

# SPDK NVMe-oF RDMA (Target & Initiator) Performance Report Release 21.01

---

**Testing Date:** February 2021

**Performed by:**

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

Maciej Wawryk ([maciejx.wawryk@intel.com](mailto:maciejx.wawryk@intel.com))

**Acknowledgments:**

James Harris ([james.r.harris@intel.com](mailto:james.r.harris@intel.com))

John Kariuki ([john.k.kariuki@intel.com](mailto:john.k.kariuki@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
Kernel & BIOS spectre-meltdown information .....	6
Introduction to SPDK NVMe-oF (Target & Initiator) .....	7
Test Case 1: SPDK NVMe-oF RDMA Target I/O core scaling .....	9
4KB Random Read Results .....	12
4KB Random Write Results .....	13
4KB Random Read-Write Results.....	14
LTO Performance Impact.....	15
Large Sequential I/O Performance .....	18
Test Case 2: SPDK NVMe-oF RDMA Initiator I/O core scaling .....	20
4KB Random Read Results .....	23
4KB Random Write Results .....	24
4KB Random 70/30 Read/Write Results .....	25
Conclusions .....	26
Test Case 3: Linux Kernel vs. SPDK NVMe-oF RDMA Latency .....	27
SPDK vs Kernel NVMe-oF RDMA Target Results .....	30
SPDK vs Kernel NVMe-oF RDMA Initiator Results .....	31
SPDK vs Kernel NVMe-oF RDMA Kernel + Initiator Results.....	32
Conclusions .....	33
Test Case 4: NVMe-oF RDMA Performance with increasing # of connections .....	34
4KB Random Read Results .....	36
4KB Random Write results.....	37
4KB Random Read-Write Results.....	38
Conclusions .....	39
Summary .....	40
List of tables.....	41
List of figures.....	43
Appendix A – Test Case 1 & 2 SPDK NVMe-oF Initiator bdev configuration.....	44
Appendix B – Test Case 3 SPDK NVMe-oF Initiator bdev configuration .....	48
Appendix C - Kernel NVMe-oF RDMA Target configuration .....	48

## ***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	<p>SuperMicro SYS-2029U-TN24R4T</p>  
CPU	<p><a href="#">Intel® Xeon® Gold 6230 Processor (27.5MB L3, 2.10 GHz)</a>            Number of cores 20 per socket, number of threads 40 per socket (Both sockets populated)            Microcode: 0x5003003</p>
Memory	<p>12 x 32GB Hynix HMA84GR7AFR4N-VK, DDR4, 2666MHz            Total of 384GB</p>
Operating System	Fedora 33
BIOS	3.4
Linux kernel version	5.8.15-300.fc33.x86_64
SPDK version	SPDK 21.01
Storage	<p><b>OS:</b> 1x 120GB Intel SSDSC2BB120G4  <b>Storage Target:</b> 16x <a href="#">Intel® SSD DC P4610™ 1.6TB</a> (FW: VDV10170)            (8 on each CPU socket)</p>
NIC	<p>2x 100GbE Mellanox ConnectX-5 NICs. Both ports connected.            1 NIC per CPU socket.</p>

## Initiator 1 Configuration

Table 2: Hardware setup configuration – Initiator system 1

Item	Description
Server Platform	<a href="#">Intel® Server System R2208WFTZSR</a>
CPU	<a href="#">Intel(R) Xeon(R) Gold 6252 CPU @ 2.10GHz (35.75MB Cache)</a> Number of cores 24 per socket, number of threads 48 per socket (Both sockets populated) Microcode: 0x5003003
Memory	6 x 32GB Micron M393A1G40EB1-CRC, DDR4, 2933MHz Total 192GBs
Operating System	Fedora 33
BIOS	02.01.0012
Linux kernel version	5.8.15-300.fc33.x86_64
SPDK version	SPDK 21.01
Storage	<b>OS:</b> 1x 240GB INTEL SSDSC2BB240G6
NIC	1x 100GbE Mellanox ConnectX-5 Ex NIC. Both ports connected to Target server. (connected to CPU socket 0)

## Initiator 2 Configuration

Table 3: Hardware setup configuration – Initiator system 2

Item	Description
Server Platform	<a href="#">Intel® Server System R2208WFTZSR</a>
CPU	<a href="#">Intel(R) Xeon(R) Gold 6252 CPU @ 2.10GHz (35.75MB Cache)</a> Number of cores 24 per socket, number of threads 48 per socket (Both sockets populated) Microcode: 0x5003003
Memory	6 x 32GB Micron M393A1G40EB1-CRC, DDR4, 2933MHz Total 192GBs
Operating System	Fedora 33
BIOS	02.01.0012
Linux kernel version	5.8.15-300.fc33.x86_64
SPDK version	SPDK 21.01
Storage	<b>OS:</b> 1x 240GB INTEL SSDSC2BB240G6
NIC	1x 100GbE Mellanox ConnectX-5 Ex NIC. Both ports connected to Target server. (connected to CPU socket 0)

## BIOS settings

Table 4: Test systems BIOS settings

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

## Kernel & BIOS spectre-meltdown information

All three server systems use kernel version available from DNF repository with default patches for spectre-meltdown issue enabled.

BIOS on all systems was updated to post spectre-meltdown versions as well.

# 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 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 sixteen Intel P4610 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 (systat).

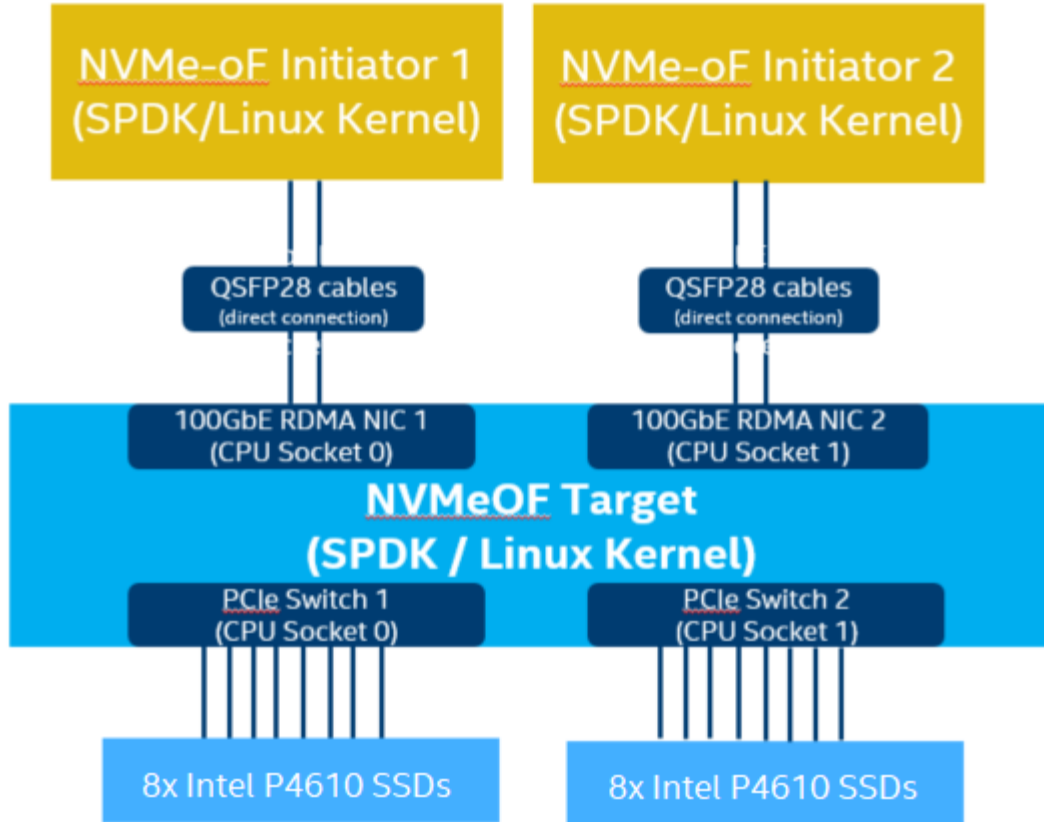


Figure 1: High-Level NVMe-oF RDMA performance testing setup



# 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 16 NVMe-oF subsystems. Each NVMe-oF subsystem ran on top of an individual bdev backed by a single Intel® SSD DC P4610 device. Each of the 2 initiators were connected to 8 individual NVMe-oF subsystems which were exposed via SPDK NVMe-oF Target over 1x 100GbE NIC. SPDK bdev FIO plugin was used to target 8 individual NVMe-oF bdevs on each of the initiators. SPDK Target Reactor Mask was configured to use up to 6 CPU cores tests while running following workloads on each initiator:

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

The table below contains more information about the test configuration. The SPDK NVMe-oF Target was configured using JSON-RPC; the table contains a sequence of commands used by `spd/scripts/rpc.py` script rather than a configuration file. 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 4KB Rand Read and 4KB Rand 70/30 Read/Write workloads to ensure that the SSDs reached their steady state where we get repeatable results. However, for the 4KB 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	<p>All the commands below were executed with <code>spd/scripts/rpc.py</code> script.</p> <p><b>Construct NVMe bdevs:</b></p> <pre>construct_nvme_bdev -t PCIe -b Nvme0 -a 0000:60:00.0 construct_nvme_bdev -t PCIe -b Nvme1 -a 0000:61:00.0 construct_nvme_bdev -t PCIe -b Nvme2 -a 0000:62:00.0 construct_nvme_bdev -t PCIe -b Nvme3 -a 0000:63:00.0 construct_nvme_bdev -t PCIe -b Nvme4 -a 0000:64:00.0 construct_nvme_bdev -t PCIe -b Nvme5 -a 0000:65:00.0 construct_nvme_bdev -t PCIe -b Nvme6 -a 0000:66:00.0 construct_nvme_bdev -t PCIe -b Nvme7 -a 0000:67:00.0 construct_nvme_bdev -t PCIe -b Nvme8 -a 0000:b5:00.0 construct_nvme_bdev -t PCIe -b Nvme9 -a 0000:b6:00.0 construct_nvme_bdev -t PCIe -b Nvme10 -a 0000:b7:00.0</pre>

	<pre> construct_nvme_bdev -t PCIe -b Nvme11 -a 0000:b8:00.0 construct_nvme_bdev -t PCIe -b Nvme12 -a 0000:b9:00.0 construct_nvme_bdev -t PCIe -b Nvme13 -a 0000:ba:00.0 construct_nvme_bdev -t PCIe -b Nvme14 -a 0000:bb:00.0 construct_nvme_bdev -t PCIe -b Nvme15 -a 0000:bc:00.0  <b>Create a RDMA transport layer:</b> nvmf_create_transport -t RDMA (creates 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: 4096 buf_cache_size: 32)  <b>Create NVMe-oF subsystems and add NVMe bdevs as namespaces:</b> 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  <b>Add listeners to NVMe-oF Subsystems:</b> 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 done </pre>
<p><b>SPDK NVMe-oF Initiator - FIO plugin configuration</b></p>	<p><b>BDEV.conf</b> See <a href="#">Appendix A</a></p> <p><b>FIO.conf</b> [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</p> <p>[filename0] filename=Nvme0n1</p>

[filename1] filename=Nvme1n1 [filename2] filename=Nvme2n1 [filename3] filename=Nvme3n1 [filename4] filename=Nvme4n1 [filename5] filename=Nvme5n1 [filename6] filename=Nvme6n1 [filename7] filename=Nvme7n1
---

## 4KB Random Read Results

Table 6: SPDK NVMe-oF RDMA Target Core Scaling results, Random Read IOPS, QD=64

# of Cores	Bandwidth (MBps)	Throughput (IOPS k)	Avg. Latency (usec)
1 core	4432.65	1134.8	901.8
2 cores	10301.77	2637.2	387.6
3 cores	15493.06	3966.2	259.8
4 cores	20757.15	5313.8	192.1
5 cores	21678.60	5549.7	184.6
6 cores	21965.38	5623.1	182.3

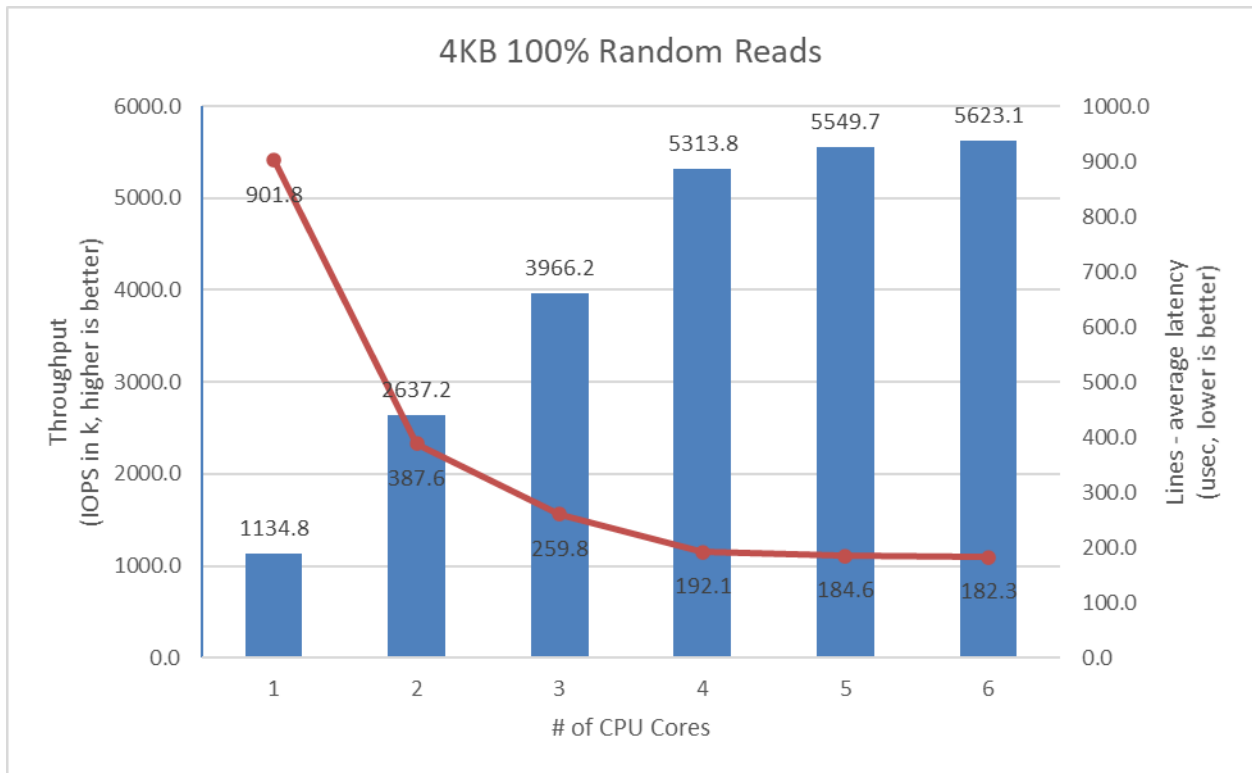


Figure 2: SPDK NVMe-oF RDMA Target I/O core scaling: IOPS vs. Latency while running 4KB 100% Random Read workload at QD = 64

## 4KB Random Write Results

Table 7: SPDK NVMe-oF RDMA Target Core Scaling results, Random Write IOPS, QD=32

# of Cores	Bandwidth (MBps)	Throughput (IOPS k)	Avg. Latency (usec)
1 core	6173.84	1580.5	321.6
2 cores	13447.33	3442.5	146.7
3 cores	18595.78	4760.5	105.8
4 cores	22021.62	5637.5	90.3
5 cores	22892.51	5860.5	86.3
6 cores	22977.95	5882.4	86.2

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.

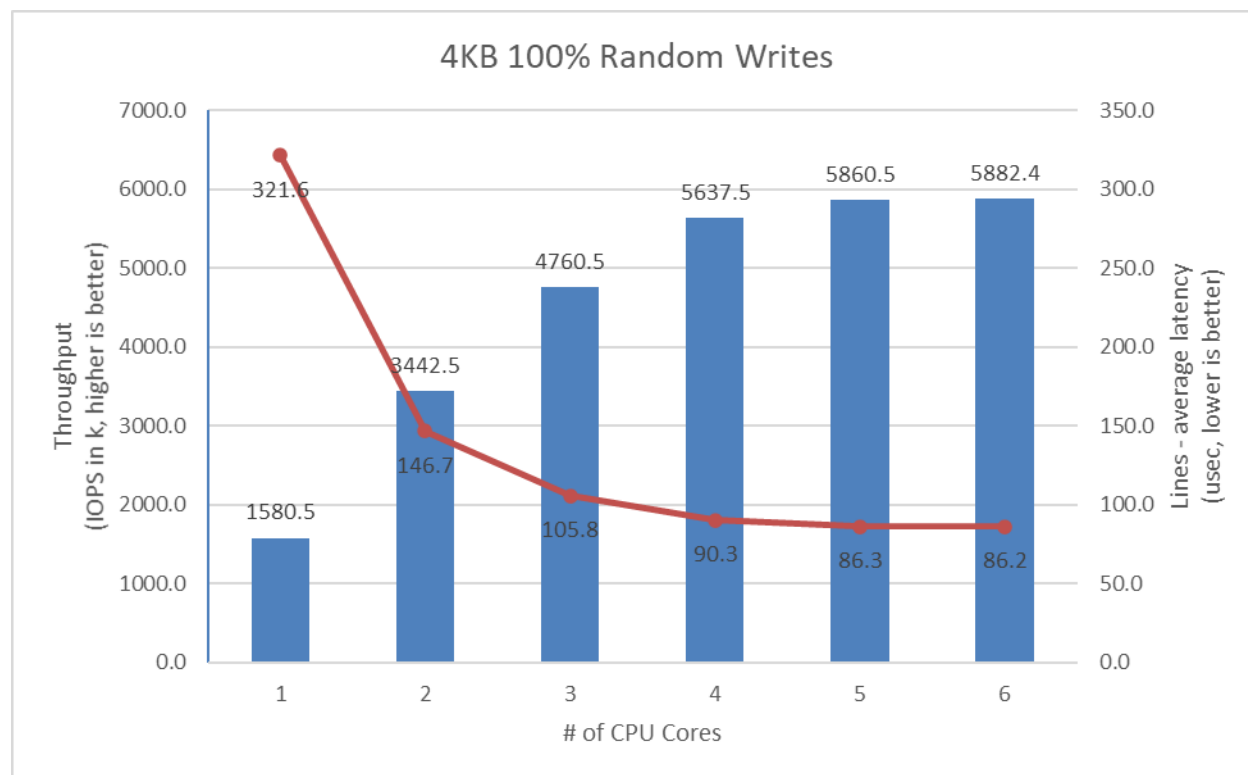


Figure 3: SPDK NVMe-oF RDMA Target I/O core scaling: IOPS vs. Latency while running 4KB 100% Random Write Workload at QD=32

## 4KB 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 (MBps)	Throughput (IOPS k)	Avg. Latency (usec)
1 core	4496.35	1151.1	1778.4
2 cores	10131.46	2593.6	787.9
3 cores	15473.67	3961.3	516.7
4 cores	20995.82	5374.9	383.0
5 cores	24043.66	6155.2	331.1
6 cores	22776.34	5830.7	353.0

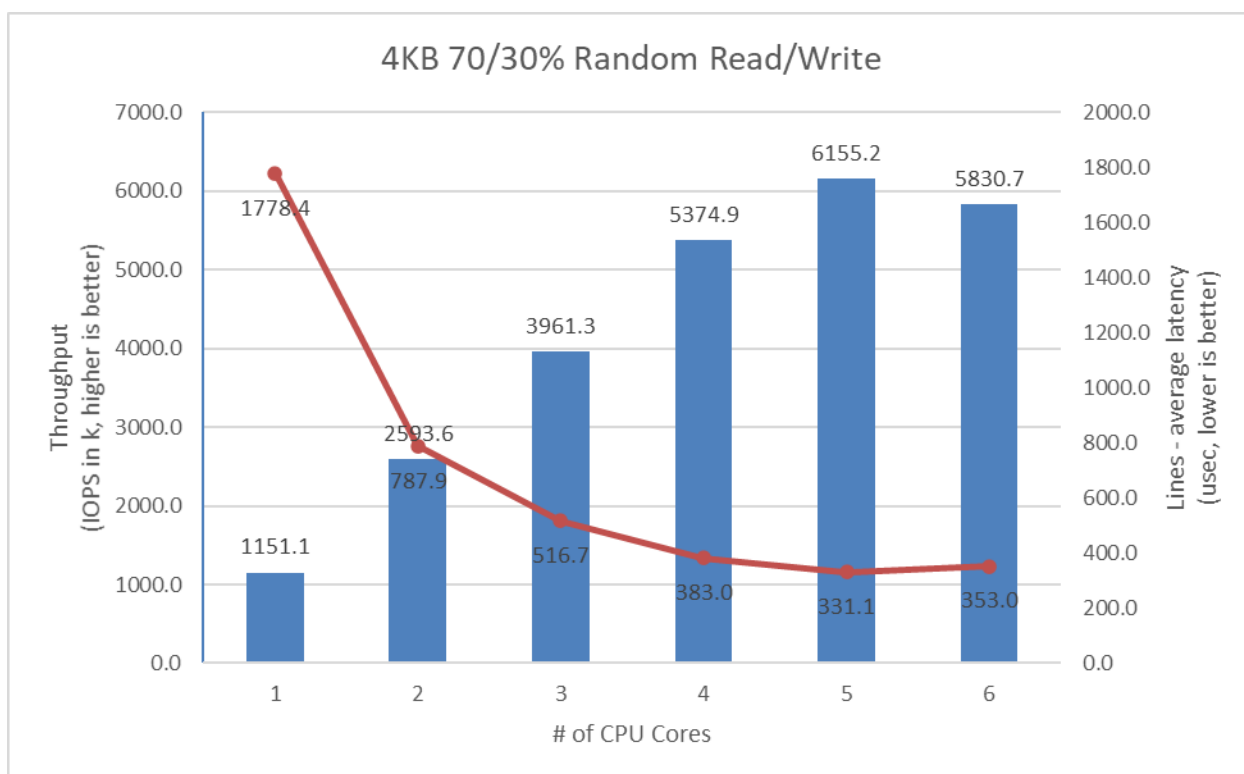


Figure 4: SPDK NVMe-oF RDMA Target I/O core scaling: IOPS vs. Latency while running 4KB Random 70/30 Read/Write workload at QD=128

## LTO Performance Impact

Below results present performance impact of using SPDK “—with-lto” build option, which enables link-time optimization.

Table 9: LTO impact for SPDK NVMe-oF RDMA Target Core Scaling performance – Bandwidth

Workload	# of Cores	Bandwidth (MBps)		
		LTO Disabled	LTO Enabled	LTO Impact
4KB Randread QD=64	1 core	4432.65	4612.36	4.05%
	2 cores	10301.77	10545.30	2.36%
	3 cores	15493.06	16495.60	6.47%
	4 cores	20757.15	20701.94	-0.27%
	5 cores	21678.60	20845.03	-3.85%
	6 cores	21965.38	23526.88	7.11%
4KB Randwrite QD=32	1 core	6173.84	6376.91	3.29%
	2 cores	13447.33	13895.37	3.33%
	3 cores	18595.78	19790.58	6.43%
	4 cores	22021.62	22822.02	3.63%
	5 cores	22892.51	23019.23	0.55%
	6 cores	22977.95	21581.69	-6.08%
4KB RandRW QD=128	1 core	4496.35	4762.72	5.92%
	2 cores	10131.46	10526.86	3.90%
	3 cores	15473.67	16172.81	4.52%
	4 cores	20995.82	21422.76	2.03%
	5 cores	24043.66	23421.93	-2.59%
	6 cores	22776.34	23041.85	1.17%

Table 10: LTO impact for SPDK NVMe-oF RDMA Target Core Scaling performance - Throughput

Workload	# of Cores	Throughput (IOPS k) – Higher is better		
		LTO Disabled	LTO Enabled	LTO Impact
4KB Randread QD=64	1 core	1134.75	1180.76	4.05%
	2 cores	2637.25	2699.59	2.36%
	3 cores	3966.22	4222.87	6.47%
	4 cores	5313.83	5299.69	-0.27%
	5 cores	5549.72	5336.32	-3.85%
	6 cores	5623.13	6022.88	7.11%
4KB Randwrite QD=32	1 core	1580.50	1632.49	3.29%
	2 cores	3442.51	3557.21	3.33%
	3 cores	4760.52	5066.39	6.43%
	4 cores	5637.53	5842.44	3.63%
	5 cores	5860.48	5892.92	0.55%
	6 cores	5882.35	5524.91	-6.08%
4KB RandRW QD=128	1 core	1151.06	1219.25	5.92%
	2 cores	2593.65	2694.87	3.90%
	3 cores	3961.25	4140.23	4.52%
	4 cores	5374.92	5484.22	2.03%
	5 cores	6155.17	5996.00	-2.59%
	6 cores	5830.73	5898.71	1.17%



Table 11: LTO impact for SPDK NVMe-oF RDMA Target Core Scaling performance – Average Latency

Workload	# of Cores	Avg. Latency (usec) – Lower is better		
		LTO Disabled	LTO Enabled	LTO Impact
4KB Randread QD=64	1 core	901.83	866.66	-3.90%
	2 cores	387.59	394.99	1.91%
	3 cores	259.75	243.26	-6.35%
	4 cores	192.12	192.94	0.43%
	5 cores	184.59	191.52	3.76%
	6 cores	182.26	170.22	-6.61%
4KB Randwrite QD=32	1 core	321.58	311.27	-3.20%
	2 cores	146.74	141.60	-3.50%
	3 cores	105.84	98.93	-6.54%
	4 cores	90.29	86.29	-4.43%
	5 cores	86.28	86.01	-0.31%
	6 cores	86.21	93.27	8.19%
4KB RandRW QD=128	1 core	1778.36	1678.95	-5.59%
	2 cores	787.94	758.66	-3.72%
	3 cores	516.73	504.18	-2.43%
	4 cores	382.99	372.74	-2.68%
	5 cores	331.09	342.94	3.58%
	6 cores	352.96	349.66	-0.94%

## Large Sequential I/O Performance

128KB 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 4KB test case in the previous part of this document. We used iodepth=8 (iodepth=1 for write workload) because higher queue depth resulted in negligible bandwidth gain and a significant increase in the latency.

Table 12: SPDK NVMe-oF RDMA Target Core Scaling results, 128KB Sequential Read IOPS, QD=4

# of Cores	Bandwidth (MBps)	Throughput (IOPS k)	Avg. Latency (usec)
<b>1 core</b>	24966.53	199.7	640.6
<b>2 cores</b>	25006.58	200.1	639.6
<b>3 cores</b>	24986.85	199.9	640.1
<b>4 cores</b>	25083.20	200.7	637.6

Table 13: SPDK NVMe-oF RDMA Target Core Scaling results, 128KB Sequential Write IOPS, QD=1

# of Cores	Bandwidth (MBps)	Throughput (IOPS k)	Avg. Latency (usec)
<b>1 core</b>	23071.24	184.6	86.4
<b>2 cores</b>	23267.77	186.1	85.7
<b>3 cores</b>	23297.32	186.4	85.6
<b>4 cores</b>	23348.46	186.8	85.4

Table 14: SPDK NVMe-oF RDMA Target Core Scaling results, 128KB Sequential 70% Read 30% Write IOPS, QD=4

# of Cores	Bandwidth (MBps)	Throughput (IOPS k)	Avg. Latency (usec)
<b>1 core</b>	24412.47	195.3	327.3
<b>2 cores</b>	24559.09	196.5	325.4
<b>3 cores</b>	24604.06	196.8	324.8
<b>4 cores</b>	24683.04	197.5	323.7

## Conclusions

1. For all tested workloads throughput scales up and latency decreases almost linearly with the addition of I/O cores up to 4 cores. Adding more CPU cores results in small performance gain as network is almost saturated.
2. For large sequential I/Os, a single CPU core Target saturated the network bandwidth. The SPDK NVMe-oF Target running on 1 core does close to 24-25 GBps of bandwidth in all tested workloads. Therefore, adding more CPU cores did not result in increased performance for these workloads because the network was the bottleneck.
3. Enabling Link-Time Optimization build option for SPDK resulted in slightly increased performance in scaling phase (up to 4 CPU cores). The efficiency increased by up to 6.5% and by about 4% on average in scaling phase. Results for 5 and 6 CPU cores are inconclusive, because in these cases network saturation was already reached.

## 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 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 8 NVMe-oF bdevs on each of the 2 initiators. The following workloads were executed on both of the initiators using fio in client-server mode.

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

Depending on the number of initiators and initiator cores, we varied the number of NVMe-oF subsystems exported by the target as shown below, instead of exposing all 16 NVMe-oF subsystems.

**1 core:** 1 initiator running on a single core connecting to 4 subsystems.

**2 cores:** 2 initiators, each running on a single core. Each initiator connected to 4 subsystems.

**3 cores:** 2 initiators, the first one running on 1 core and other one running on 2 cores. Initiator 1 connected to 4 subsystems and initiator 2 connected to 8 subsystems.

**4 cores:** 2 initiators, both running on 2 cores each. Both initiators connected to 8 subsystems.

This was done to avoid a situation where single initiator would connect to all 16 target subsystems, which resulted in unnatural latency spike in the 1 CPU core test case.

*Table 15: 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><b>BDEV.conf</b> See <a href="#">Appendix A</a></p> <p><b>FIO.conf</b> [global] ioengine=/tmp/spdk/examples/bdev/fio_plugin/fio_plugin spdk_json_conf=/tmp/spdk/bdev.conf thread=1 group_reporting=1 direct=1</p> <p>norandommap=1 rw=randrw</p>

	<pre> rwmixread={100, 70, 0} bs=4k iodepth={32, 64, 128, 256} time_based=1 ramp_time=60 runtime=300  {filename_section}  <b>FIO.conf filename section for 1 &amp; 2 CPU test run</b> [filename0] filename=Nvme0n1 filename=Nvme1n1 filename=Nvme2n1 filename=Nvme3n1  <b>FIO.conf filename section for 3 &amp; 4 CPU test run</b> [filename0] filename=Nvme0n1 filename=Nvme1n1 filename=Nvme2n1 filename=Nvme3n1 [filename1] filename=Nvme4n1 filename=Nvme5n1 filename=Nvme6n1 filename=Nvme7n1 </pre>
<p><b>SPDK NVMe-oF Initiator 2 - FIO plugin configuration</b></p>	<p><b>BDEV.conf</b> See <a href="#">Appendix A</a></p> <p><b>FIO.conf</b> Similar as Initiator 1. The only differences are were in the filename section which are shown below.</p> <p><b>FIO.conf filename section for 1 CPU test run</b> N/A (only Initiator 1 is used)</p> <p><b>FIO.conf filename section for 2 CPU test run</b> [filename0] filename=Nvme0n1 filename=Nvme1n1 filename=Nvme2n1 filename=Nvme3n1</p> <p><b>FIO.conf filename section for 3 CPU test run</b> [filename0] filename=Nvme0n1 filename=Nvme1n1 filename=Nvme2n1 filename=Nvme3n1</p> <p><b>FIO.conf filename section for 4 CPU test run</b> [filename0] filename=Nvme0n1 filename=Nvme1n1</p>

	filename=Nvme2n1 filename=Nvme3n1 [filename1] filename=Nvme4n1 filename=Nvme5n1 filename=Nvme6n1 filename=Nvme7n1
--	---

## 4KB Random Read Results

Table 16: SPDK NVMe-oF RDMA Initiator Core Scaling results, 4KB Random Read IOPS, QD=64, SPDK Target 6 CPU Cores

# of Initiator CPU Cores	Bandwidth (MBps)	Throughput (IOPS k)	Avg. Latency (usec)
1 core	5926.98	1517.30	547.61
2 cores	16582.37	4245.08	213.97
3 cores	19407.25	4968.25	195.91
4 cores	22200.55	5683.34	179.47
5 cores	22334.73	5717.69	178.58
6 cores	20899.58	5350.29	189.79
7 cores	22460.29	5749.83	177.54
8 cores	22373.40	5727.59	178.87

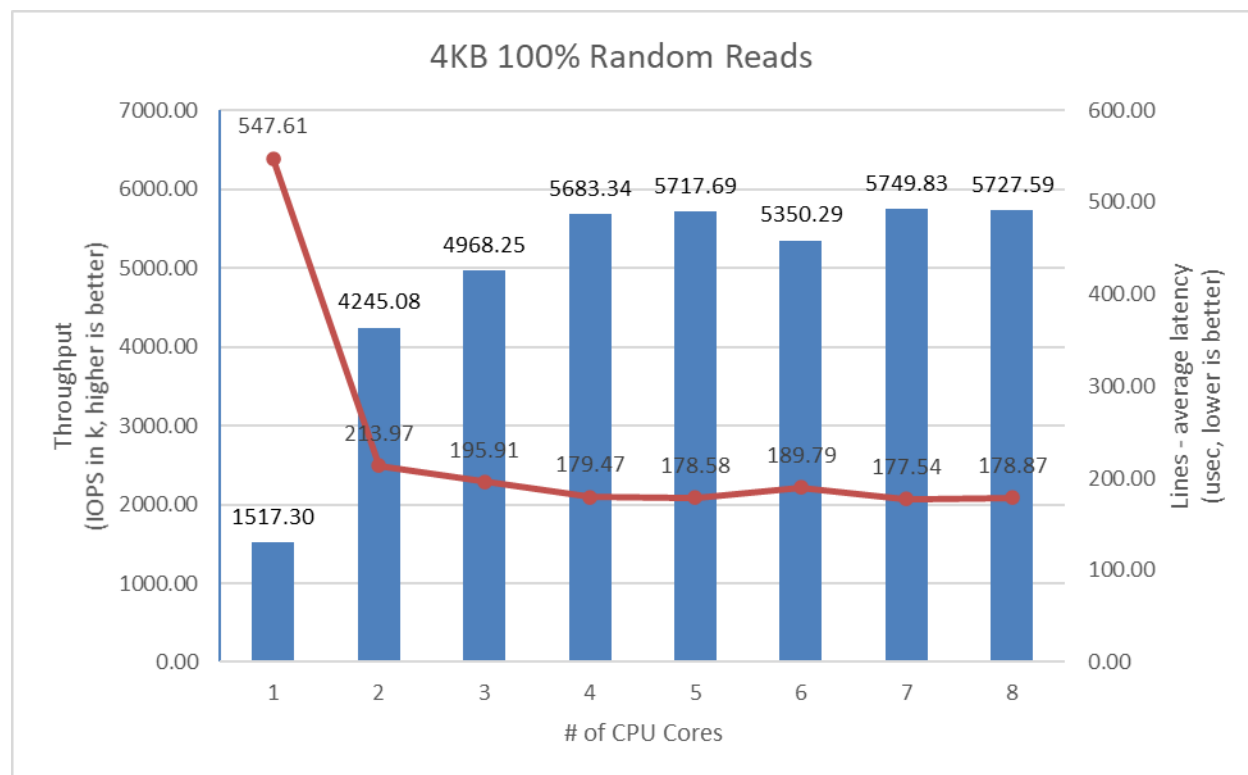


Figure 5: SPDK NVMe-oF RDMA Initiator I/O core scaling: IOPS vs. Latency while running 4KB 100% Random Read QD=64 workload

## 4KB 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 17: SPDK NVMe-oF RDMA Initiator Core Scaling results, 4KB Random Write IOPS, QD=32, SPDK Target 6 CPU Cores

# of Initiator CPU Cores	Bandwidth (MBps)	Throughput (IOPS k)	Avg. Latency (usec)
1 core	6339.92	1623.02	251.80
2 cores	14152.85	3623.13	113.95
3 cores	18507.67	4737.96	99.78
4 cores	22825.76	5843.39	84.08
5 cores	22697.79	5810.63	85.79
6 cores	22931.45	5870.45	85.70
7 cores	22708.62	5813.40	86.82
8 cores	21420.00	5483.52	94.53

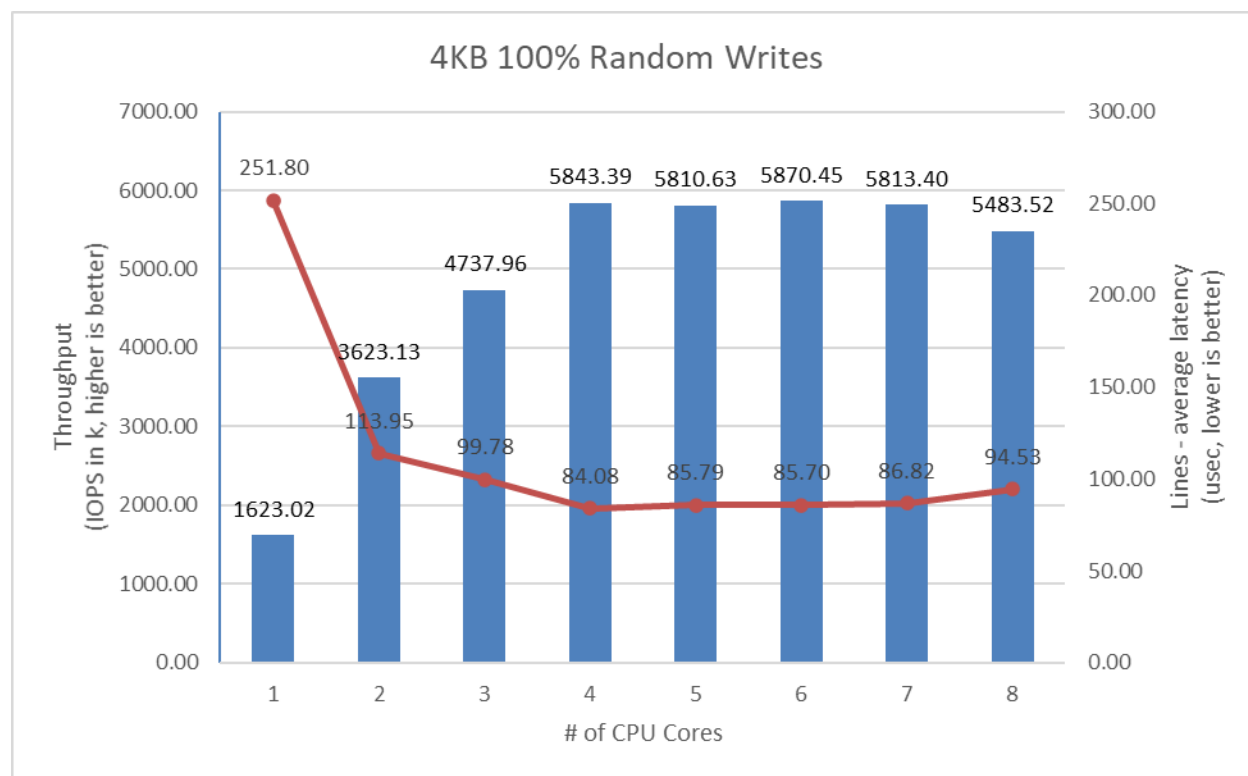


Figure 6: SPDK NVMe-oF RDMA Initiator I/O core scaling: IOPS vs. Latency while running 4KB 100% Random Write Workload at QD=32



## 4KB Random 70/30 Read/Write Results

Table 18: SPDK NVMe-oF RDMA Initiator Core Scaling results, 4KB Random 70%/30% Read/Write IOPS, QD=128, SPDK Target 6 CPU Cores

# of Initiator CPU Cores	Bandwidth (MBps)	Throughput (IOPS k)	Avg. Latency (usec)
1 core	4160.22	1065.01	1544.48
2 cores	11601.16	2969.89	572.17
3 cores	16996.30	4351.04	473.59
4 cores	21728.94	5562.60	372.12
5 cores	20708.98	5301.49	387.68
6 cores	21618.48	5534.32	367.68
7 cores	21825.31	5587.27	368.91
8 cores	22008.02	5634.04	367.28

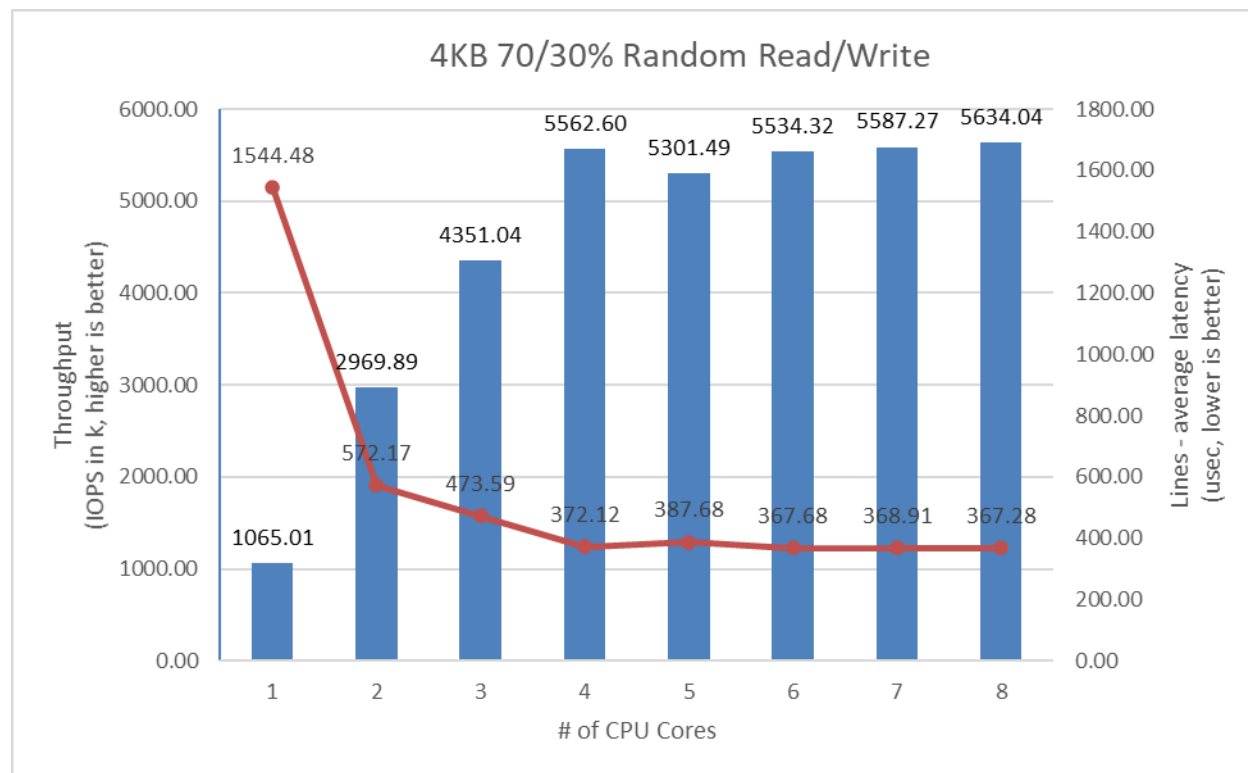


Figure 7: SPDK NVMe-oF RDMA Initiator I/O core scaling: IOPS vs. Latency while running 4KB Random 70% Read 30% Write Workload at QD=128

## Conclusions

1. Random Read workload scaling was not linear when increasing the number of initiator CPU cores.
2. Random Write and Random Read-Write IOPS scaling was linear or close to linear as we increased the number of initiator cores up to 4 CPU cores.
3. For all workloads beyond 4 CPU cores there was no significant change in IOPS because the 200 Gbps network was saturated.
4. Peak performance of 5.7 million IOPS was reached at 5 Initiator CPU cores for random read workload, 5.8 million at 4 CPU cores for random write workload and 5.6 million at 8 CPU cores for mixed random read/write workload.

## 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 containing a *Null Block Device*. The null block device (bdev) was chosen as the backend block device to eliminate the media latency during these tests.

Table 19: Linux Kernel vs. SPDK NVMe-oF RDMA Latency test configuration

Item	Description
<b>Test Case</b>	Linux Kernel vs. SPDK NVMe-oF RDMA Latency
<b>Test configuration</b>	
<b>SPDK NVMe-oF Target configuration</b>	<p>The following commands are executed with <code>spdk/scripts/rpc.py</code> script to configure the SPDK NVMe-oF target.</p> <pre> nvmf_create_transport -t RDMA (creates 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: 4096 buf_cache_size: 32)  construct_null_bdev Nvme0n1 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 Nvme0n1 nvmf_subsystem_add_listener nqn.2018-09.io.spdk:cnode1 -t rdma -f ipv4 -s 4420 -a 20.0.0.1                     </pre>
<b>Kernel NVMe-oF Target configuration</b>	<p>The following target configuration file loaded using <code>nvmet-cli</code> tool.</p> <pre> {   "ports": [     {       "addr": {         "adrfam": "ipv4",         "traddr": "20.0.0.1",         "trsvcid": "4420",         "trtype": "rdma"       },       "portid": 1,       "referrals": [],       "subsystems": [                     </pre>

	<pre> "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>
<b>FIO configuration</b>	
<p><b>SPDK NVMe-oF Initiator FIO plugin configuration</b></p>	<p><b>BDEV.conf</b> See <a href="#">Appendix B</a>.</p> <p><b>FIO.conf</b> [global] ioengine=/tmp/spdk/examples/bdev/fio_plugin/fio_plugin spdk_json_conf=/tmp/spdk/bdev.conf thread=1 group_reporting=1 direct=1</p> <p>norandommap=1 rw=randrw rwmixread={100, 70, 0} bs=4k iodepth=1 time_based=1 ramp_time=60 runtime=300</p> <p>[filename0] filename=NvmeOn1</p>
<p><b>Kernel initiator configuration</b></p>	<p><b>Device config</b> The following configuration was performed using nvme-cli tool. modprobe nvme-fabrics nvme connect -n nqn.2018-09.io.spdk:cnode1 -t rdma -a 20.0.0.1 -s 4420</p> <p><b>FIO.conf</b></p>

```
[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
```

## 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.

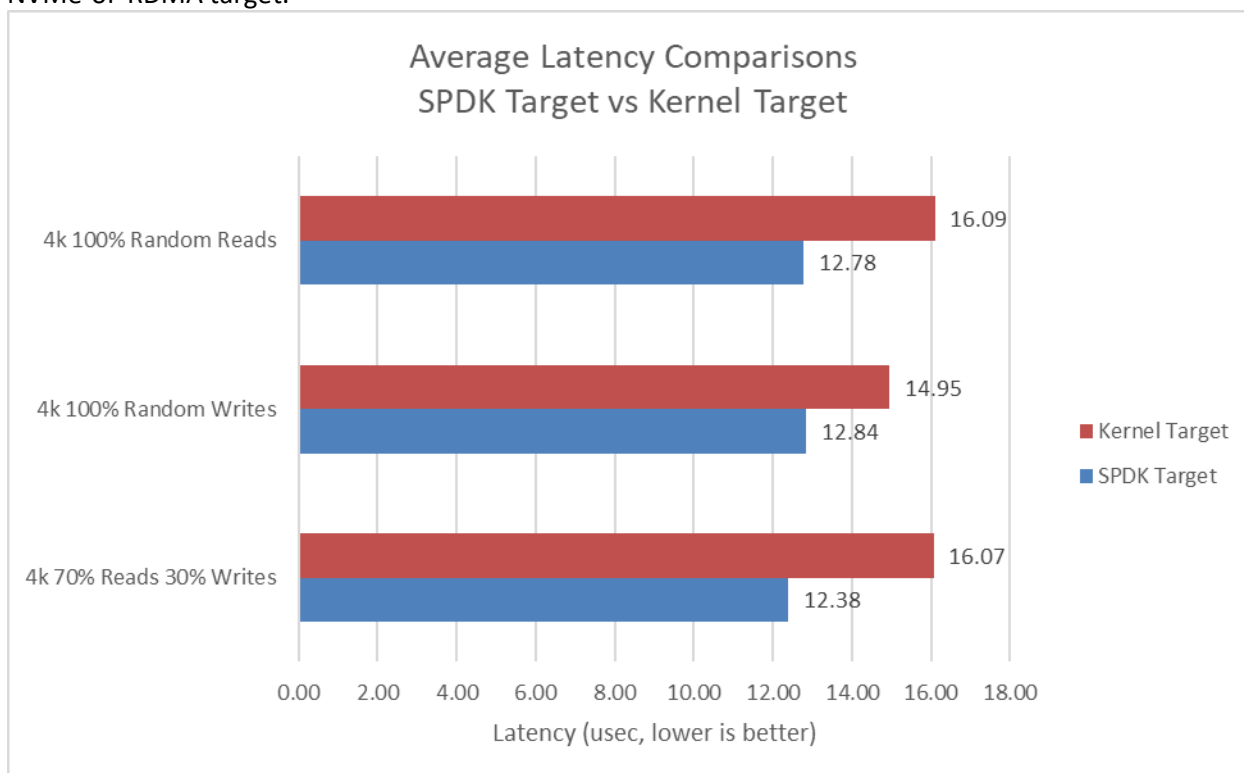


Figure 8: SPDK vs. Kernel NVMe-oF RDMA Target average I/O latency for various workloads run using the Kernel Initiator

Table 20: 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)
<b>4KB 100% Random Read</b>	12.78	76877	14.7	23.9	30.2	134.8
<b>4KB 100% Random Write</b>	12.84	76136	19.6	26.5	34.2	134.1
<b>4KB 70/30% Random Read/Write</b>	12.38	79067	11.3	22.2	26.8	136.5

Table 21: 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)
<b>4KB 100% Random Read</b>	16.09	61200	16.1	28.4	30.6	49.7
<b>4KB 100% Random Write</b>	14.95	65785	15.3	23.9	29.3	45.3
<b>4KB 70/30% Random Read/Write</b>	16.07	61276	16.2	26.3	30.3	44.4

## 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.

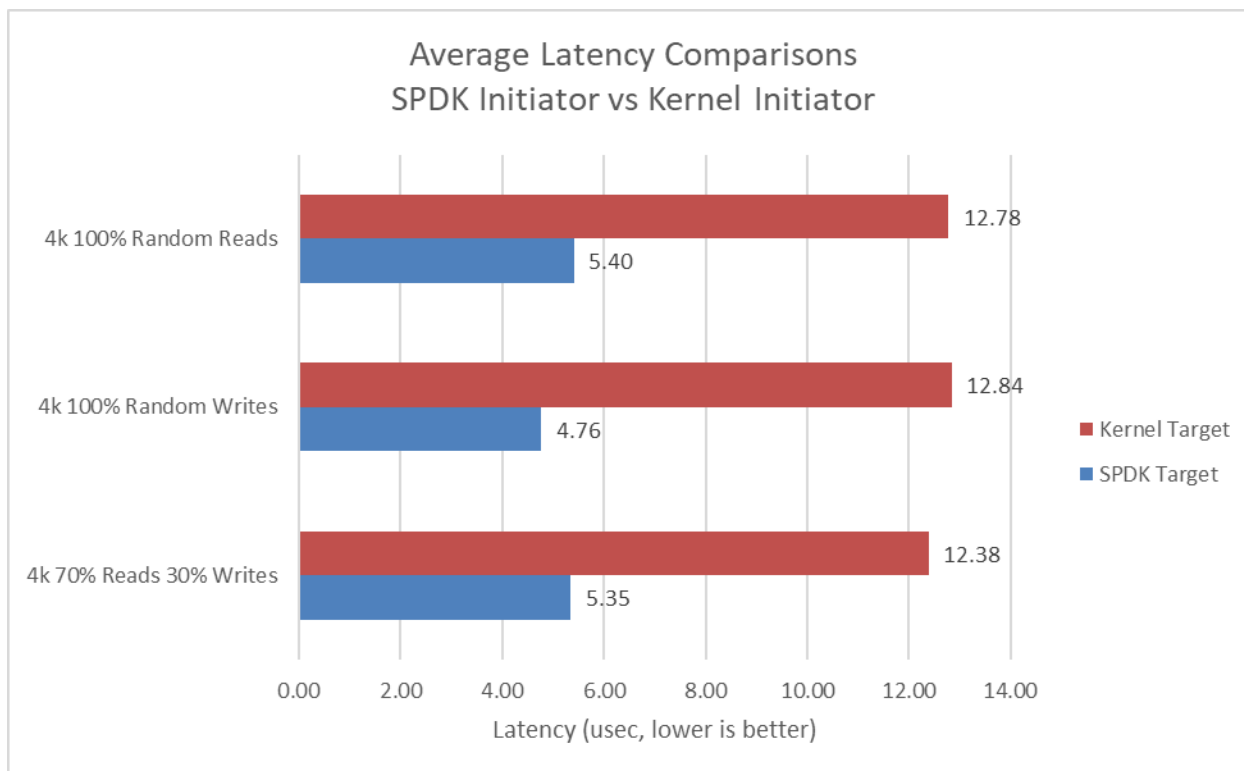


Figure 9: SPDK vs. Kernel NVMe-oF RDMA Initiator average I/O latency for various workloads against SPDK Target

Table 22: 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)
<b>4KB 100% Random Read</b>	5.40	179439	5.7	10.8	15.6	50.3
<b>4KB 100% Random Write</b>	4.76	203456	5.0	10.9	14.6	45.7
<b>4KB 70/30% Random Read/Write</b>	5.35	180618	5.5	16.2	21.2	50.0

Table 23: 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)
<b>4KB 100% Random Read</b>	12.78	76877	14.7	23.9	30.2	134.8
<b>4KB 100% Random Write</b>	12.84	76136	19.6	26.5	34.2	134.1
<b>4KB 70/30% Random Read/Write</b>	12.38	79067	11.3	22.2	26.8	136.5

## 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.

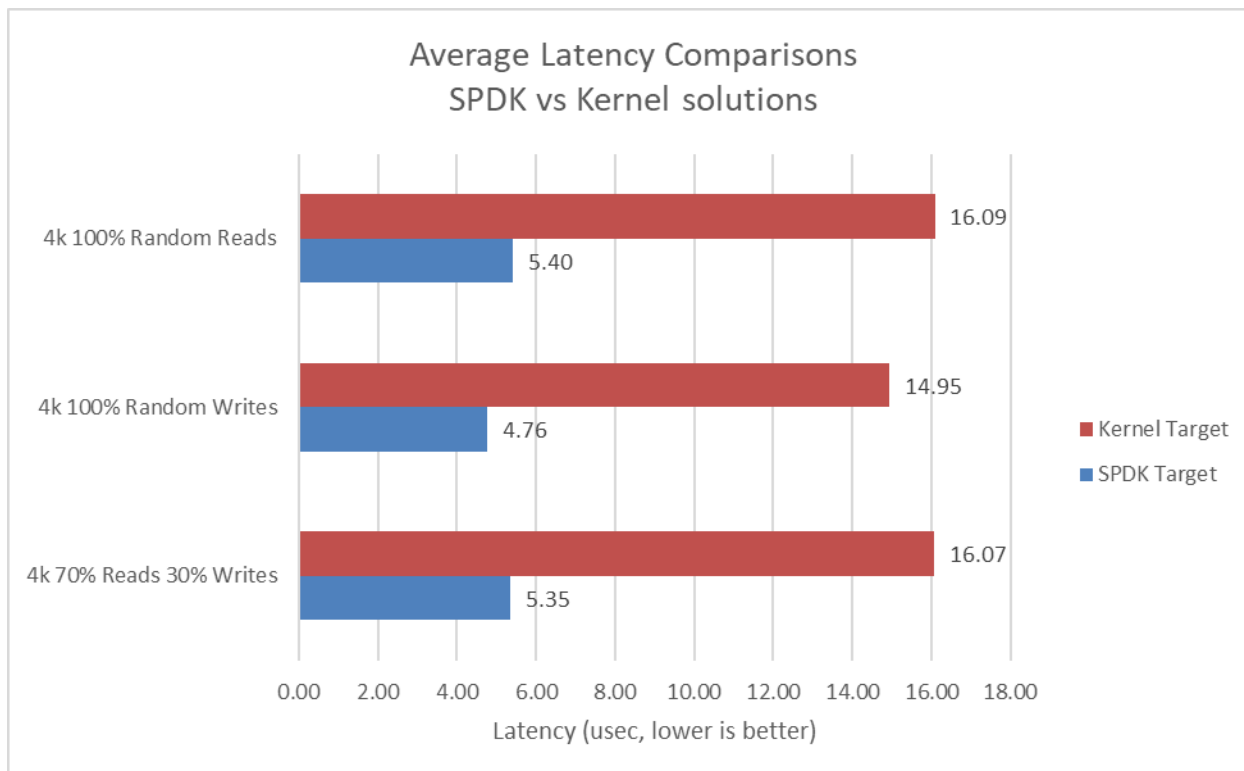


Figure 10: SPDK vs. Kernel NVMe-oF RDMA solutions average I/O Latency for various workloads against SPDK Target

Table 24: 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)
<b>4KB 100% Random Read</b>	5.40	179439	5.7	10.8	15.6	50.3
<b>4KB 100% Random Write</b>	4.76	203456	5.0	10.9	14.6	45.7
<b>4KB 70/30% Random Read/Write</b>	5.35	180618	5.5	16.2	21.2	50.0

Table 25: 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)
<b>4KB 100% Random Read</b>	16.09	61200	16.1	28.4	30.6	49.7
<b>4KB 100% Random Write</b>	14.95	65785	15.3	23.9	29.3	45.3
<b>4KB 70/30% Random Read/Write</b>	16.07	61276	16.2	26.3	30.3	44.4



## 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 3.69 usec vs. the Linux Kernel NVMe-oF target used in Fedora 33 5.8.15 setup. This is entirely software overhead, therefore, using the SPDK NVMe-oF target reduces the NVMe-oF software overhead by approximately 23% vs. the Linux Kernel NVMe-oF target.
2. The SPDK NVMe-oF initiator reduces the NVMe-oF software overhead by up to 2.6x times vs. the Linux Kernel NVMe-oF Initiator for the RDMA transport.
3. The SPDK NVMe-oF RDMA target and initiator reduced the average latency by up to 68% vs. the Linux Kernel NVMe-oF RDMA target and initiator.

## 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 connections (or I/O queue pairs) per NVMe-oF subsystem were 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 (systat package in Linux). The SPDK NVMe-oF RDMA Target was configured to run on 4 CPU cores, export 16 NVMe-oF subsystems (1 per Intel P4610) and 2 initiators were used both running the I/O workloads below to 8 separate subsystems using Kernel NVMe-oF RDMA initiator.

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

Table 26: NVMe-oF RDMA Performance with increasing number of connections test configuration

Item	Description
<b>Test Case</b>	NVMe-oF RDMA Target performance with increasing # of connections
<b>SPDK NVMe-oF Target configuration</b>	Same as in Test Case #1, using 4 CPU cores.
<b>Kernel NVMe-oF Target configuration</b>	Target configuration file loaded using nvmet-cli tool. For detailed configuration file contents please see <a href="#">Appendix C</a> .
<b>Kernel NVMe-oF Initiator #1</b>	<p><b>Device config</b> 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 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 nvme connect -n nqn.2018-09.io.spdk:cnode8 -t rdma -a 20.0.1.1 -s 4420</pre>
<b>Kernel NVMe-oF Initiator #2</b>	<p><b>Device config</b> Performed using nvme-cli tool.</p> <pre>modprobe nvme-fabrics 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 10.0.0.1 -s 4420</pre>

	<pre>nvme connect -n nqn.2018-09.io.spdk:cnode13 -t rdma -a 10.0.1.1 -s 4420 nvme connect -n nqn.2018-09.io.spdk:cnode14 -t rdma -a 10.0.1.1 -s 4420 nvme connect -n nqn.2018-09.io.spdk:cnode15 -t rdma -a 10.0.1.1 -s 4420 nvme connect -n nqn.2018-09.io.spdk:cnode16 -t rdma -a 10.0.1.1 -s 4420</pre>
<p><b>FIO configuration (used on both initiators)</b></p>	<pre><b>FIO.conf</b> [global] ioengine=libaio thread=1 group_reporting=1 direct=1  norandommap=1 rw=randrw rwmixread={100, 70, 0} bs=4k iodepth={8, 16, 32, 64, 128} time_based=1 ramp_time=60 runtime=300 numjobs={1, 4, 16, 32}  [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 [filename8] filename=/dev/nvme7n1</pre>

The SPDK NVMe-oF Target was configured to use 4 CPU cores, whereas 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. For the case of Kernel NVMe-oF target, total CPU cores was calculated from % CPU utilization which was measured using sar utility in Linux.

## 4KB Random Read Results

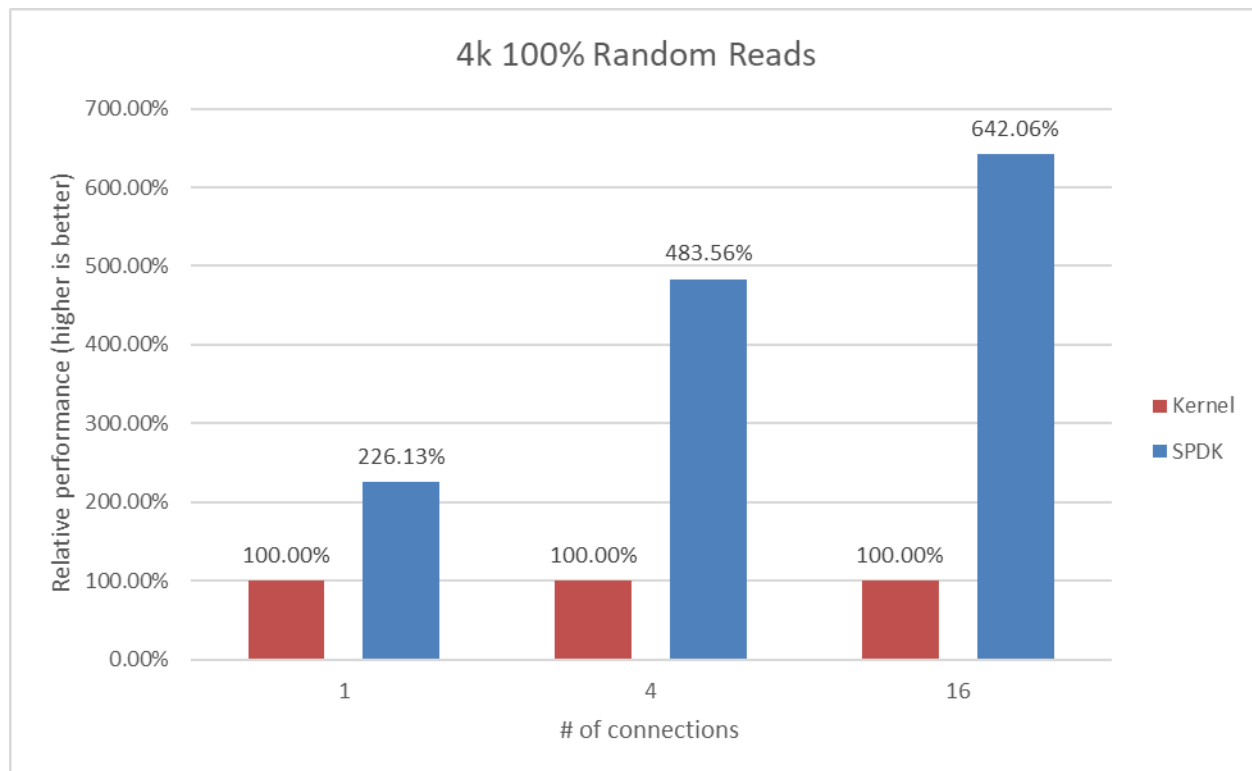


Figure 11: Relative Efficiency Comparison of Linux Kernel vs. SPDK NVMe-oF RDMA Target IOPS/Core for 4KB 100% Random Reads QD=64, using the Kernel Initiator

Table 27: Linux Kernel NVMe-oF RDMA Target: 4KB 100% Random Reads, QD=64

Connections per subsystem	Bandwidth (MBps)	IOPS (k)	Avg. Latency (usec)	# CPU Cores
<b>1</b>	16442.01	4209.2	243.1	10.2
<b>4</b>	15344.70	3928.2	260.7	21.4
<b>16</b>	14193.84	3633.6	281.1	30.7

Table 28: SPDK NVMe-oF RDMA Target: 4KB 100% Random Reads, QD=64

Connections per subsystem	Bandwidth (MBps)	IOPS (k)	Avg. Latency (usec)	# CPU Cores
<b>1</b>	14595.82	3736.5	274.7	4.0
<b>4</b>	13897.39	3557.7	287.8	4.0
<b>16</b>	11878.66	3040.9	336.1	4.0

## 4KB Random Write results

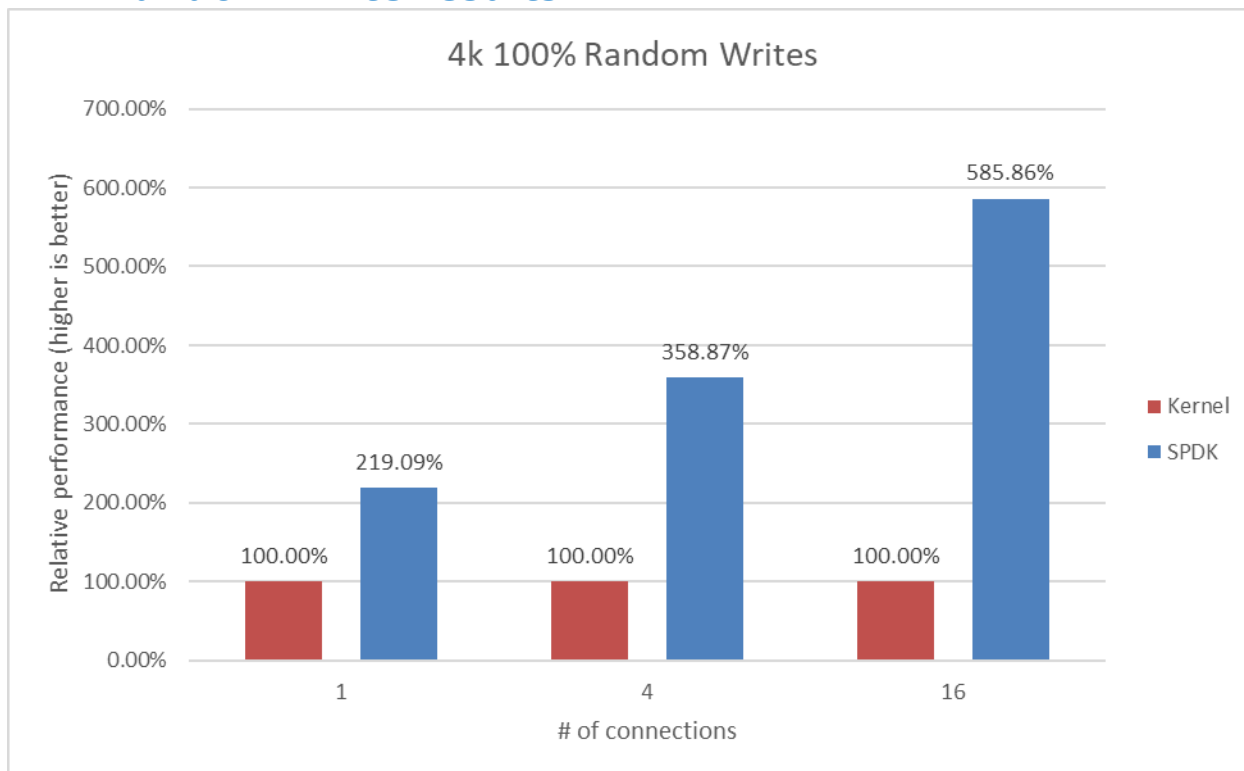


Figure 12: Relative Efficiency Comparison of Linux Kernel vs. SPDK NVMe-oF RDMA Target IOPS/Core for 4KB 100% Random Writes QD=64 workload, using Kernel Initiators

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

Table 29: Linux Kernel NVMe-oF RDMA Target: 4KB 100% Random Writes, QD=64

Connections per subsystem	Bandwidth (MBps)	IOPS (k)	Avg. Latency (usec)	# CPU Cores
1	22377.10	5728.5	178.6	11.1
4	19471.61	4984.7	212.5	17.2
16	13832.29	3541.1	288.5	24.0

Table 30: SPDK NVMe-oF RDMA Target: 4KB 100% Random Writes, QD=64

Connections per subsystem	Bandwidth (MBps)	IOPS (k)	Avg. Latency (usec)	# CPU Cores
1	17709.01	4533.5	226.3	4.0
4	16242.82	4158.2	246.0	4.0
16	13496.49	3455.1	295.8	4.0

## 4KB Random Read-Write Results

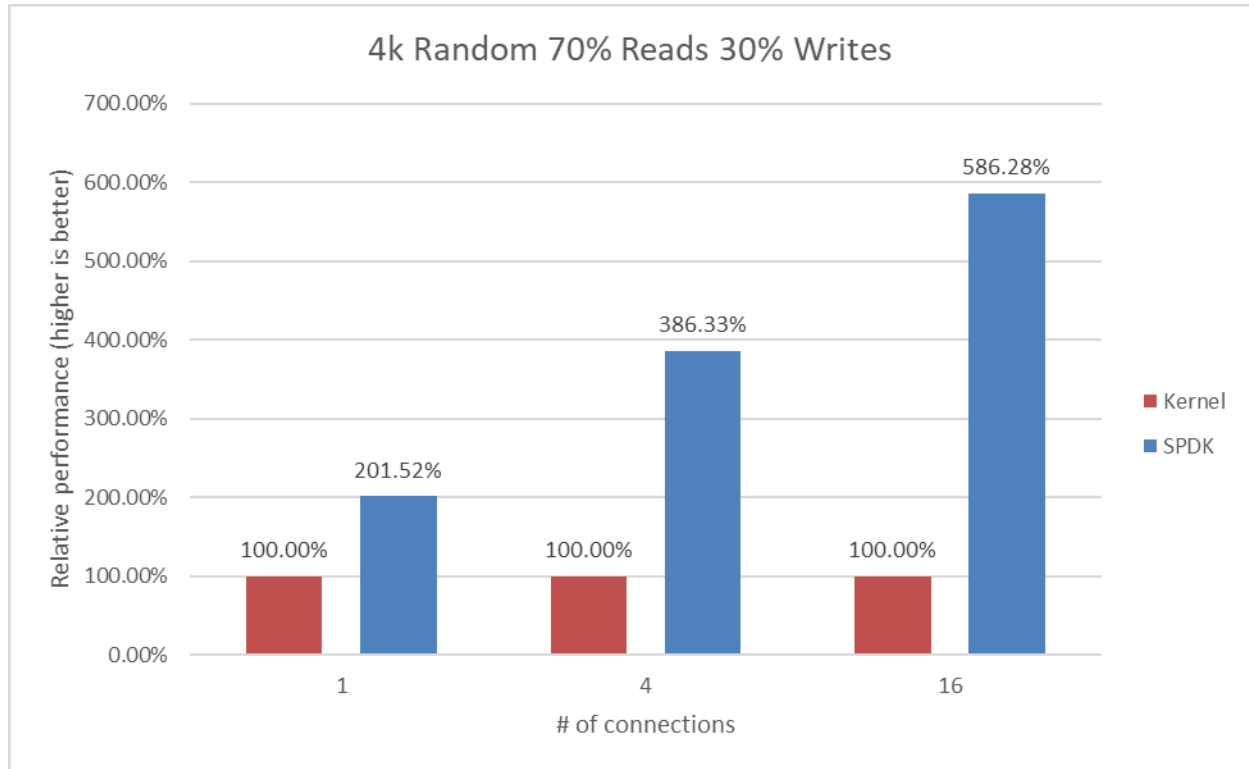


Figure 13: Relative Efficiency Comparison of Linux Kernel vs. SPDK NVMe-oF RDMA Target IOPS/Core for 4KB Random 70% Reads 30% Writes QD=64 Workload, using Kernel Initiators

Table 31: Linux Kernel NVMe-oF RDMA Target: 4KB 70% Random Read 30% Random Write, QD=64

Connections per subsystem	Bandwidth (MBps)	IOPS (k)	Avg. Latency (usec)	# CPU Cores
<b>1</b>	17199.80	4403.1	232.4	10.2
<b>4</b>	18687.08	4783.9	213.7	21.1
<b>16</b>	16583.12	4245.3	240.4	33.2

Table 32: SPDK NVMe-oF RDMA Target: 4KB 70% Random Read 30% Random Write, QD=64

Connections per subsystem	Bandwidth (MBps)	IOPS (k)	Avg. Latency (usec)	# CPU Cores
<b>1</b>	13527.96	3463.2	296.0	4.0
<b>4</b>	13714.01	3510.8	291.8	4.0
<b>16</b>	11717.53	2999.7	340.6	4.0

## Conclusions

1. When the SPDK NVMe-oF Target was configured with 4 CPU cores the performance peaked at 1 connection per subsystem for all workloads. Increasing the number of connections per subsystem beyond these values resulted in higher average latency and lower IOPS.
2. The performance for the Linux Kernel NVMe-oF Target peaked at 1 connection per subsystem for all workloads. Increasing the number of connections only increases the latency and CPU utilization.
3. The SPDK NVMe-oF target shows up to 6.4x 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 last report, throughput scales up and latency decreases almost linearly with the scaling of SPDK NVMe-oF target and initiator 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 2.6x times lower than Kernel initiator.

The SPDK NVMe-oF Target performed up to 6.4 times better w.r.t IOPS/core than Linux Kernel NVMe-oF target while running 4KB 100% random read workload with increasing number of 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=64 .....	12
Table 7: SPDK NVMe-oF RDMA Target Core Scaling results, Random Write IOPS, QD=32 .....	13
Table 8: SPDK NVMe-oF RDMA Target Core Scaling results, Random Read/Write 70%/30% IOPS, QD=128 .....	14
Table 9: LTO impact for SPDK NVMe-oF RDMA Target Core Scaling performance – Bandwidth .....	15
Table 10: LTO impact for SPDK NVMe-oF RDMA Target Core Scaling performance - Throughput.....	16
Table 11: LTO impact for SPDK NVMe-oF RDMA Target Core Scaling performance – Average Latency	17
Table 12: SPDK NVMe-oF RDMA Target Core Scaling results, 128KB Sequential Read IOPS, QD=4 ....	18
Table 13: SPDK NVMe-oF RDMA Target Core Scaling results, 128KB Sequential Write IOPS, QD=1 ...	18
Table 14: SPDK NVMe-oF RDMA Target Core Scaling results, 128KB Sequential 70% Read 30% Write IOPS, QD=4 .....	18
Table 15: SPDK NVMe-oF RDMA Initiator Core Scaling test configuration .....	20
Table 16: SPDK NVMe-oF RDMA Initiator Core Scaling results, 4KB Random Read IOPS, QD=64, SPDK Target 6 CPU Cores .....	23
Table 17: SPDK NVMe-oF RDMA Initiator Core Scaling results, 4KB Random Write IOPS, QD=32, SPDK Target 6 CPU Cores .....	24
Table 18: SPDK NVMe-oF RDMA Initiator Core Scaling results, 4KB Random 70%/30% Read/Write IOPS, QD=128, SPDK Target 6 CPU Cores.....	25
Table 19: Linux Kernel vs. SPDK NVMe-oF RDMA Latency test configuration .....	27
Table 20: SPDK NVMe-oF RDMA Target Latency and IOPS at QD=1, Null Block Device.....	30
Table 21: Linux Kernel NVMe-oF RDMA Target Latency and IOPS at QD=1, Null Block Device.....	30
Table 22: SPDK NVMe-oF RDMA Initiator Latency and IOPS at QD=1, Null Block Device.....	31
Table 23: Linux Kernel NVMe-oF RDMA Initiator Latency and IOPS at QD=1, Null Block Device.....	31
Table 24: SPDK NVMe-oF RDMA Latency and IOPS at QD=1, Null Block Device.....	32
Table 25: Linux Kernel NVMe-oF RDMA Latency and IOPS at QD=1, Null Block Device .....	32
Table 26: NVMe-oF RDMA Performance with increasing number of connections test configuration .....	34
Table 27: Linux Kernel NVMe-oF RDMA Target: 4KB 100% Random Reads, QD=64 .....	36

*Table 28: SPDK NVMe-oF RDMA Target: 4KB 100% Random Reads, QD=64.....36*

*Table 29: Linux Kernel NVMe-oF RDMA Target: 4KB 100% Random Writes, QD=64 .....37*

*Table 30: SPDK NVMe-oF RDMA Target: 4KB 100% Random Writes, QD=64 .....37*

*Table 31: Linux Kernel NVMe-oF RDMA Target: 4KB 70% Random Read 30% Random Write, QD=6438*

*Table 32: SPDK NVMe-oF RDMA Target: 4KB 70% Random Read 30% Random Write, QD=64 .....38*

## 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 4KB 100% Random Read workload at QD = 64 .....	12
Figure 3: SPDK NVMe-oF RDMA Target I/O core scaling: IOPS vs. Latency while running 4KB 100% Random Write Workload at QD=32 .....	13
Figure 4: SPDK NVMe-oF RDMA Target I/O core scaling: IOPS vs. Latency while running 4KB 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 4KB 100% Random Read QD=64 workload .....	23
Figure 6: SPDK NVMe-oF RDMA Initiator I/O core scaling: IOPS vs. Latency while running 4KB 100% Random Write Workload at QD=32 .....	24
Figure 7: SPDK NVMe-oF RDMA Initiator I/O core scaling: IOPS vs. Latency while running 4KB Random 70% Read 30% Write Workload at QD=128 .....	25
Figure 8: SPDK vs. Kernel NVMe-oF RDMA Target average I/O latency for various workloads run using the Kernel Initiator.....	30
Figure 9: SPDK vs. Kernel NVMe-oF RDMA Initiator average I/O latency for various workloads against SPDK Target.....	31
Figure 10: SPDK vs. Kernel NVMe-oF RDMA solutions average I/O Latency for various workloads against SPDK Target.....	32
Figure 11: Relative Efficiency Comparison of Linux Kernel vs. SPDK NVMe-oF RDMA Target IOPS/Core for 4KB 100% Random Reads QD=64, using the Kernel Initiator.....	36
Figure 12: Relative Efficiency Comparison of Linux Kernel vs. SPDK NVMe-oF RDMA Target IOPS/Core for 4KB 100% Random Writes QD=64 workload, using Kernel Initiators.....	37
Figure 13: Relative Efficiency Comparison of Linux Kernel vs. SPDK NVMe-oF RDMA Target IOPS/Core for 4KB Random 70% Reads 30% Writes QD=64 Workload, using Kernel Initiators .....	38

## 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"
    }
}
]
}
}
```

## 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"
    ]
  }
}
```

```

    },
    {
      "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-5c7b23c7a6eb"
      },
      "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"
},
]
```



```
"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](http://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.

§