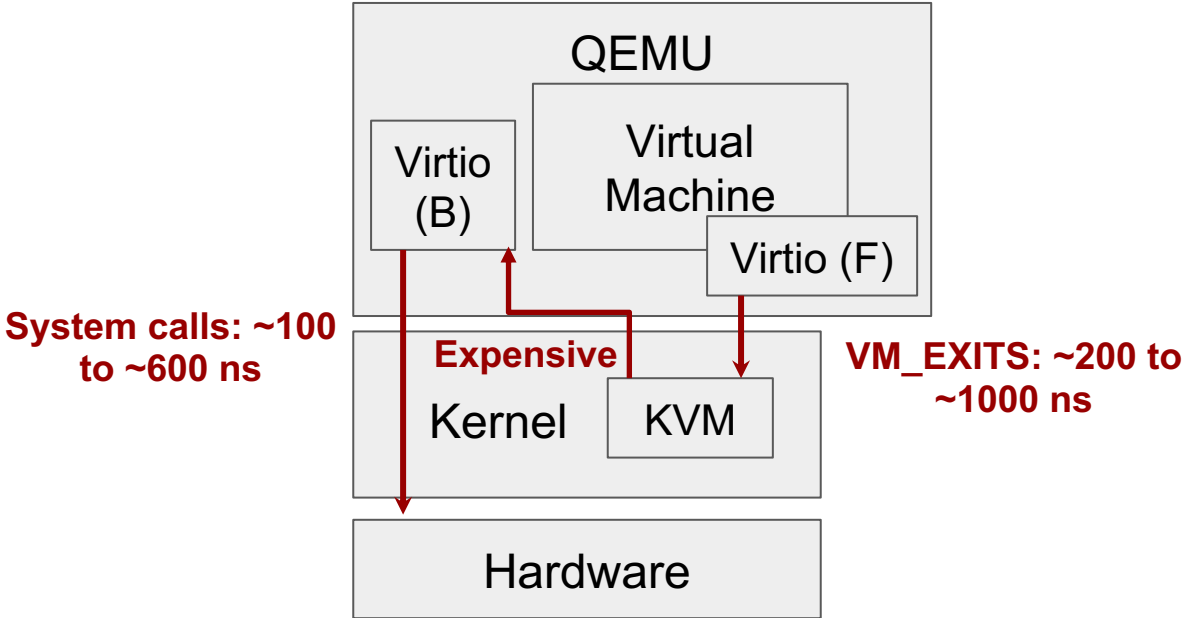4. Userspace drivers

**NUTANIX**  **vfio-user**

McGill UNIVERSITY

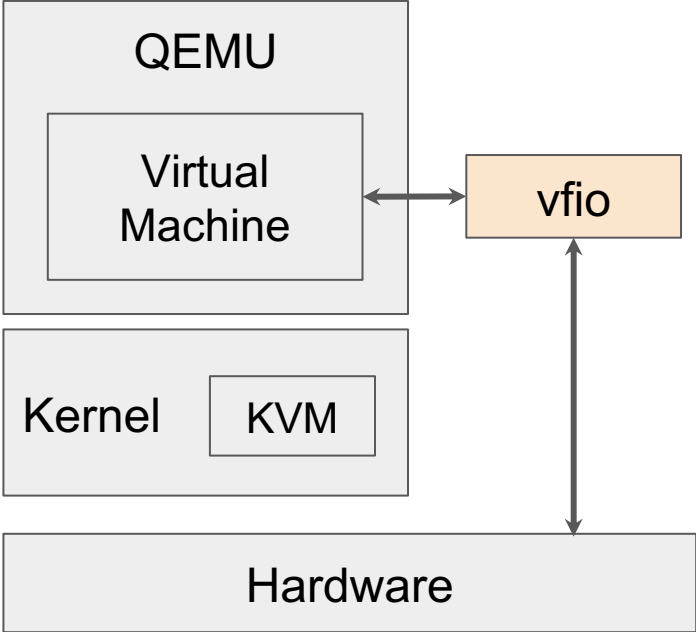# Can **vfio-user** be used as an alternative to current VM storage?

# Virtualization and QEMU/KVM

# Context switches are expensive



QEMU
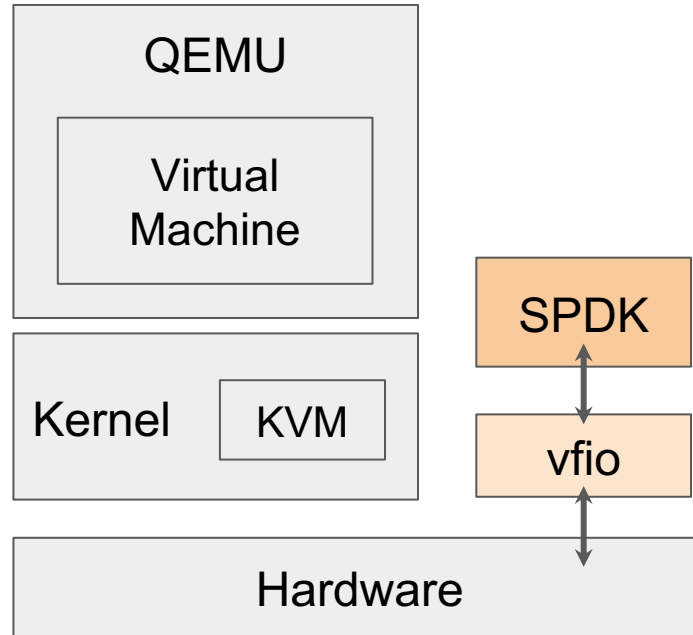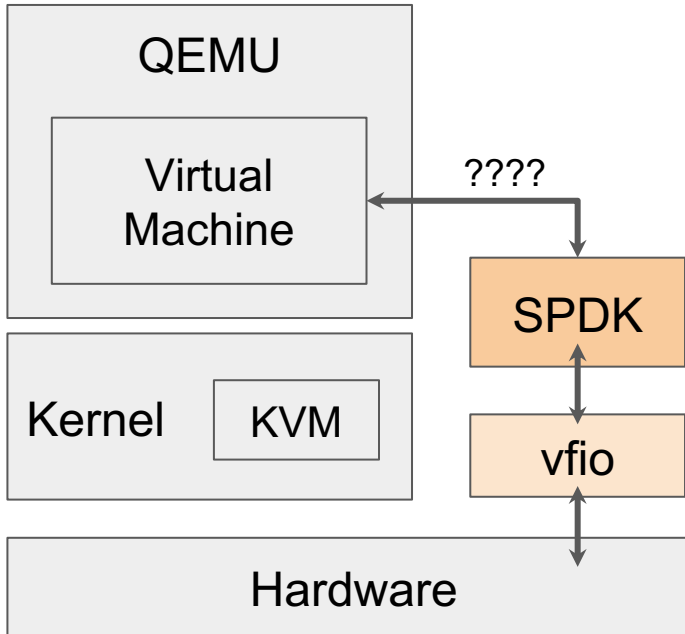
Virtio (B)

Virtual Machine

Virtio (F)

**System calls: ~100 to ~600 ns**

**Expensive**

Kernel

KVM

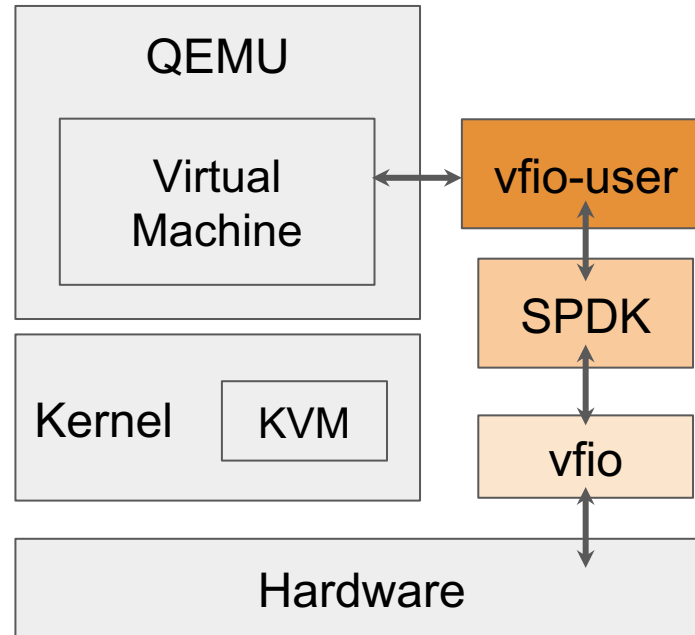**VM_EXITS: ~200 to ~1000 ns**

Hardware

# Userspace hardware access

# Abstracting hardware access

# How do we connect after abstracting?
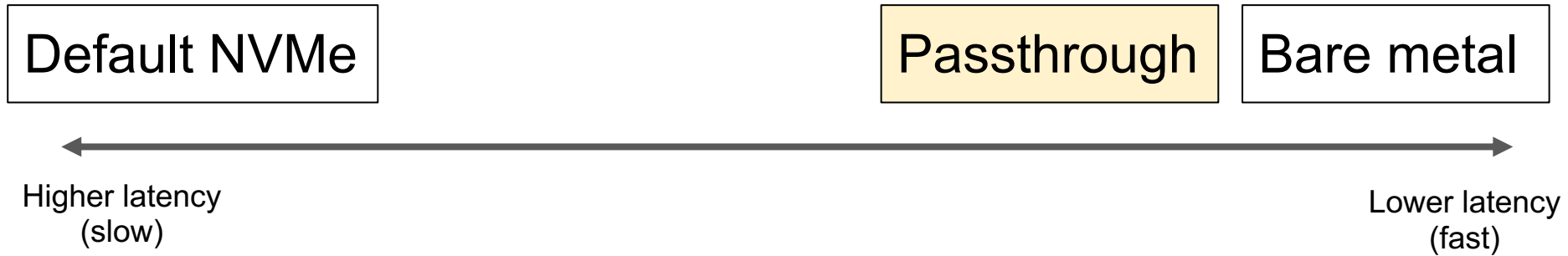
# Vfio-user virtualizes hardware over a channel
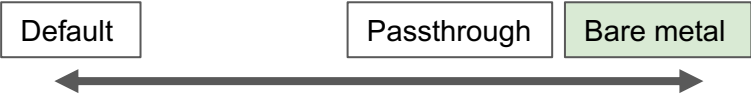
# I want my process to see an NVMe

| Default NVMe | | Passthrough | Bare metal |

Higher latency
(slow)
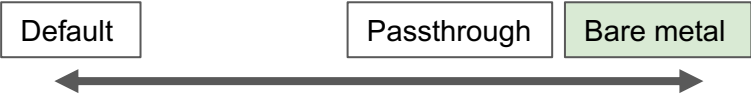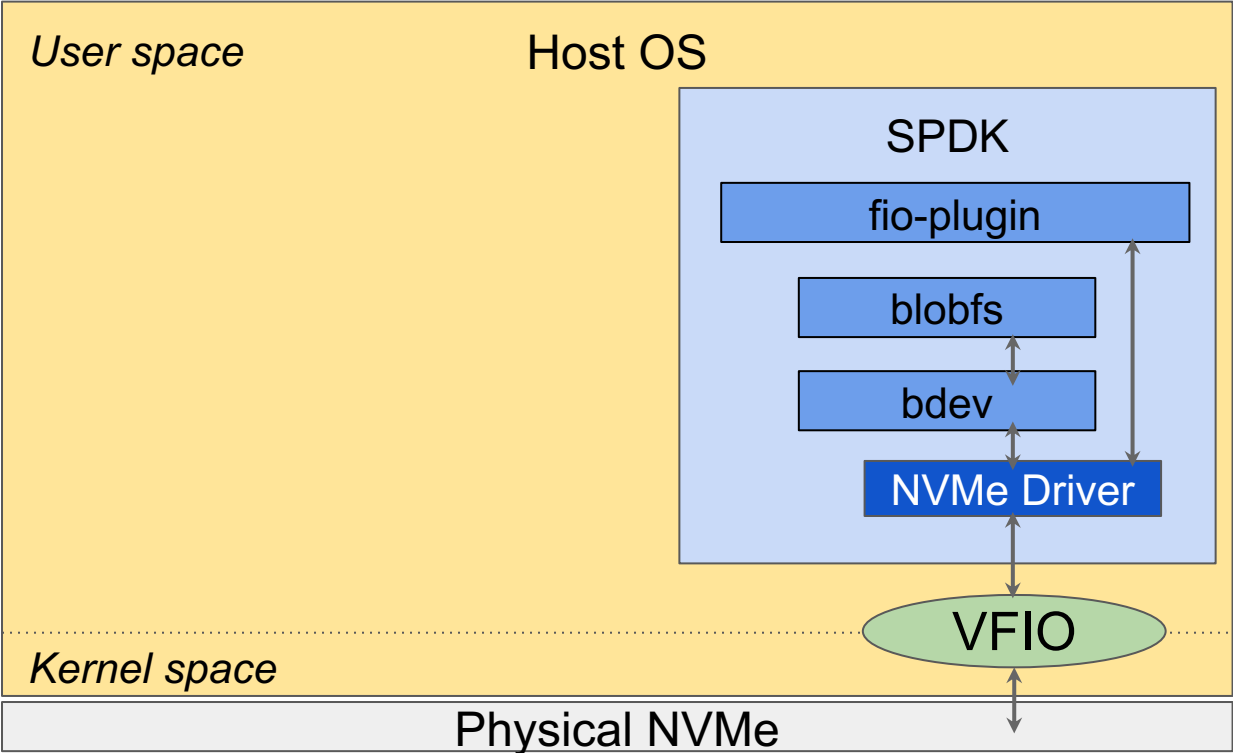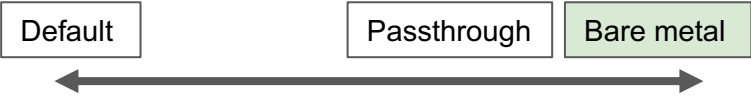
Lower latency
(fast)

# Fast and simple

| Default NVMe | | Passthrough | Bare metal |

Higher latency
(slow)

Lower latency
(fast)

# Fast, involves some configuration

| Default NVMe |   | Passthrough | Bare metal |
|:---:|:---:|:---:|:---:|

← Higher latency (slow) · · · · · · · · · · · · → Lower latency (fast)

# Slow, requires one command-line flag

| Default NVMe | | Passthrough | Bare metal |

Higher latency
(slow)

Lower latency
(fast)

# Where does vfio-user fit?



Default NVMe

Passthrough

Bare metal

Higher latency
(slow)

Lower latency
(fast)

# Bare metal configuration

*User space*          Host OS

*Kernel space*

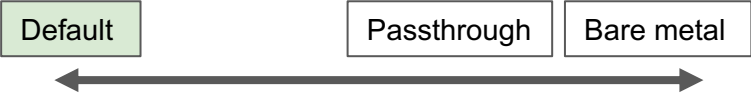VFIO

Physical NVMe

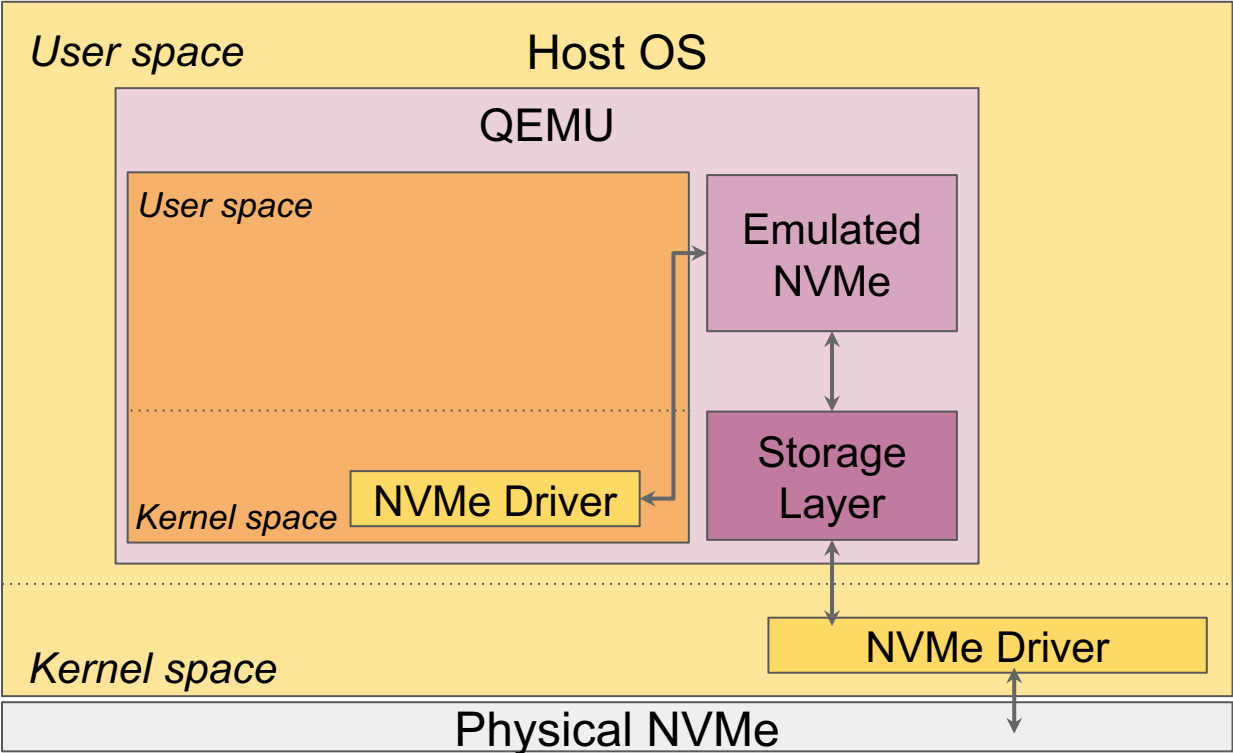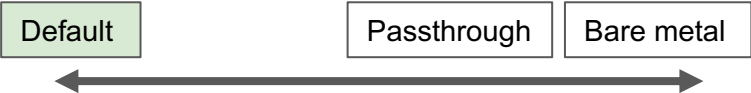# Bare metal configuration

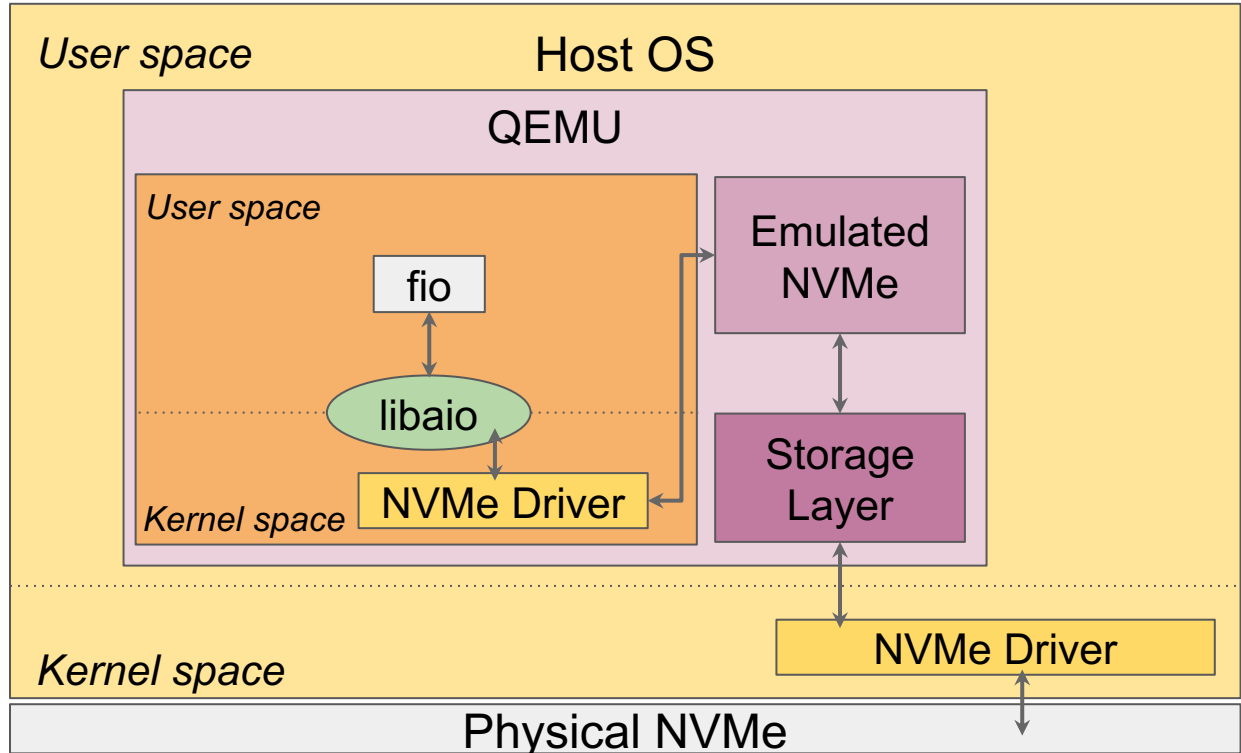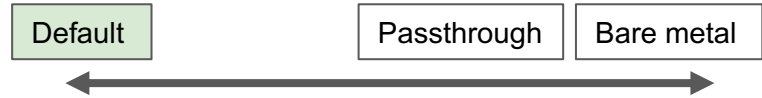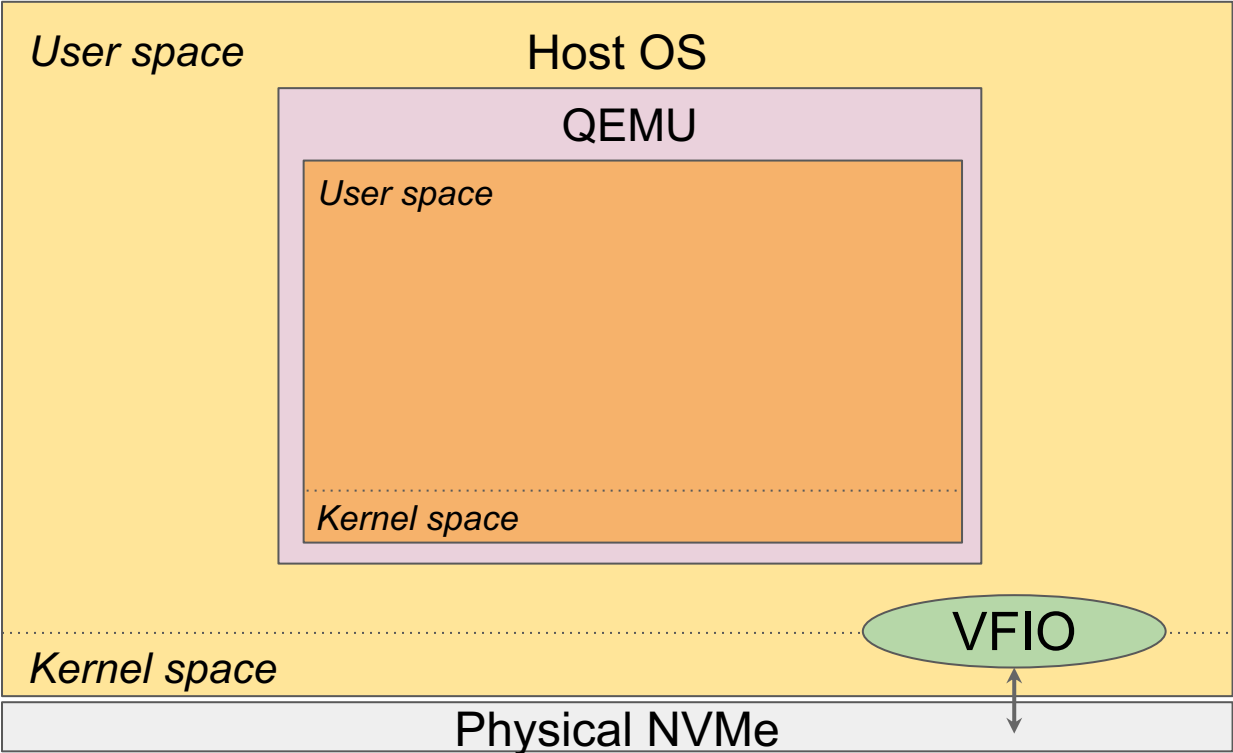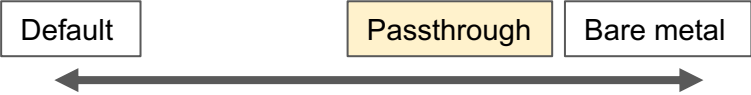# Bare metal configuration

# Bare metal configuration

# Default NVMe configuration
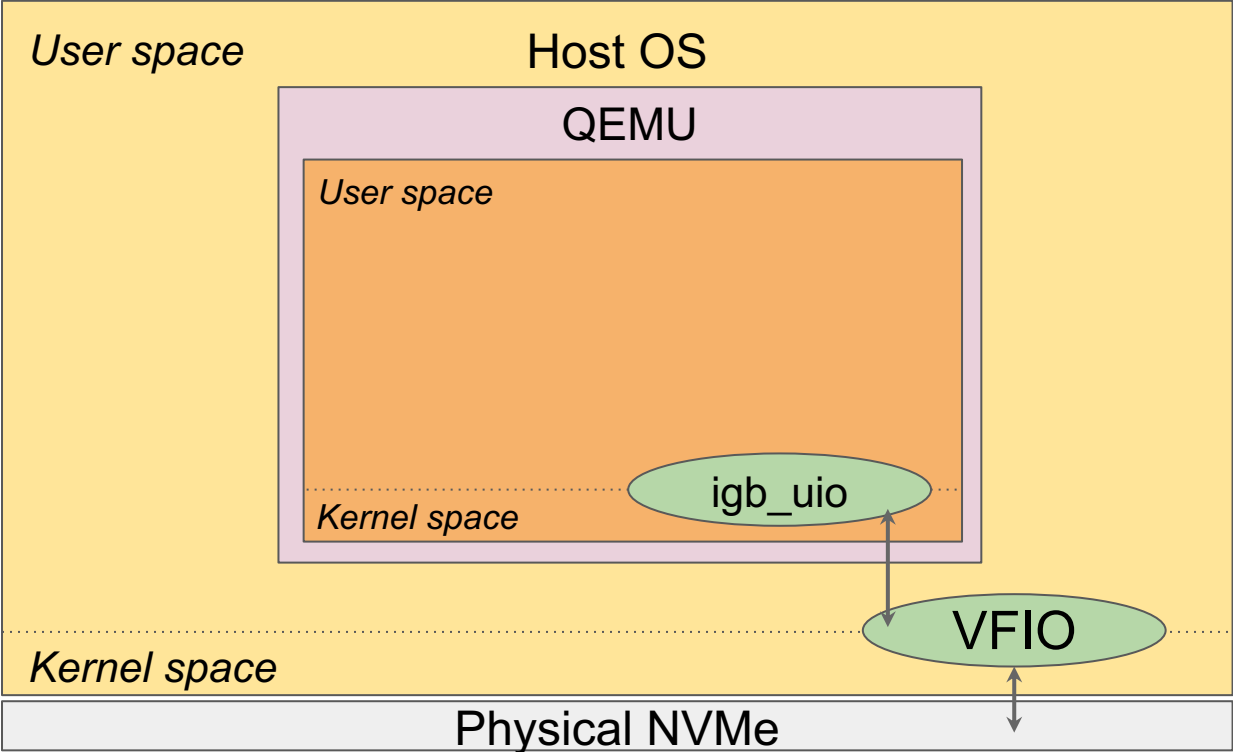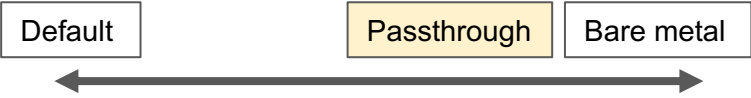
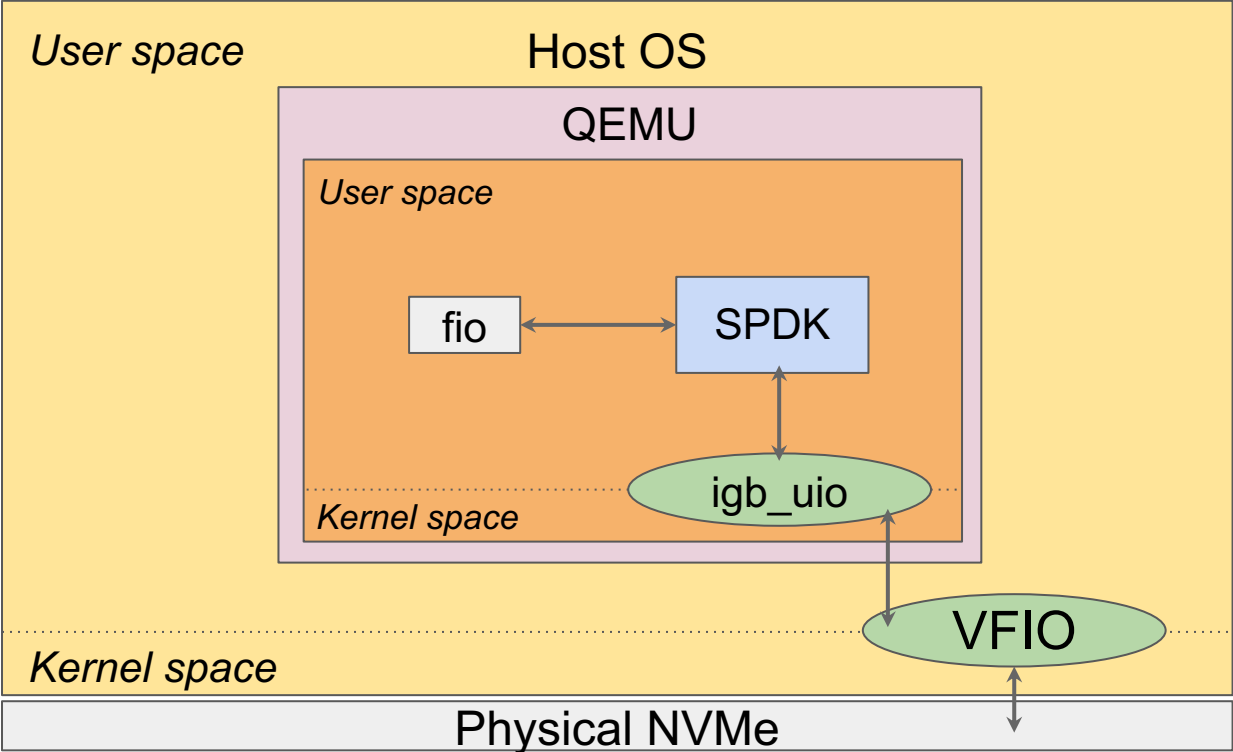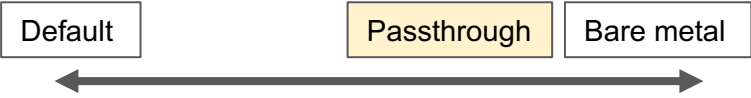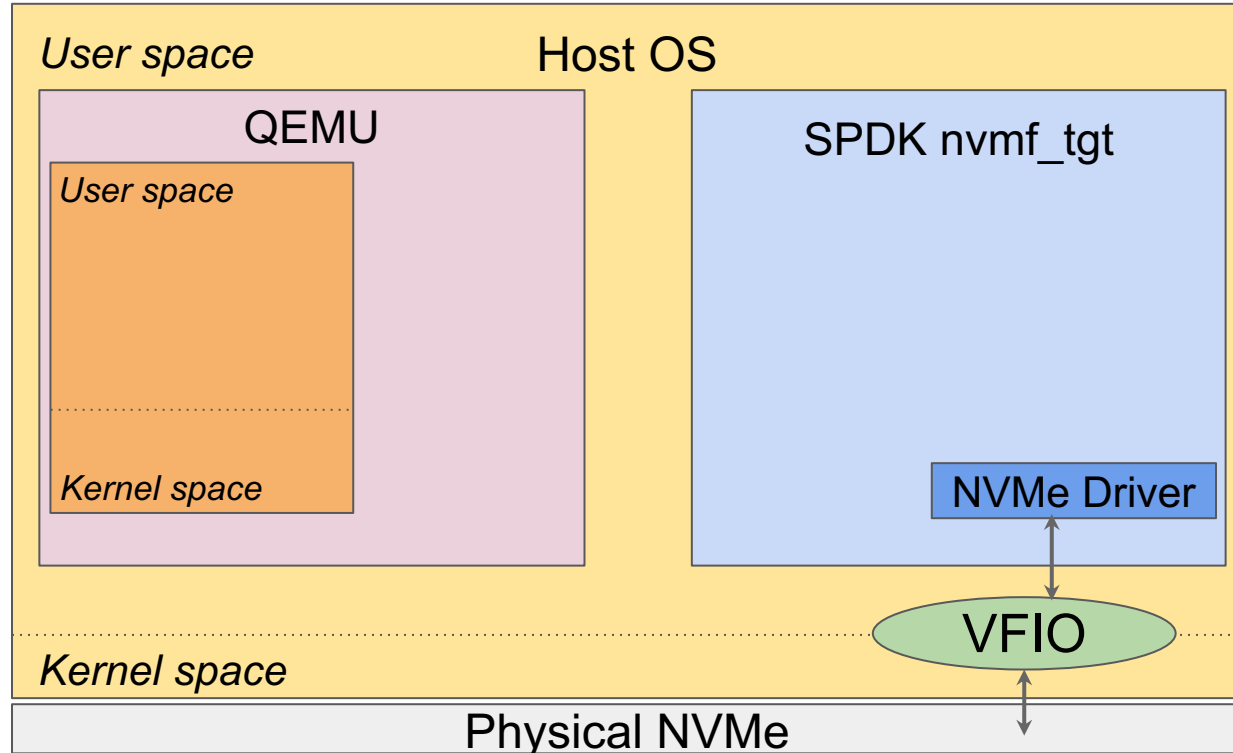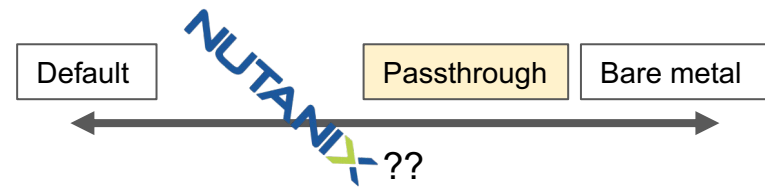# Default NVMe configuration

# Default NVMe configuration

# Default NVMe configuration

# Passthrough configuration

# Passthrough configuration

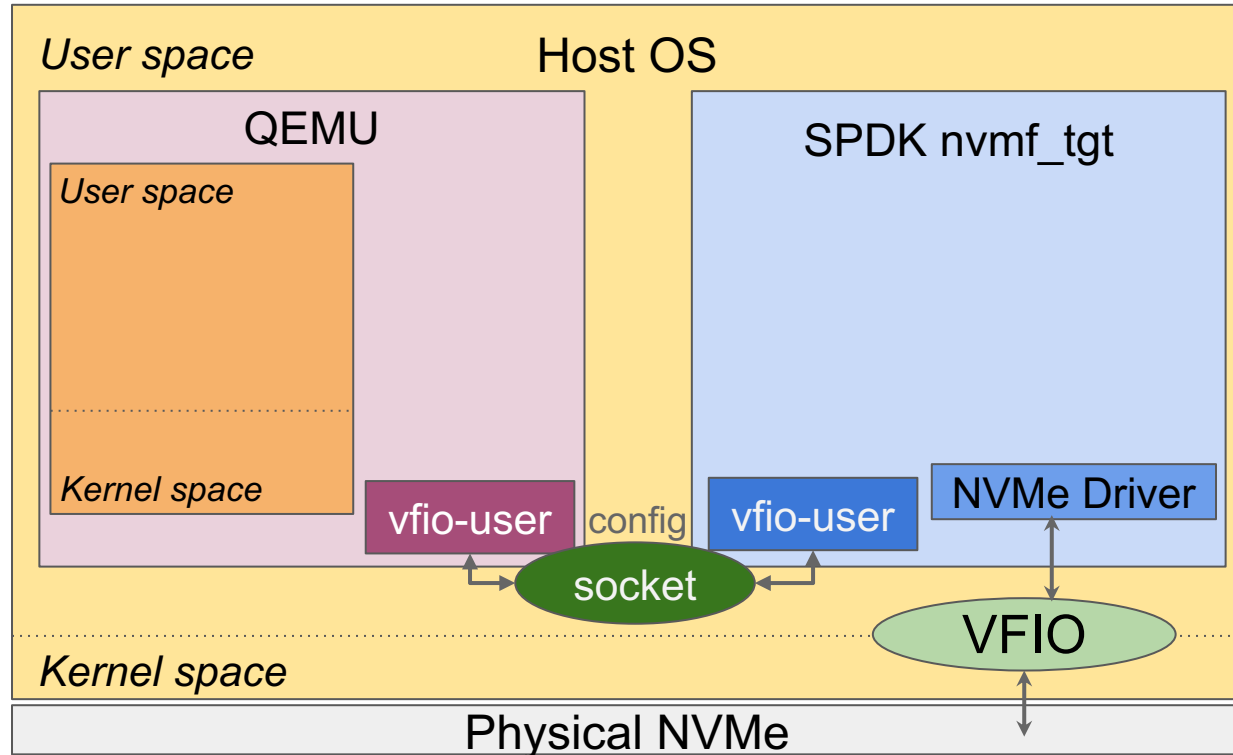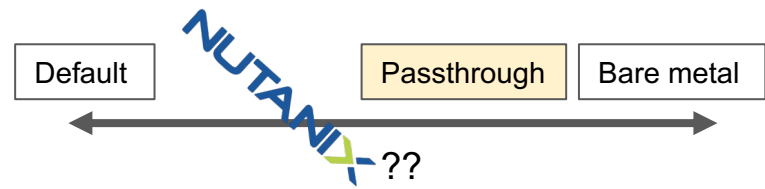# Passthrough configuration

27

# Vfio-user configuration

**Host OS**

*User space*

QEMU
- *User space*
- *Kernel space*

SPDK nvmf_tgt

NVMe Driver

VFIO

*Kernel space*

Physical NVMe

# Vfio-user configuration

Host OS

*User space*

QEMU

*User space*

*Kernel space*

vfio-user

config
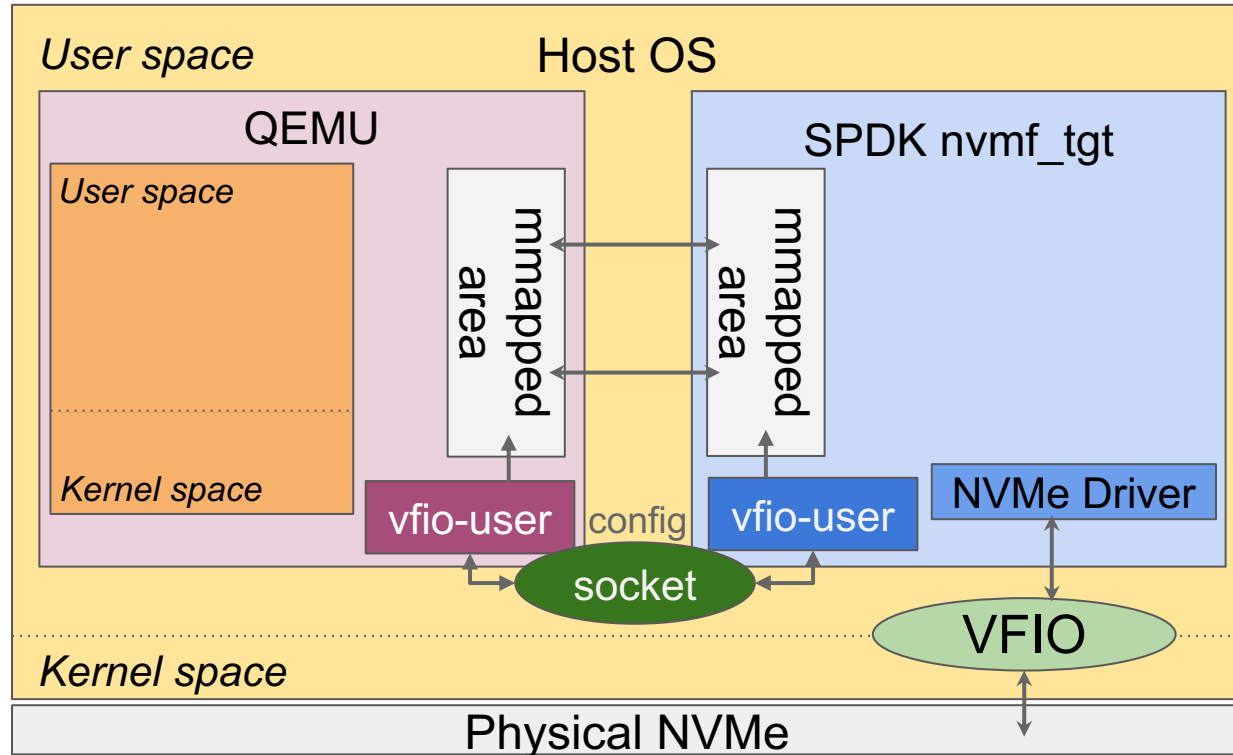
socket
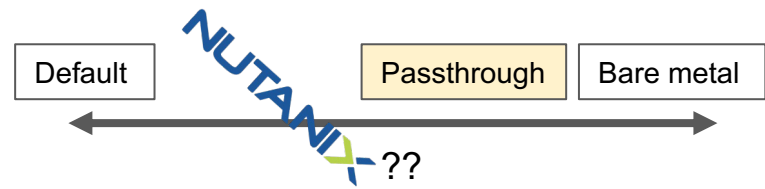
SPDK nvmf_tgt

vfio-user

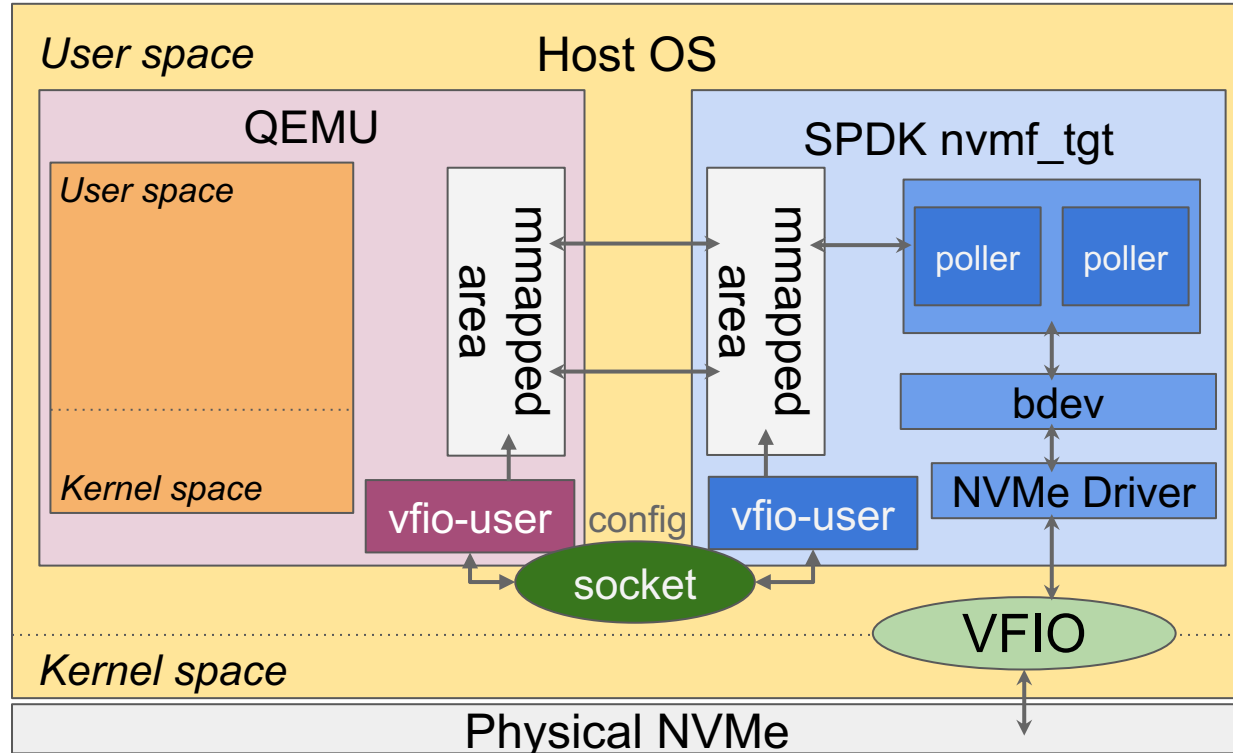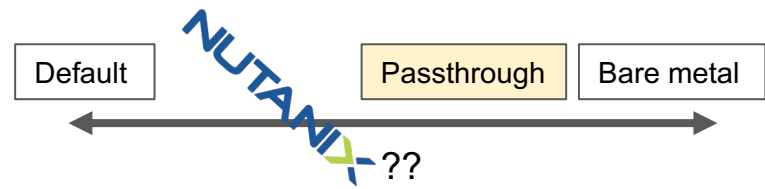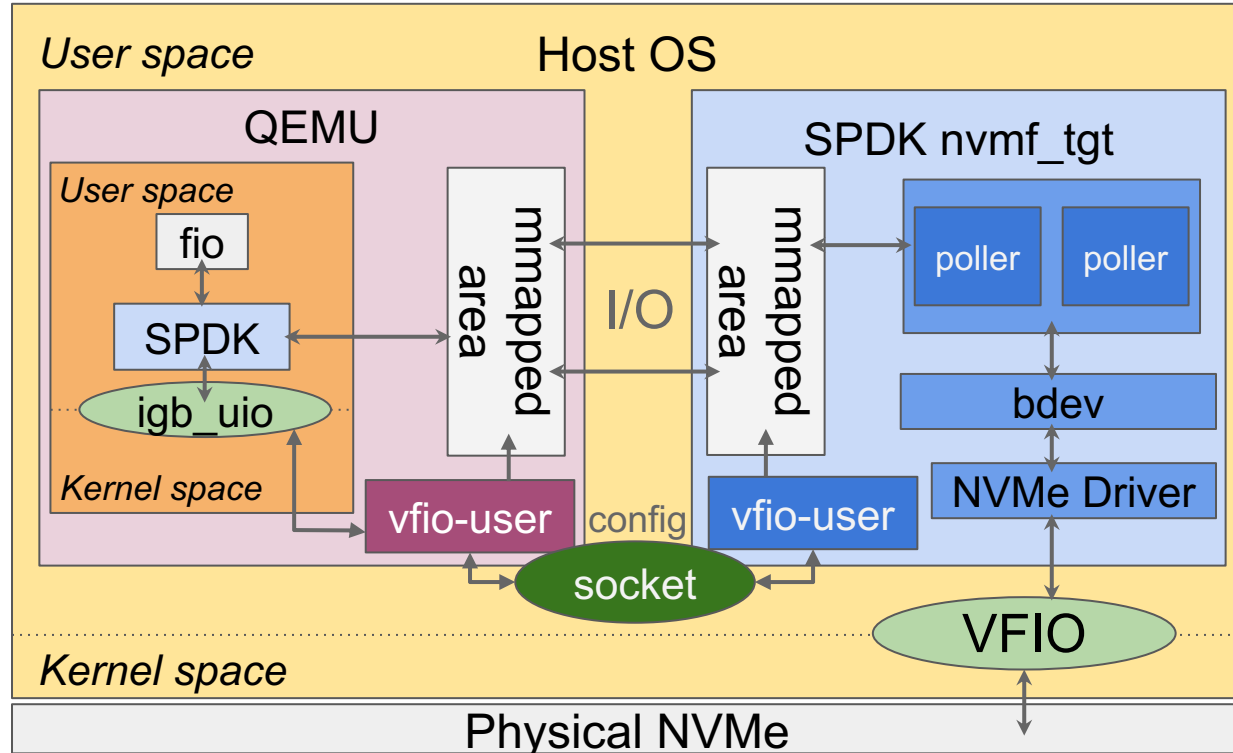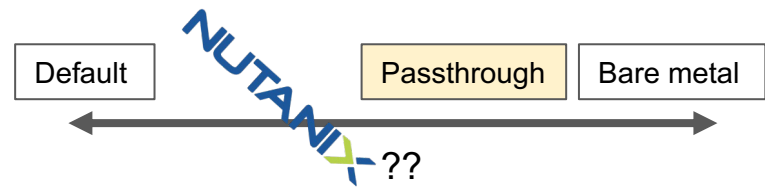NVMe Driver

VFIO

*Kernel space*

Physical NVMe

# Vfio-user configuration

# Vfio-user configuration

# Vfio-user configuration

# Experimental evaluation

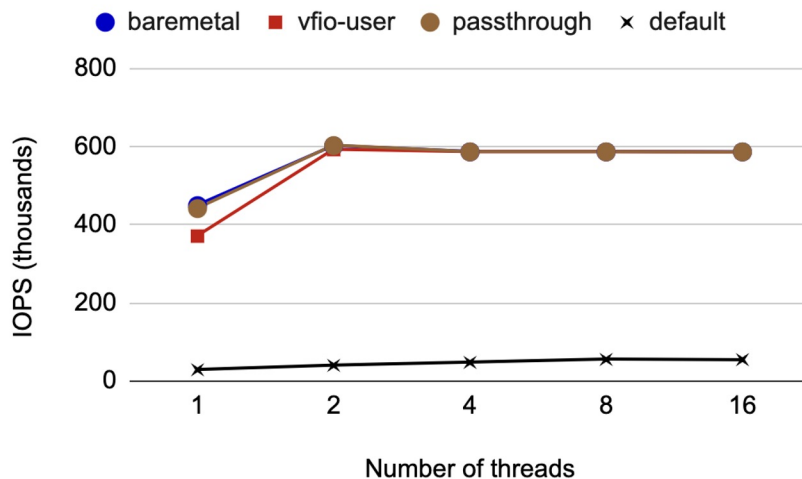| Hardware Specifications | |
|---|---|
| CPU | 36 Core Xeon Gold 6240L @ 2.40 GHz |
| Memory | 768 GiB 3200MHz DDR4 DIMM |
| SSD | 375GiB Dell Express Flash NVMe P4800X |

# Results: fio random reads

# Results: fio random writes

# Results: RocksDB benchmarks



Random reads (100M keys)

Random writes (100M keys)

baremetal   vfio-user
passthrough   default

# Layer-by-layer latency measurement

# Layer-by-layer latency measurement

# Layer-by-layer latency measurement
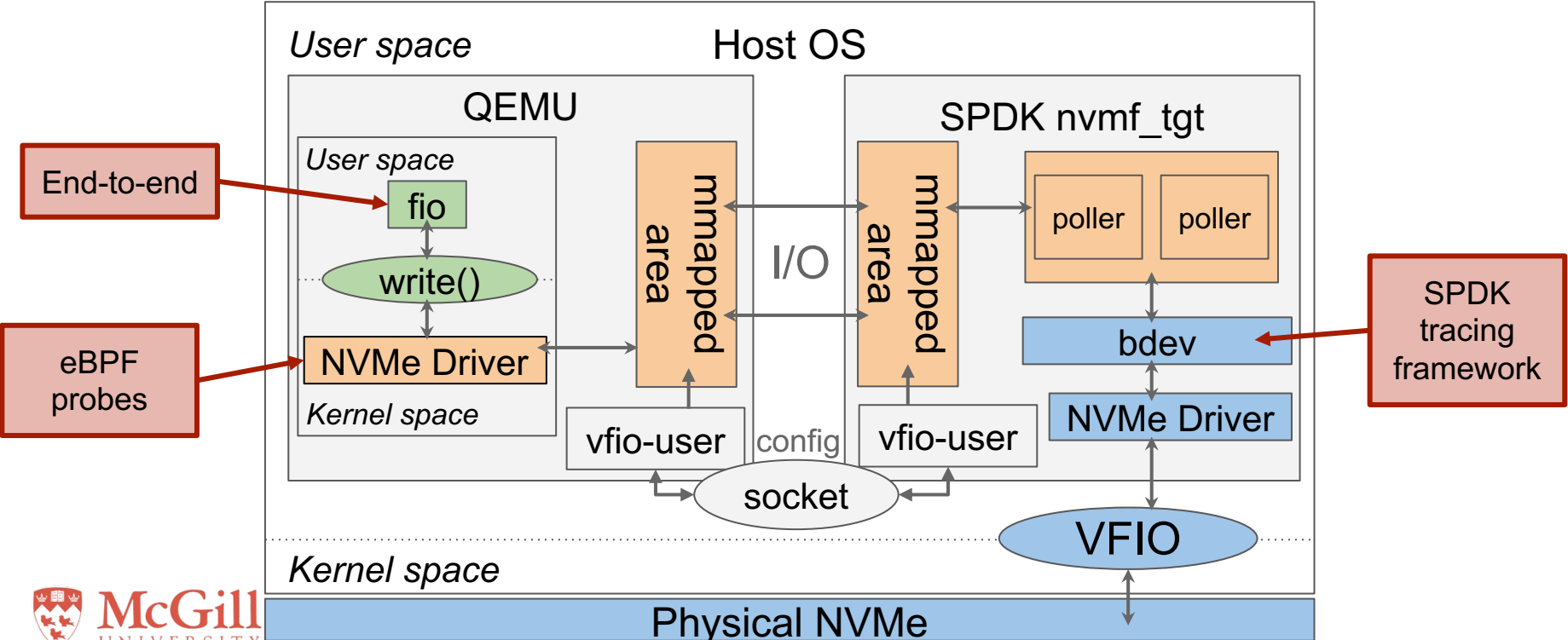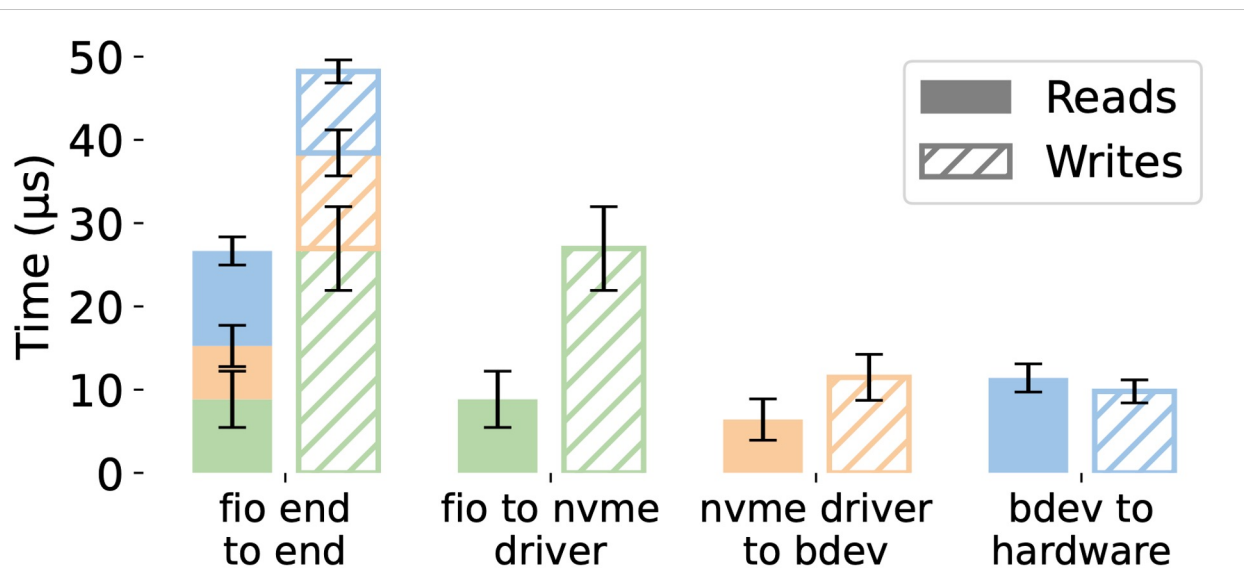
# Layer-by-layer latency measurement

# Layer-by-layer latency measurements

# Where does vfio-user fit again?



Default NVMe

?

Passthrough    Bare metal

Higher latency
(slow)

Lower latency
(fast)

# Right about next to Passthrough



Default NVMe

Passthrough

Bare metal

Higher latency
(slow)

Lower latency
(fast)

# Conclusion

- **Vfio-user** appears to have comparable performance to passthrough
- Could be viable for VM storage
- See more benchmarks and analysis in our paper!

Visit the DISCS Lab:
https://discslab.cs.mcgill.ca

On the job market!
js@rolon.co