



containercon

CHINA 中国



THINK OPEN

开放性思维

# Accelerating NVMe I/Os in Virtual Machines via SPDK vhost

Ziye Yang, Changpeng Liu  
Senior software Engineer  
Intel

# Notices & Disclaimers

Intel technologies' features and benefits depend on system configuration and may require enabled hardware, software or service activation. Performance varies depending on system configuration.

No computer system can be absolutely secure.

Tests document performance of components on a particular test, in specific systems. Differences in hardware, software, or configuration will affect actual performance. For more complete information about performance and benchmark results, visit <http://www.intel.com/benchmarks>.

Software and workloads used in performance tests may have been optimized for performance only on Intel microprocessors. Performance tests, such as SYSmark and MobileMark, are measured using specific computer systems, components, software, operations and functions. Any change to any of those factors may cause the results to vary. You should consult other information and performance tests to assist you in fully evaluating your contemplated purchases, including the performance of that product when combined with other products. For more complete information visit <http://www.intel.com/benchmarks>.

Benchmark results were obtained prior to implementation of recent software patches and firmware updates intended to address exploits referred to as "Spectre" and "Meltdown." Implementation of these updates may make these results inapplicable to your device or system.

Intel® Advanced Vector Extensions (Intel® AVX)\* provides higher throughput to certain processor operations. Due to varying processor power characteristics, utilizing AVX instructions may cause a) some parts to operate at less than the rated frequency and b) some parts with Intel® Turbo Boost Technology 2.0 to not achieve any or maximum turbo frequencies. Performance varies depending on hardware, software, and system configuration and you can learn more at <http://www.intel.com/go/turbo>.

Intel's compilers may or may not optimize to the same degree for non-Intel microprocessors for optimizations that are not unique to Intel microprocessors. These optimizations include SSE2, SSE3, and SSSE3 instruction sets and other optimizations. Intel does not guarantee the availability, functionality, or effectiveness of any optimization on microprocessors not manufactured by Intel. Microprocessor-dependent optimizations in this product are intended for use with Intel microprocessors. Certain optimizations not specific to Intel microarchitecture are reserved for Intel microprocessors. Please refer to the applicable product User and Reference Guides for more information regarding the specific instruction sets covered by this notice.

Cost reduction scenarios described are intended as examples of how a given Intel-based product, in the specified circumstances and configurations, may affect future costs and provide cost savings. Circumstances will vary. Intel does not guarantee any costs or cost reduction.

Intel does not control or audit third-party benchmark data or the web sites referenced in this document. You should visit the referenced web site and confirm whether referenced data are accurate.

© 2018 Intel Corporation.

Intel, the Intel logo, and Intel Xeon are trademarks of Intel Corporation in the U.S. and/or other countries.

\*Other names and brands may be claimed as property of others.

# Agenda

- Background
- SPDK vhost solution
- Experiments
- Conclusion



LINUXCON

containercon



CLOUDOPEN

CHINA 中国

THINK OPEN

开放性思维

# Background

# NVMe & virtualization

- NVMe specification enables highly optimized drives (e.g., NVMe SSD)
  - For example, multiple I/O queues allows lockless submission from CPU cores in parallel
- However, even the best kernel mode drivers have non-trivial software overhead
  - Long I/O stack in kernel with resource contention
- Virtualization adds additional overhead
  - Long I/O stack in both guest OS kernel and host OS kernel
  - Context switch overhead (e.g., VM\_EXIT caused by I/O interrupt in guest OS)

# What is in QEMU's solution?

- The solution in QEMU to virtualize NVMe device:
  - Virtio virtualization
  - NVMe controller virtualization
  - Hardware assisted virtualization
- Virtio virtualization
  - Virtio SCSI/block Controllers
- NVMe controller virtualization
  - QEMU emulated NVMe Device (file based NVMe backend)
  - QEMU NVMe Block Driver based on VFIO (exclusive access by QEMU)

# Background: What is in QEMU

Guest VM  
(Linux\*, Windows\*, FreeBSD\*, etc.)

virtio front-end drivers

virtqueue

virtio back-end drivers

device emulation

Hypervisor (i.e. QEMU/KVM)

- Paravirtualized driver specification
- Common mechanisms and layouts for device discovery, I/O queues, etc.
- virtio device types include:
  - virtio-net
  - virtio-blk
  - virtio-scsi
  - virtio-gpu
  - virtio-rng
  - virtio-crypto

# Accelerate virtio via vhost target

Guest VM  
(Linux\*, Windows\*, FreeBSD\*, etc.)

virtio front-end drivers

virtqueue

virtio back-end drivers

Device emulation

vhost

Hypervisor (i.e. QEMU/KVM)

vhost

vhost target  
(kernel or userspace)

- Separate process for I/O processing
- vhost protocol for communicating guest VM parameters
  - memory
  - number of virtqueues
  - virtqueue locations





containercon



CHINA 中国

THINK OPEN

开放性思维

# SPDK vhost solution

# What is SPDK?

## Storage Performance Development Kit



### ***Intel® Platform Storage Reference Architecture***

- Optimized for *Intel platform* characteristics
- Open source building blocks (BSD licensed)
- Available via [github.com/spdk](https://github.com/spdk) or [spdk.io](https://spdk.io)



### ***Scalable and Efficient Software Ingredients***

- User space, lockless, polled-mode components
- Up to millions of IOPS per core
- Designed for Intel Optane™ technology latencies

# SPDK architecture

18.01 Release

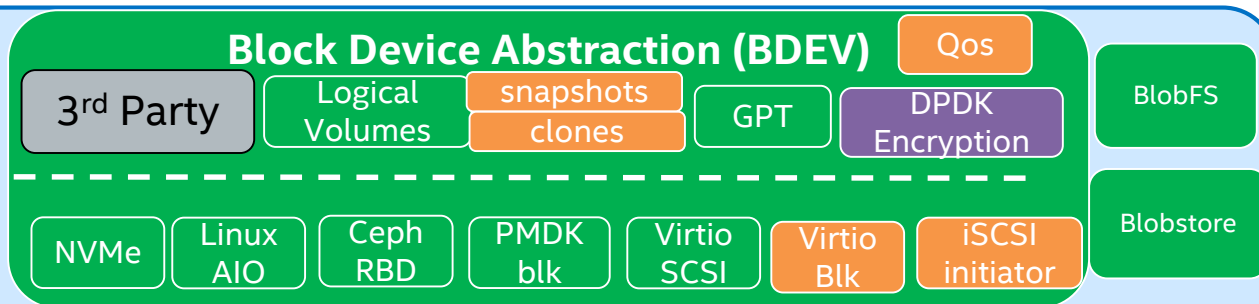
18.04 Release

18.07 Release

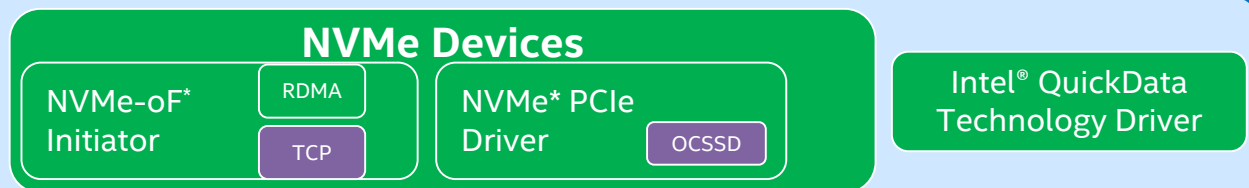
## Storage Protocols



## Storage Services



## Drivers



## Integration

Cinder

VPP TCP/IP

RocksDB

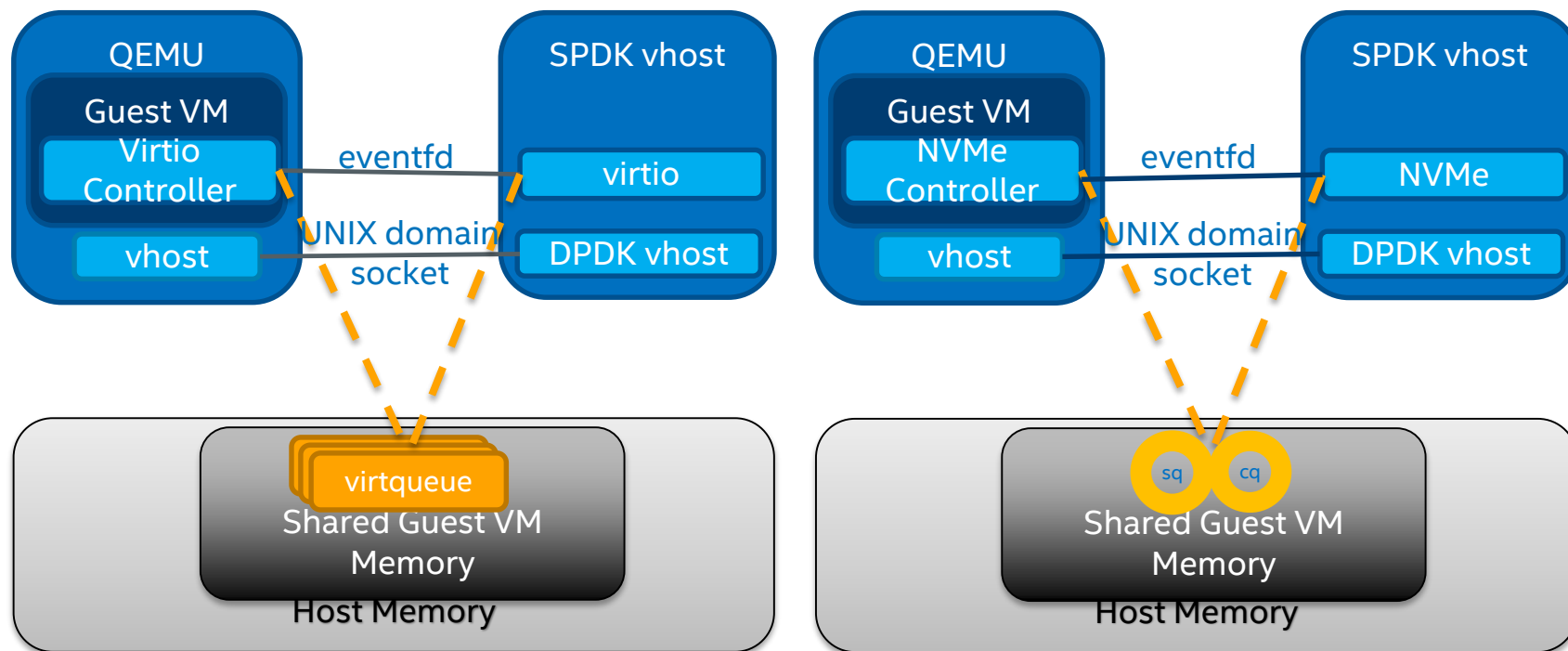
Ceph

QEMU

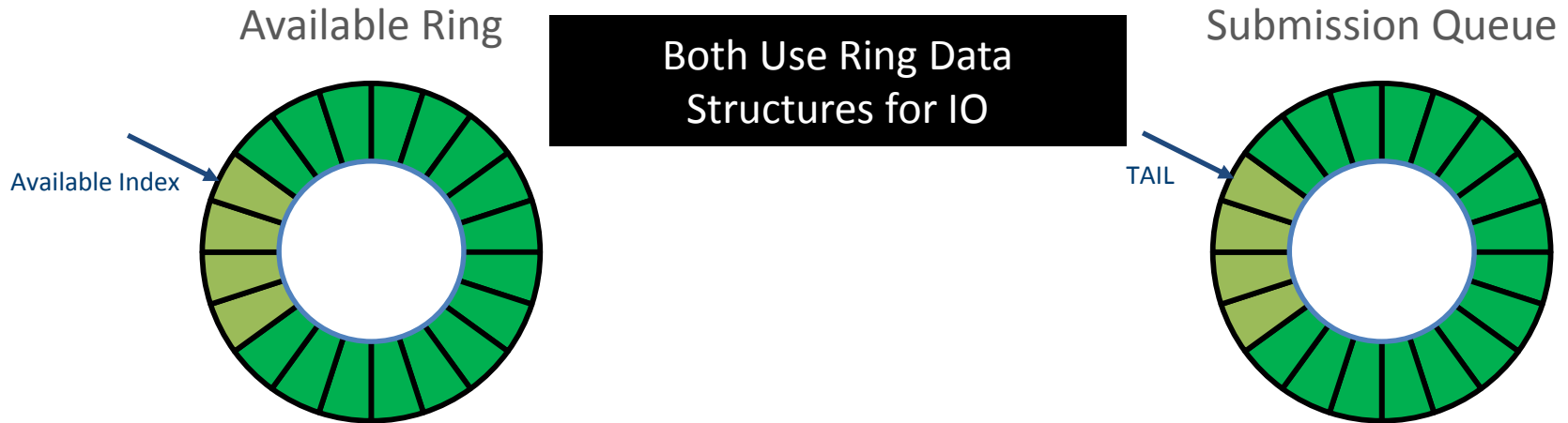
## Core

Application Framework

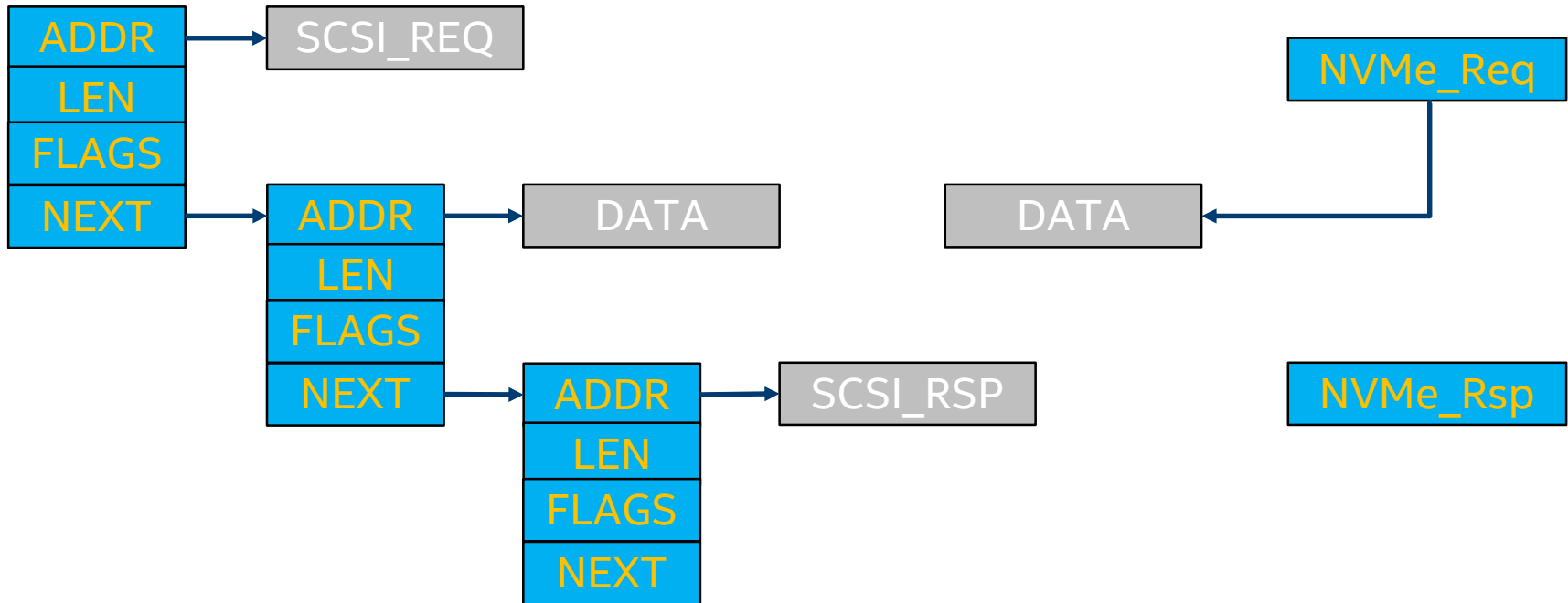
# Combine virtio and NVMe to inform a uniform SPDK vhost solution



# Virtio VS NVMe



# Virtio-SCSI and NVMe protocol format comparison



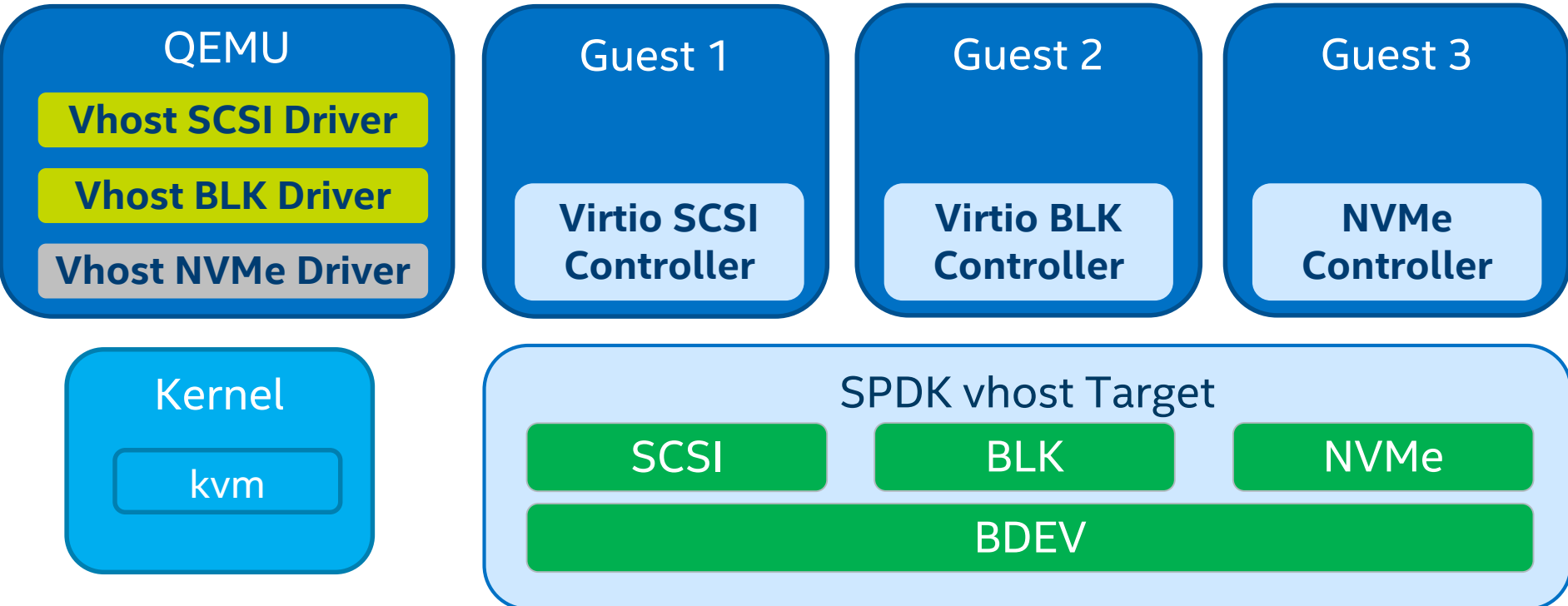
$(16 * 3 + \text{SCSI\_Req} + \text{SCSI\_Rsp} + \text{Data})$  Bytes

$(\text{NVMe\_Req} + \text{Data} + \text{NVMe\_Rsp})$  Bytes

# SPDK vhost architecture

QEMU Released

Separate Patch for QEMU



# Comparison of known solutions

Solution / Usage	QEMU Emulated NVMe device	QEMU VFIO Based solution	SPDK Vhost-SCSI	SPDK Vhost-BLK	SPDK Vhost-NVMe
Guest OS driver Interface	NVMe	NVMe	Virtio SCSI	Virtio BLK	NVMe
Backend Device sharing	Y	N	Y	Y	Y
Application Transparent support	Y	Y	Y	N (e.g., Command set is very small )	Y
Live Migration support	Y	N	Y	Y	N
VFIO dependency	N	Y	N	N	N
QEMU Change	No modification	Upstream is done	Upstream is done	Upstream is done	Upstream is in process





LINUXCON

containercon

CLLOUDOPEN

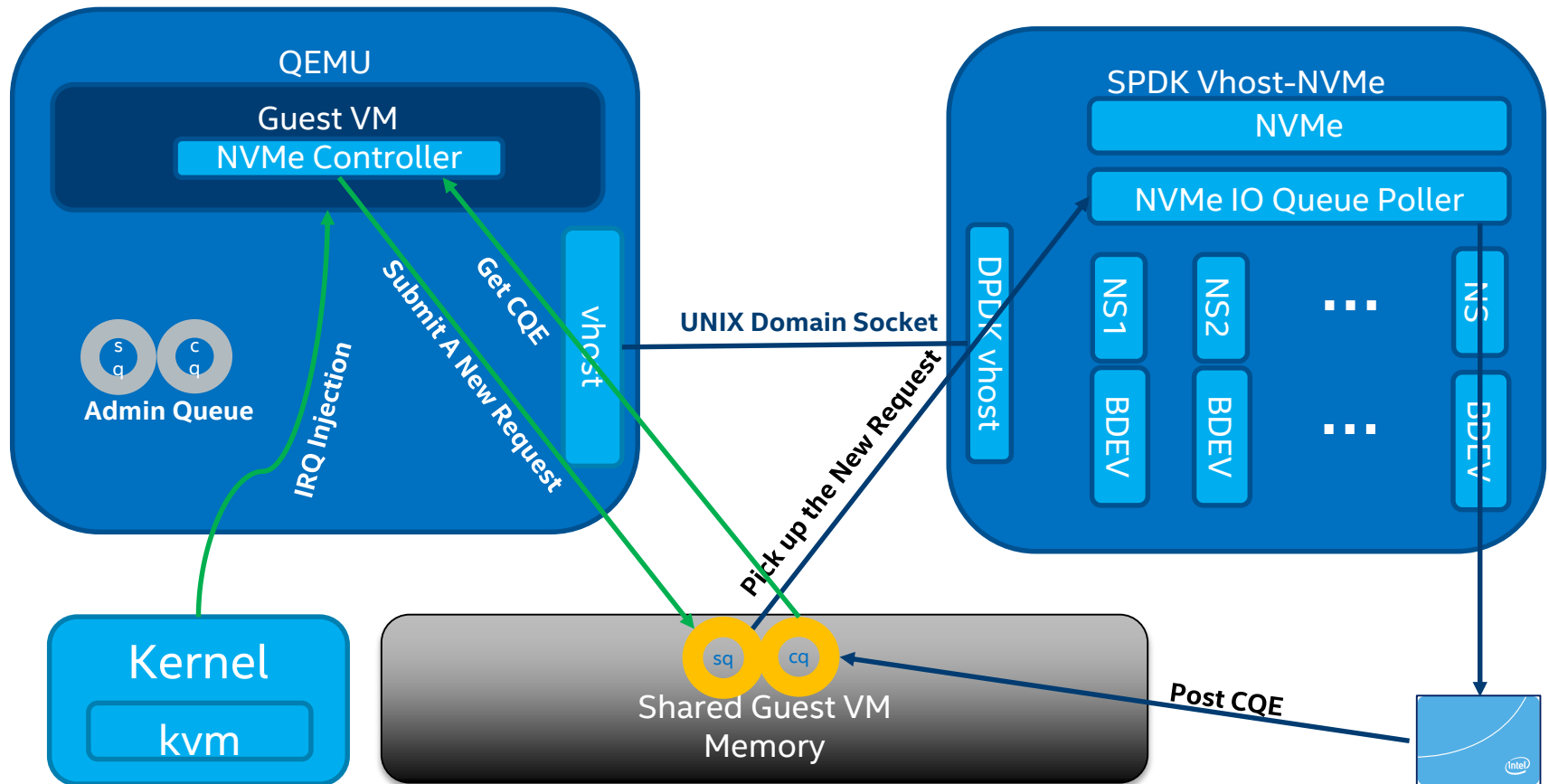
CHINA 中国

THINK OPEN

开放性思维

# SPDK vhost NVMe implementation details

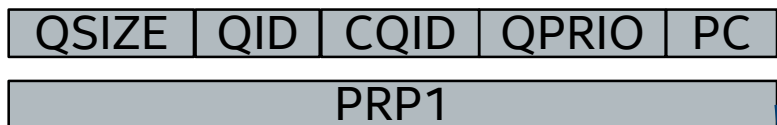
# vhost NVMe implementation details



# Create io queue

Guest: Create IO Queue

SPDK: Start to Create IO Queue



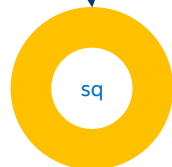
Guest: Submit to Admin, Write DB

SPDK: Memory Translation

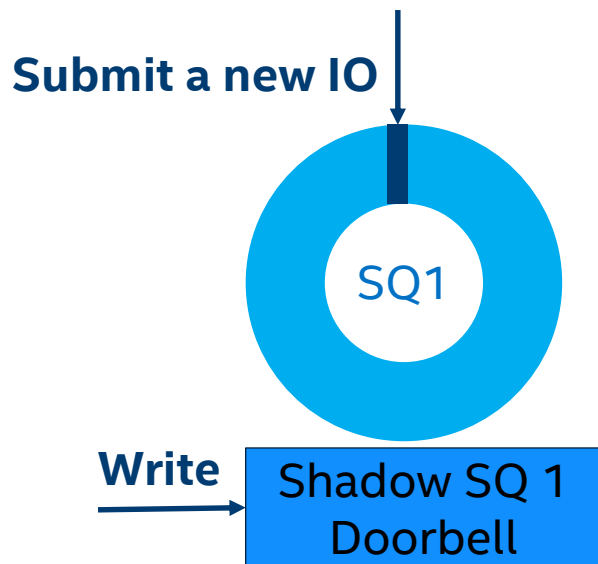
QEMU: Pick up Admin Command

SPDK: Both Guest and SPDK see same IO Queue now

QEMU: Send via Domain Socket



# New feature to address guest NVMe performance issue



MMIO Writes happened, which will cause VM\_EXIT

**NVMe 1.3 New Feature: Optional Admin Command support for Doorbell Buffer Config, only used for emulated NVMe controllers, Guest can update shadow doorbell buffer instead of submission queue's doorbell registers**

# Shadow doorbell buffer

Start	End	Description
00h	03h	Submission Queue 0 Tail Doorbell or Eventidx (Admin)
04h	07h	Completion Queue 0 Head Doorbell or Eventidx (Admin)
08h	0Bh	Submission Queue 1 Tail Doorbell or Eventidx
0Ch	0Fh	Completion Queue 1 Head Doorbell or Eventidx

Command	Description
PRP1	Shadow doorbell memory address, updated by Guest NVMe Driver
PRP2	Eventidx memory address, updated by SPDK vhost target



LINUXCON

containercon



CHINA 中国

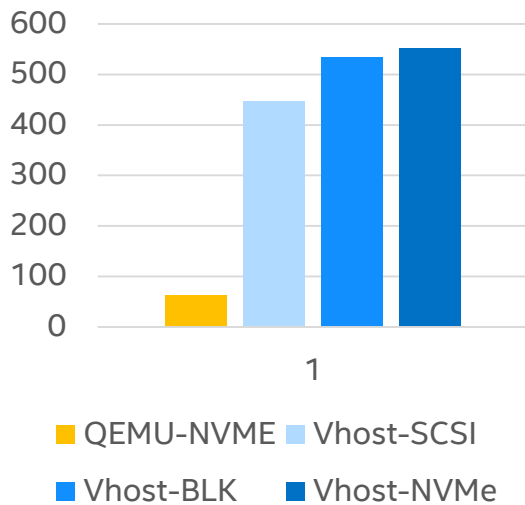
THINK OPEN

开放性思维

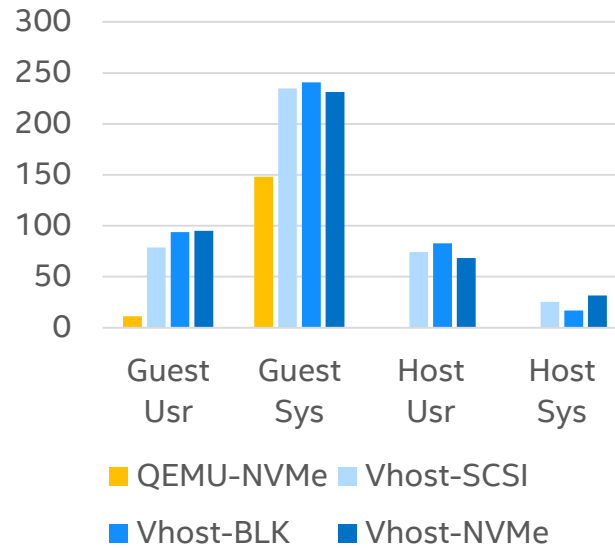
# Experiments

# 1 VM with 1 NVMe SSD

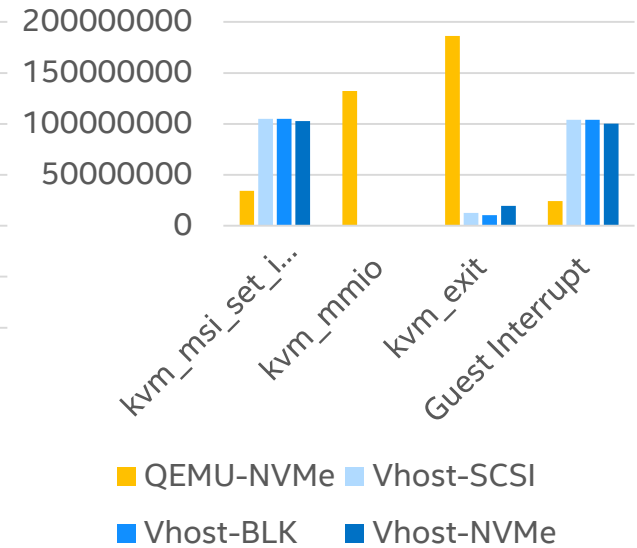
### IOPS (K)



### CPU Utilization (%)

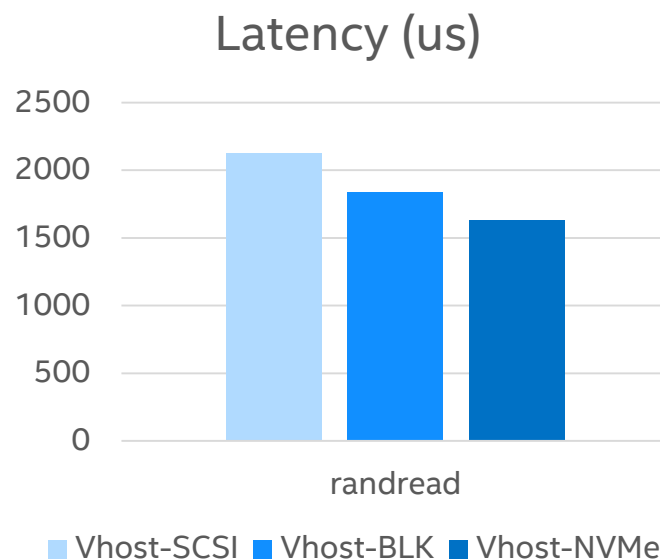
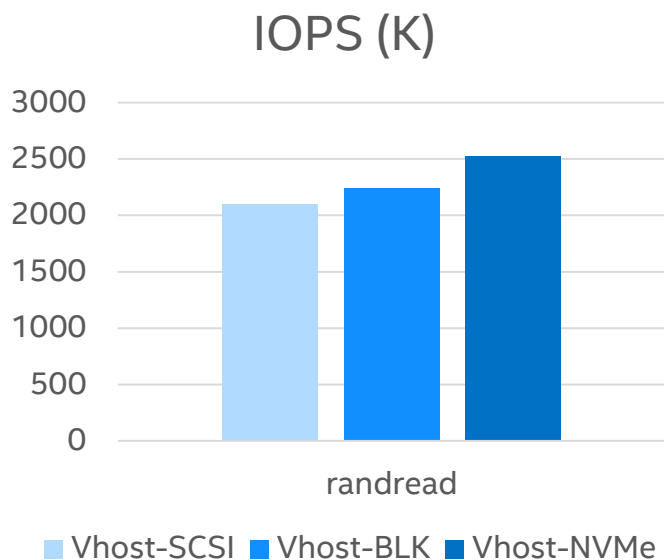


### KVM Events



System Configuration: 2 \* Intel Xeon E5 2699v4 @ 2.2GHz; 128GB, 2667 DDR4, 6 memory Channels; SSD: Intel Optane™ P4800X, FW: E2010324, 375GiB; Bios: HT disabled, Turbo disabled; OS: Fedora 25, kernel 4.16.0. 1 VM, VM config : 4 vcpu 4GB memory, 4 IO queues; VM OS: Fedora 27, kernel 4.16.5-200, blk-mq enabled; Software: QEMU-2.12.0 with SPDK Vhost-NVMe driver patch, IO distribution: 1 vhost-cores for SPDK, FIO 3.3, io depth=32, numjobs=4, direct=1, block size=4k, total tested data size=400GiB

# 8 VMs with 4 NVMe SSDs



- Linux kernel NVMe driver will poll completion queue when submitting a new request, which can help to decrease interrupt numbers and vm\_exit events.

System Configuration: 2 \* Intel Xeon E5 2699v4 @ 2.2GHz; 256GB, 2667 DDR4, 6 memory Channels; SSD: Intel DC P4510, FW: VDV10110, 2TiB; BIOS: HT disabled, Turbo disabled; Host OS: CentOS 7, kernel 4.16.7. 8 VMs, VM config : 4 vcpu 4GB memory, 4 IO queues; Guest OS: Fedora 27, kernel 4.16.5-200, blk-mq enabled; Software: QEMU-2.12.0 with SPDK Vhost-NVMe driver patch, IO distribution: 2 vhost-cores for SPDK, FIO 3.3, io depth=128, numjobs=4, direct=1, block size=4k, runtime=300s, ramp\_time=60s; SSDs well preconditioned with 2 hours randwrites before randread tests.





LINUXCON

containercon

CLLOUDOPEN

CHINA 中国

THINK OPEN

开放性思维

# Conclusion

# Conclusion & Future work

- **Conclusion**
  - In this presentation, we introduce SPDK vhost solution(i.e., SCSI/Blk/NVMe) to accelerate NVMe I/Os in virtual machines
- **Future work**
  - VM live migration support for the whole SPDK vhost solution(i.e., vhost SCSI/BLK/NVMe)
  - Upstream QEMU vhost driver.
- **Promotion**
  - Welcome to evaluate & use SPDK vhost target !
  - Welcome to contribute to SPDK community !



LINUXCON

containercon



CLOUDOPEN

CHINA 中国

THINK OPEN

开放性思维

Q & A



LINUXCON

containercon



CLOUDOPEN

CHINA 中国

THINK OPEN

开放性思维