

SPDK NVMe-oF RDMA (Target & Initiator)

Performance Report

Release 23.01

Testing Date: February 2023

Performed by:

Karol Latecki ([karol.latecki@intel.com](mailto:karel.latecki@intel.com))

Jaroslaw Chachulski (jaroslawx.chachulski@intel.com)

Acknowledgments:

James Harris (james.r.harris@intel.com)

Krzysztof Karas (krzysztof.karas@intel.com)

Contents

Contents	2
Audience and Purpose	3
Test setup	4
Target Configuration	4
Initiator 1 Configuration	5
Initiator 2 Configuration	5
BIOS settings	6
SPDK Build Options	6
Introduction to SPDK NVMe-oF (Target & Initiator)	7
Test Case 1: SPDK NVMe-oF RDMA Target I/O core scaling	9
4KiB Random Read Results	12
4KiB Random Write Results	13
Large Sequential I/O Performance	15
Conclusions	16
Test Case 2: SPDK NVMe-oF RDMA Initiator I/O core scaling	17
4KiB Random Read Results	19
4KiB Random Write Results	20
4KiB Random 70/30 Read/Write Results	21
Conclusions	22
Test Case 3: Linux Kernel vs. SPDK NVMe-oF RDMA Latency	23
SPDK vs Kernel NVMe-oF RDMA Target Results	26
SPDK vs Kernel NVMe-oF RDMA Initiator Results	27
SPDK vs Kernel NVMe-oF RDMA Kernel + Initiator Results	28
Conclusions	29
Test Case 4: NVMe-oF RDMA Performance with increasing # of connections	30
4KiB Random Read Results	32
4KiB Random Write results	33
4KiB Random Read-Write Results	34
Conclusions	35
Summary	36
List of tables	37
List of figures	39
Appendix A – Test Case 1 & 2 SPDK NVMe-oF Initiator bdev configuration	40
Appendix B – Test Case 3 SPDK NVMe-oF Initiator bdev configuration	44
Appendix C - Kernel NVMe-oF RDMA Target configuration	44

Audience and Purpose

This report is intended for people who are interested in evaluating SPDK NVMe-oF (Target & Initiator) performance as compared to the Linux Kernel NVMe-oF (Target & Initiator). This report contains performance and efficiency of the SPDK vs. Linux Kernel NVMe-oF Target and Initiator for the RDMA transport only.

The purpose of report is not to imply a single “correct” approach, but rather to provide a baseline of well-tested configurations and procedures that produce repeatable results. This report can also be viewed as information regarding best known method/practice when performance testing SPDK NVMe-oF Target and Initiator components.

Test setup

Target Configuration

Table 1: Hardware setup configuration – Target system

Item	Description
Server Platform	<u>SuperMicro® Ultra SuperServer SYS-220U-TNR</u> 
Motherboard	Server board X12DPU-6
CPU	2 CPU sockets, Intel(R) Xeon(R) Gold 6348 CPU @ 2.60GHz Number of cores 28 per socket, number of threads 56 per socket Both sockets populated Microcode: 0xd000375
Memory	16 x 32GB SK Hynix HMA84GR7DJR4N-XN, DDR4, 3200MHz Total of 512GB
Operating System	Fedora 37
BIOS	1.4a
Linux kernel version	6.0.18-300.fc37.x86_64 Spectre-meltdown mitigations enabled
SPDK version	SPDK 23.01
Storage	OS: 1x 250GB Crucial CT250MX500SSD1 Storage Target: 14x Kioxia® KCM61VUL3T20 3.2TBs (FW: 0105) (6 on CPU NUMA Node 0, 8 on CPU NUMA Node 1)
NIC	4x 100GbE Intel(R) Ethernet Network Adapter E810-CQDA2. Single port connected on each NIC. ice driver version: 1.10.1.2.2 NVM FW version: v4.20 2 NICs per CPU socket. Protocol: RoCEv2

Initiator 1 Configuration

Table 2: Hardware setup configuration – Initiator system 1

Item	Description
Server Platform	Intel® Server System M50CYP2UR208
CPU	Intel® Xeon® Gold 6348 Processor @ 2.60GHz (42MB Cache) Number of cores 28 per socket, number of threads 56 per socket (Both sockets populated) Microcode: 0xd000375
Memory	16 x 32GB Micron 36ASF4G72PZ-3G2J3, DDR4, 3200MHz Total 512GBs
Operating System	Fedora 37
BIOS	SE5C620.86B.01.01.0007.2210270543
Linux kernel version	6.0.18-300.fc37.x86_64 Spectre-meltdown mitigations enabled
SPDK version	SPDK 23.01
Storage	OS: 1x 250GB Crucial CT250MX500SSD1
NIC	2x 100GbE Intel(R) Ethernet Network Adapter E810-CQDA2. Single port connected on each NIC. ice driver version: 1.10.1.2.2 NVM FW version: v4.20 Protocol: RoCEv2

Initiator 2 Configuration

Table 3: Hardware setup configuration – Initiator system 2

Item	Description
Server Platform	Intel® Server System M50CYP2UR208
CPU	Intel® Xeon® Gold 6348 Processor @ 2.60GHz (42MB Cache) Number of cores 28 per socket, number of threads 56 per socket (Both sockets populated) Microcode: 0xd000375
Memory	16 x 32GB Micron 36ASF4G72PZ-3G2J3, DDR4, 3200MHz Total 512GBs
Operating System	Fedora 37
BIOS	SE5C620.86B.01.01.0007.2210270543
Linux kernel version	6.0.18-300.fc37.x86_64 Spectre-meltdown mitigations enabled
SPDK version	SPDK 23.01
Storage	OS: 1x 250GB Crucial CT250MX500SSD1
NIC	2x 100GbE Intel(R) Ethernet Network Adapter E810-CQDA2. Single port connected on each NIC. ice driver version: 1.10.1.2.2 NVM FW version: v4.20

Protocol: RoCEv2

BIOS settings

Table 4: Test systems BIOS settings

Item	Description
BIOS <i>(Applied to all 3 systems)</i>	Hyper threading Enabled CPU Power and Performance Policy: <ul style="list-style-type: none">• “Extreme Performance” for Target• “Performance” for Initiators CPU C-state No Limit CPU P-state Enabled Enhanced Intel® SpeedStep® Tech Enabled Turbo Boost Enabled

SPDK Build Options

All measurements included in this report document were done with SPDK build with “—enable-lto” option enabled. Link time optimization allows better SPDK performance thanks to code optimization done by inlining functions across compilation units, which in turn results in reduced function call overhead.

Introduction to SPDK NVMe-oF (Target & Initiator)

The NVMe over Fabrics (NVMe-oF) protocol extends the parallelism and efficiencies of the NVMe Express* (NVMe) block protocol over network fabrics such as RDMA (iWARP, RoCE, InfiniBand™), Fibre Channel and TCP. SPDK provides both a user-space NVMe-oF target and initiator that extends the software efficiencies of the rest of the SPDK stack over the network. The SPDK NVMe-oF target uses the SPDK user-space, polled-mode NVMe driver to submit and complete I/O requests to NVMe devices which reduces the software processing overhead. Likewise, it pins connections to CPU cores to avoid synchronization and cache thrashing so that the data for those connections is kept as close to the CPU cache as possible.

The SPDK NVMe-oF target and initiator uses the Infiniband/RDMA verbs API to access an RDMA-capable NIC. These should work on all flavors of RDMA transports but for the purpose of this report document are tested against RoCEv2. Similar to the SPDK NVMe driver, SPDK provides a user-space, lockless, polled-mode NVMe-oF initiator. The host system uses the initiator to establish a connection and submit I/O requests to an NVMe subsystem within an NVMe-oF target. NVMe subsystems contain namespaces, each of which maps to a single block device exposed via SPDK's bdev layer. SPDK's bdev layer is a block device abstraction layer and general-purpose block storage stack akin to what is found in many operating systems. Using the bdev interface completely decouples the storage media from the front-end protocol used to access storage. Users can build their own virtual bdevs that provide complex storage services and integrate them with the SPDK NVMe-oF target with no additional code changes. There can be many subsystems within an NVMe-oF target and each subsystem may hold many namespaces. Subsystems and namespaces can be configured dynamically via a JSON-RPC interface.

Figure 1 shows a high-level schematic of the systems used for testing in the rest of this report. The set up consists of three individual systems (two used as initiators and one used as the target). The NVMe-oF target is connected to both initiator systems point-to-point using QSFP28 cables without any switches. The target system has fourteen Kioxia® KCM61VUL3T20 SSDs which were used as block devices for NVMe-oF subsystems 100GbE Intel® E810-CQDA2 NICs connected to provide up to 200GbE of network bandwidth. Each Initiator system has two Intel® E810-CQDA2 100GbE NICs connected directly to the target without any switch.

One goal of this report was to make clear the advantages and disadvantages inherent to the design of the SPDK NVMe-oF components. These components are written using techniques such as run-to completion, polling, and asynchronous I/O. The report covers four real-world use cases.

For performance benchmarking the fio tool is used with two storage engines:

- 1) Linux Kernel libaio engine
- 2) SPDK bdev engine

Performance numbers reported are aggregate I/O per second, average latency, and CPU utilization as a percentage for various scenarios. Aggregate I/O per second and average latency data is reported from fio and CPU utilization was collected using sar (sysstat).

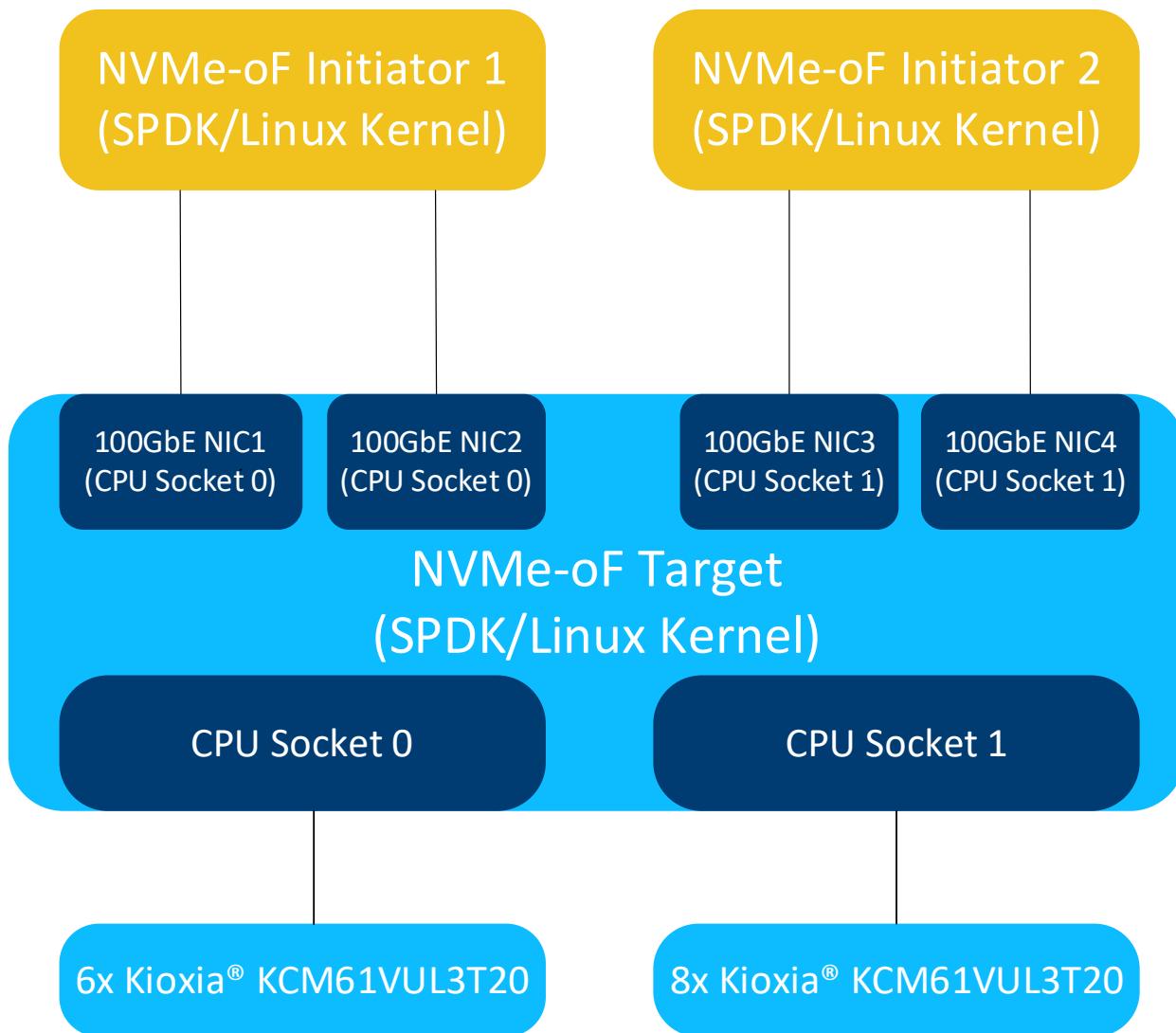


Figure 1: High-Level NVMe-oF RDMA performance testing setup

Test Case 1: SPDK NVMe-oF RDMA Target I/O core scaling

This test case was designed to demonstrate how the SPDK NVMe-oF target throughput in IOPS (I/O per second) scales when additional CPU cores are added to the SPDK NVMe-oF target application.

The SPDK NVMe-oF RDMA target was configured to run with 14 NVMe-oF subsystems. Each NVMe-oF subsystem ran on top of an individual NVMe bdev backed by a single Kioxia KCM61VUL3T20 NVMe drive. Each of the 2 host systems was connected to 7 NVMe-oF subsystems, which were exported by the SPDK NVMe-oF Target over 2 x 100GbE link. The SPDK bdev fio plugin was used to target 7 NVMe-oF bdevs on each of the host. SPDK Target Reactor Mask was configured to use up to 10 CPU cores tests while running following workloads on each initiator:

- 4KiB 100% Random Read
- 4KiB 100% Random Write
- 4KiB Random 70% Read 30% Write

Below table contains information about the test configuration in form of a sequence of commands used by spdk/scripts/rpc.py script to configure the SPDK NVMe-oF Target. The SPDK NVMe-oF Initiator (bdev_fio_plugin) still uses plain configuration files.

Each workload was run three times at each CPU count and the reported results are the average of the 3 runs. We preconditioned the SSDs once before running the 4KiB Rand Read and 4KiB Rand 70/30 Read/Write workloads to ensure that the SSDs reached their steady state where we get repeatable results. However, for the 4KiB Rand Write workload we didn't precondition the NVMe devices to ensure workload saturated the network rather than being limited to the steady state performance of the SSDs which is much lower than the available network bandwidth.

Table 5: SPDK NVMe-oF RDMA Target Core Scaling test configuration

Item	Description
Test Case	SPDK NVMe-oF Target I/O core scaling
SPDK NVMe-oF Target configuration	All the commands below were executed with spdk/scripts/rpc.py script. Construct NVMe bdevs: bdev_nvme_attach_controller -t PCIe -b Nvme0 -a 0000:17:00.0 bdev_nvme_attach_controller -t PCIe -b Nvme1 -a 0000:18:00.0 bdev_nvme_attach_controller -t PCIe -b Nvme2 -a 0000:65:00.0 bdev_nvme_attach_controller -t PCIe -b Nvme3 -a 0000:66:00.0 bdev_nvme_attach_controller -t PCIe -b Nvme4 -a 0000:67:00.0 bdev_nvme_attach_controller -t PCIe -b Nvme5 -a 0000:68:00.0 bdev_nvme_attach_controller -t PCIe -b Nvme6 -a 0000:98:00.0 bdev_nvme_attach_controller -t PCIe -b Nvme7 -a 0000:99:00.0 bdev_nvme_attach_controller -t PCIe -b Nvme8 -a 0000:9a:00.0 bdev_nvme_attach_controller -t PCIe -b Nvme9 -a 0000:9b:00.0 bdev_nvme_attach_controller -t PCIe -b Nvme10 -a 0000:e3:00.0 bdev_nvme_attach_controller -t PCIe -b Nvme11 -a 0000:e4:00.0 bdev_nvme_attach_controller -t PCIe -b Nvme12 -a 0000:e5:00.0 bdev_nvme_attach_controller -t PCIe -b Nvme13 -a 0000:e6:00.0

	<pre> Create RDMA transport layer: nvmf_create_transport -t RDMA -n 8192 { trtype: "RDMA" max_queue_depth: 128 max_qpairs_per_ctrlr: 64 in_capsule_data_size: 4096 max_io_size: 131072 io_unit_size: 8192 max_aq_depth: 128 num_shared_buffers: 8192 buf_cache_size: 32 } Create NVMe-oF subsystems and add NVMe bdevs as namespaces: for i in \$(seq 1 16); do nvmf_subsystem_create nqn.2018-09.io.spdk:cnode\${i} -s SPDK00\${i} -a -m 8 nvmf_subsystem_add_ns nqn.2018-09.io.spdk:cnode\${i} Nvme\$((i-1))n1 done Add listeners to NVMe-oF Subsystems: i=1 ips=(20.0.0.1 20.0.1.1 10.0.0.1 10.0.1.1) for ip in \${ips[@]}; do for j in \$(seq 1 4); do nvmf_subsystem_add_listener nqn.2018-09.io.spdk:cnode\${i} -t rdma \ -f ipv4 -s 4420 -a \${ip} ((i++)) done done </pre>
SPDK NVMe-oF Initiator - fio plugin configuration	<p>BDEV.conf See Appendix A</p> <p>fio.conf</p> <pre> [global] ioengine=/tmp/spdk/examples/bdev/fio_plugin/fio_plugin spdk_json_conf=/tmp/spdk/bdev.conf thread=1 group_reporting=1 direct=1 norandommap=1 rw=randrw rwmixread={100, 70, 0} bs=4k iodepth={1, 8, 16, 32} time_based=1 ramp_time=60 runtime=300 [filename0] filename=NvmeOn1 [filename1] filename=Nvme1n1 [filename2] filename=Nvme2n1 [filename3] filename=Nvme3n1 [filename4] filename=Nvme4n1 [filename5] </pre>

	filename=Nvme5n1 [filename6] filename=Nvme6n1
--	---

4KiB Random Read Results

Table 6: SPDK NVMe-oF RDMA Target Core Scaling results, Random Read IOPS, QD=192

# of Cores	Bandwidth (MiBps)	Throughput (IOPS k)	Avg. Latency (usec)
1 core	6543.53	1675.1	1604.7
2 cores	15052.49	3853.4	697.2
3 cores	24489.33	6269.3	428.4
4 cores	30858.24	7899.7	340.1
5 cores	32999.30	8447.8	317.9
6 cores	32847.03	8408.8	319.3
8 cores	32909.85	8424.9	318.7
10 cores	33056.34	8462.4	317.3

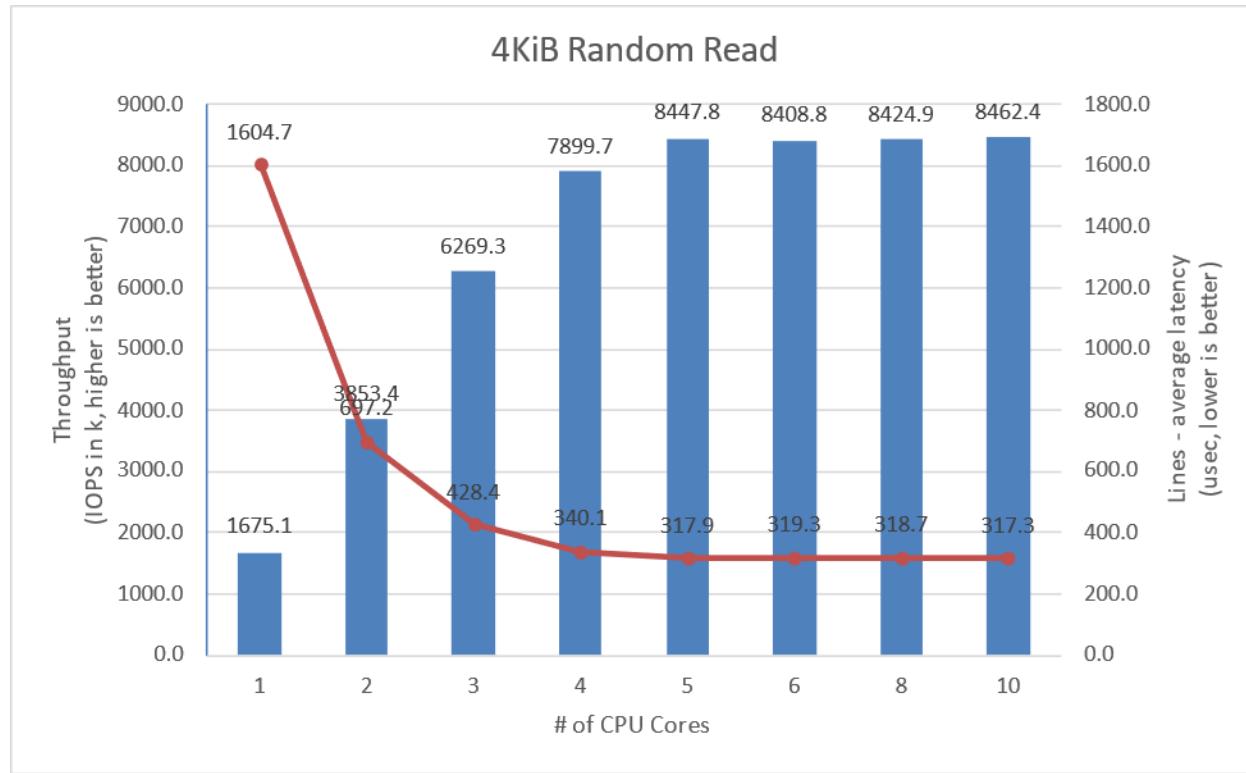


Figure 2: SPDK NVMe-oF RDMA Target I/O core scaling: IOPS vs. Latency while running 4KiB 100% Random Read workload at QD = 192

4KiB Random Write Results

Table 7: SPDK NVMe-oF RDMA Target Core Scaling results, Random Write IOPS, QD=64

# of Cores	Bandwidth (MiBps)	Throughput (IOPS k)	Avg. Latency (usec)
1 core	8879.68	2273.2	392.8
2 cores	20009.44	5122.4	173.1
3 cores	30908.68	7912.6	111.1
4 cores	34634.20	8866.4	99.6
5 cores	34207.84	8757.2	101.5
6 cores	33560.50	8591.5	103.7
8 cores	33150.78	8486.6	105.1
10 cores	33196.89	8498.4	105.0

Note that the SSDs were not preconditioned for the 4K random write workload because that would limit the workload performance to the SSDs steady state performance.

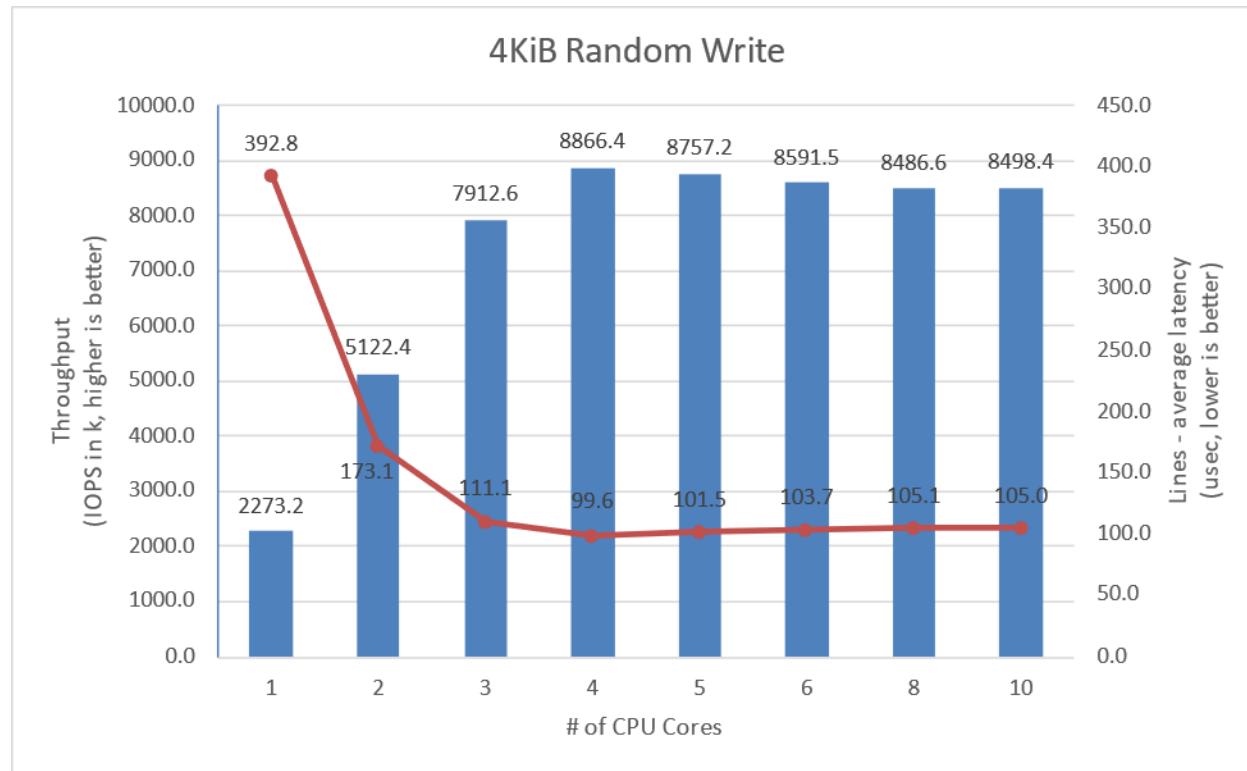


Figure 3: SPDK NVMe-oF RDMA Target I/O core scaling: IOPS vs. Latency while running 4KiB 100% Random Write Workload at QD=64

4KiB Random Read-Write Results

Table 8: SPDK NVMe-oF RDMA Target Core Scaling results, Random Read/Write 70%/30% IOPS, QD=128

# of Cores	Bandwidth (MiBps)	Throughput (IOPS k)	Avg. Latency (usec)
1 core	6664.57	1706.1	1049.9
2 cores	15465.58	3959.2	451.9
3 cores	25013.37	6403.4	279.2
4 cores	32816.74	8401.1	213.3
5 cores	37609.12	9627.9	185.5
6 cores	37917.42	9706.9	184.0
8 cores	38462.25	9846.3	181.4
10 cores	38780.18	9927.7	179.9



Figure 4: SPDK NVMe-oF RDMA Target I/O core scaling: IOPS vs. Latency while running 4KiB Random 70/30 Read/Write workload at QD=128

Large Sequential I/O Performance

128KiB block size I/O tests were performed with sequential I/O workloads at queue depth 8. The rest of the fio configuration is similar to the 4KiB test case in the previous part of this document. We used iodepth=4 and iodepth=8 because higher queue depth resulted in negligible bandwidth gain and a significant increase in the latency.

Table 9: SPDK NVMe-oF RDMA Target Core Scaling results, 128KiB Sequential Read IOPS, QD=4

# of Cores	Bandwidth (MiBps)	Throughput (IOPS k)	Avg. Latency (usec)
1 core	40986.94	327.9	170.6
2 cores	41076.10	328.6	170.2
3 cores	41080.95	328.6	170.2
4 cores	41065.13	328.5	170.3

Table 10: SPDK NVMe-oF RDMA Target Core Scaling results, 128KiB Sequential Write IOPS, QD=4

# of Cores	Bandwidth (MiBps)	Throughput (IOPS k)	Avg. Latency (usec)
1 core	43387.77	347.1	161.0
2 cores	43413.42	347.3	160.9
3 cores	43432.50	347.5	160.9
4 cores	43355.91	346.8	161.1

Table 11: SPDK NVMe-oF RDMA Target Core Scaling results, 128KiB Sequential 70% Read 30% Write IOPS, QD=8

# of Cores	Bandwidth (MiBps)	Throughput (IOPS k)	Avg. Latency (usec)
1 core	40511.93	324.1	345.2
2 cores	40578.42	324.6	344.6
3 cores	40563.80	324.5	344.8
4 cores	40590.03	324.7	344.5

Conclusions

1. 100% 4KiB Random Read workload throughput scales up and latency decreases linearly with the addition of I/O cores up to 4 cores. Further increasing the number to 5 CPU cores results in non-linear improvement. Beyond this point increasing the number of Target IO Cores does not affect performance which peaks at 8.4 million IOPS. Maximum bandwidth measured by fio is about 276 Gbps leaving network link not saturated.
2. 100% 4KiB Write Read workload throughput scales up and latency decreases linearly with the addition of I/O cores up to 3 cores reaching 7.9 million IOPS. Further increasing the number of CPU cores results in non-linear or no performance improvement. Peak performance is reached with 5 CPU cores and 8.7 million IOPS. Network link is not saturated.
3. 70/30% 4KiB Random Read/Write throughput scales up and latency decreases linearly with the addition of I/O cores up to 4 cores reaching 8.4 million IOPS. Further increasing the number of CPU cores results in non-linear or no performance improvement. Peak performance is reached with 10 CPU cores and 9.92 million IOPS. Network link is not saturated.
4. For large sequential I/Os, a single CPU core Target saturated the network bandwidth only for 128k Write workload. Increasing IO CPUs and Queue Depth for other workloads did not affect performance and network was not fully saturated.

Test Case 2: SPDK NVMe-oF RDMA Initiator I/O core scaling

This test case was designed to demonstrate how the SPDK NVMe-oF initiator throughput in IOPS (I/O per second) scales when additional CPU cores are added to the SPDK NVMe-oF initiator.

The test setup for this test case is slightly different than the set up described in [introduction chapter](#), as we used just a single SPDK NVMe-oF RDMA Initiator. The Initiator was connected to Target server with two 100 Gbps network links, using single port from two separate network interface cards.

The SPDK NVMe-oF RDMA Target was configured using 6 cores; all the other configurations are similar to test case 1. The SPDK bdev fio plugin was used to target 14 individual NVMe-oF subsystems exported by the Target. The number of CPU threads used by the fio process was managed by setting the fio job sections and numjobs parameter and ranged from 1 to 8 CPUs. For detailed fio job configuration see table below.

- 4KiB 100% Random Read
- 4KiB 100% Random Write
- 4KiB Random 70% Read 30% Write

It is important to note that fio io depth parameter values presented in the table below are actual queue depths used for each of the connected subsystem. These values were calculated in test based on number of fio job sections, numjobs parameter and the number of “filename” targets grouped in each of the fio job sections.

Table 12: SPDK NVMe-oF RDMA Initiator Core Scaling test configuration

Item	Description
Test Case	SPDK NVMe-oF RDMA Initiator I/O core scaling
SPDK NVMe-oF Target configuration	Same as in Test Case #1, using 6 CPU cores.
SPDK NVMe-oF Initiator 1 - fio plugin configuration	<p>BDEV.conf See appendix B.</p> <p>fio.conf</p> <p>For 1 CPU initiator configuration:</p> <pre>[global] ioengine=/tmp/spdk/examples/bdev/fio_plugin/fio_plugin spdk_conf=/tmp/spdk/bdev.conf thread=1 group_reporting=1 direct=1 norandommap=1 rw=randrw rwmixread={100, 70, 0} bs=4k iodepth={1,32, 64, 128, 192, 256} time_based=1</pre>

	<pre>ramp_time=60 runtime=300 numjobs=1 [filename0] filename=Nvme0n1 filename=Nvme1n1 filename=Nvme2n1 filename=Nvme3n1 filename=Nvme4n1 filename=Nvme5n1 filename=Nvme6n1 filename=Nvme7n1 filename=Nvme8n1 filename=Nvme9n1 filename=Nvme10n1 filename=Nvme11n1 filename=Nvme12n1 filename=Nvme13n1</pre>
	<p>fio.conf</p> <p>For CPU > 1 (up to N=8) initiator configuration “filename=NvmeXn1” are evenly spread across fio job threads:</p> <pre>[global] ioengine=/tmp/spdk/examples/bdev/fio_plugin/fio_plugin spdk_conf=/tmp/spdk/bdev.conf thread=1 group_reporting=1 direct=1 norandommap=1 rw=randrw rwmixread={100, 70, 0} bs=4k iodepth={1,32, 64, 128, 192, 256} time_based=1 ramp_time=60 runtime=300 numjobs=X [filename0] filename=Nvme0n1 filename=Nvme1n1 [filename1] filename=Nvme2n1 filename=Nvme3n1 [...] [filename N-1] filename=Nvme10n1 filename=Nvme11n1 [filename N] filename=Nvme12n1 filename=Nvme13n1</pre>

4KiB Random Read Results

Table 13: SPDK NVMe-oF RDMA Initiator Core Scaling results, 4KiB Random Read IOPS, QD=64, SPDK Target 6 CPU Cores

# of Initiator CPU Cores	Bandwidth (MiBps)	Throughput (IOPS k)	Avg. Latency (usec)
1 core	5467.89	1399.8	526.3
2 cores	15638.31	4003.4	195.8
3 cores	19520.90	4997.3	163.7
4 cores	21318.20	5457.5	160.6
5 cores	21296.50	5451.9	162.6
6 cores	21258.48	5442.2	164.0
7 cores	21260.71	5442.7	164.2
8 cores	21264.12	5443.6	164.1



Figure 5: SPDK NVMe-oF RDMA Initiator I/O core scaling: IOPS vs. Latency while running 4KiB 100% Random Read QD=64 workload

4KiB Random Write Results

Note: The SSDs were not pre-conditioned before running the 100% Random Write test cases. This allowed the throughput to scale to the 2x 100GbE network bandwidth when testing with 3 and 4 CPU cores rather than limiting the workload performance to the storage bottleneck (which is approx. 3.2M IOPS).

Table 14: SPDK NVMe-oF RDMA Initiator Core Scaling results, 4KiB Random Write IOPS, QD=64, SPDK Target 6 CPU Cores

# of Initiator CPU Cores	Bandwidth (MiBps)	Throughput (IOPS k)	Avg. Latency (usec)
1 core	5408.13	1384.5	526.1
2 cores	12610.73	3228.3	231.3
3 cores	16676.03	4269.1	193.3
4 cores	18237.55	4668.8	191.3
5 cores	18801.28	4813.1	187.7
6 cores	20285.57	5193.1	172.0
7 cores	20789.94	5322.2	167.9
8 cores	21120.82	5406.9	165.1

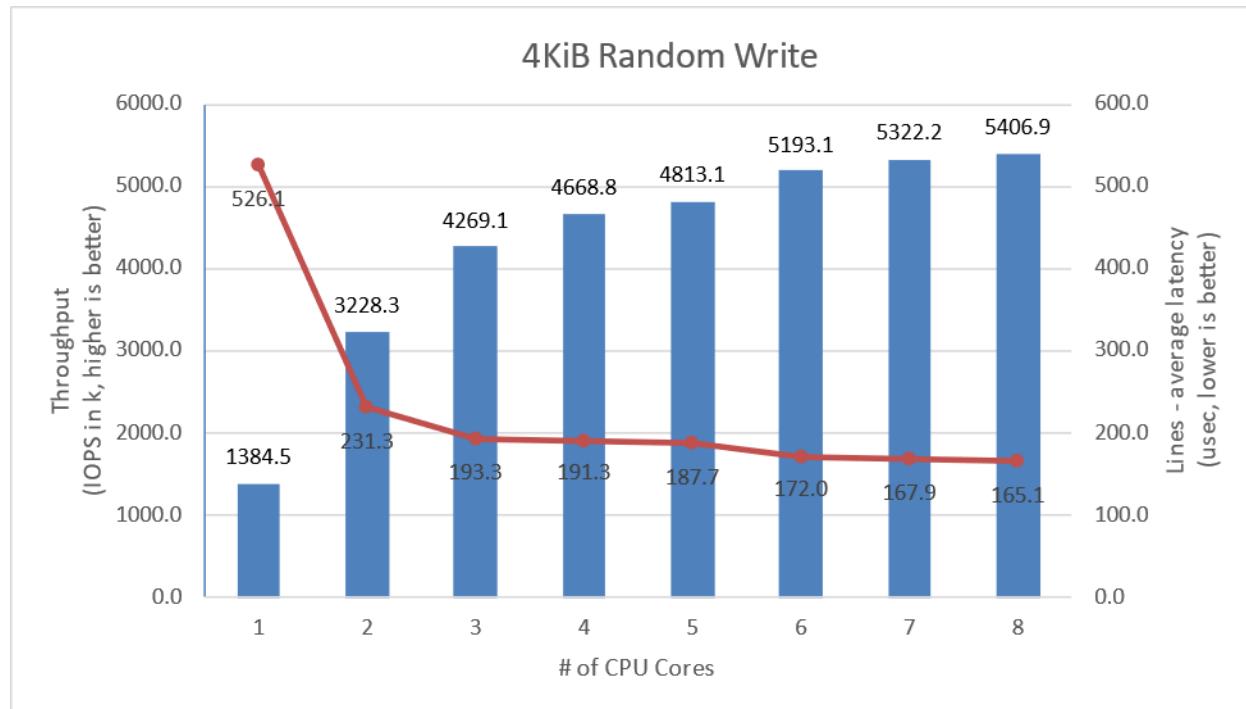


Figure 6: SPDK NVMe-oF RDMA Initiator I/O core scaling: IOPS vs. Latency while running 4KiB 100% Random Write Workload at QD=64

4KiB Random 70/30 Read/Write Results

Table 15: SPDK NVMe-oF RDMA Initiator Core Scaling results, 4KiB Random 70%/30% Read/Write IOPS, QD=64, SPDK Target 6 CPU Cores

# of Initiator CPU Cores	Bandwidth (MiBps)	Throughput (IOPS k)	Avg. Latency (usec)
1 core	5378.20	1376.8	527.4
2 cores	14709.77	3765.7	203.6
3 cores	18426.78	4717.3	169.7
4 cores	23573.89	6034.9	141.7
5 cores	25486.38	6524.5	132.9
6 cores	27214.68	6967.0	126.2
7 cores	27148.87	6950.1	128.0
8 cores	27971.65	7160.7	124.2

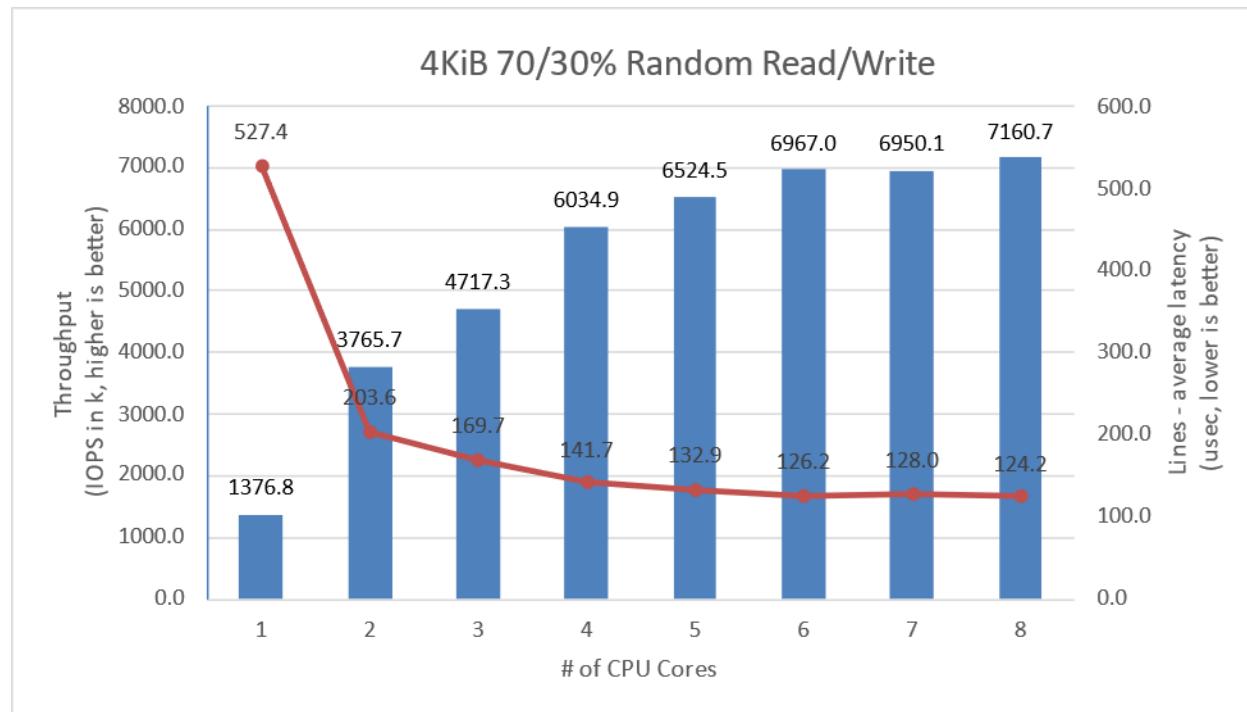


Figure 7: SPDK NVMe-oF RDMA Initiator I/O core scaling: IOPS vs. Latency while running 4KiB Random 70% Read 30% Write Workload at QD=64

Conclusions

1. Random Read workload scaling was not linear when increasing the number of initiator CPU cores. Peak performance of about 5.4 million IOPS was reached with SPDK NVMe-oF RDMA Initiator using 4 CPU cores, saturating 200 GbE link between Target and Initiator.
2. Random Write workload scaling was not linear when increasing the number of initiator CPU cores. Peak performance of 5.4 million IOPS was reached with SPDK NVMe-oF RDMA Initiator using 8 CPU cores, saturating Target NVMe drives 4k Random Write processing capability.
3. Random Read-Write scaling was not linear when increasing the number of initiator CPU cores. Performance peaked at 7.16 million IOPS with the SPDK NVMe-oF RDMA initiator using 8 CPU cores. Throughput exceeding 200 GbE link between Target and Initiator can be explained by the fact that mixed read-write workload can be considered a bi-directional network traffic, not being strictly affected by 200 GbE “one-way” network limit.

Test Case 3: Linux Kernel vs. SPDK NVMe-oF RDMA Latency

This test case was designed to understand latency characteristics of SPDK NVMe-oF RDMA Target and Initiator vs. the Linux Kernel NVMe-oF RDMA Target and Initiator implementations on a single NVMe-oF subsystem. The average I/O latency and p99 latency was compared between SPDK NVMe-oF (Target/Initiator) vs. Linux Kernel (Target/Initiator). Both SPDK and Kernel NVMe-oF Targets were configured to run on a single core, with a single NVMe-oF subsystem backed by a Null Block Device. The null block device (bdev) was chosen as the backend block device to eliminate the media latency during these tests.

Kernel NVMe-oF Initiator disclaimer:

For establishing Kernel NVMe-oF RDMA Initiator connections “nvme-cli” tool was used. While performing benchmark tests two issues were encountered:

- It was not possible to establish connection and create a NVMe block device on Initiator side with poll queues enabled ([link](#)). Using “nvme-cli” with “—nr-poll-queues” parameter present resulted in “Kernel Oops” to be generated. Because of this issue the fio workload for Kernel Initiator connection was configured to use “libaio” engine.
- Attempts to establish connection with default number of IO queues (which is equal to number of CPU cores on Initiator system) resulted in connection timeouts. To work around this issue “—nr-io-queues=32” was added to nvme-cli command. This does not affect the results in this test case as only a single connection with very small queue depth is tested.

Table 16: Linux Kernel vs. SPDK NVMe-oF RDMA Latency test configuration

Item	Description
Test Case	Linux Kernel vs. SPDK NVMe-oF RDMA Latency
SPDK NVMe-oF Target configuration	<p>Test configuration</p> <p>The following commands are executed with spdk/scripts/rpc.py script to configure the SPDK NVMe-oF target.</p> <pre>nvmf_create_transport -t RDMA (create RDMA transport layer with default values: trtype: "RDMA" max_queue_depth: 128 max_qpairs_per_ctrlr: 64 in_capsule_data_size: 4096 max_io_size: 131072 io_unit_size: 8192 max_aq_depth: 128 num_shared_buffers: 8192 buf_cache_size: 32) bdev_null_create NvmeOn1 10240 4096 nvmf_subsystem_create nqn.2018-09.io.spdk:cnode1 -s SPDK001 -a -m 8 nvmf_subsystem_add_ns nqn.2018-09.io.spdk:cnode1 NvmeOn1 nvmf_subsystem_add_listener nqn.2018-09.io.spdk:cnode1 -t rdma -f ipv4 -s 4420 -a 20.0.0.1</pre>
Kernel NVMe-oF Target configuration	The following target configuration file loaded using nvmet-cl tool. {

	<pre> "ports": [{ "addr": { "adrfam": "ipv4", "traddr": "20.0.0.1", "trsvcid": "4420", "trtype": "rdma" }, "portid": 1, "referrals": [], "subsystems": ["nqn.2018-09.io.spdk:cnode1"] }], "hosts": [], "subsystems": [{ "allowed_hosts": [], "attr": { "allow_any_host": "1", "version": "1.3" }, "namespaces": [{ "device": { "path": "/dev/nullb0", "uuid": "621e25d2-8334-4c1a-8532-b6454390b8f9" }, "enable": 1, "nsid": 1 }], "nqn": "nqn.2018-09.io.spdk:cnode1" }] } </pre>
fio configuration	
SPDK NVMe-oF Initiator fio plugin configuration	<p>BDEV.conf See Appendix B.</p> <p>fio.conf</p> <pre> [global] ioengine=/tmp/spdk/examples/bdev/fio_plugin/fio_plugin spdk_json_conf=/tmp/spdk/bdev.conf thread=1 group_reporting=1 direct=1 norandommap=1 rw=randrw rwmixread={100, 70, 0} bs=4k iodepth=1 time_based=1 ramp_time=60 runtime=300 [filename0] filename=Nvme0n1 </pre>
Kernel initiator	Device config

configuration	<p>The following configuration was performed using nvme-cli tool.</p> <pre>modprobe nvme-fabrics nvme connect -n nqn.2018-09.io.spdk:cnode1 -t rdma -a 20.0.0.1 -s 4420</pre> <p>fio.conf</p> <pre>[global] ioengine=libaio thread=1 group_reporting=1 direct=1 norandommap=1 rw=randrw rwmixread={100, 70, 0} bs=4k iodepth=1 time_based=1 ramp_time=60 runtime=300 [filename0] filename=/dev/nvme0n1</pre>
----------------------	---

SPDK vs Kernel NVMe-oF RDMA Target Results

The following data was collected using the Linux Kernel initiator against both SPDK and Linux Kernel NVMe-oF RDMA target.

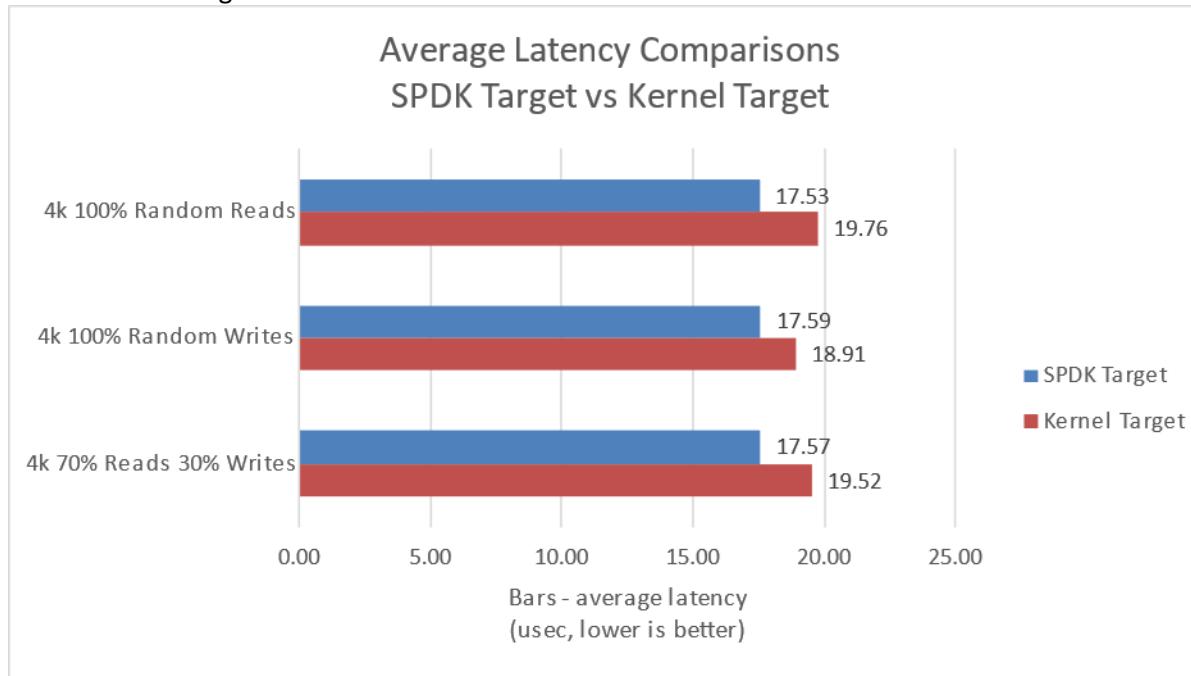


Figure 8: SPDK vs. Kernel NVMe-oF RDMA Target average I/O latency for various workloads run using the Kernel Initiator

Table 17: SPDK NVMe-oF RDMA Target Latency and IOPS at QD=1, Null Block Device

Access Pattern	Average Latency (usec)	IOPS	p99 (usec)	p99.9 (usec)	p99.99 (usec)	p99.999 (usec)
4KiB 100% Random Read	17.53	56173	20.6	24.4	99.8	389.8
4KiB 100% Random Write	17.59	56020	19.7	23.9	99.8	392.5
4KiB 70/30% Random Read/Write	17.57	55939	20.6	24.5	100.6	391.6

Table 18: Linux Kernel NVMe-oF RDMA Target Latency and IOPS at QD=1, Null Block Device

Access Pattern	Average Latency (usec)	IOPS	p99 (usec)	p99.9 (usec)	p99.99 (usec)	p99.999 (usec)
4KiB 100% Random Read	19.76	49932	21.6	25.8	34.4	207.2
4KiB 100% Random Write	18.91	52140	25.6	27.8	36.6	207.9
4KiB 70/30% Random Read/Write	19.52	50432	21.8	26.0	31.8	207.3

SPDK vs Kernel NVMe-oF RDMA Initiator Results

This following data was collected using the Linux Kernel and SPDK NVMe-oF RDMA initiator against an SPDK NVMe-oF RDMA target.

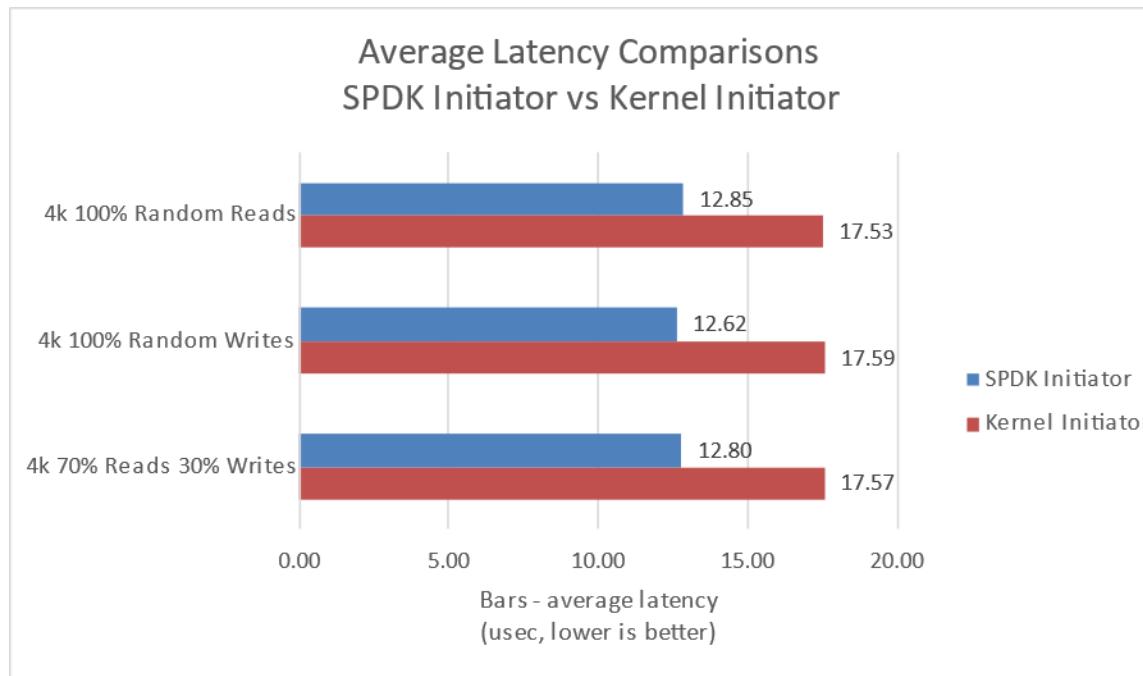


Figure 9: SPDK vs. Kernel NVMe-oF RDMA Initiator average I/O latency for various workloads against SPDK Target

Table 19: SPDK NVMe-oF RDMA Initiator Latency and IOPS at QD=1, Null Block Device

Access Pattern	Average Latency (usec)	IOPS	p99 (usec)	p99.9 (usec)	p99.99 (usec)	p99.999 (usec)
4KiB 100% Random Read	12.85	76838	16.4	21.1	43.1	374.8
4KiB 100% Random Write	12.62	78136	16.1	21.6	41.6	366.6
4KiB 70/30% Random Read/Write	12.80	77053	15.5	18.6	40.1	315.5

Table 20: Linux Kernel NVMe-oF RDMA Initiator Latency and IOPS at QD=1, Null Block Device

Access Pattern	Average Latency (usec)	IOPS	p99 (usec)	p99.9 (usec)	p99.99 (usec)	p99.999 (usec)
4KiB 100% Random Read	17.53	56173	20.6	24.4	99.8	389.8
4KiB 100% Random Write	17.59	56020	19.7	23.9	99.8	392.5
4KiB 70/30% Random Read/Write	17.57	55939	20.6	24.5	100.6	391.6

SPDK vs Kernel NVMe-oF RDMA Kernel + Initiator Results

Following data was collected using SPDK Target with SPDK Initiator and Linux Target with Linux Initiator.

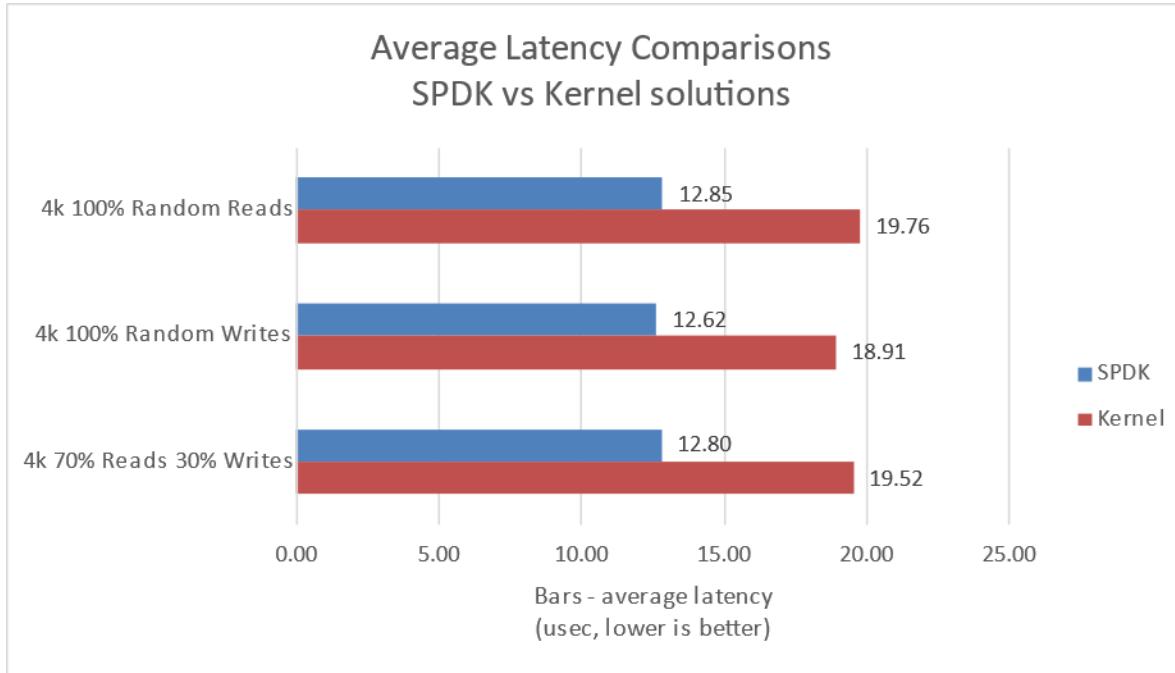


Figure 10: SPDK vs. Kernel NVMe-oF RDMA solutions average I/O Latency for various workloads against SPDK Target

Table 21: SPDK NVMe-oF RDMA Latency and IOPS at QD=1, Null Block Device

Access Pattern	Average Latency (usec)	IOPS	p99 (usec)	p99.9 (usec)	p99.99 (usec)	p99.999 (usec)
4KiB 100% Random Read	12.85	76838	16.4	21.1	43.1	374.8
4KiB 100% Random Write	12.62	78136	16.1	21.6	41.6	366.6
4KiB 70/30% Random Read/Write	12.80	77053	15.5	18.6	40.1	315.5

Table 22: Linux Kernel NVMe-oF RDMA Latency and IOPS at QD=1, Null Block Device

Access Pattern	Average Latency (usec)	IOPS	p99 (usec)	p99.9 (usec)	p99.99 (usec)	p99.999 (usec)
4KiB 100% Random Read	19.76	49932	21.6	25.8	34.4	207.2
4KiB 100% Random Write	18.91	52140	25.6	27.8	36.6	207.9
4KiB 70/30% Random Read/Write	19.52	50432	21.8	26.0	31.8	207.3

Conclusions

1. For the RDMA transport, the SPDK NVMe-oF Target reduces the NVMe-oF average round trip I/O latency (reads/writes) by up to 2.23 usec vs. the Linux Kernel NVMe-oF target used in Fedora 37 6.0.18 setup. This is entirely software overhead, therefore, using the SPDK NVMe-oF target reduces the NVMe-oF software overhead by approximately 11.28% vs. the Linux Kernel NVMe-oF target.
2. The SPDK NVMe-oF Initiator reduces the NVMe-oF software overhead by up to 4.96 usec vs. the Linux Kernel NVMe-oF Initiator for the RDMA transport, which is approximately 28.22% of Linux Kernel NVMe-oF Initiator overhead.

Test Case 4: NVMe-oF RDMA Performance with increasing # of connections

This test case was designed to demonstrate the throughput and latency of the SPDK NVMe-oF RDMA Target vs. Linux Kernel NVMe-oF RDMA Target under increasing number of connections per subsystem. The number of active connections (or I/O queue pairs) per NVMe-oF subsystem was varied, we measured the aggregated IOPS and number of CPU cores used by each target. The number of CPU cores metric was calculated from %CPU utilization measured using sar (sysstat package in Linux). The SPDK NVMe-oF RDMA Target was configured to run on 8 CPU cores, export 14 NVMe-oF subsystems (1 per Kioxia NVMe SSD) and 2 initiators were used both running the I/O workloads below to 7 separate subsystems using Kernel NVMe-oF RDMA initiator.

- 4KiB 100% Random Read
- 4KiB 100% Random Write
- 4KiB Random 70% Read 30% Write

Kernel NVMe-oF Initiator disclaimer:

For establishing Kernel NVMe-oF RDMA Initiator connections “nvme-cli” tool was used. While performing benchmark tests two issues were encountered:

- It was not possible to establish connection and create a NVMe block device on Initiator side with poll queues enabled ([link](#)). Using “nvme-cli” with “—nr-poll-queues” parameter present resulted in “Kernel Oops” to be generated. Because of this issue the fio workload for Kernel Initiator connection was configured to use “libaio” engine.
- Attempts to establish connection with default number of IO queues (which is equal to number of CPU cores on Initiator system) resulted in connection timeouts. To work around this issue “—nr-io-queues=32” was added to nvme-cli command. This does not affect the results in this test case as connections are not scaled beyond 16 per NVMe-oF subsystem, thus maximum number of used IO queues per NVMe-oF subsystem is 16.

Table 23: NVMe-oF RDMA Performance with increasing number of connections test configuration

Item	Description
Test Case	NVMe-oF RDMA Target performance with increasing # of connections
SPDK NVMe-oF Target configuration	Same as in Test Case #1, using 8 CPU cores.
Kernel NVMe-oF Target configuration	Target configuration file loaded using nvmet-clt tool. For detailed configuration file contents please see Appendix C .
Kernel NVMe-oF Initiator #1	Device config Performed using nvme-cli tool. modprobe nvme-fabrics

	<pre>nvme connect -n nqn.2018-09.io.spdk:cnode1 -t rdma -a 20.0.0.1 -s 4420 nvme connect -n nqn.2018-09.io.spdk:cnode2 -t rdma -a 20.0.0.1 -s 4420 nvme connect -n nqn.2018-09.io.spdk:cnode3 -t rdma -a 20.0.0.1 -s 4420 nvme connect -n nqn.2018-09.io.spdk:cnode4 -t rdma -a 20.0.0.1 -s 4420 nvme connect -n nqn.2018-09.io.spdk:cnode5 -t rdma -a 20.0.1.1 -s 4420 nvme connect -n nqn.2018-09.io.spdk:cnode6 -t rdma -a 20.0.1.1 -s 4420 nvme connect -n nqn.2018-09.io.spdk:cnode7 -t rdma -a 20.0.1.1 -s 4420</pre>
Kernel NVMe-oF Initiator #2	<p>Device config Performed using nvme-cli tool.</p> <pre>modprobe nvme-fabrics nvme connect -n nqn.2018-09.io.spdk:cnode8 -t rdma -a 10.0.0.1 -s 4420 nvme connect -n nqn.2018-09.io.spdk:cnode9 -t rdma -a 10.0.0.1 -s 4420 nvme connect -n nqn.2018-09.io.spdk:cnode10 -t rdma -a 10.0.0.1 -s 4420 nvme connect -n nqn.2018-09.io.spdk:cnode11 -t rdma -a 10.0.0.1 -s 4420 nvme connect -n nqn.2018-09.io.spdk:cnode12 -t rdma -a 10.0.1.1 -s 4420 nvme connect -n nqn.2018-09.io.spdk:cnode13 -t rdma -a 10.0.1.1 -s 4420</pre>
fio configuration (used on both initiators)	<pre>fio.conf [global] ioengine=libaio thread=1 group_reporting=1 direct=1 norandommap=1 rw=randrw rwmixread={100, 70, 0} bs=4k iodepth={32, 64, 128, 192} time_based=1 ramp_time=60 runtime=300 numjobs={1, 4, 8, 12, 16} [filename1] filename=/dev/nvme0n1 [filename2] filename=/dev/nvme1n1 [filename3] filename=/dev/nvme2n1 [filename4] filename=/dev/nvme3n1 [filename5] filename=/dev/nvme4n1 [filename6] filename=/dev/nvme5n1 [filename7] filename=/dev/nvme6n1</pre>

The SPDK NVMe-oF Target was configured to use 8 CPU cores for Random Read and Random Read/Write workloads and 4 CPU cores for Random Write workloads. We did not limit the number of CPU cores for the Linux Kernel NVMe-oF target. The graph below shows the relative efficiency in terms of IOPS/core which was calculated by dividing the total aggregate IOPS by the total CPU cores used while running that specific workload.

4KiB Random Read Results

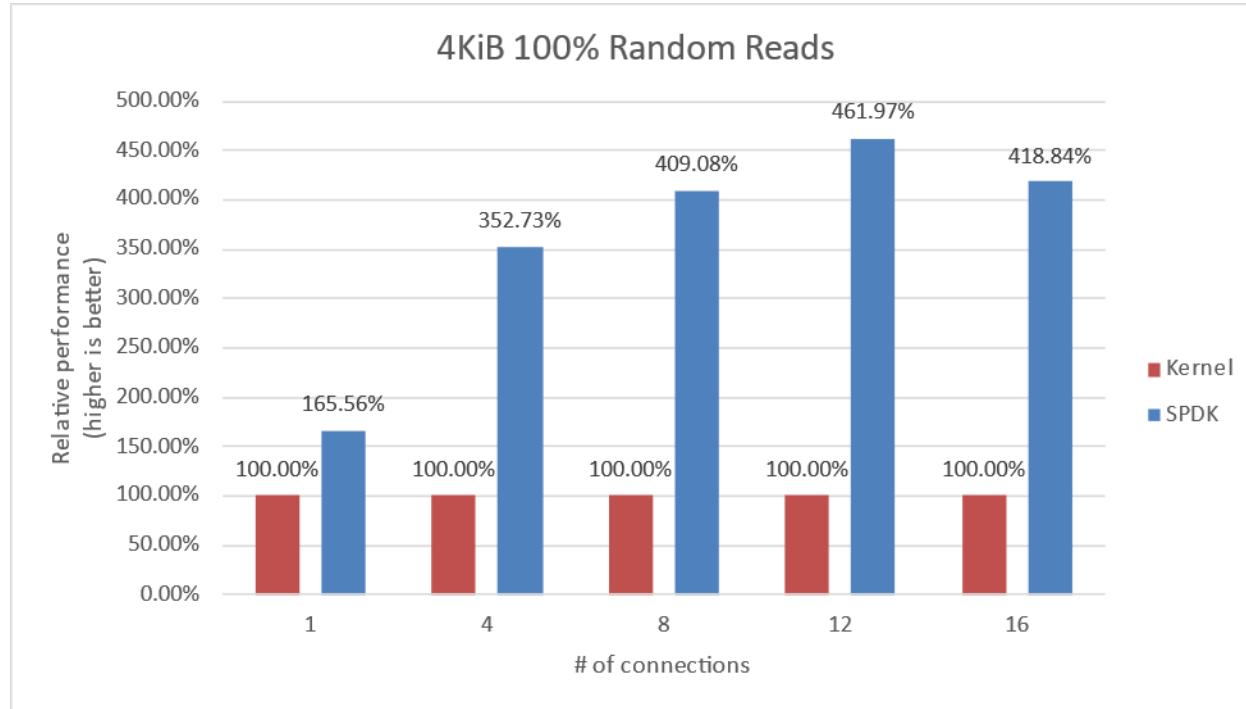


Figure 11: Relative Efficiency Comparison of Linux Kernel vs. SPDK NVMe-oF RDMA Target IOPS/Core for 4KiB 100% Random Reads QD=192, using the Kernel Initiator

Table 24: Linux Kernel NVMe-oF RDMA Target: 4KiB 100% Random Reads, QD=192

Connections per subsystem	Bandwidth (MiBps)	IOPS (k)	Avg. Latency (usec)	# CPU Cores
1	20583.53	5269.4	510.8	12.9
4	40218.82	10296.0	260.9	30.3
8	38902.64	9959.1	269.8	34.0
12	38214.51	9782.9	274.4	33.8
16	38894.27	9956.9	269.7	35.0

Table 25: SPDK NVMe-oF RDMA Target: 4KiB 100% Random Reads, QD=192

Connections per subsystem	Bandwidth (MiBps)	IOPS (k)	Avg. Latency (usec)	# CPU Cores
1	21352.79	5466.3	492.6	8.1
4	42485.74	10876.3	246.9	9.1
8	42439.25	10864.4	247.2	9.1
12	42204.92	10804.4	248.5	8.1
16	42051.68	10765.2	249.3	9.0

4KiB Random Write results

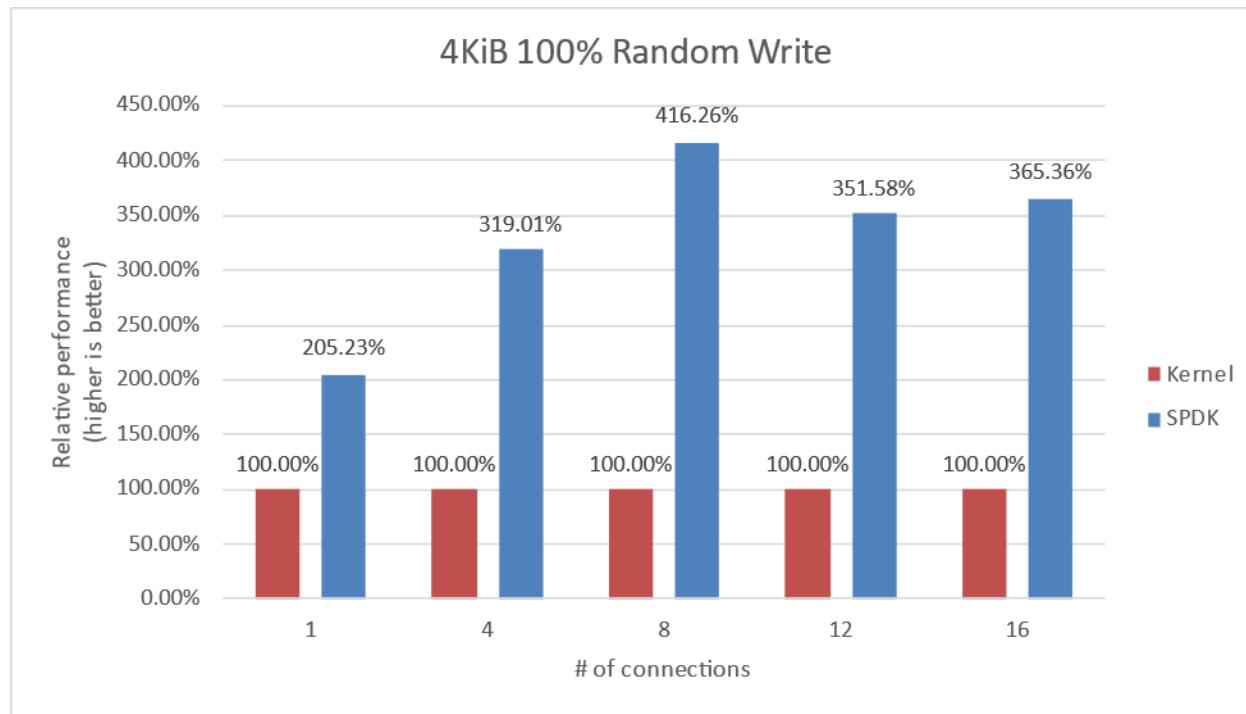


Figure 12: Relative Efficiency Comparison of Linux Kernel vs. SPDK NVMe-oF RDMA Target IOPS/Core for 4KiB 100% Random Writes QD=128 workload, using Kernel Initiators

Note: The SSDs were not pre-conditioned before running 100% Random Write I/O test.

Table 26: Linux Kernel NVMe-oF RDMA Target: 4KiB 100% Random Writes, QD=128

Connections per subsystem	Bandwidth (MiBps)	IOPS (k)	Avg. Latency (usec)	# CPU Cores
1	21479.85	5498.8	330.0	10.2
4	26529.28	6791.5	271.5	15.7
8	28236.90	7228.6	258.0	17.7
12	29505.17	7553.3	255.0	19.2
16	29698.12	7602.7	245.1	20.8

Table 27: SPDK NVMe-oF RDMA Target: 4KiB 100% Random Writes, QD=128

Connections per subsystem	Bandwidth (MiBps)	IOPS (k)	Avg. Latency (usec)	# CPU Cores
1	21717.27	5559.6	325.9	5.0
4	27297.14	6988.1	261.6	5.1
8	27207.55	6965.1	262.5	4.1
12	27305.00	6990.1	268.3	5.1
16	26446.66	6770.3	266.8	5.1

4KiB Random Read-Write Results



Figure 13: Relative Efficiency Comparison of Linux Kernel vs. SPDK NVMe-oF RDMA Target IOPS/Core for 4KiB Random 70% Reads 30% Writes QD=192 Workload, using Kernel Initiators

Table 28: Linux Kernel NVMe-oF RDMA Target: 4KiB 70% Random Read 30% Random Write, QD=192

Connections per subsystem	Bandwidth (MiBps)	IOPS (k)	Avg. Latency (usec)	# CPU Cores
1	20666.14	5290.5	510.1	11.2
4	41186.29	10543.7	254.7	29.5
8	40172.07	10284.0	261.1	33.1
12	39139.23	10019.6	268.0	32.5
16	39781.78	10184.1	263.5	34.0

Table 29: SPDK NVMe-oF RDMA Target: 4KiB 70% Random Read 30% Random Write, QD=192

Connections per subsystem	Bandwidth (MiBps)	IOPS (k)	Avg. Latency (usec)	# CPU Cores
1	21500.67	5504.2	489.2	9.0
4	41633.58	10658.2	252.1	9.0
8	40070.74	10258.1	261.7	9.0
12	40273.00	10309.9	260.3	9.0
16	40153.35	10279.2	261.1	9.1

Conclusions

1. SPDK NVMe-oF Target performance peaked at 4 connections for all workloads. Increasing the number of connections per subsystem beyond these values did not result in significant changes to the IOPS or latency.
2. The performance for the Linux Kernel NVMe-oF Target peaked at 4 connections per subsystem for all Random Read and 70/30 Random Read/Write workloads. Increasing the number of connections for these workloads increases the latency and CPU utilization and lowers IOPS. For Random Write workload performance peaked at 16 connections per subsystem.
3. The SPDK NVMe-oF target shows up to 4.62x more IOPS/Core relative to the Linux Kernel NVMe-oF target as the number of connections per subsystem increased.

Summary

This report showcased performance results with SPDK NVMe-oF RDMA Target and Initiator under various test cases, including:

- I/O core scaling
- Average I/O latency
- Performance with increasing number of connections per subsystems

It compared performance results while running the Linux Kernel NVMe-oF RDMA (Target/Initiator) against the accelerated polled-mode driven SPDK NVMe-oF RDMA (Target/Initiator) implementation. Like in the previous reports, throughput scales up and latency decreases almost linearly with the scaling of SPDK NVMe-oF Target cores.

It was also observed that the SPDK NVMe-oF Target average latency is up to 2.23 usec lower than Kernel when testing against null bdev based backend. The advantage of SPDK is even greater when comparing NVMe-oF Initiators: the SPDK NVMe-oF RDMA average latency is lower by up to 28.22% lower than Kernel initiator.

The SPDK NVMe-oF Target performed up to 4.55 times better w.r.t IOPS/core than Linux Kernel NVMe-oF target while running 4KiB 100% Random workloads with increasing number of active connections per NVMe-oF subsystem.

This report provides information regarding methodologies and practices while benchmarking NVMe-oF using SPDK, as well as the Linux Kernel. It should be noted that the performance data showcased in this report is based on specific hardware and software configurations and that performance results may vary depending on different hardware and software configurations.

List of tables

Table 1: Hardware setup configuration – Target system.....	4
Table 2: Hardware setup configuration – Initiator system 1	5
Table 3: Hardware setup configuration – Initiator system 2	5
Table 4: Test systems BIOS settings	6
Table 5: SPDK NVMe-oF RDMA Target Core Scaling test configuration	9
Table 6: SPDK NVMe-oF RDMA Target Core Scaling results, Random Read IOPS, QD=192	12
Table 7: SPDK NVMe-oF RDMA Target Core Scaling results, Random Write IOPS, QD=64	13
Table 8: SPDK NVMe-oF RDMA Target Core Scaling results, Random Read/Write 70%/30% IOPS, QD=128	14
Table 9: SPDK NVMe-oF RDMA Target Core Scaling results, 128KiB Sequential Read IOPS, QD=4	15
Table 10: SPDK NVMe-oF RDMA Target Core Scaling results, 128KiB Sequential Write IOPS, QD=4	15
Table 11: SPDK NVMe-oF RDMA Target Core Scaling results, 128KiB Sequential 70% Read 30% Write IOPS, QD=8	15
Table 12: SPDK NVMe-oF RDMA Initiator Core Scaling test configuration	17
Table 13: SPDK NVMe-oF RDMA Initiator Core Scaling results, 4KiB Random Read IOPS, QD=64, SPDK Target 6 CPU Cores	19
Table 14: SPDK NVMe-oF RDMA Initiator Core Scaling results, 4KiB Random Write IOPS, QD=64, SPDK Target 6 CPU Cores	20
Table 15: SPDK NVMe-oF RDMA Initiator Core Scaling results, 4KiB Random 70%/30% Read/Write IOPS, QD=64, SPDK Target 6 CPU Cores	21
Table 16: Linux Kernel vs. SPDK NVMe-oF RDMA Latency test configuration	23
Table 17: SPDK NVMe-oF RDMA Target Latency and IOPS at QD=1, Null Block Device.....	26
Table 18: Linux Kernel NVMe-oF RDMA Target Latency and IOPS at QD=1, Null Block Device	26
Table 19: SPDK NVMe-oF RDMA Initiator Latency and IOPS at QD=1, Null Block Device.....	27
Table 20: Linux Kernel NVMe-oF RDMA Initiator Latency and IOPS at QD=1, Null Block Device	27
Table 21: SPDK NVMe-oF RDMA Latency and IOPS at QD=1, Null Block Device.....	28
Table 22: Linux Kernel NVMe-oF RDMA Latency and IOPS at QD=1, Null Block Device	28
Table 23: NVMe-oF RDMA Performance with increasing number of connections test configuration	30
Table 24: Linux Kernel NVMe-oF RDMA Target: 4KiB 100% Random Reads, QD=192	32
Table 25: SPDK NVMe-oF RDMA Target: 4KiB 100% Random Reads, QD=192	32
Table 26: Linux Kernel NVMe-oF RDMA Target: 4KiB 100% Random Writes, QD=128.....	33
Table 27: SPDK NVMe-oF RDMA Target: 4KiB 100% Random Writes, QD=128.....	33

Table 28: Linux Kernel NVMe-oF RDMA Target: 4KiB 70% Random Read 30% Random Write, QD=19234

Table 29: SPDK NVMe-oF RDMA Target: 4KiB 70% Random Read 30% Random Write, QD=19234

List of figures

Figure 1: High-Level NVMe-oF RDMA performance testing setup	8
Figure 2: SPDK NVMe-oF RDMA Target I/O core scaling: IOPS vs. Latency while running 4KiB 100% Random Read workload at QD = 192.....	12
Figure 3: SPDK NVMe-oF RDMA Target I/O core scaling: IOPS vs. Latency while running 4KiB 100% Random Write Workload at QD=64	13
Figure 4: SPDK NVMe-oF RDMA Target I/O core scaling: IOPS vs. Latency while running 4KiB Random 70/30 Read/Write workload at QD=128	14
Figure 5: SPDK NVMe-oF RDMA Initiator I/O core scaling: IOPS vs. Latency while running 4KiB 100% Random Read QD=64 workload	19
Figure 6: SPDK NVMe-oF RDMA Initiator I/O core scaling: IOPS vs. Latency while running 4KiB 100% Random Write Workload at QD=64	20
Figure 7: SPDK NVMe-oF RDMA Initiator I/O core scaling: IOPS vs. Latency while running 4KiB Random 70% Read 30% Write Workload at QD=64.....	21
Figure 8: SPDK vs. Kernel NVMe-oF RDMA Target average I/O latency for various workloads run using the Kernel Initiator.....	26
Figure 9: SPDK vs. Kernel NVMe-oF RDMA Initiator average I/O latency for various workloads against SPDK Target.....	27
Figure 10: SPDK vs. Kernel NVMe-oF RDMA solutions average I/O Latency for various workloads against SPDK Target.....	28
Figure 11: Relative Efficiency Comparison of Linux Kernel vs. SPDK NVMe-oF RDMA Target IOPS/Core for 4KiB 100% Random Reads QD=192, using the Kernel Initiator	32
Figure 12: Relative Efficiency Comparison of Linux Kernel vs. SPDK NVMe-oF RDMA Target IOPS/Core for 4KiB 100% Random Writes QD=128 workload, using Kernel Initiators	33
Figure 13: Relative Efficiency Comparison of Linux Kernel vs. SPDK NVMe-oF RDMA Target IOPS/Core for 4KiB Random 70% Reads 30% Writes QD=192 Workload, using Kernel Initiators	34

Appendix A – Test Case 1 & 2 SPDK NVMe-oF Initiator bdev configuration

Initiator system 1

```
{  
    "subsystems": [  
        {  
            "subsystem": "bdev",  
            "config": [  
                {  
                    "method": "bdev_nvme_attach_controller",  
                    "params": {  
                        "name": "Nvme0",  
                        "trtype": "rdma",  
                        "traddr": "20.0.0.1",  
                        "trsvcid": "4420",  
                        "subnqn": "nqn.2018-09.io.spdk:cnode0",  
                        "adrifam": "IPv4"  
                    }  
                },  
                {  
                    "method": "bdev_nvme_attach_controller",  
                    "params": {  
                        "name": "Nvme1",  
                        "trtype": "rdma",  
                        "traddr": "20.0.0.1",  
                        "trsvcid": "4420",  
                        "subnqn": "nqn.2018-09.io.spdk:cnode1",  
                        "adrifam": "IPv4"  
                    }  
                },  
                {  
                    "method": "bdev_nvme_attach_controller",  
                    "params": {  
                        "name": "Nvme2",  
                        "trtype": "rdma",  
                        "traddr": "20.0.0.1",  
                        "trsvcid": "4420",  
                        "subnqn": "nqn.2018-09.io.spdk:cnode2",  
                        "adrifam": "IPv4"  
                    }  
                },  
                {  
                    "method": "bdev_nvme_attach_controller",  
                    "params": {  
                        "name": "Nvme3",  
                        "trtype": "rdma",  
                        "traddr": "20.0.0.1",  
                        "trsvcid": "4420",  
                        "subnqn": "nqn.2018-09.io.spdk:cnode3",  
                        "adrifam": "IPv4"  
                    }  
                }  
            ]  
        }  
    ]  
}
```

```
        },
        {
            "method": "bdev_nvme_attach_controller",
            "params": {
                "name": "Nvme4",
                "trtype": "rdma",
                "traddr": "20.0.1.1",
                "trsvcid": "4420",
                "subnqn": "nqn.2018-09.io.spdk:cnode4",
                "adrfam": "IPv4"
            }
        },
        {
            "method": "bdev_nvme_attach_controller",
            "params": {
                "name": "Nvme5",
                "trtype": "rdma",
                "traddr": "20.0.1.1",
                "trsvcid": "4420",
                "subnqn": "nqn.2018-09.io.spdk:cnode5",
                "adrfam": "IPv4"
            }
        },
        {
            "method": "bdev_nvme_attach_controller",
            "params": {
                "name": "Nvme6",
                "trtype": "rdma",
                "traddr": "20.0.1.1",
                "trsvcid": "4420",
                "subnqn": "nqn.2018-09.io.spdk:cnode6",
                "adrfam": "IPv4"
            }
        },
        {
            "method": "bdev_nvme_attach_controller",
            "params": {
                "name": "Nvme7",
                "trtype": "rdma",
                "traddr": "20.0.1.1",
                "trsvcid": "4420",
                "subnqn": "nqn.2018-09.io.spdk:cnode7",
                "adrfam": "IPv4"
            }
        }
    ]
}
```

Initiator system 2

```
{
    "subsystems": [
        {
            "subsystem": "bdev",
            "config": [
                {

```

```

        "method": "bdev_nvme_attach_controller",
        "params": {
            "name": "Nvme0",
            "trtype": "rdma",
            "traddr": "10.0.0.1",
            "trsvcid": "4420",
            "subnqn": "nqn.2018-09.io.spdk:cnode0",
            "adrifam": "IPv4"
        }
    },
    {
        "method": "bdev_nvme_attach_controller",
        "params": {
            "name": "Nvme1",
            "trtype": "rdma",
            "traddr": "10.0.0.1",
            "trsvcid": "4420",
            "subnqn": "nqn.2018-09.io.spdk:cnode1",
            "adrifam": "IPv4"
        }
    },
    {
        "method": "bdev_nvme_attach_controller",
        "params": {
            "name": "Nvme2",
            "trtype": "rdma",
            "traddr": "10.0.0.1",
            "trsvcid": "4420",
            "subnqn": "nqn.2018-09.io.spdk:cnode2",
            "adrifam": "IPv4"
        }
    },
    {
        "method": "bdev_nvme_attach_controller",
        "params": {
            "name": "Nvme3",
            "trtype": "rdma",
            "traddr": "10.0.0.1",
            "trsvcid": "4420",
            "subnqn": "nqn.2018-09.io.spdk:cnode3",
            "adrifam": "IPv4"
        }
    },
    {
        "method": "bdev_nvme_attach_controller",
        "params": {
            "name": "Nvme4",
            "trtype": "rdma",
            "traddr": "10.0.1.1",
            "trsvcid": "4420",
            "subnqn": "nqn.2018-09.io.spdk:cnode4",
            "adrifam": "IPv4"
        }
    },
    {
        "method": "bdev_nvme_attach_controller",
        "params": {
    
```

```
        "name": "Nvme5",
        "trtype": "rdma",
        "traddr": "10.0.1.1",
        "trsvcid": "4420",
        "subnqn": "nqn.2018-09.io.spdk:cnode5",
        "adrifam": "IPv4"
    }
},
{
    "method": "bdev_nvme_attach_controller",
    "params": {
        "name": "Nvme6",
        "trtype": "rdma",
        "traddr": "10.0.1.1",
        "trsvcid": "4420",
        "subnqn": "nqn.2018-09.io.spdk:cnode6",
        "adrifam": "IPv4"
    }
},
{
    "method": "bdev_nvme_attach_controller",
    "params": {
        "name": "Nvme7",
        "trtype": "rdma",
        "traddr": "10.0.1.1",
        "trsvcid": "4420",
        "subnqn": "nqn.2018-09.io.spdk:cnode7",
        "adrifam": "IPv4"
    }
}
]
}
```

Appendix B – Test Case 3 SPDK NVMe-oF Initiator bdev configuration

```
{  
    "subsystems": [  
        {  
            "subsystem": "bdev",  
            "config": [  
                {  
                    "method": "bdev_nvme_attach_controller",  
                    "params": {  
                        "name": "Nvme0",  
                        "trtype": "tcp",  
                        "traddr": "20.0.0.1",  
                        "trsvcid": "4420",  
                        "subnqn": "nqn.2018-09.io.spdk:cnode0",  
                        "adrifam": "IPv4"  
                    }  
                }  
            ]  
        }  
    ]  
}
```

Appendix C - Kernel NVMe-oF RDMA Target configuration

Example Linux Kernel NVMe-oF RDMA Target configuration for Test Case 4.

```
{  
    "ports": [  
        {  
            "addr": {  
                "adrifam": "ipv4",  
                "traddr": "20.0.0.1",  
                "trsvcid": "4420",  
                "trtype": "rdma"  
            },  
            "portid": 1,  
            "referrals": [],  
            "subsystems": [  
                "nqn.2018-09.io.spdk:cnode1"  
            ]  
        },  
        {  
            "addr": {  
                "adrifam": "ipv4",  
                "traddr": "20.0.0.1",  
                "trsvcid": "4421",  
                "trtype": "rdma"  
            },  
            "portid": 2,  
            "referrals": [],  
            "subsystems": [  
                "nqn.2018-09.io.spdk:cnode2"  
            ]  
        }  
    ]  
}
```

```
        "nqn.2018-09.io.spdk:cnode2"
    ],
},
{
    "addr": {
        "adrfam": "ipv4",
        "traddr": "20.0.0.1",
        "trsvid": "4422",
        "trtype": "rdma"
    },
    "portid": 3,
    "referrals": [],
    "subsystems": [
        "nqn.2018-09.io.spdk:cnode3"
    ]
},
{
    "addr": {
        "adrfam": "ipv4",
        "traddr": "20.0.0.1",
        "trsvid": "4423",
        "trtype": "rdma"
    },
    "portid": 4,
    "referrals": [],
    "subsystems": [
        "nqn.2018-09.io.spdk:cnode4"
    ]
},
{
    "addr": {
        "adrfam": "ipv4",
        "traddr": "20.0.1.1",
        "trsvid": "4424",
        "trtype": "rdma"
    },
    "portid": 5,
    "referrals": [],
    "subsystems": [
        "nqn.2018-09.io.spdk:cnode5"
    ]
},
{
    "addr": {
        "adrfam": "ipv4",
        "traddr": "20.0.1.1",
        "trsvid": "4425",
        "trtype": "rdma"
    },
    "portid": 6,
    "referrals": [],
    "subsystems": [
        "nqn.2018-09.io.spdk:cnode6"
    ]
},
{
    "addr": {
```

```

        "adrfam": "ipv4",
        "traddr": "20.0.1.1",
        "trsvcid": "4426",
        "trtype": "rdma"
    },
    "portid": 7,
    "referrals": [],
    "subsystems": [
        "nqn.2018-09.io.spdk:cnode7"
    ]
},
{
    "addr": {
        "adrfam": "ipv4",
        "traddr": "20.0.1.1",
        "trsvcid": "4427",
        "trtype": "rdma"
    },
    "portid": 8,
    "referrals": [],
    "subsystems": [
        "nqn.2018-09.io.spdk:cnode8"
    ]
},
{
    "addr": {
        "adrfam": "ipv4",
        "traddr": "10.0.0.1",
        "trsvcid": "4428",
        "trtype": "rdma"
    },
    "portid": 9,
    "referrals": [],
    "subsystems": [
        "nqn.2018-09.io.spdk:cnode9"
    ]
},
{
    "addr": {
        "adrfam": "ipv4",
        "traddr": "10.0.0.1",
        "trsvcid": "4429",
        "trtype": "rdma"
    },
    "portid": 10,
    "referrals": [],
    "subsystems": [
        "nqn.2018-09.io.spdk:cnode10"
    ]
},
{
    "addr": {
        "adrfam": "ipv4",
        "traddr": "10.0.0.1",
        "trsvcid": "4430",
        "trtype": "rdma"
    },

```

```
"portid": 11,
"referrals": [],
"subsystems": [
    "nqn.2018-09.io.spdk:cnode11"
]
},
{
    "addr": {
        "adrfam": "ipv4",
        "traddr": "10.0.0.1",
        "trsvid": "4431",
        "trtype": "rdma"
    },
    "portid": 12,
    "referrals": [],
    "subsystems": [
        "nqn.2018-09.io.spdk:cnode12"
    ]
},
{
    "addr": {
        "adrfam": "ipv4",
        "traddr": "10.0.1.1",
        "trsvid": "4432",
        "trtype": "rdma"
    },
    "portid": 13,
    "referrals": [],
    "subsystems": [
        "nqn.2018-09.io.spdk:cnode13"
    ]
},
{
    "addr": {
        "adrfam": "ipv4",
        "traddr": "10.0.1.1",
        "trsvid": "4433",
        "trtype": "rdma"
    },
    "portid": 14,
    "referrals": [],
    "subsystems": [
        "nqn.2018-09.io.spdk:cnode14"
    ]
},
{
    "addr": {
        "adrfam": "ipv4",
        "traddr": "10.0.1.1",
        "trsvid": "4434",
        "trtype": "rdma"
    },
    "portid": 15,
    "referrals": [],
    "subsystems": [
        "nqn.2018-09.io.spdk:cnode15"
    ]
}
```

```

        },
        {
            "addr": {
                "adrfam": "ipv4",
                "traddr": "10.0.1.1",
                "trsvcid": "4435",
                "trtype": "rdma"
            },
            "portid": 16,
            "referrals": [],
            "subsystems": [
                "nqn.2018-09.io.spdk:cnode16"
            ]
        }
    ],
    "hosts": [],
    "subsystems": [
        {
            "allowed_hosts": [],
            "attr": {
                "allow_any_host": "1",
                "version": "1.3"
            },
            "namespaces": [
                {
                    "device": {
                        "path": "/dev/nvme0n1",
                        "uuid": "b53be81d-6f5c-4768-b3bd-203614d8cf20"
                    },
                    "enable": 1,
                    "nsid": 1
                }
            ],
            "nqn": "nqn.2018-09.io.spdk:cnode1"
        },
        {
            "allowed_hosts": [],
            "attr": {
                "allow_any_host": "1",
                "version": "1.3"
            },
            "namespaces": [
                {
                    "device": {
                        "path": "/dev/nvme1n1",
                        "uuid": "12fcf584-9c45-4b6b-abc9-63a763455cf7"
                    },
                    "enable": 1,
                    "nsid": 2
                }
            ],
            "nqn": "nqn.2018-09.io.spdk:cnode2"
        },
        {
            "allowed_hosts": [],
            "attr": {
                "allow_any_host": "1",

```

```
        "version": "1.3"
    },
    "namespaces": [
        {
            "device": {
                "path": "/dev/nvme2n1",
                "uuid": "ceae8569-69e9-4831-8661-90725bdf768d"
            },
            "enable": 1,
            "nsid": 3
        }
    ],
    "nqn": "nqn.2018-09.io.spdk:cnode3"
},
{
    "allowed_hosts": [],
    "attr": {
        "allow_any_host": "1",
        "version": "1.3"
    },
    "namespaces": [
        {
            "device": {
                "path": "/dev/nvme3n1",
                "uuid": "39f36db4-2cd5-4f69-b37d-1192111d52a6"
            },
            "enable": 1,
            "nsid": 4
        }
    ],
    "nqn": "nqn.2018-09.io.spdk:cnode4"
},
{
    "allowed_hosts": [],
    "attr": {
        "allow_any_host": "1",
        "version": "1.3"
    },
    "namespaces": [
        {
            "device": {
                "path": "/dev/nvme4n1",
                "uuid": "984aed55-90ed-4517-ae36-d3afb92dd41f"
            },
            "enable": 1,
            "nsid": 5
        }
    ],
    "nqn": "nqn.2018-09.io.spdk:cnode5"
},
{
    "allowed_hosts": [],
    "attr": {
        "allow_any_host": "1",
        "version": "1.3"
    },
    "namespaces": [
```

```
{
    "device": {
        "path": "/dev/nvme5n1",
        "uuid": "d6d16e74-378d-40ad-83e7-b8d8af3d06a6"
    },
    "enable": 1,
    "nsid": 6
},
],
"nqn": "nqn.2018-09.io.spdk:cnode6"
},
{
    "allowed_hosts": [],
    "attr": {
        "allow_any_host": "1",
        "version": "1.3"
    },
    "namespaces": [
        {
            "device": {
                "path": "/dev/nvme6n1",
                "uuid": "a65dc00e-d35c-4647-9db6-c2a8d90db5e8"
            },
            "enable": 1,
            "nsid": 7
        }
    ],
    "nqn": "nqn.2018-09.io.spdk:cnode7"
},
{
    "allowed_hosts": [],
    "attr": {
        "allow_any_host": "1",
        "version": "1.3"
    },
    "namespaces": [
        {
            "device": {
                "path": "/dev/nvme7n1",
                "uuid": "1b242cb7-8e47-4079-a233-83e2cd47c86c"
            },
            "enable": 1,
            "nsid": 8
        }
    ],
    "nqn": "nqn.2018-09.io.spdk:cnode8"
},
{
    "allowed_hosts": [],
    "attr": {
        "allow_any_host": "1",
        "version": "1.3"
    },
    "namespaces": [
        {
            "device": {
                "path": "/dev/nvme8n1",

```

```
        "uuid": "f12bb0c9-a2c6-4eef-a94f-5c4887bbf77f"
    },
    "enable": 1,
    "nsid": 9
}
],
"nqn": "nqn.2018-09.io.spdk:cnode9"
},
{
    "allowed_hosts": [],
    "attr": {
        "allow_any_host": "1",
        "version": "1.3"
    },
    "namespaces": [
        {
            "device": {
                "path": "/dev/nvme9n1",
                "uuid": "40fae536-227b-47d2-bd74-8ab76ec7603b"
            },
            "enable": 1,
            "nsid": 10
        }
    ],
    "nqn": "nqn.2018-09.io.spdk:cnode10"
},
{
    "allowed_hosts": [],
    "attr": {
        "allow_any_host": "1",
        "version": "1.3"
    },
    "namespaces": [
        {
            "device": {
                "path": "/dev/nvme10n1",
                "uuid": "b9756b86-263a-41cf-a68c-5cfb23c7a6eb"
            },
            "enable": 1,
            "nsid": 11
        }
    ],
    "nqn": "nqn.2018-09.io.spdk:cnode11"
},
{
    "allowed_hosts": [],
    "attr": {
        "allow_any_host": "1",
        "version": "1.3"
    },
    "namespaces": [
        {
            "device": {
                "path": "/dev/nvme11n1",
                "uuid": "9d7e74cc-97f3-40fb-8e90-f2d02b5fff4c"
            },
            "enable": 1,
```

```

        "nsid": 12
    }
],
"nqn": "nqn.2018-09.io.spdk:cnode12"
},
{
    "allowed_hosts": [],
    "attr": {
        "allow_any_host": "1",
        "version": "1.3"
    },
    "namespaces": [
        {
            "device": {
                "path": "/dev/nvme12n1",
                "uuid": "d3f4017b-4f7d-454d-94a9-ea75ffc7436d"
            },
            "enable": 1,
            "nsid": 13
        }
    ],
    "nqn": "nqn.2018-09.io.spdk:cnode13"
},
{
    "allowed_hosts": [],
    "attr": {
        "allow_any_host": "1",
        "version": "1.3"
    },
    "namespaces": [
        {
            "device": {
                "path": "/dev/nvme13n1",
                "uuid": "6b9a65a3-6557-4713-8bad-57d9c5cb17a9"
            },
            "enable": 1,
            "nsid": 14
        }
    ],
    "nqn": "nqn.2018-09.io.spdk:cnode14"
},
{
    "allowed_hosts": [],
    "attr": {
        "allow_any_host": "1",
        "version": "1.3"
    },
    "namespaces": [
        {
            "device": {
                "path": "/dev/nvme14n1",
                "uuid": "ed69ba4d-8727-43bd-894a-7b08ade4f1b1"
            },
            "enable": 1,
            "nsid": 15
        }
    ],

```

```
"nqn": "nqn.2018-09.io.spdk:cnode15"
},
{
  "allowed_hosts": [],
  "attr": {
    "allow_any_host": "1",
    "version": "1.3"
  },
  "namespaces": [
    {
      "device": {
        "path": "/dev/nvme15n1",
        "uuid": "5b8e9af4-0ab4-47fb-968f-b13e4b607f4e"
      },
      "enable": 1,
      "nsid": 16
    }
  ],
  "nqn": "nqn.2018-09.io.spdk:cnode16"
}
]
```

Notices & Disclaimers

Performance varies by use, configuration and other factors. Learn more at www.intel.com/PerformanceIndex.

Performance results are based on testing as of dates shown in configurations and may not reflect all publicly available updates.

Your costs and results may vary.

No product or component can be absolutely secure.

Intel technologies may require enabled hardware, software or service activation.

© Intel Corporation. Intel, the Intel logo, and other Intel marks are trademarks of Intel Corporation or its subsidiaries. Other names and brands may be claimed as the property of others.

§