



OPEN SOURCE SUMMIT

China 2019





SPDK BASED USER SPACE NVME OVER TCP TRANSPORT SOLUTION

Presenters: Ziye Yang

Company: Intel

Notices and 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.

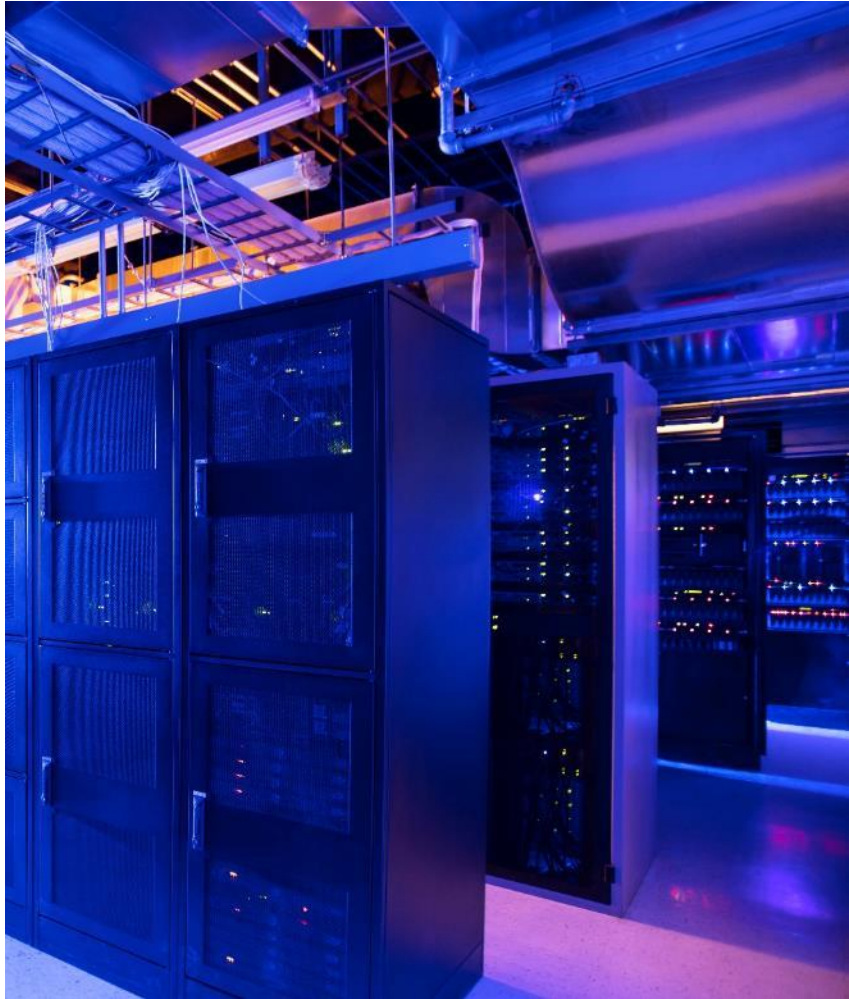
© 2018 Intel Corporation.

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.

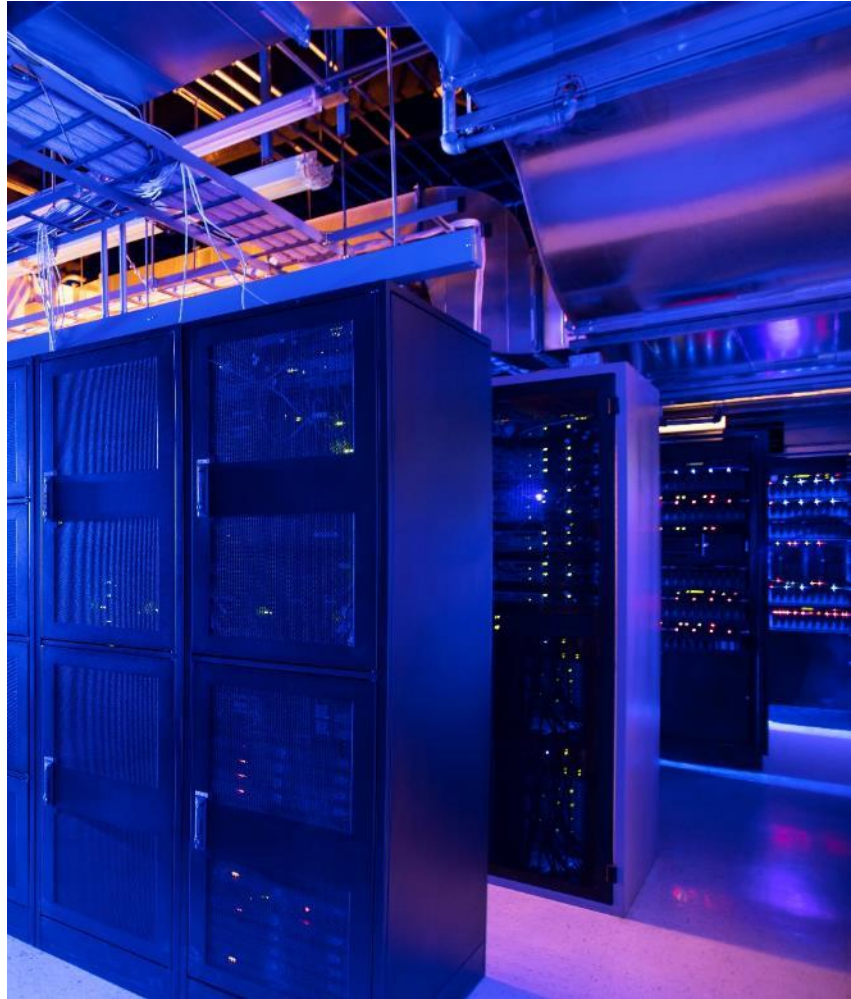
AGENDA

- SPDK NVMe-oF development history & status
- SPDK TCP transport introduction
- Conclusion

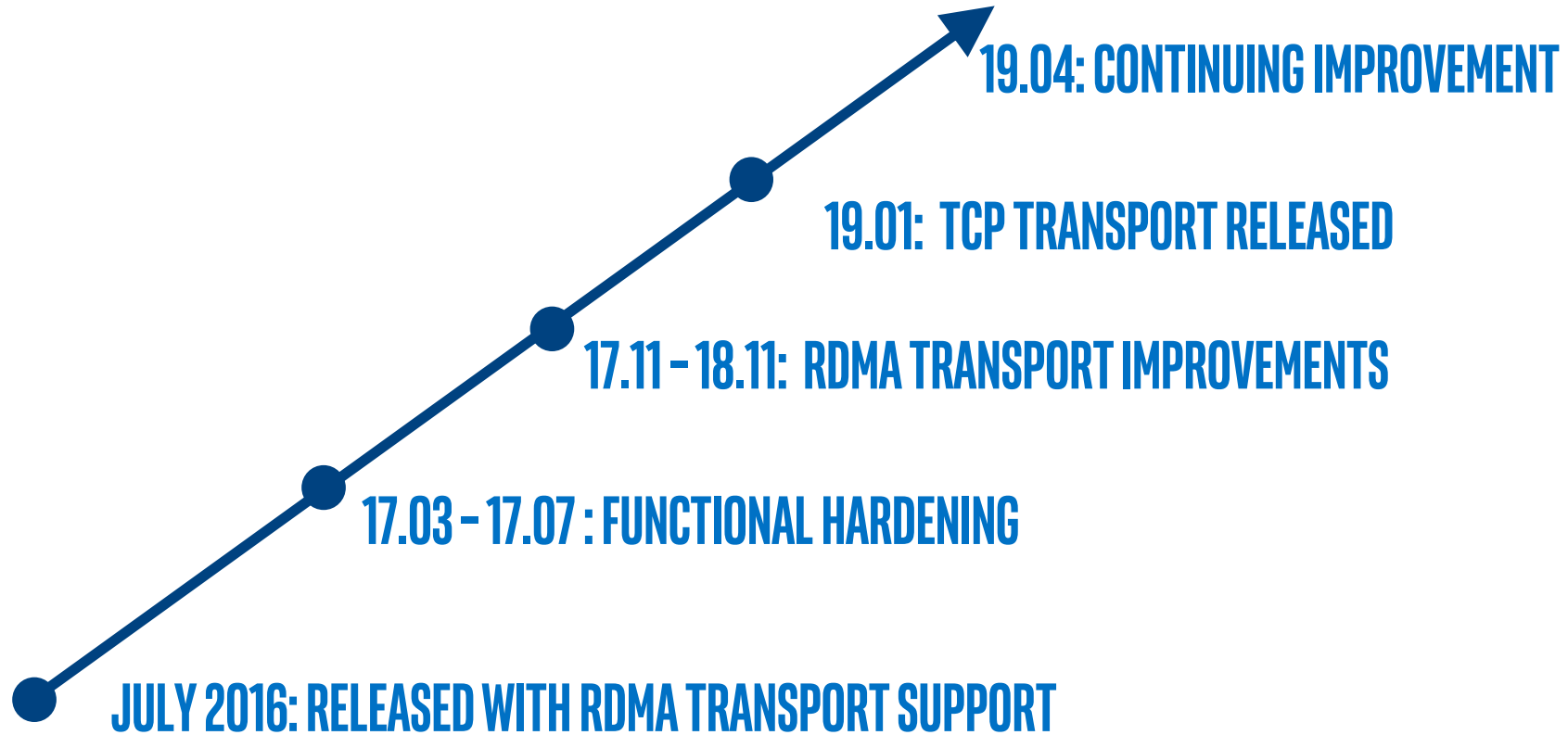


AGENDA

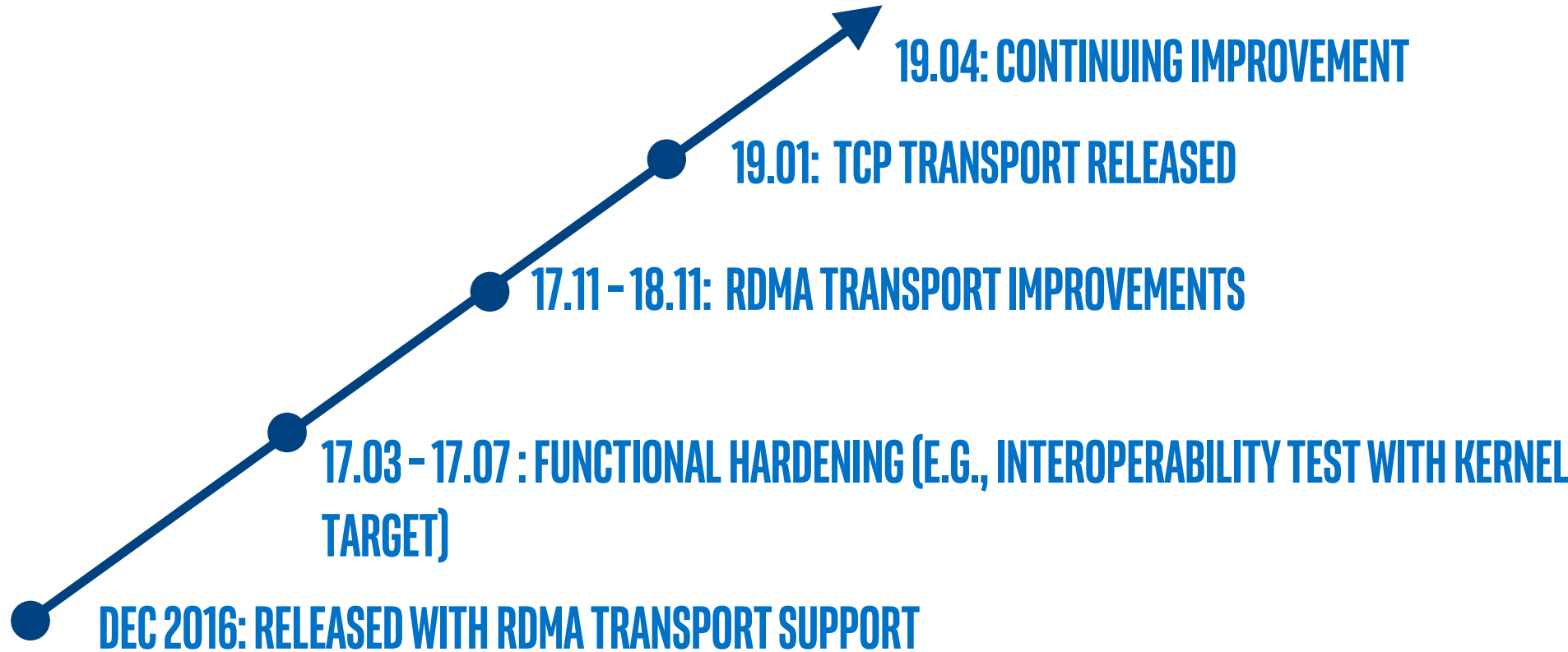
- **SPDK NVMe-oF development history & status**
- SPDK TCP transport introduction
- Conclusion



SPDK NVMe-oF Target Timeline



SPDK NVMe-oF Host Timeline

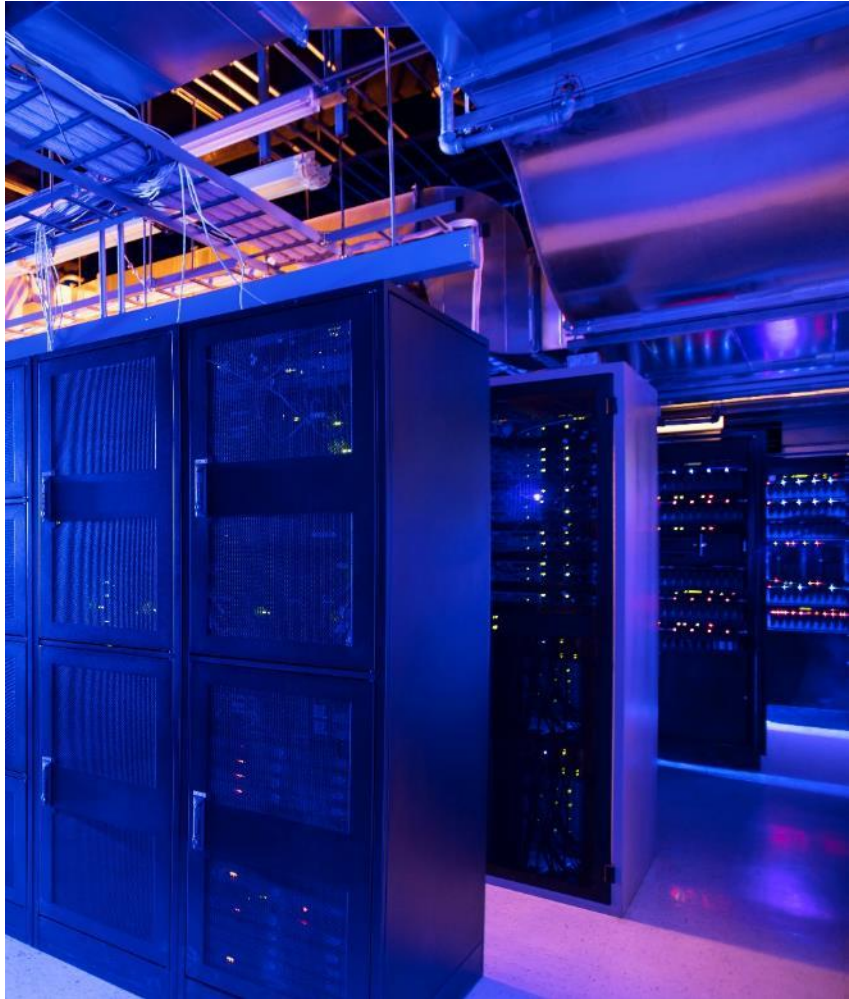


SPDK NVMe-oF target design highlights

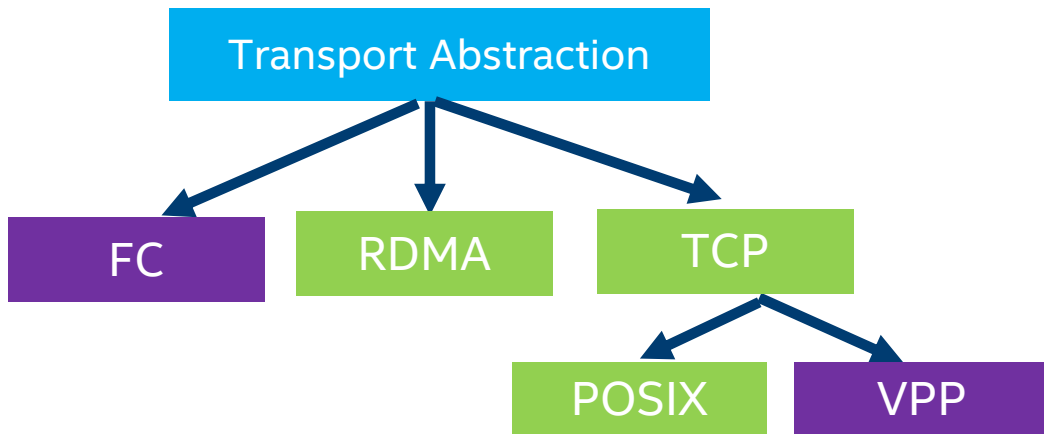
NVMe* over Fabrics Target Features	Performance Benefit
Utilizes user space NVM Express* (NVMe) Polled Mode Driver	Reduced overhead per NVMe I/O
Group polling on each SPDK thread (binding on CPU core) for multiple transports	No interrupt overhead
Connections pinned to dedicated SPDK thread	No synchronization overhead
Asynchronous NVMe CMD handling in whole life cycle	No locks in NVMe CMD data handling path

AGENDA

- SPDK NVMe-oF Development history & status
- **SPDK TCP transport introduction**
- Conclusion



General design and implementation



- Follow the SPDK transport abstraction:

- Host side code:
`lib/nvme/nvme_tcp.c`
- Target side code:
`lib/nvmf/tcp.c`



Already released



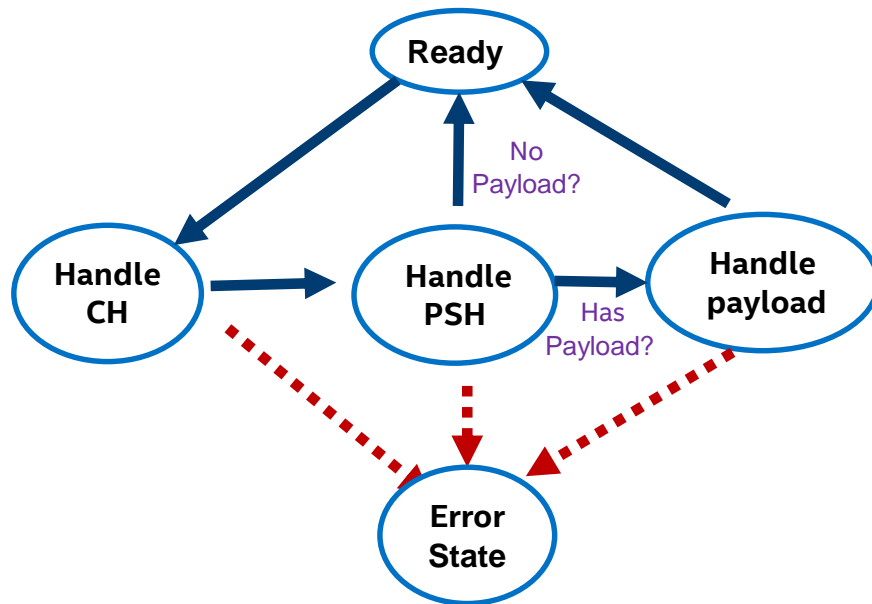
Inprogress

Performance design consideration for TCP transport in target side

Ingredients	Methodology
Design framework	Follow the general SPDK NVMe-oF framework (e.g., polling group)
TCP connection optimization	Use the SPDK encapsulated Socket API (preparing for integrating other stack, e.g., VPP)
NVMe/TCP PDU handling	Use state machine to track
NVMe/TCP request life time cycle	Use state machine to track (Purpose: Easy to debug and good for further performance improvement)

TCP PDU Receiving handling for each connection

```
enum nvme_tcp_pdu_rcv_state {  
    /* Ready to wait PDU */  
    NVME_TCP_PDU_RECV_STATE_AWAIT_PDU_READY,  
  
    /* Active tpair waiting for any PDU common header */  
    NVME_TCP_PDU_RECV_STATE_AWAIT_PDU_CH,  
  
    /* Active tpair waiting for any PDU specific header */  
    NVME_TCP_PDU_RECV_STATE_AWAIT_PDU_PSH,  
  
    /* Active tpair waiting for payload */  
    NVME_TCP_PDU_RECV_STATE_AWAIT_PDU_PAYLOAD,  
  
    /* Active tpair does not wait for payload */  
    NVME_TCP_PDU_RECV_STATE_ERROR,  
};
```



Error Path→

SPDK NVMe-oF TCP request life cycle of each connection in target side

```
/* spdk nvme related structure */
enum spdk_nvme_tcp_req_state {

    /* The request is not currently in use */
    TCP_REQUEST_STATE_FREE = 0,

    /* Initial state when request first received */
    TCP_REQUEST_STATE_NEW,

    /* The request is queued until a data buffer is available. */
    TCP_REQUEST_STATE_NEED_BUFFER,

    /* The request is currently transferring data from the host to the controller. */
    TCP_REQUEST_STATE_TRANSFERRING_HOST_TO_CONTROLLER,

    /* The request is ready to execute at the block device */
    TCP_REQUEST_STATE_READY_TO_EXECUTE,

    /* The request is currently executing at the block device */
    TCP_REQUEST_STATE_EXECUTING,

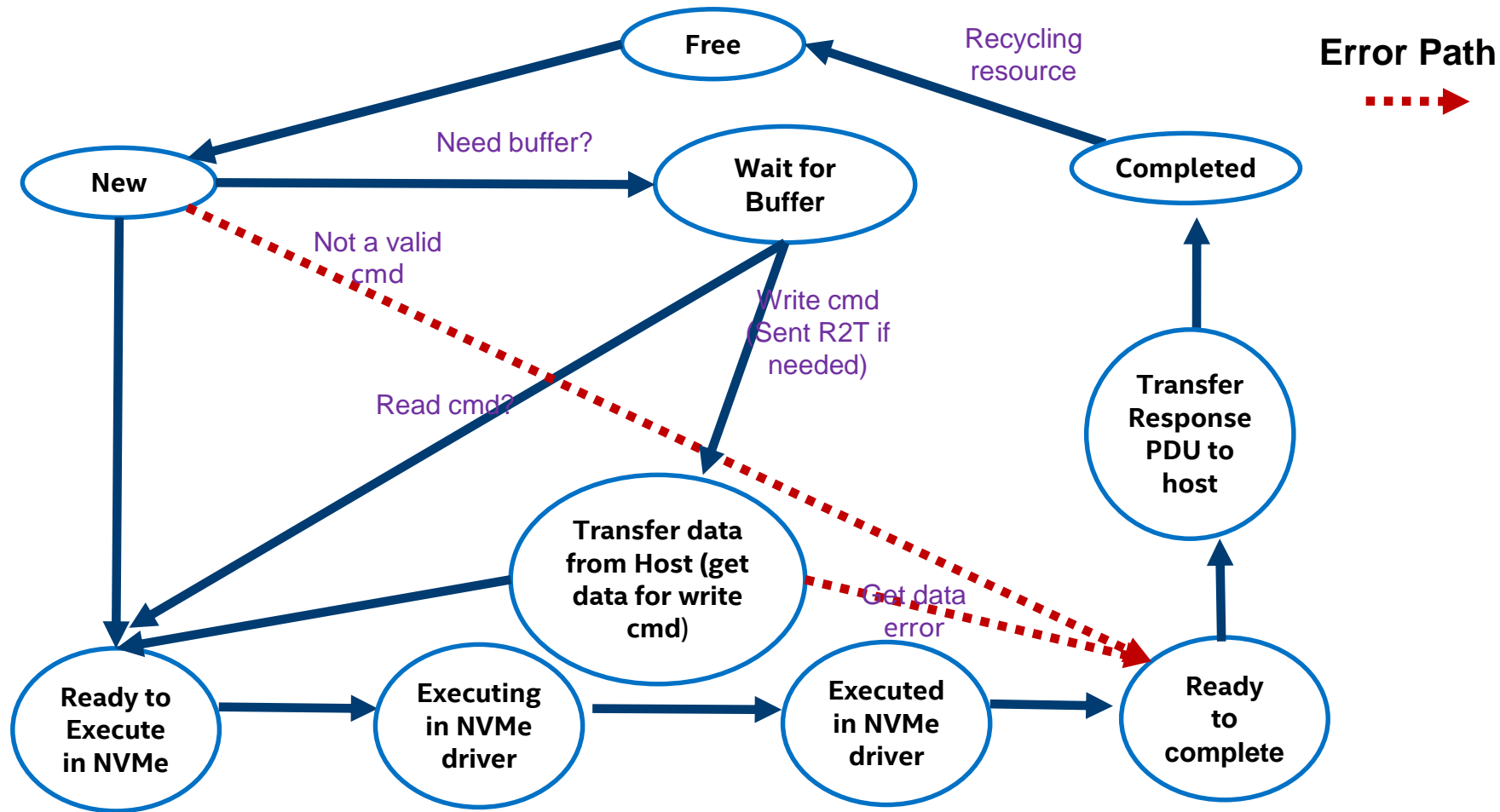
    /* The request finished executing at the block device */
    TCP_REQUEST_STATE_EXECUTED,

    /* The request is ready to send a completion */
    TCP_REQUEST_STATE_READY_TO_COMPLETE,

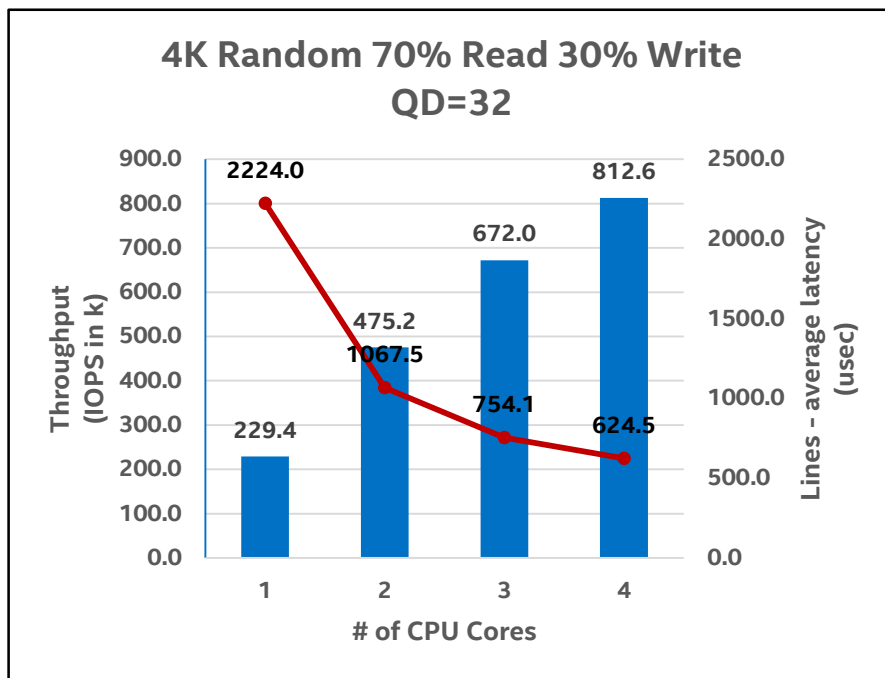
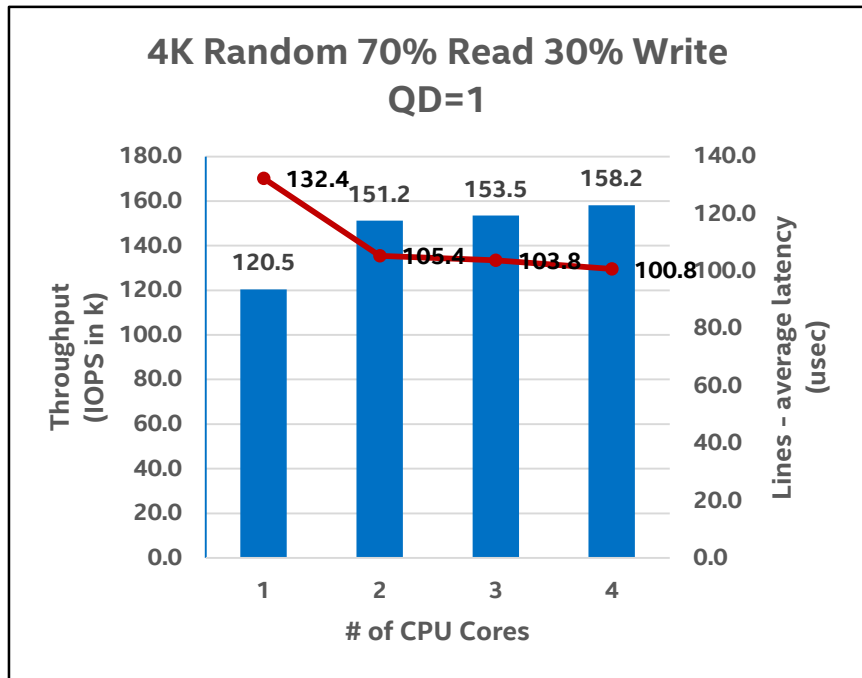
    /* The request is currently transferring final pdus from the controller to the host. */
    TCP_REQUEST_STATE_TRANSFERRING_CONTROLLER_TO_HOST,

    /* The request completed and can be marked free. */
    TCP_REQUEST_STATE_COMPLETED,

    /* Terminator */
    TCP_REQUEST_NUM_STATES,
};
```

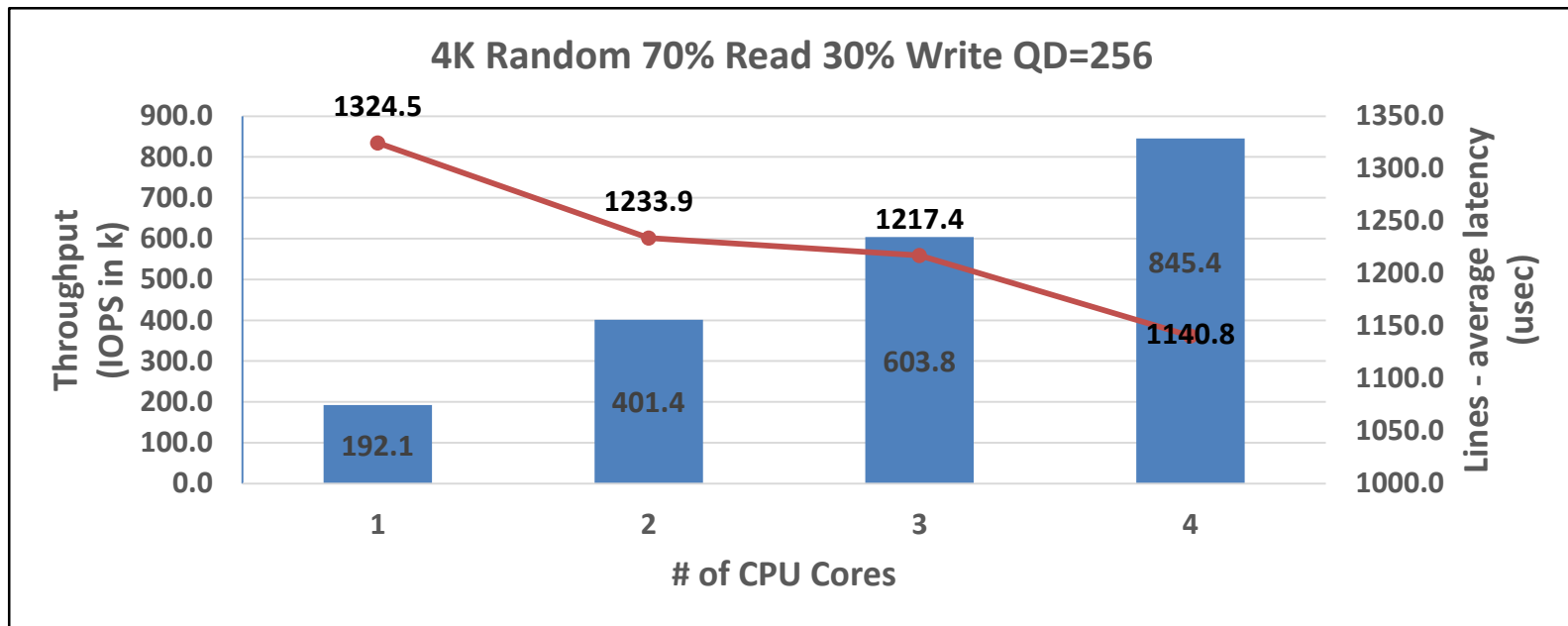


SPDK TARGET SIDE (TCP TRANSPORT): I/O SCALING



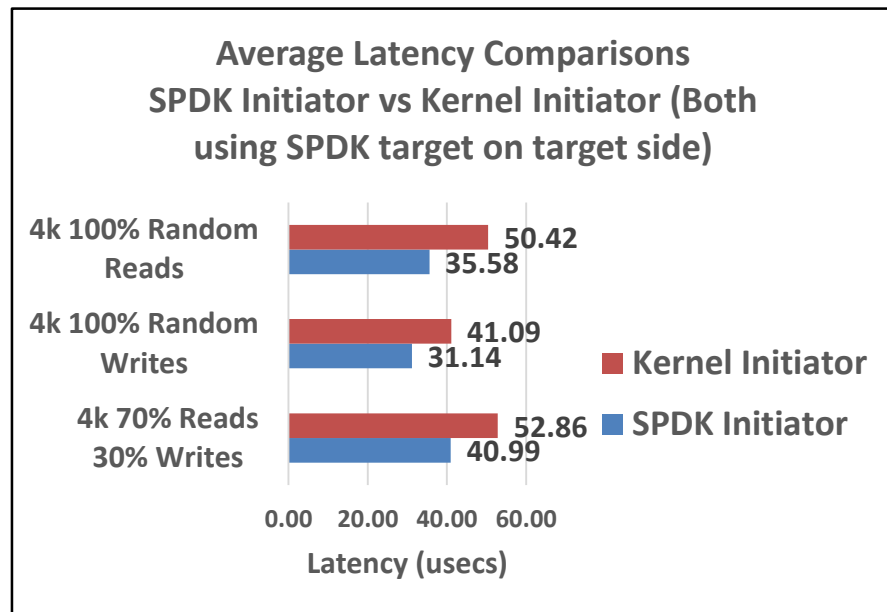
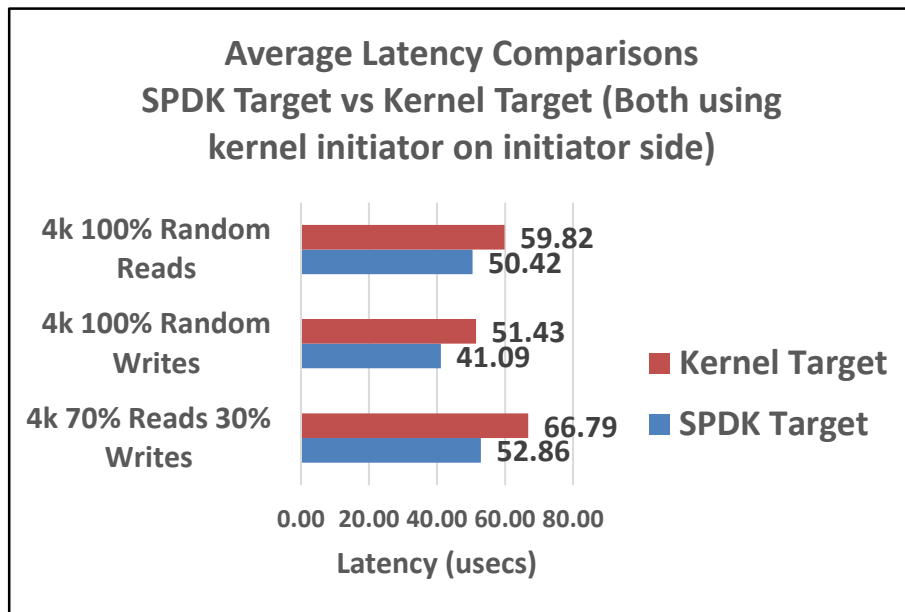
System configuration: (1) Target: server platform: SuperMicro SYS2029U-TN24R4T; 2x Intel® Xeon® Platinum 8180 CPU @ 2.50 GHz, Intel® Speed Step enabled, Intel® Turbo Boost Technology enabled, 4x 2GB DDR4 2666 MT/s, 1 DIMM per channel; 2x 100GbE Mellanox ConnectX-5 NICs; Fedora 28, Linux kernel 5.05, SPDK 19.01.1; 6x Intel® P4600TM P4600x 2.0TB; (2) initiator: Server platform: SuperMicro SYS-2028U TN24R4T+; 44x Intel(R) Xeon(R) CPU E5-2699 v4 @ 2.20GHz (HT off); 1x 100GbE Mellanox ConnectX-4 NIC; Fedora 28, Linux kernel 5.05, SPDK 19.0.1. (3): Fio ver: fio-3.3; Fio workload: blocksize=4k, iodepth=1, iodepth_batch=128, iodepth_low=256, ioengine=libaio or SPDK bdev engine, size=10G, ramp_time=0, run_time=300, group_reporting, thread, direct=1, rw=read/rw/randread/randwrite/randrw

SPDK HOST SIDE (TCP TRANSPORT): I/O SCALING



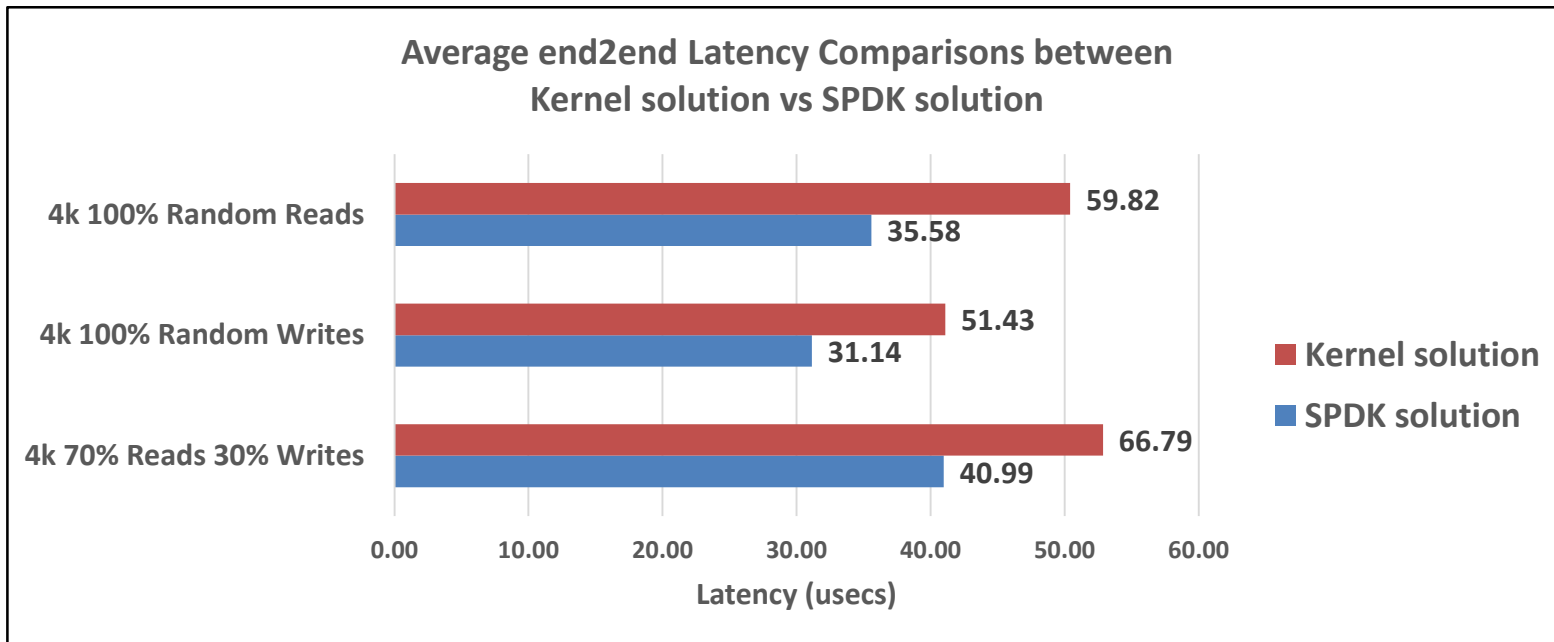
System configuration: (1) Target: server platform: SuperMicro SYS2029U-TN24R4T; 2x Intel® Xeon® Platinum 8180 CPU @ 2.50 GHz, Intel® Speed Step enabled, Intel® Turbo Boost Technology enabled, 4x 2GB DDR4 2666 MT/s, 1 DIMM per channel; 2x 100GbE Mellanox ConnectX-5 NICs; Fedora 28, Linux kernel 5.05, SPDK 19.01.1; 6x Intel® P4600TM P4600x 2.0TB; (2) initiator: Server platform: SuperMicro SYS-2028U TN24R4T+; 44x Intel(R) Xeon(R) CPU E5-2699 v4 @ 2.20GHz (HT off); 1x 100GbE Mellanox ConnectX-4 NIC; Fedora 28, Linux kernel 5.05, SPDK 19.0.1. (3): Fio ver: fio-3.3; Fio workload: blocksize=4k, iodepth=1, iodepth_batch=128, iodepth_low=256, ioengine=libaio or SPDK bdev engine, size=10G, ramp_time=0, run_time=300, group_reporting, thread, direct=1, rw=read/write/rw/randread/randwrite/randrw

LATENCY COMPARISON BETWEEN SPDK AND KERNEL (NULL BDEV IS USED)



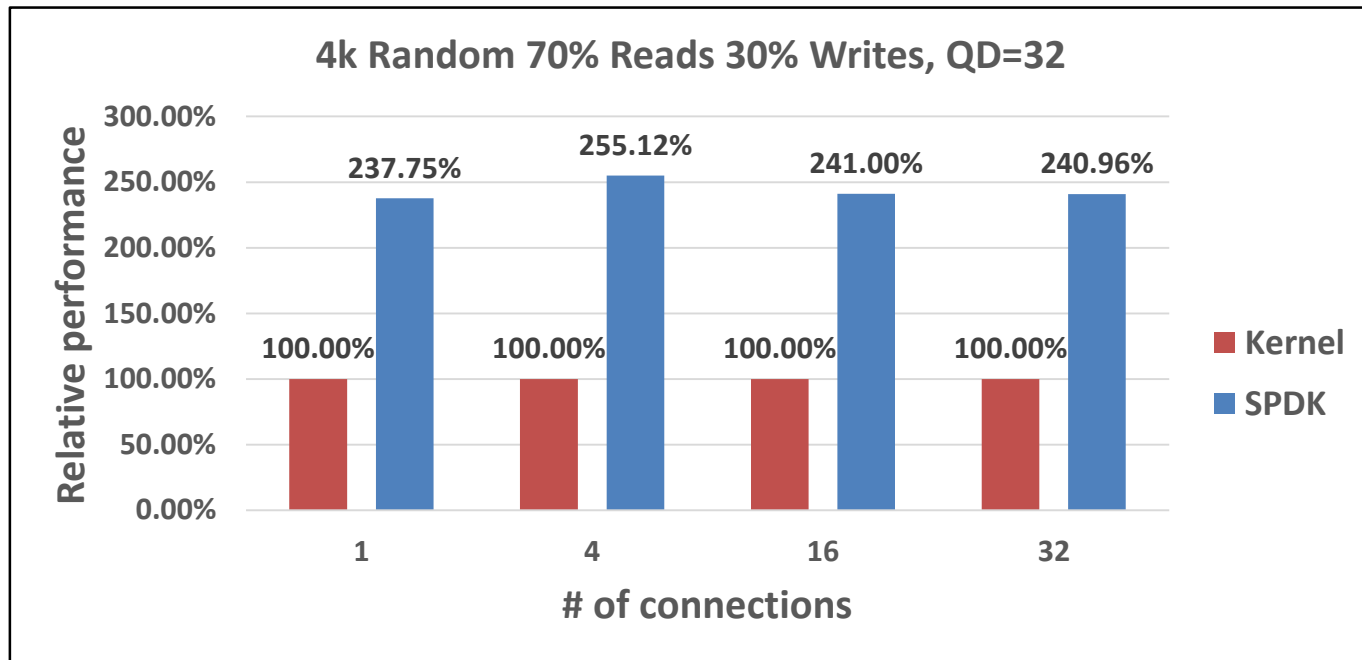
System configuration: (1) Target: server platform: SuperMicro SYS2029U-TN24R4T; 2x Intel® Xeon® Platinum 8180 CPU @ 2.50 GHz, Intel® Speed Step enabled, Intel® Turbo Boost Technology enabled, 4x 2GB DDR4 2666 MT/s, 1 DIMM per channel; 2x 100GbE Mellanox ConnectX-5 NICs; Fedora 28, Linux kernel 5.05, SPDK 19.01.1; 6x Intel® P4600TM P4600x 2.0TB; (2) initiator: Server platform: SuperMicro SYS-2028U TN24R4T+; 44x Intel(R) Xeon(R) CPU E5-2699 v4 @ 2.20GHz (HT off); 1x 100GbE Mellanox ConnectX-4 NIC; Fedora 28, Linux kernel 5.05, SPDK 19.0.1. (3) : Fio ver: fio-3.3; Fio workload: blocksize=4k, iodepth=1, iodepth_batch=128, iodepth_low=256, ioengine=libaio or SPDK bdev engine, size=10G, ramp_time=0, run_time=300, group_reporting, thread, direct=1, rw=read/write/rw/randread/randwrite/randrw

LATENCY COMPARISON BETWEEN SPDK AND KERNEL (NULL BDEV IS USED)



System configuration: (1) Target: server platform: SuperMicro SYS2029U-TN24R4T; 2x Intel® Xeon® Platinum 8180 CPU @ 2.50 GHz, Intel® Speed Step enabled, Intel® Turbo Boost Technology enabled, 4x 2GB DDR4 2666 MT/s, 1 DIMM per channel; 2x 100GbE Mellanox ConnectX-5 NICs; Fedora 28, Linux kernel 5.05, SPDK 19.01.1; 6x Intel® P4600TM P4600x 2.0TB; (2) initiator: Server platform: SuperMicro SYS-2028U TN24R4T+; 44x Intel(R) Xeon(R) CPU E5-2699 v4 @ 2.20GHz (HT off); 1x 100GbE Mellanox ConnectX-4 NIC; Fedora 28, Linux kernel 5.05, SPDK 19.0.1. (3) : Fio ver: fio-3.3; Fio workload: blocksize=4k, iodepth=1, iodepth_batch=128, iodepth_low=256, ioengine=libaio or SPDK bdev engine, size=10G, ramp_time=0, run_time=300, group_reporting, thread, direct=1, rw=read/write/rw/randread/randwrite/randrw

IOPS/CORE COMPARISON BETWEEN SPDK AND KERNEL ON TARGET SIDE



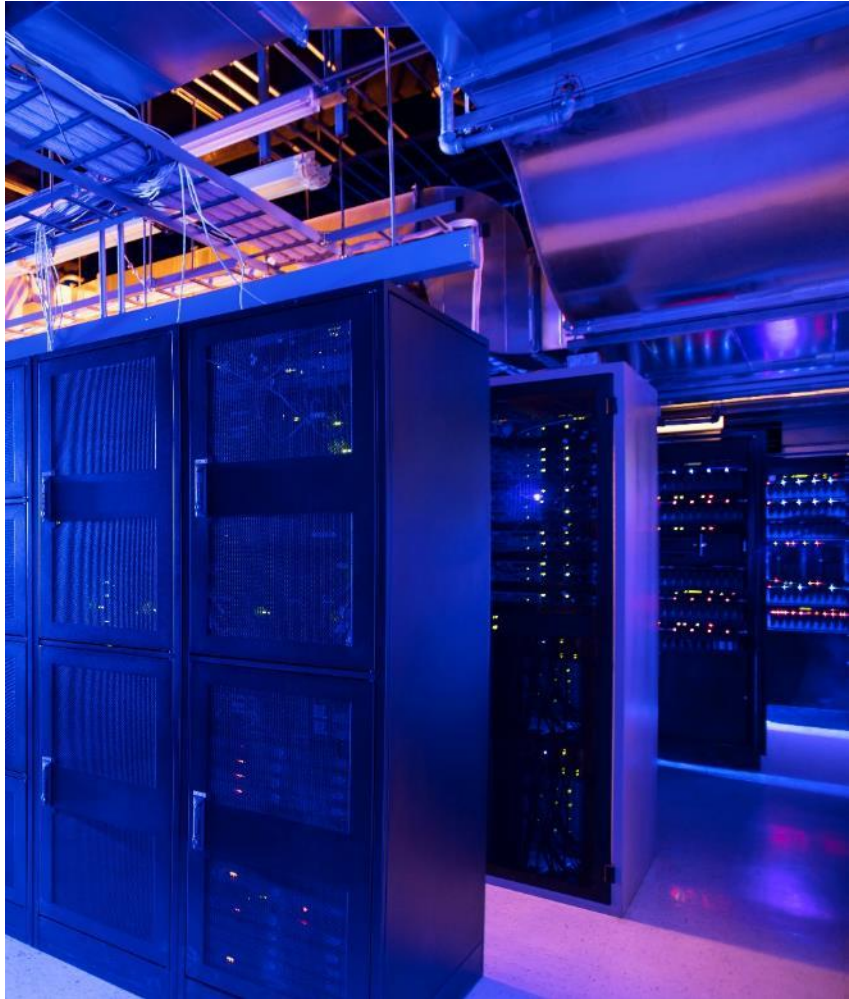
System configuration: (1) Target: server platform: SuperMicro SYS2029U-TN24R4T; 2x Intel® Xeon® Platinum 8180 CPU @ 2.50 GHz, Intel® Speed Step enabled, Intel® Turbo Boost Technology enabled, 4x 2GB DDR4 2666 MT/s, 1 DIMM per channel; 2x 100GbE Mellanox ConnectX-5 NICs; Fedora 28, Linux kernel 5.05, SPDK 19.01.1; 6x Intel® P4600TM P4600x 2.0TB; (2) initiator: Server platform: SuperMicro SYS-2028U TN24R4T+; 44x Intel(R) Xeon(R) CPU E5-2699 v4 @ 2.20GHz (HT off); 1x 100GbE Mellanox ConnectX-4 NIC; Fedora 28, Linux kernel 5.05, SPDK 19.0.1. (3): Fio ver: fio-3.3; Fio workload: blocksize=4k, iodepth=1, iodepth_batch=128, iodepth_low=256, ioengine=libaio or SPDK bdev engine, size=10G, ramp_time=0, run_time=300, group_reporting, thread, direct=1, rw=read/write/rw/randread/randwrite/randrw

Further development plan

- Continue enhancing the functionality
 - Including the compatible test with Linux kernel solution.
- Performance tuning
- Integration with third party software
 - Deep integration with user space stack: VPP + DPDK
- Leveraging hardware features
 - Use existing hardware features of NICs for performance improvement, e.g., VMA from Mellanox's NIC; load balance features (ADQ) from Intel's 100Gbit NIC.
 - Figuring out offloading methods with hardware, e.g., FPGA, Smart NIC, and etc.

AGENDA

- SPDK NVMe-oF Development history & status
- SPDK TCP transport introduction
- **Conclusion**



Conclusion

- SPDK NVMe-oF solution is well adopted by the industry. In this presentation, followings are introduced, i.e.,
 - The development status of SPDK NVMe-oF solution
 - SPDK TCP transport development status.
- Further development
 - Continue following the NVMe-oF spec and adding more features.
 - Continue performance enhancements and integration with other solutions.
- Call for activity in community
 - Welcome to bug submission, idea discussion and patch submission for NVMe-oF



Q&A

