



SPDK China Summit 2018

# SPDK PROGRAMMING FRAMEWORK AND NVME-OF OPTIMIZATION

Ziye Yang

Senior Software Engineer

Network Platforms Group, Intel Corporation

# Agenda

- SPDK programming framework
- Accelerated NVMe-oF via SPDK
- Conclusion

# Agenda

- **SPDK programming framework**
- Accelerated NVMe-oF via SPDK
- Conclusion

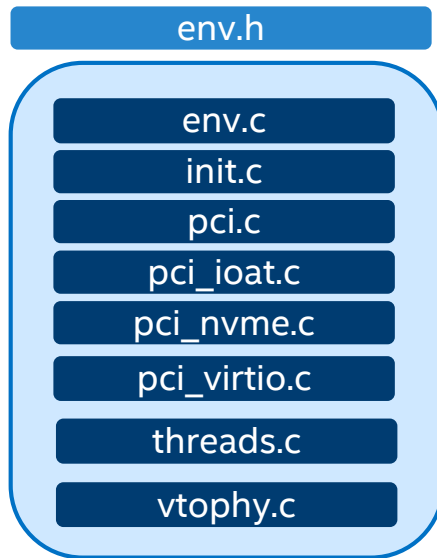
# SPDK ENVIRONMENT ABSTRACTION

# WHY AN ENVIRONMENT ABSTRACTION?

**FLEXIBILITY FOR USER**

# ENVIRONMENT ABSTRACTION

- Memory allocation (pinned for DMA) and address translation
- PCI enumeration and resource mapping
- Thread startup (pinned to cores)
- Lock-free ring and memory pool data structures



# ENVIRONMENT ABSTRACTION

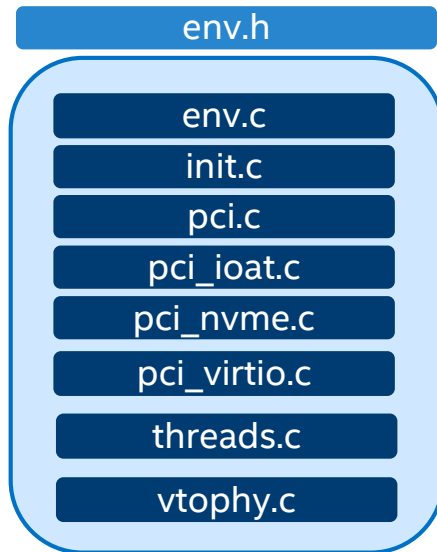
Configurable:

```
./configure --with-env=...
```

Interface defined in `spdk/env.h`

Default implementation uses **DPDK**  
(`lib/env_dpdk`)

**FLEXIBILITY: DECOUPLING AND DPDK ENHANCEMENTS**



# APPLICATION FRAMEWORK



**HOW DO WE COMBINE SPDK COMPONENTS?**

**THE SPDK APP FRAMEWORK PROVIDES THE GLUE**

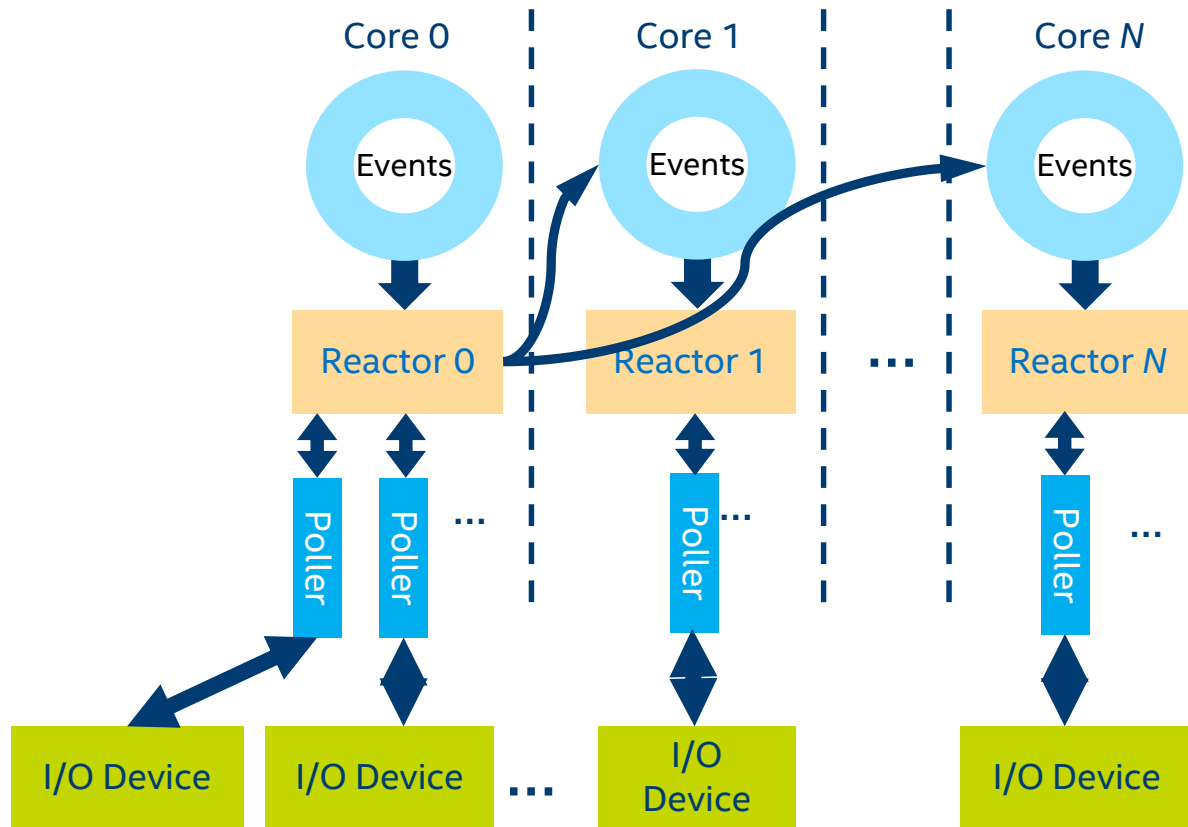
# APP FRAMEWORK COMPONENTS

REACTOR

POLLER

EVENT

I/O CHANNEL



# POLLER

Essentially a “task” running on a reactor

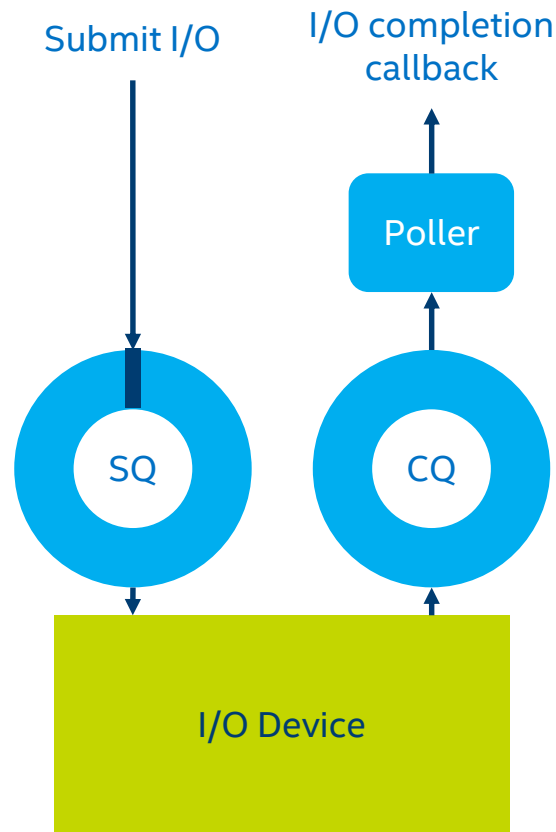
Primarily checks hardware for async events

Can run periodically on a timer

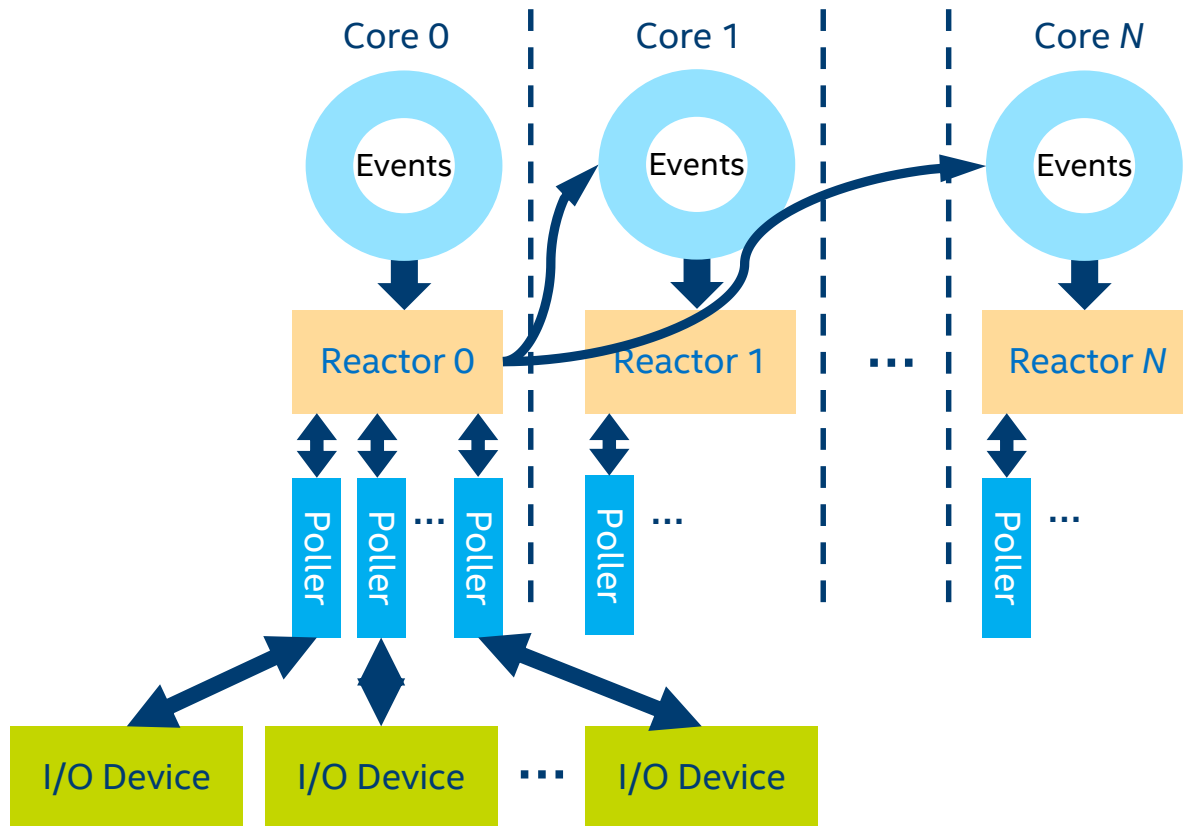
Example: poll completion queue

Callback runs to completion on reactor thread

Completion handler may send an event



# EVENT



# EVENT

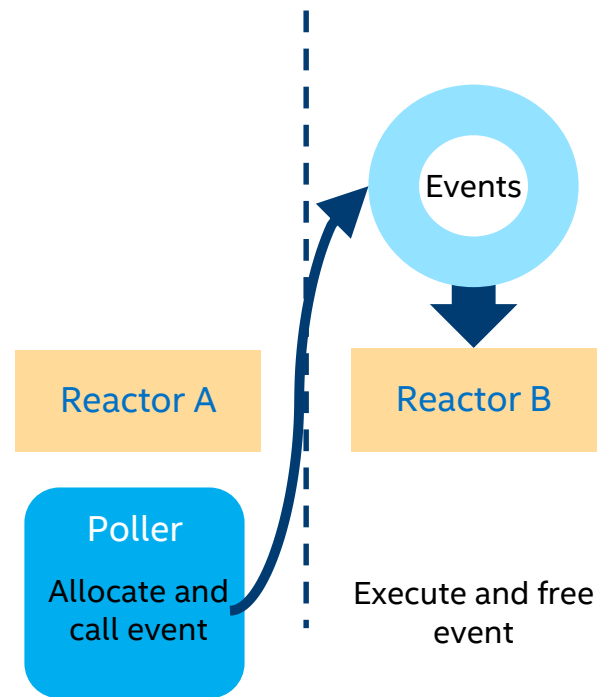
Cross-thread communication

Function pointer + arguments

One-shot message passed between reactors

Multi-producer/single-consumer ring

Runs to completion on reactor thread



# I/O CHANNEL

Abstracts hardware I/O queues

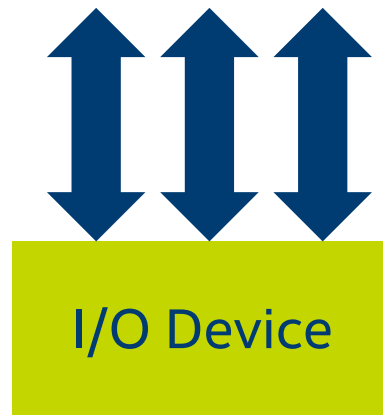
Register I/O devices

Create I/O channel per thread/device combination

Provides hooks for driver resource allocation

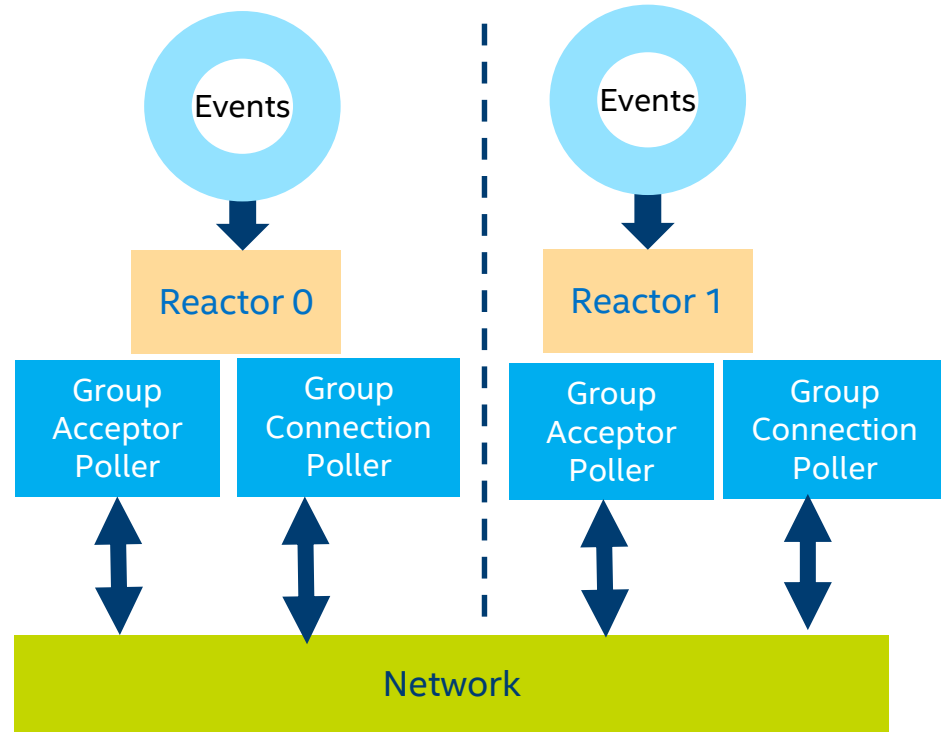
I/O channel creation drives poller creation

Pervasive in SPDK



# NVME OVER FABRICS TARGET EXAMPLE

- `nvme_tgt_advance_state`
  - `spdk_nvme_parse_conf` (listen on transport)
  - NVMe-oF tgt I/O channel creation:  
`spdk_nvme_tgt_create`
  - Group data poller creation in each core:  
Trigger the `create_cb` (`spdk_nvme_tgt_create_poll_group`) of I/O channel, then we will have `spdk_nvme_poll_group_poll` in each core
  - Group Acceptor network poller creation:  
`spdk_nvme_tgt_accept` will be used to connect events in each core





# NVME OVER FABRICS TARGET EXAMPLE

- Group Acceptor network poller handles connect events
- New qpair (connection) is allocated to different cores via Round Robin manner. Asynchronous message passing is used, then `spdk_nvme_poll_group_add` is called.
- I/O request arrives over network, and handled by the group poller in the designated core.
- I/O submitted to storage
- Storage device poller checks completions
- Response sent

**ALL ASYNCHRONOUS WORK IS DRIVEN BY POLLERS**

# Agenda

- SPDK programming framework
- Accelerated NVMe-oF via SPDK
- Conclusion

# SPDK NVMe-oF Components

## NVMe over Fabrics Target

- Released July 2016 (with spec)
- **Hardening:**
  - Intel test infrastructure
  - Discovery simplification
  - Correctness & kernel interop
- **Performance improvements:**
  - Read latency improvement
  - Scalability validation (up to 150Gbps)
  - Event Framework enhancements
  - Multiple connection performance improvement (e.g., group transport polling,)

## NVMe over Fabrics Host (Initiator)

- New component added in Dec 2016
- Performance improvements
  - Eliminate copy: now true zero-copy
  - SGL (single SGL element)

# SPDK NVMe-oF transport work

## Existing work: RDMA transport

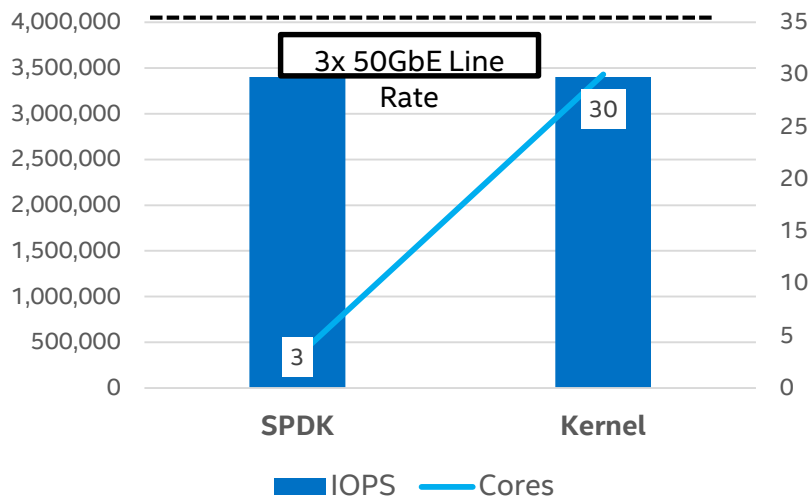
- **DPDK components used which is encapsulated in libspdk\_env\_dpdk.a, e.g.,**
  - PCI device management
  - CPU/thread scheduling
  - Memory management (e.g., lock free rings)
  - Log management

## Upcoming work: TCP transport

- Kernel based TCP transport
- VPP/DPDK based user space TCP transport
  - Use DPDK Ethernet PMDs
  - Use user space TCP/IP stack (e.g., VPP)

# NVMe-oF Target Throughput Performance (RDMA)

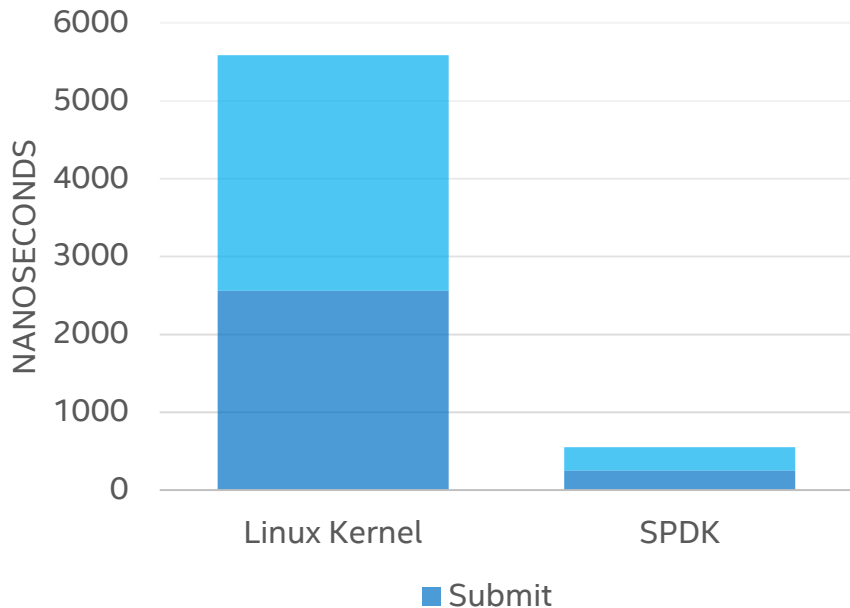
SPDK vs. Kernel NVMe-oF I/O Efficiency



NVMe* over Fabrics Target Features	Realized Benefit
Utilizes NVM Express* (NVMe) Polled Mode Driver	Reduced overhead per NVMe I/O
RDMA Queue Pair group Polling	No interrupt overhead
Connections pinned to CPU cores	No synchronization overhead

SPDK reduces NVMe over Fabrics software overhead up to 10x!

# NVM Express\* Driver Software Overhead

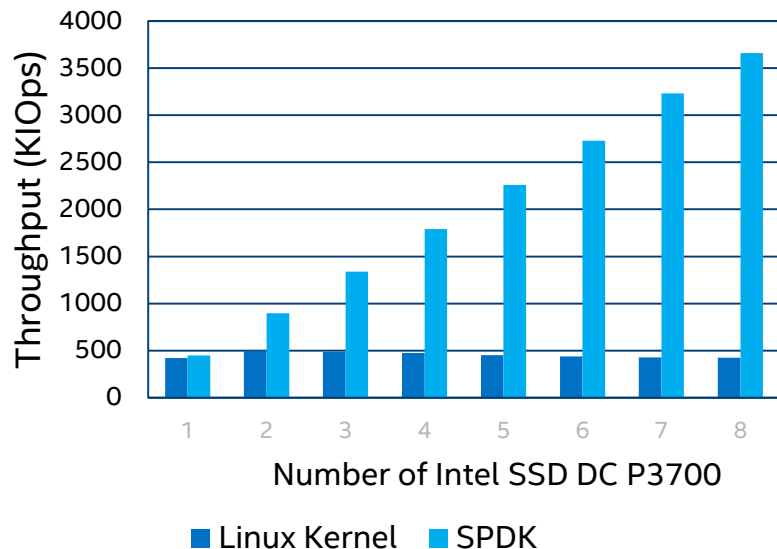


Kernel Source of Overhead	SPDK Approach
Interrupts	Asynchronous Polled Mode
Synchronization	Lockless
System Calls	User Space Hardware Access
DMA Mapping	Hugepages
Generic Block Layer	Specific for Flash Latencies

SPDK reduces NVM Express\* (NVMe) software overhead up to 10x!

# NVM Express\* Driver Throughput Scalability

I/O Performance on  
Single Intel® Xeon® core

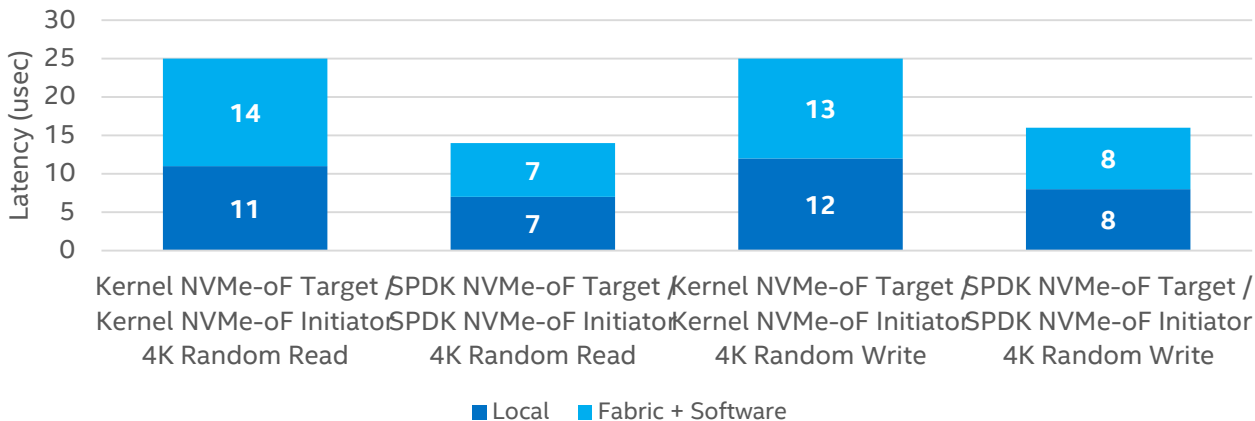


- Systems with multiple NVM Express\* (NVMe) SSDs capable of millions of I/O per second
- Results in many cores of software overhead with kernel-based interrupt-driven driver model
- SPDK enables:
  - more CPU cycles for storage services
  - lower I/O latency

**SPDK saturates 8 NVMe SSDs with a single CPU core!**

# SPDK Host + Target vs. Kernel Host + Target

Avg. I/O Round Trip Time  
Kernel vs. SPDK NVMe-oF Stacks  
Coldstream, Perf, qd=1



SPDK reduces Optane NVMe-oF latency by 44%, write latency by 36%!



# Agenda

- SPDK programming framework
- Accelerated NVMe-oF via SPDK
- **Conclusion**

# Conclusion

- In this presentation, we introduce
  - SPDK library
  - The accelerated NVMe-oF target built from SPDK library
- SPDK proves to be useful to accelerate storage applications equipped with NVMe based devices
- Call for action:
  - Welcome to use SPDK in storage area (similar as using DPDK in network ) and contribute into SPDK community.

**Q&A**

# 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.

