intel.

# SPDK NVMe-oF RDMA (Target & Initiator) Performance Report

# Release 24.05

## Mellanox ConnectX-5 version

**Testing Date:** June 2024

**Performed by:**

Karol Latecki (karol.latecki@intel.com)

Jaroslaw Chachulski (jaroslawx.chachulski@intel.com)

**Acknowledgments:**

Krzysztof Karas (krzysztof.karas@intel.com)

# *Contents*

# *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 | SuperMicro® Ultra SuperServer SYS-220U-TNR |
| Motherboard | Server board X12DPU-6 |
| CPU | 2 CPU sockets, Intel(R) Xeon(R) Gold 6348 CPU @ 2.60GHz<br><br>Number of cores: 28 per socket, number of threads: 56 per socket<br>Both sockets populated<br>Microcode: 0xd0003d1 |
| Memory | 16 x 32GB SK Hynix HMA84GR7DJR4N-XN, DDR4, 3200MHz<br>Total of 512GB |
| Operating System | Fedora 37 |
| BIOS | 1.8 |
| Linux kernel version | 6.0.18-300.fc37.x86_64<br>Spectre-meltdown mitigations enabled |
| SPDK version | SPDK 24.05 |
| Storage | **OS:** 1x 250GB Crucial CT250MX500SSD1<br><br>**Storage Target**:<br>14x Kioxia® KCM61VUL3T20 3.2TBs (FW: 0105) (6 on CPU NUMA Node 0, 8 on CPU NUMA Node 1) |
| NIC | 4x 100GbE Mellanox ConnectX-5 NICs. Both ports connected.<br>2 NICs per CPU socket. |

# 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)<br>Number of cores: 28 per socket, number of threads: 56 per socket (Both sockets populated)<br>Microcode: 0xd0003b9 |
| Memory | 16 x 32GB Micron 36ASF4G72PZ-3G2J3, DDR4, 3200MHz<br>Total 512GBs |
| Operating System | Fedora 37 |
| BIOS | SE5C620.86B.01.01.0009.2311021928 |
| Linux kernel version | 6.0.18-300.fc37.x86_64<br>Spectre-meltdown mitigations enabled |
| SPDK version | SPDK 24.05 |
| Storage | **OS:** 1x 250GB Crucial CT250MX500SSD1 |
| NIC | 2x 100GbE Mellanox ConnectX-5 Ex NIC. Single port on each NIC connected to Target server. 1 NIC per CPU socket. |

# 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)<br>Number of cores: 28 per socket, number of threads: 56 per socket (Both sockets populated)<br>Microcode: 0xd0003b9 |
| Memory | 16 x 32GB Micron 36ASF4G72PZ-3G2J3, DDR4, 3200MHz<br>Total 512GBs |
| Operating System | Fedora 37 |
| BIOS | SE5C620.86B.01.01.0009.2311021928 |
| Linux kernel version | 6.0.18-300.fc37.x86_64<br>Spectre-meltdown mitigations enabled |
| SPDK version | SPDK 24.05 |
| Storage | **OS:** 1x 250GB Crucial CT250MX500SSD1 |
| NIC | 2x 100GbE Mellanox ConnectX-5 Ex NIC. Single port on each NIC connected to Target server. 1 NIC per CPU socket. |

# BIOS settings

*Table 4: Test systems BIOS settings*

| Item | Description |
|---|---|
| **BIOS**<br>*(Applied to all 3 systems)* | Hyper threading Enabled<br>CPU Power and Performance Policy:<br>    •   "Extreme Performance" for Target<br>    •   "Performance" for Initiators<br>CPU C-state No Limit<br>CPU P-state Enabled<br>Enhanced Intel® SpeedStep® Tech Enabled<br>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 NVM Express* (NVMe) block protocol over network fabrics such as RDMA (iWARP, RoCE, InfiniBand™), Fiber 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 use the Infiniband/RDMA verbs API to access an RDMA-capable NIC. These should work on all flavors of RDMA transports but are currently 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 and two 100GbE Mellanox ConnectX-5 NICs connected to provide up to 200GbE of network bandwidth. Each Initiator system has one Mellanox ConnectX-5 Ex 100GbE NIC 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).
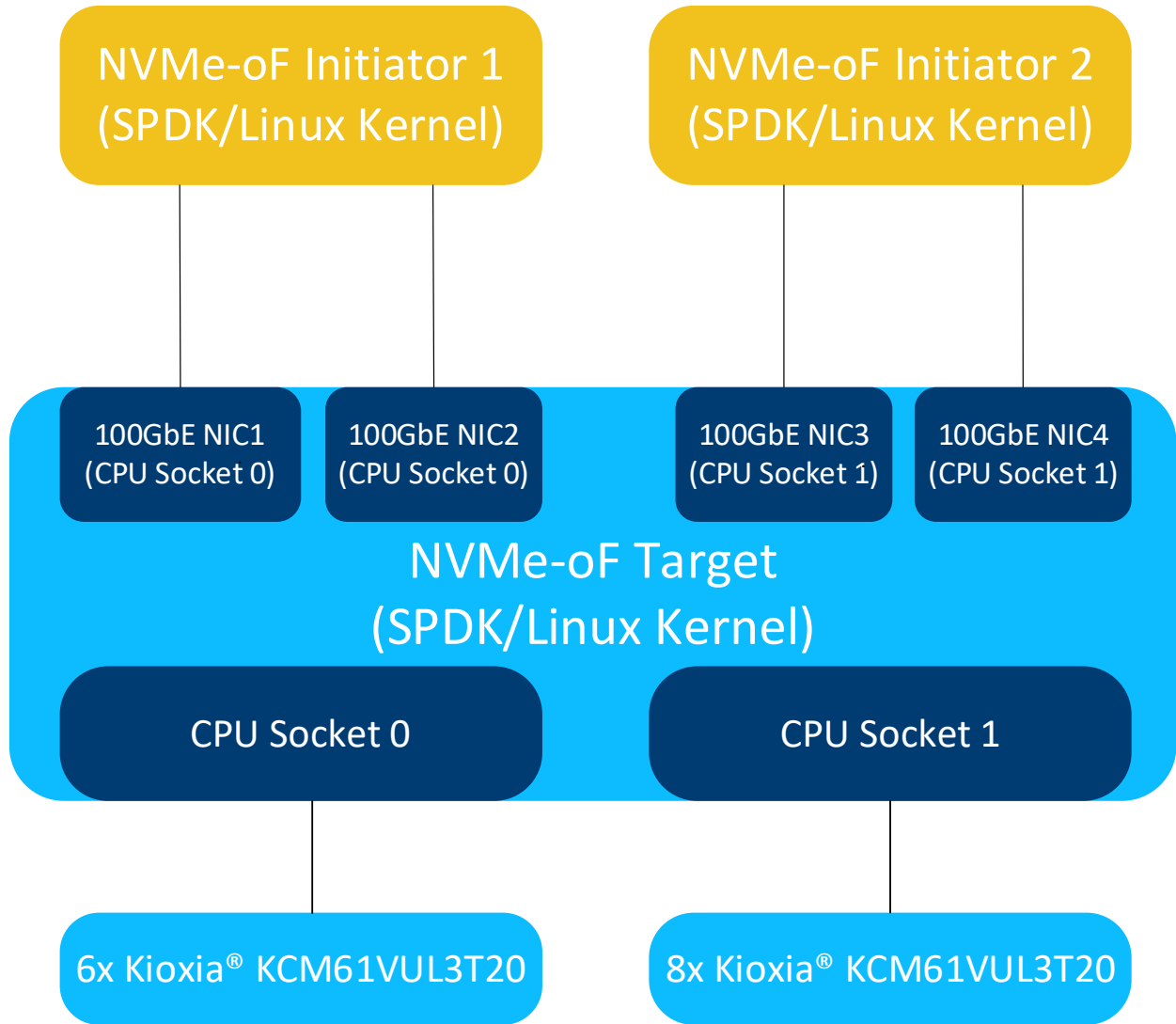
*Figure 1: High-Level NVMe-oF RDMA performance testing setup*

**intel.**

# *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 NIC. The SPDK bdev fio plugin was used to target 7 NVMe-oF bdevs on each of the host. The SPDK Target was configured to use up to 10 CPU cores. We ran the 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 Random Read and 4KiB Random 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.<br><br>**Construct NVMe bdevs:**<br>bdev_nvme_attach_controller -t PCIe -b Nvme0 -a 0000:17:00.0<br>bdev_nvme_attach_controller -t PCIe -b Nvme1 -a 0000:18:00.0<br>bdev_nvme_attach_controller -t PCIe -b Nvme2 -a 0000:65:00.0<br>bdev_nvme_attach_controller -t PCIe -b Nvme3 -a 0000:66:00.0<br>bdev_nvme_attach_controller -t PCIe -b Nvme4 -a 0000:67:00.0<br>bdev_nvme_attach_controller -t PCIe -b Nvme5 -a 0000:68:00.0<br>bdev_nvme_attach_controller -t PCIe -b Nvme6 -a 0000:98:00.0<br>bdev_nvme_attach_controller -t PCIe -b Nvme7 -a 0000:99:00.0<br>bdev_nvme_attach_controller -t PCIe -b Nvme8 -a 0000:9a:00.0<br>bdev_nvme_attach_controller -t PCIe -b Nvme9 -a 0000:9b:00.0<br>bdev_nvme_attach_controller -t PCIe -b Nvme10 -a 0000:e3:00.0<br>bdev_nvme_attach_controller -t PCIe -b Nvme11 -a 0000:e4:00.0<br>bdev_nvme_attach_controller -t PCIe -b Nvme12 -a 0000:e5:00.0<br>bdev_nvme_attach_controller -t PCIe -b Nvme13 -a 0000:e6:00.0 |

| | |
|---|---|
| | **Create RDMA transport layer:**<br>nvmf_create_transport -t RDMA -n 8192<br>{<br>    trtype: "RDMA"<br>    max_queue_depth: 128<br>    max_qpairs_per_ctrlr: 64<br>    in_capsule_data_size: 4096<br>    max_io_size: 131072<br>    io_unit_size: 8192<br>    max_aq_depth: 128<br>    num_shared_buffers: 8192<br>    buf_cache_size: 32<br>}<br><br>**Create NVMe-oF subsystems and add NVMe bdevs as namespaces:**<br>for i in $(seq 1 16); do<br>    nvmf_subsystem_create nqn.2018-09.io.spdk:cnode${i} -s SPDK00${i} -a -m 8<br>    nvmf_subsystem_add_ns nqn.2018-09.io.spdk:cnode${i} Nvme$((i-1))n1<br>done<br><br>**Add listeners to NVMe-oF Subsystems:**<br>i=1<br>ips=(20.0.0.1 20.0.1.1 10.0.0.1 10.0.1.1)<br>for ip in ${ips[@]}; do<br>    for j in $(seq 1 4); do<br>        nvmf_subsystem_add_listener nqn.2018-09.io.spdk:cnode${i} -t rdma \\<br>                          -f ipv4 -s 4420 -a ${ip}<br>        ((i++))<br>    done<br>done |
| **SPDK NVMe-oF Initiator - fio plugin configuration** | **BDEV.conf**<br>  See [Appendix A](#)<br><br>**fio.conf**<br>[global]<br>ioengine=/tmp/spdk/examples/bdev/fio_plugin/fio_plugin<br>spdk_json_conf=/tmp/spdk/bdev.conf<br>thread=1<br>group_reporting=1<br>direct=1<br>norandommap=1<br>rw=randrw<br>rwmixread={100, 70, 0}<br>bs=4k<br>iodepth={1, 64, 128, 192, 256}<br>time_based=1<br>ramp_time=60<br>runtime=300<br><br>[filename0]<br>filename=Nvme0n1<br>[filename1]<br>filename=Nvme1n1<br>[filename2]<br>filename=Nvme2n1<br>[filename3]<br>filename=Nvme3n1<br>[filename4]<br>filename=Nvme4n1<br>[filename5] |

| | filename=Nvme5n1<br>[filename6]<br>filename=Nvme6n1 | |
|---|---|---|

# 4KiB Random Read Results

*Table 6: SPDK NVMe-oF RDMA Target Core Scaling results, Random Read IOPS, QD=128*

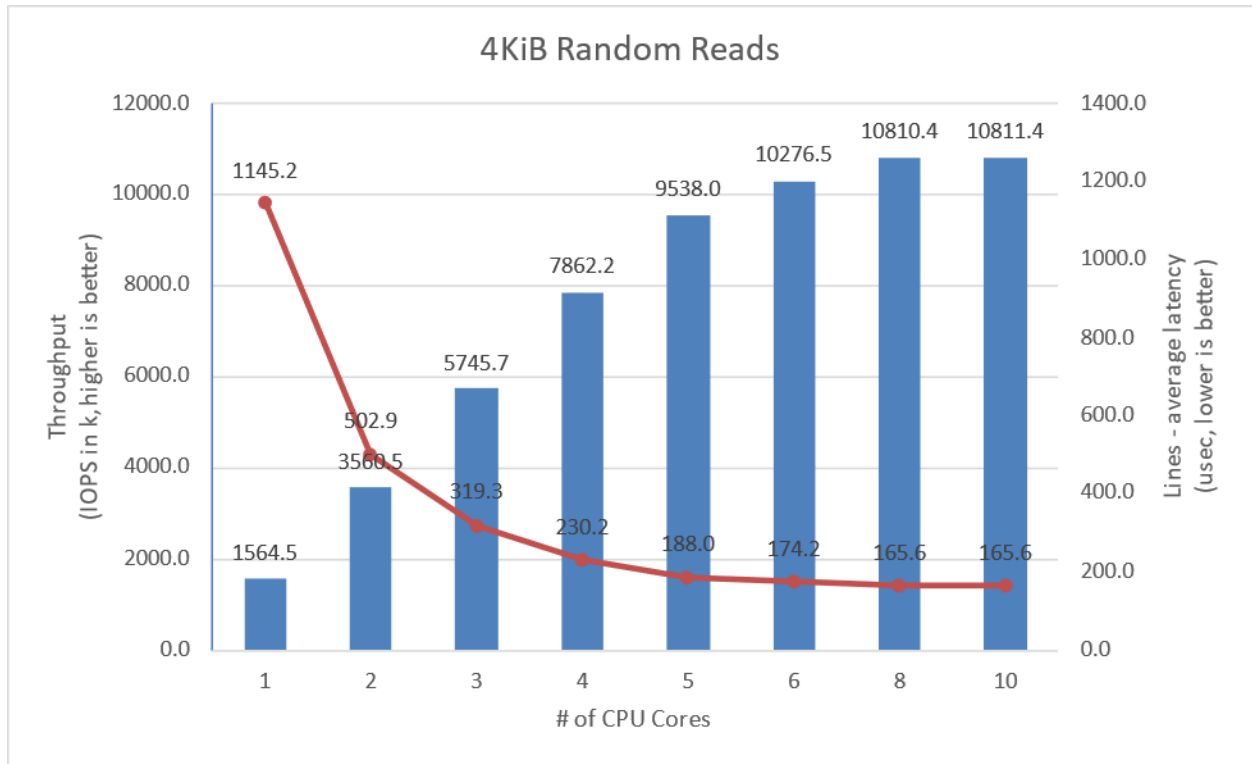| # of Cores | Bandwidth (MiBps) | Throughput (IOPS k) | Avg. Latency (usec) |
|---|---|---|---|
| **1 core** | 6111.18 | 1564.5 | 1145.2 |
| **2 cores** | 13908.15 | 3560.5 | 502.9 |
| **3 cores** | 22444.09 | 5745.7 | 319.3 |
| **4 cores** | 30711.77 | 7862.2 | 230.2 |
| **5 cores** | 37257.97 | 9538.0 | 188.0 |
| **6 cores** | 40142.72 | 10276.5 | 174.2 |
| **8 cores** | 42227.96 | 10810.4 | 165.6 |
| **10 cores** | 42231.96 | 10811.4 | 165.6 |



*Figure 2: SPDK NVMe-oF RDMA Target I/O core scaling: IOPS vs. Latency while running 4KiB 100% Random Read workload at QD=128*

# 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 | 8181.81 | 2094.5 | 426.0 |
| 2 cores | 19075.22 | 4883.3 | 181.0 |
| 3 cores | 29145.51 | 7461.2 | 121.4 |
| 4 cores | 33741.26 | 8637.8 | 106.0 |
| 5 cores | 31810.81 | 8143.6 | 110.0 |
| 6 cores | 40859.20 | 10460.0 | 84.5 |
| 8 cores | 40954.82 | 10484.4 | 84.8 |
| 10 cores | 41059.12 | 10511.1 | 84.7 |

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.
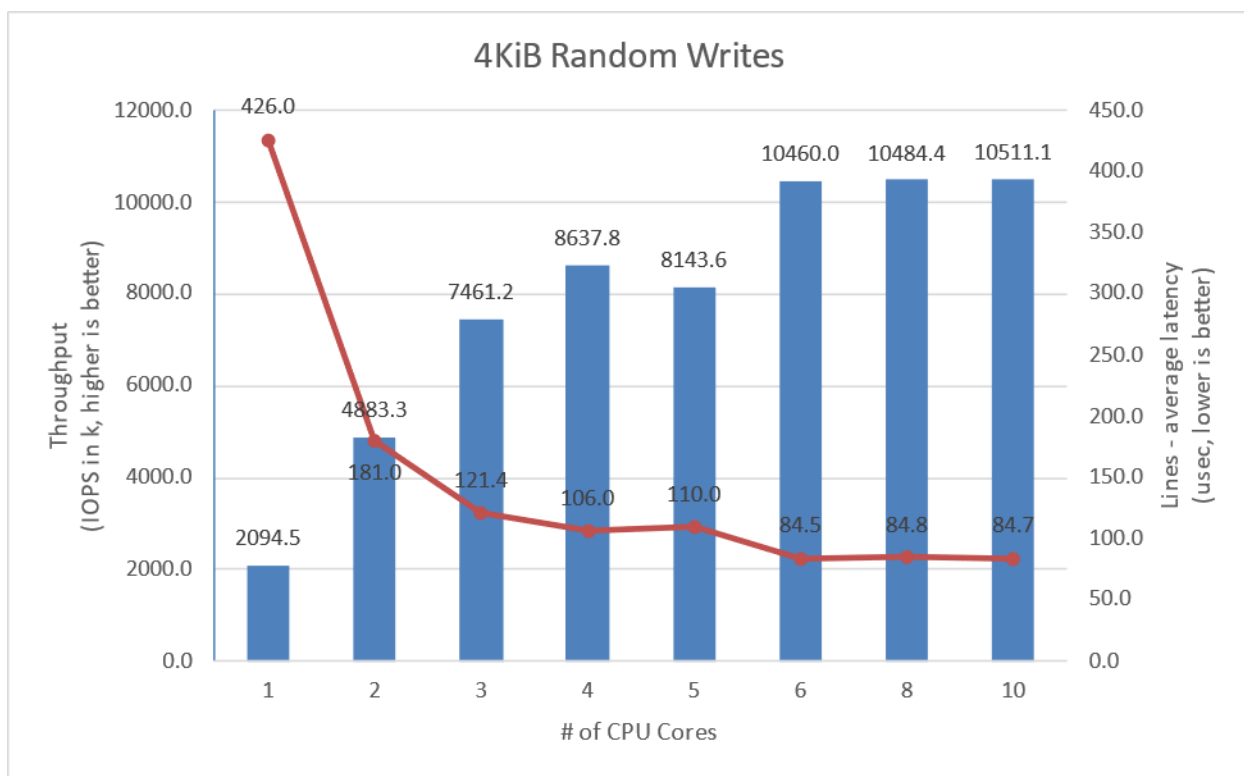


*Figure 3: SPDK NVMe-oF RDMA Target I/O core scaling: IOPS vs. Latency while running 4KiB 100% Random Write Workload at QD=64*

# 4KiB Random Read-Write Results

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

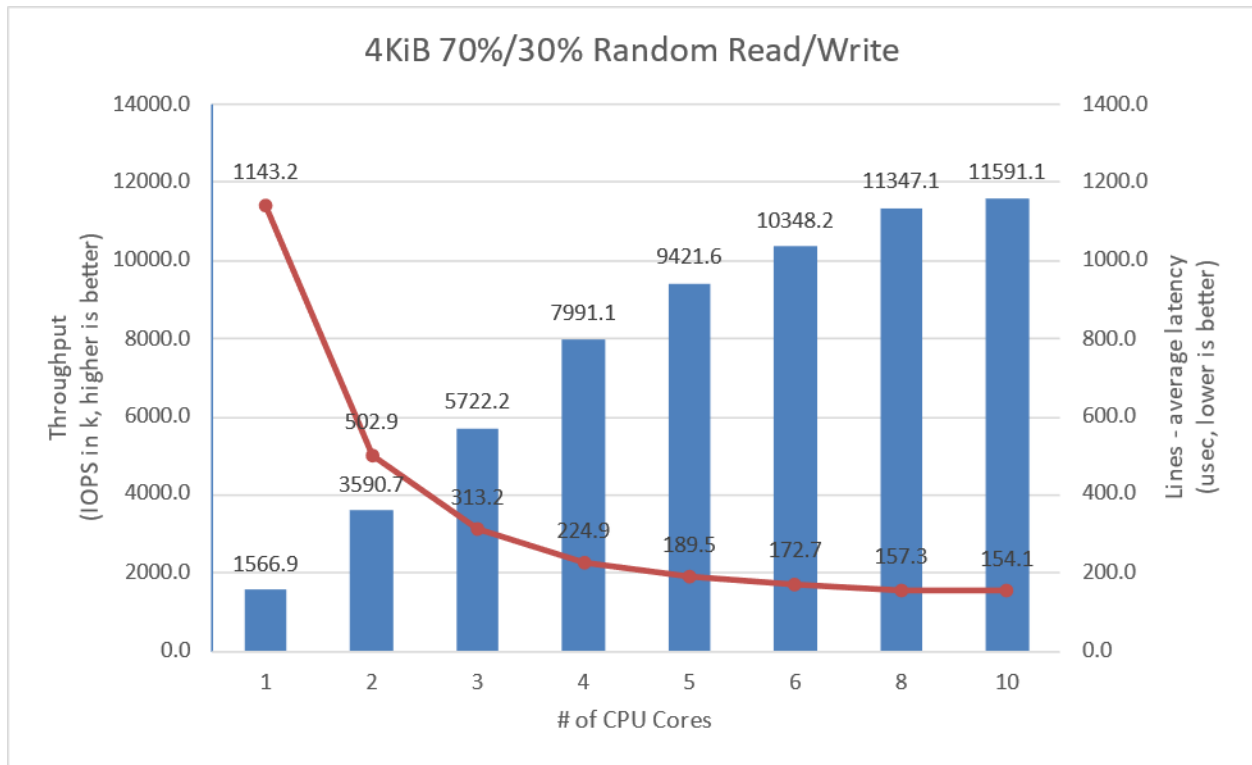| # of Cores | Bandwidth (MiBps) | Throughput (IOPS k) | Avg. Latency (usec) |
|------------|-------------------|---------------------|---------------------|
| 1 core | 6120.87 | 1566.9 | 1143.2 |
| 2 cores | 14026.07 | 3590.7 | 502.9 |
| 3 cores | 22352.19 | 5722.2 | 313.2 |
| 4 cores | 31215.41 | 7991.1 | 224.9 |
| 5 cores | 36803.26 | 9421.6 | 189.5 |
| 6 cores | 40422.82 | 10348.2 | 172.7 |
| 8 cores | 44324.64 | 11347.1 | 157.3 |
| 10 cores | 45277.72 | 11591.1 | 154.1 |



*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 4 and 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 | 43337.00 | 346.7 | 161.4 |
| 2 cores | 43542.69 | 348.3 | 160.6 |
| 3 cores | 43580.40 | 348.6 | 160.5 |
| 4 cores | 43606.98 | 348.9 | 160.4 |

*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 | 43971.37 | 351.8 | 159.0 |
| 2 cores | 43994.18 | 352.0 | 158.9 |
| 3 cores | 43996.86 | 352.0 | 158.9 |
| 4 cores | 43999.60 | 352.0 | 158.9 |

*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 | 49076.23 | 392.6 | 285.0 |
| 2 cores | 49388.10 | 395.1 | 283.1 |
| 3 cores | 49513.11 | 396.1 | 282.4 |
| 4 cores | 49553.00 | 396.4 | 282.2 |

# Conclusions

1.  For 100% 4KiB Random Read and 70/30% 4KiB Random Read/Write I/O workloads, throughput scales up almost linearly with the addition of I/O cores up to 4 cores. Adding more CPU cores results in non-linear performance gain until reaching network saturation at 8 CPU cores.

2.  For 100% 4KiB Random Write workload throughput scales up almost linearly with addition of up to 4 CPU cores and increasing CPU core number beyond 6 results in insignificant performance gain. Latency decreases non-linearly with addition of I/O cores.

3.  For large sequential I/Os, a single CPU core Target saturated the network bandwidth. Therefore, adding more CPU cores did not result in increased performance for these workloads because the network was saturated.

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

This test case demonstrates how the throughput of the SPDK NVMe-oF initiator, measured in IOPS (I/O per second), scales with the addition of more CPU cores 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.

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 iodepth 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** | **BDEV.conf**<br>See appendix B.<br><br>**fio.conf**<br>**For 1 CPU initiator configuration:**<br>[global]<br>ioengine=/tmp/spdk/examples/bdev/fio_plugin/fio_plugin<br>spdk_conf=/tmp/spdk/bdev.conf<br>thread=1<br>group_reporting=1<br>direct=1<br><br>norandommap=1<br>rw=randrw<br>rwmixread={100, 70, 0}<br>bs=4k<br>iodepth={1,32, 64, 128, 192}<br>time_based=1 |

```
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
```

**fio.conf**
**For CPU > 1 (up to N=8) initiator configuration "filename=NvmeXn1" are evenly spread across fio job threads:**

```
[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}
time_based=1
ramp_time=60
runtime=300
numjobs=1

[filename0]
filename=Nvme0n1
filename=Nvme1n1

[filename1]
filename=Nvme2n1
filename=Nvme3n1

[…]

[filename N-1]
filename=Nvme10n1
filename=Nvme11n1

[filename N]
filename=Nvme12n1
filename=Nvme13n1
```

# 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*

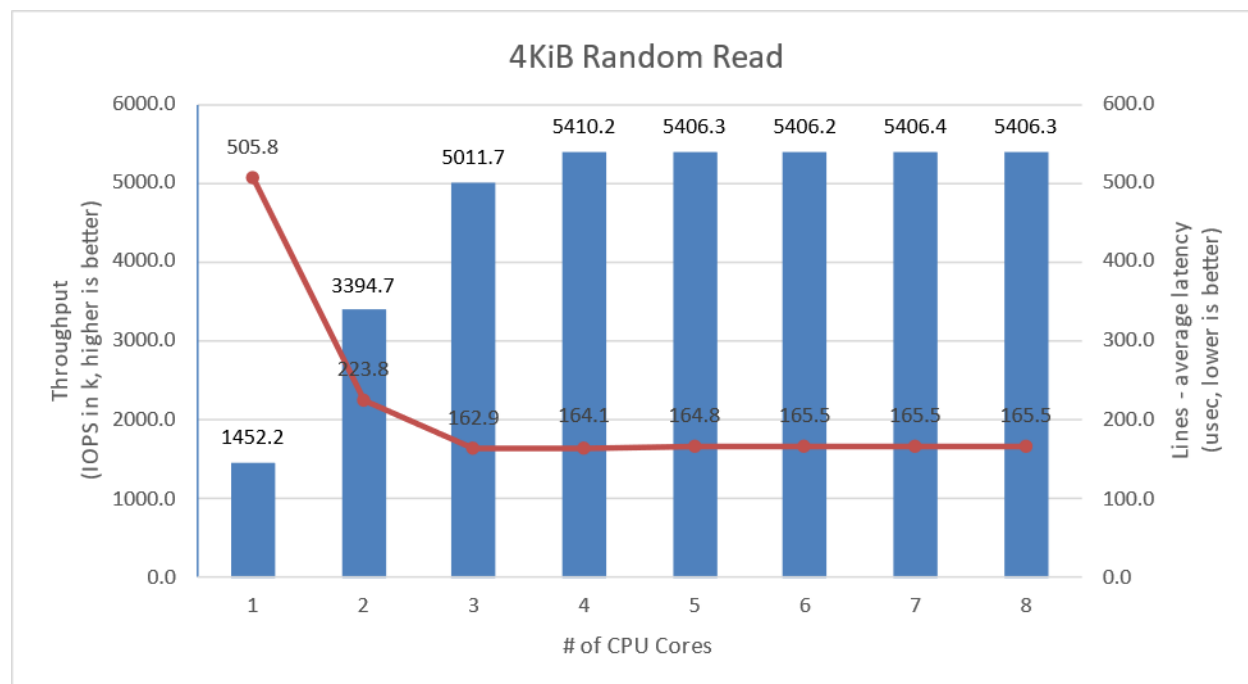| # of Initiator CPU Cores | Bandwidth (MiBps) | Throughput (IOPS k) | Avg. Latency (usec) |
|---|---|---|---|
| **1 core** | 5672.74 | 1452.2 | 505.8 |
| **2 cores** | 13260.56 | 3394.7 | 223.8 |
| **3 cores** | 19577.10 | 5011.7 | 162.9 |
| **4 cores** | 21133.60 | 5410.2 | 164.1 |
| **5 cores** | 21118.39 | 5406.3 | 164.8 |
| **6 cores** | 21117.97 | 5406.2 | 165.5 |
| **7 cores** | 21118.60 | 5406.4 | 165.5 |
| **8 cores** | 21118.31 | 5406.3 | 165.5 |



*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 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*

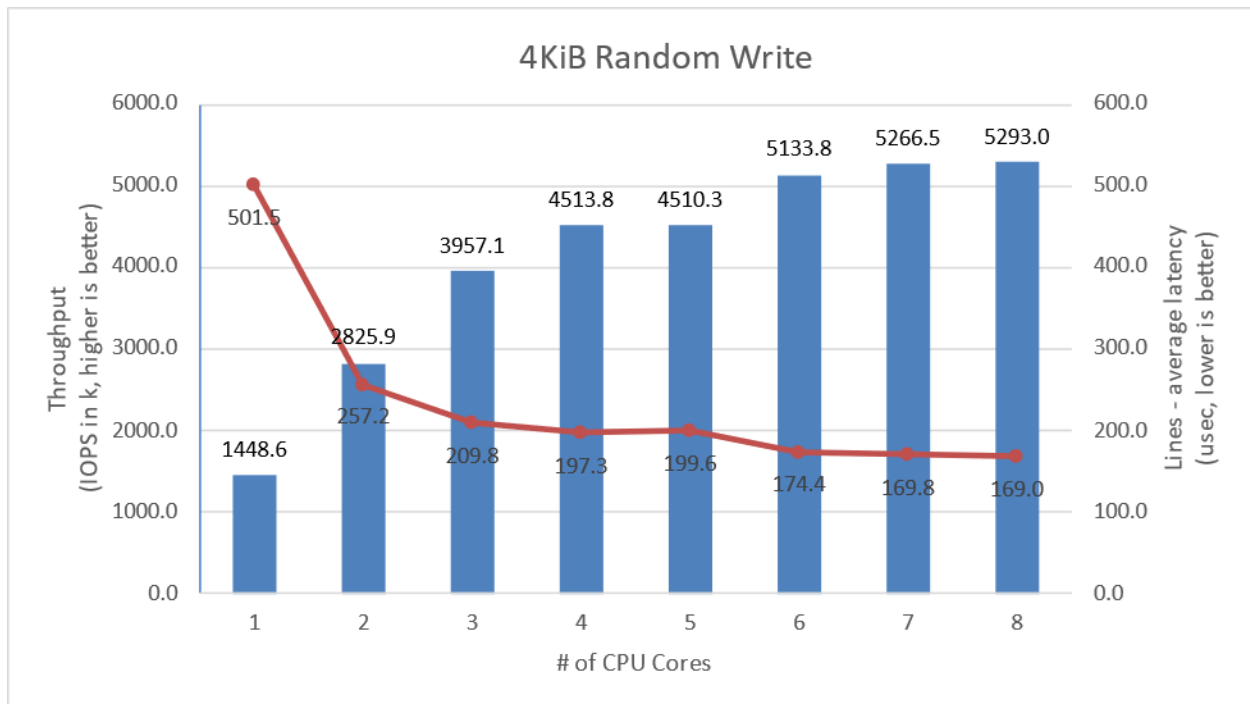| # of Initiator CPU Cores | Bandwidth (MiBps) | Throughput (IOPS k) | Avg. Latency (usec) |
|---|---|---|---|
| 1 core | 5658.74 | 1448.6 | 501.5 |
| 2 cores | 11038.65 | 2825.9 | 257.2 |
| 3 cores | 15457.44 | 3957.1 | 209.8 |
| 4 cores | 17632.22 | 4513.8 | 197.3 |
| 5 cores | 17618.54 | 4510.3 | 199.6 |
| 6 cores | 20053.85 | 5133.8 | 174.4 |
| 7 cores | 20572.36 | 5266.5 | 169.8 |
| 8 cores | 20675.69 | 5293.0 | 169.0 |



*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*

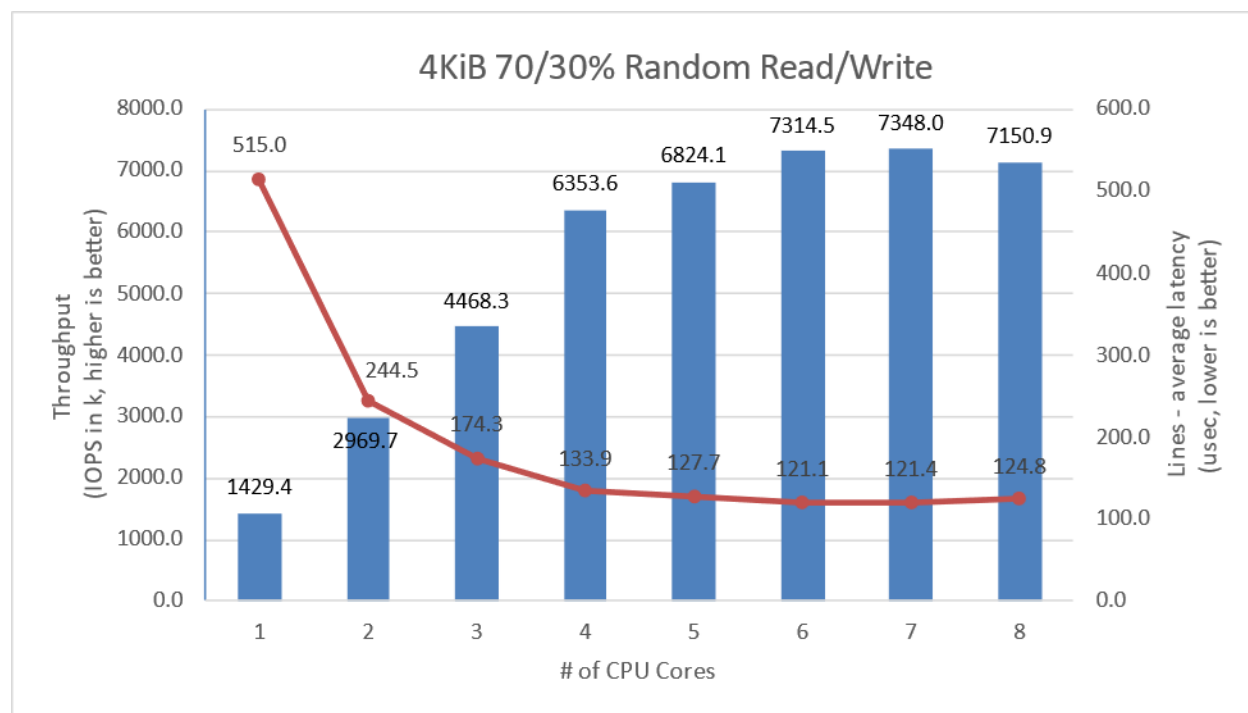| # of Initiator CPU Cores | Bandwidth (MiBps) | Throughput (IOPS k) | Avg. Latency (usec) |
|---|---|---|---|
| 1 core | 5583.42 | 1429.4 | 515.0 |
| 2 cores | 11600.25 | 2969.7 | 244.5 |
| 3 cores | 17454.38 | 4468.3 | 174.3 |
| 4 cores | 24818.91 | 6353.6 | 133.9 |
| 5 cores | 26656.63 | 6824.1 | 127.7 |
| 6 cores | 28572.24 | 7314.5 | 121.1 |
| 7 cores | 28703.02 | 7348.0 | 121.4 |
| 8 cores | 27933.32 | 7150.9 | 124.8 |



*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 close to linear when increasing the number of initiator CPUs up to 3 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 5.3 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 linear as we increased the number of initiator CPU cores up to 3. Beyond that scaling became non-linear. Performance peaked at around 7.35 million IOPS with the SPDK NVMe-oF RDMA initiator using 7 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 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 |
| **Test configuration** | |
| **SPDK NVMe-oF Target configuration** | The following commands are executed with spdk/scripts/rpc.py script to configure the SPDK NVMe-oF target.<br><br>nvmf_create_transport -t RDMA<br>(creates RDMA transport layer with default values:<br>trtype: "RDMA"<br>max_queue_depth: 128<br>max_qpairs_per_ctrlr: 64<br>in_capsule_data_size: 4096<br>max_io_size: 131072<br>io_unit_size: 8192<br>max_aq_depth: 128<br>num_shared_buffers: 8192<br>buf_cache_size: 32)<br><br>bdev_null_create Nvme0n1 10240 4096<br>nvmf_subsystem_create nqn.2018-09.io.spdk:cnode1 -s SPDK001 -a -m 8<br>nvmf_subsystem_add_ns nqn.2018-09.io.spdk:cnode1 Nvme0n1<br>nvmf_subsystem_add_listener nqn.2018-09.io.spdk:cnode1 -t rdma -f ipv4 -s 4420 -a 20.0.0.1 |

| Kernel NVMe-oF Target configuration | The following target configuration file loaded using nvmet-cli tool. |
|---|---|
| | {<br>  "ports": [<br>    {<br>      "addr": {<br>        "adrfam": "ipv4",<br>        "traddr": "20.0.0.1",<br>        "trsvcid": "4420",<br>        "trtype": "rdma"<br>      },<br>      "portid": 1,<br>      "referrals": [],<br>      "subsystems": [<br>        "nqn.2018-09.io.spdk:cnode1"<br>      ]<br>    }<br>  ],<br>  "hosts": [],<br>  "subsystems": [<br>    {<br>      "allowed_hosts": [],<br>      "attr": {<br>        "allow_any_host": "1",<br>        "version": "1.3"<br>      },<br>      "namespaces": [<br>        {<br>          "device": {<br>            "path": "/dev/nullb0",<br>            "uuid": "621e25d2-8334-4c1a-8532-b6454390b8f9"<br>          },<br>          "enable": 1,<br>          "nsid": 1<br>        }<br>      ],<br>      "nqn": "nqn.2018-09.io.spdk:cnode1"<br>    }<br>  ]<br>} |

| fio configuration | |
|---|---|
| SPDK NVMe-oF Initiator fio plugin configuration | **BDEV.conf**<br>See Appendix B.<br><br>**fio.conf**<br>[global]<br>ioengine=/tmp/spdk/examples/bdev/fio_plugin/fio_plugin<br>spdk_json_conf=/tmp/spdk/bdev.conf<br>thread=1<br>group_reporting=1<br>direct=1<br><br>norandommap=1<br>rw=randrw<br>rwmixread={100, 70, 0}<br>bs=4k<br>iodepth=1<br>time_based=1<br>ramp_time=60<br>runtime=300 |

| | [filename0]<br>filename=Nvme0n1 | |
|---|---|---|
| **Kernel initiator configuration** | **Device config**<br>The following configuration was performed using nvme-cli tool.<br>modprobe nvme-fabrics<br>nvme connect –n nqn.2018-09.io.spdk:cnode1 –t rdma –a 20.0.0.1 –s 4420<br><br>**fio.conf**<br>[global]<br>ioengine=libaio<br>thread=1<br>group_reporting=1<br>direct=1<br><br>norandommap=1<br>rw=randrw<br>rwmixread={100, 70, 0}<br>bs=4k<br>iodepth=1<br>time_based=1<br>ramp_time=60<br>runtime=300<br><br>[filename0]<br>filename=/dev/nvme0n1 | |

# SPDK vs Kernel NVMe-oF RDMA Target Results

This 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** | 9.39 | 103784 | 9.8 | 11.5 | 22.7 | 210.6 |
| **4KiB 100% Random Write** | 7.42 | 130364 | 7.5 | 10.1 | 17.3 | 173.7 |
| **4KiB 70/30% Random Read/Write** | 8.89 | 108984 | 9.7 | 13.7 | 21.3 | 202.2 |

*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** | 12.10 | 80983 | 12.6 | 13.7 | 18.6 | 100.9 |
| **4KiB 100% Random Write** | 11.12 | 87904 | 12.1 | 12.9 | 16.3 | 99.8 |
| **4KiB 70/30% Random Read/Write** | 11.81 | 82700 | 12.3 | 13.4 | 18.0 | 100.4 |

# 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.
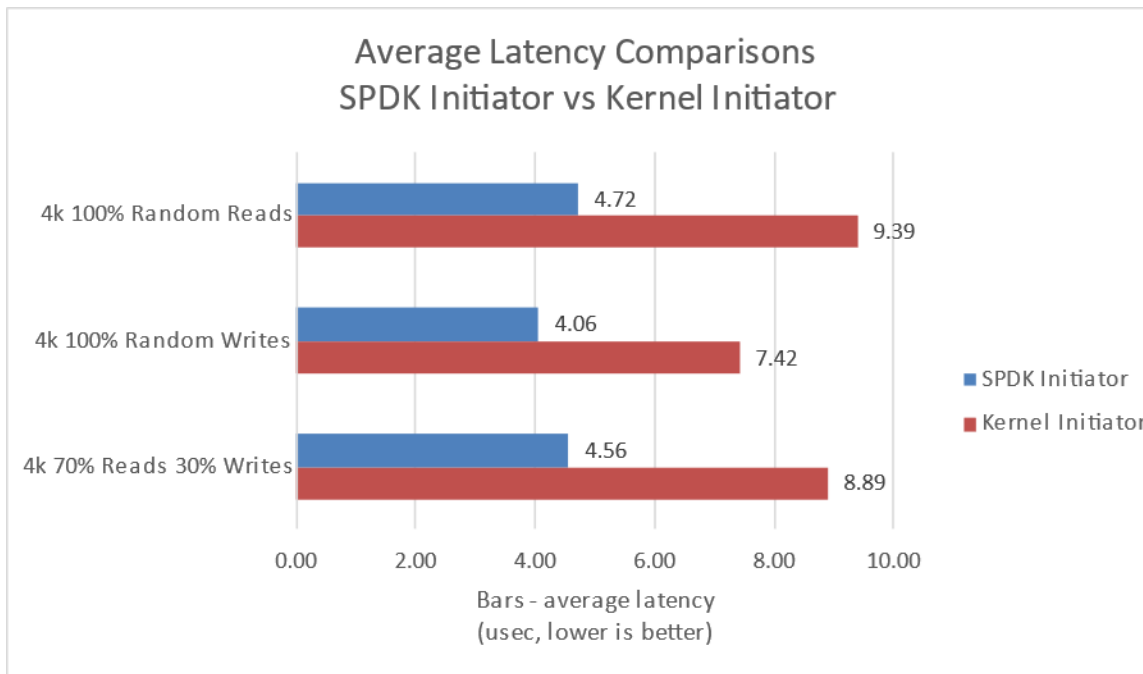


*Figure 9: SPDK vs. Kernel NVMe-oF RDMA Initiator average I/O latency for various workloads against SPDK Target*

*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 | 4.72 | 204989 | 4.9 | 6.4 | 10.3 | 149.2 |
| 4KiB 100% Random Write | 4.06 | 236857 | 4.0 | 5.8 | 9.4 | 102.6 |
| 4KiB 70/30% Random Read/Write | 4.56 | 211242 | 4.7 | 6.4 | 10.4 | 138.6 |

*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 | 9.39 | 103784 | 9.8 | 11.5 | 22.7 | 210.6 |
| 4KiB 100% Random Write | 7.42 | 130364 | 7.5 | 10.1 | 17.3 | 173.7 |
| 4KiB 70/30% Random Read/Write | 8.89 | 108984 | 9.7 | 13.7 | 21.3 | 202.2 |

# 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.
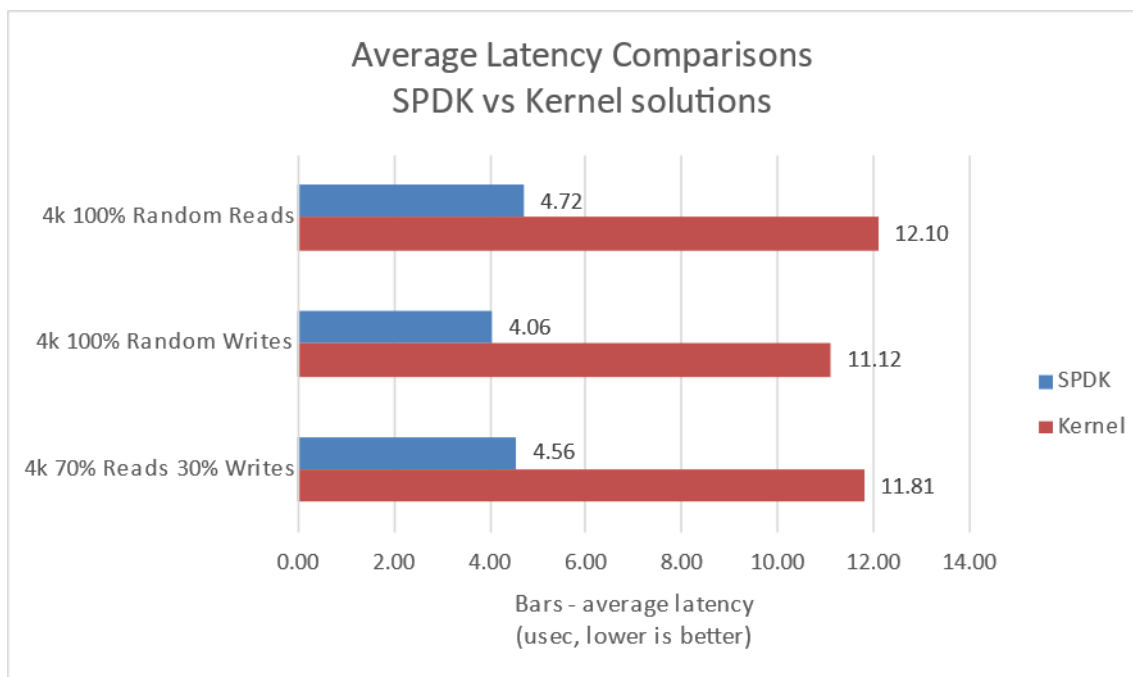


*Figure 10: SPDK vs. Kernel NVMe-oF RDMA solutions average I/O Latency for various workloads against SPDK Target*

*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 | 4.72 | 204989 | 4.9 | 6.4 | 10.3 | 149.2 |
| 4KiB 100% Random Write | 4.06 | 236857 | 4.0 | 5.8 | 9.4 | 102.6 |
| 4KiB 70/30% Random Read/Write | 4.56 | 211242 | 4.7 | 6.4 | 10.4 | 138.6 |

*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 | 12.10 | 80983 | 12.6 | 13.7 | 18.6 | 100.9 |
| 4KiB 100% Random Write | 11.12 | 87904 | 12.1 | 12.9 | 16.3 | 99.8 |
| 4KiB 70/30% Random Read/Write | 11.81 | 82700 | 12.3 | 13.4 | 18.0 | 100.4 |

# 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 3.7 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 33.29% vs. the Linux Kernel NVMe-oF target.

2.  The SPDK NVMe-oF Initiator reduces the NVMe-oF software overhead by up to 4.67 usec vs. the Linux Kernel NVMe-oF Initiator for the RDMA transport, which is approximately 49.73% lower than the Linux Kernel NVMe-oF Initiator overhead.

3.  **The SPDK NVMe-oF TCP target and initiator reduced the average latency by up to 7.38 usec vs.** the Linux Kernel NVMe-oF target and initiator, which eliminates up to 61% NVMe-oF software 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-cli tool.<br>For detailed configuration file contents please see Appendix C. |
| **Kernel NVMe-oF Initiator #1** | **Device config**<br>Performed using nvme-cli tool.<br><br>modprobe nvme-fabrics |

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

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.
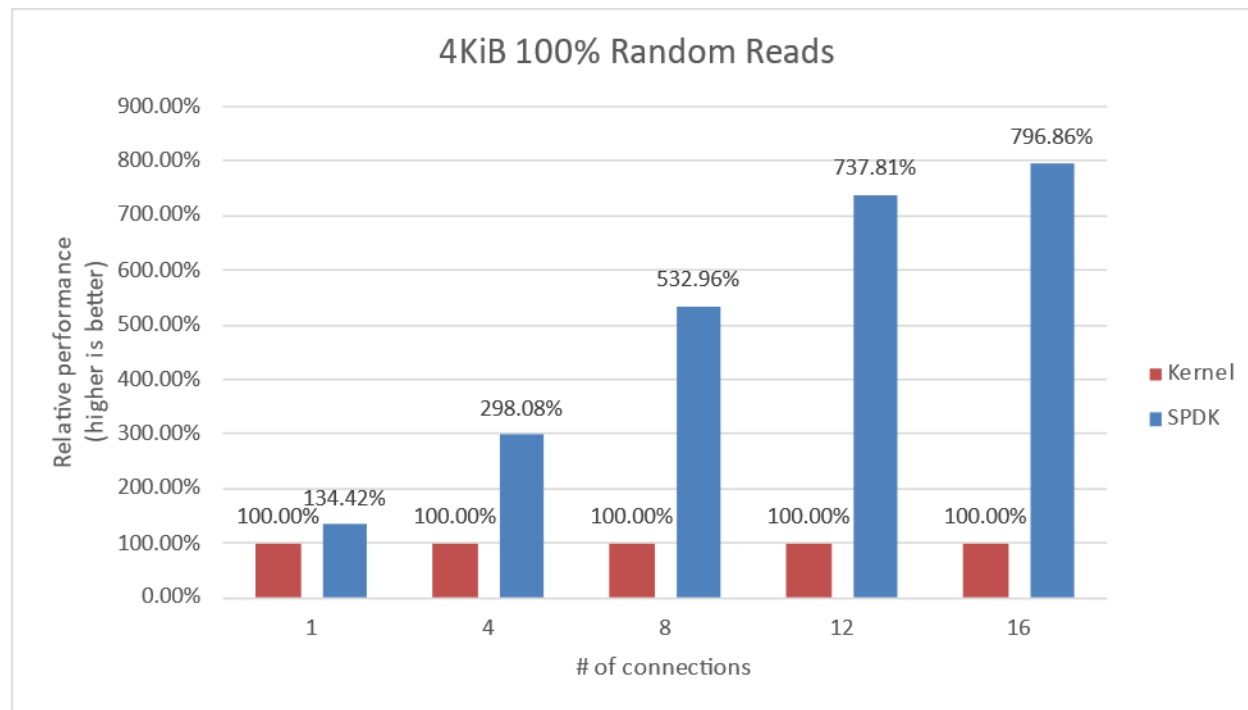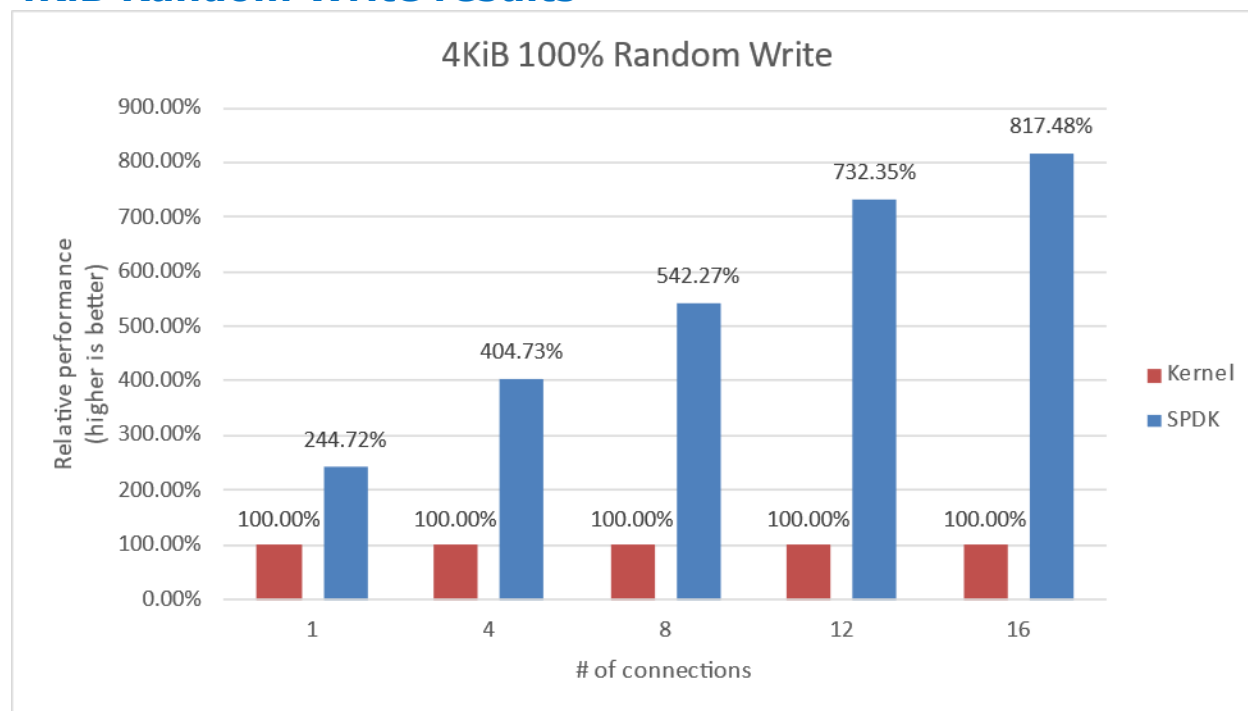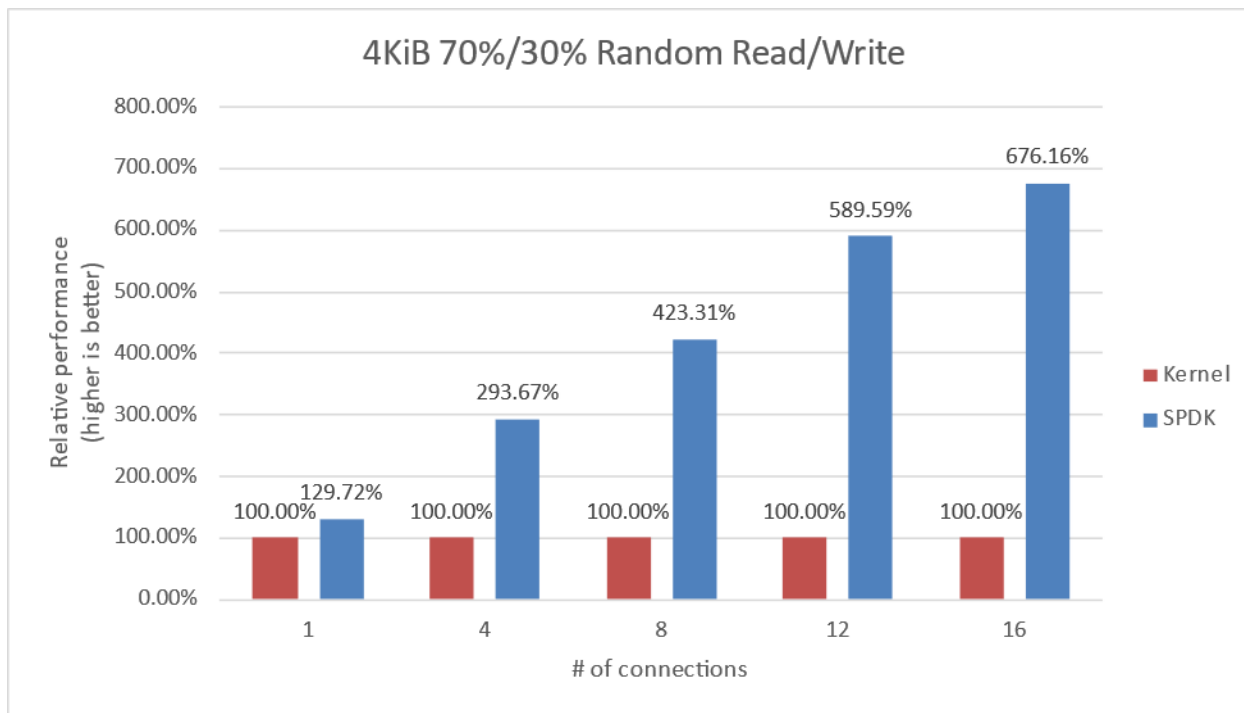
# 4KiB Random Read Results



*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 | 24119.26 | 6174.5 | 435.5 | 10.2 |
| 4 | 42061.10 | 10767.6 | 249.4 | 24.0 |
| 8 | 31495.24 | 8062.8 | 343.5 | 33.0 |
| 12 | 24411.88 | 6249.4 | 438.0 | 36.8 |
| 16 | 23071.07 | 5906.2 | 456.8 | 37.7 |

*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 | 25576.51 | 6547.6 | 410.3 | 8.1 |
| 4 | 42186.41 | 10799.7 | 248.7 | 8.1 |
| 8 | 41037.39 | 10505.6 | 255.6 | 8.1 |
| 12 | 39550.80 | 10125.0 | 265.1 | 8.1 |
| 16 | 39393.34 | 10084.7 | 266.2 | 8.1 |

# 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 | 23750.41 | 6080.1 | 299.9 | 9.5 |
| 4 | 27827.28 | 7123.8 | 262.5 | 17.2 |
| 8 | 26499.62 | 6783.9 | 270.1 | 23.6 |
| 12 | 25475.99 | 6521.8 | 287.0 | 30.7 |
| 16 | 24096.45 | 6168.7 | 292.0 | 32.7 |

*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 | 24800.84 | 6349.0 | 286.5 | 4.1 |
| 4 | 26727.07 | 6842.1 | 268.1 | 4.1 |
| 8 | 24792.08 | 6346.8 | 285.5 | 4.1 |
| 12 | 24760.09 | 6338.6 | 293.5 | 4.1 |
| 16 | 24537.18 | 6281.5 | 286.3 | 4.1 |

# 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 | 24003.22 | 6144.8 | 437.5 | 10.0 |
| 4 | 43039.67 | 11018.1 | 243.7 | 24.2 |
| 8 | 37845.56 | 9688.5 | 277.3 | 34.3 |
| 12 | 29249.75 | 7487.9 | 360.0 | 37.3 |
| 16 | 26603.94 | 6810.6 | 397.8 | 38.5 |

*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 | 25228.42 | 6458.5 | 416.0 | 8.1 |
| 4 | 42156.72 | 10792.1 | 248.9 | 8.1 |
| 8 | 37682.79 | 9646.8 | 278.3 | 8.1 |
| 12 | 37316.57 | 9553.0 | 281.0 | 8.1 |
| 16 | 37703.59 | 9652.1 | 278.1 | 8.1 |

34

# Conclusions

1. When the SPDK NVMe-oF Target was configured with 8 CPU cores the performance peaked at 4 connections for all 4KiB 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 connection per subsystem for all workloads. Further Increasing the number of connections also increased the latency and CPU utilization, while lowering IOPS.

3. The SPDK NVMe-oF target shows up to 8.17x 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 3.7 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 by up to 49.73% lower than Kernel initiator.

The SPDK NVMe-oF Target performed up to 8.17 times better w.r.t IOPS/core than Linux Kernel NVMe-oF target while running 4KiB 100% Random Write 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*

# *List of figures*

## *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",
            "adrfam": "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",
            "adrfam": "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",
            "adrfam": "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",
            "adrfam": "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",
        "adrfam": "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",
        "adrfam": "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",
        "adrfam": "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",
        "adrfam": "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",
        "adrfam": "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",
                    "adrfam": "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",
                    "adrfam": "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",
                    "adrfam": "IPv4"
                }
            }
        ]
    }
]
}
```

SPDK NVMe-oF RDMA Performance Report (Mellanox ConnectX-5)
Release 24.05

## *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",
            "adrfam": "IPv4"
          }
        }
      ]
    }
  ]
}
```

## *Appendix C - Kernel NVMe-oF RDMA Target configuration*

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

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

```
            "nqn.2018-09.io.spdk:cnode2"
        ]
    },
    {
        "addr": {
            "adrfam": "ipv4",
            "traddr": "20.0.0.1",
            "trsvcid": "4422",
            "trtype": "rdma"
        },
        "portid": 3,
        "referrals": [],
        "subsystems": [
            "nqn.2018-09.io.spdk:cnode3"
        ]
    },
    {
        "addr": {
            "adrfam": "ipv4",
            "traddr": "20.0.0.1",
            "trsvcid": "4423",
            "trtype": "rdma"
        },
        "portid": 4,
        "referrals": [],
        "subsystems": [
            "nqn.2018-09.io.spdk:cnode4"
        ]
    },
    {
        "addr": {
            "adrfam": "ipv4",
            "traddr": "20.0.1.1",
            "trsvcid": "4424",
            "trtype": "rdma"
        },
        "portid": 5,
        "referrals": [],
        "subsystems": [
            "nqn.2018-09.io.spdk:cnode5"
        ]
    },
    {
        "addr": {
            "adrfam": "ipv4",
            "traddr": "20.0.1.1",
            "trsvcid": "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",
        "trsvcid": "4431",
        "trtype": "rdma"
      },
      "portid": 12,
      "referrals": [],
      "subsystems": [
        "nqn.2018-09.io.spdk:cnode12"
      ]
    },
    {
      "addr": {
        "adrfam": "ipv4",
        "traddr": "10.0.1.1",
        "trsvcid": "4432",
        "trtype": "rdma"
      },
      "portid": 13,
      "referrals": [],
      "subsystems": [
        "nqn.2018-09.io.spdk:cnode13"
      ]
    },
    {
      "addr": {
        "adrfam": "ipv4",
        "traddr": "10.0.1.1",
        "trsvcid": "4433",
        "trtype": "rdma"
      },
      "portid": 14,
      "referrals": [],
      "subsystems": [
        "nqn.2018-09.io.spdk:cnode14"
      ]
    },
    {
      "addr": {
        "adrfam": "ipv4",
        "traddr": "10.0.1.1",
        "trsvcid": "4434",
        "trtype": "rdma"
      },
      "portid": 15,
      "referrals": [],
      "subsystems": [
        "nqn.2018-09.io.spdk:cnode15"
      ]
```

```
            "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.

§