

# SPDK NVMe-oF RDMA (Target & Initiator) Performance Report Release 20.04

---

**Testing Date:** June 2020

**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:**

John Kariuki ([john.k.kariuki@intel.com](mailto:john.k.kariuki@intel.com))

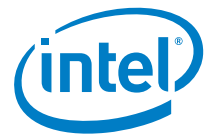
James Harris ([james.r.harris@intel.com](mailto:james.r.harris@intel.com))



# 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 Impact .....	15
Conclusions .....	16
Test Case 2: SPDK NVMe-oF RDMA Initiator I/O core scaling .....	19
4KB Random Read Results .....	22
4KB Random Write Results .....	23
4KB Random 70/30 Read/Write Results .....	24
Conclusions .....	25
Test Case 3: Linux Kernel vs. SPDK NVMe-oF RDMA Latency .....	26
SPDK vs Kernel NVMe-oF RDMA Target Results .....	29
Conclusions .....	30
SPDK vs Kernel NVMe-oF RDMA Initiator Results .....	31
Conclusions .....	32
Test Case 4: NVMe-oF RDMA Performance with increasing # of connections .....	33
4KB Random Read Results .....	35
4KB Random Write results.....	36
4KB Random Read-Write Results.....	37
Conclusions .....	38
Summary .....	39
Appendix A.....	40



## ***Audience and Purpose***

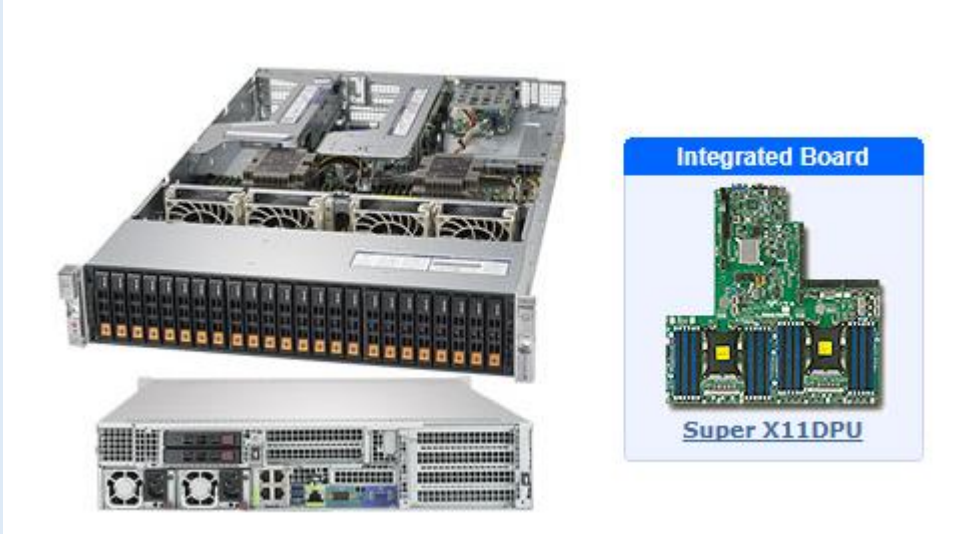
---

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

Item	Description
Server Platform	SuperMicro SYS-2029U-TN24R4T 
CPU	<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)
Memory	10 x 32GB Hynix HMA84GR7AFR4N-VK, DDR4, 2666MHz Total of 320GB
Operating System	Fedora 30
BIOS	3.1a
Linux kernel version	5.4.14-100.fc30
SPDK version	SPDK 20.04
Storage	<b>OS:</b> 1x 120GB Intel SSDSC2BB120G4 <b>Storage Target:</b> 16x <a href="#">Intel® SSD DC P4610™ 1.6TB</a> (FW: QDV10190) (8 on each CPU socket)
NIC	2x 100GbE Mellanox ConnectX-5 NICs. Both ports connected. 1 NIC per CPU socket.



## Initiator 1 Configuration

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)
Memory	6 x 32GB Micron M393A1G40EB1-CRC, DDR4, 2933MHz Total 192GBs
Operating System	Fedora 30
BIOS	02.01.0008 03/19/2019
Linux kernel version	5.4.14-100.fc30
SPDK version	SPDK 20.04
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

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)
Memory	6 x 32GB Micron M393A1G40EB1-CRC, DDR4, 2933MHz Total 192GBs
Operating System	02.01.0008 03/19/2019
BIOS	3.1 06/08/2018
Linux kernel version	5.4.14-100.fc30
SPDK version	SPDK 20.04
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

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 Fedora 5.4.14-100.fc30 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, TCP and Intel® Omni-Path. 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, iWARP, and Omni-Path NICs. 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).

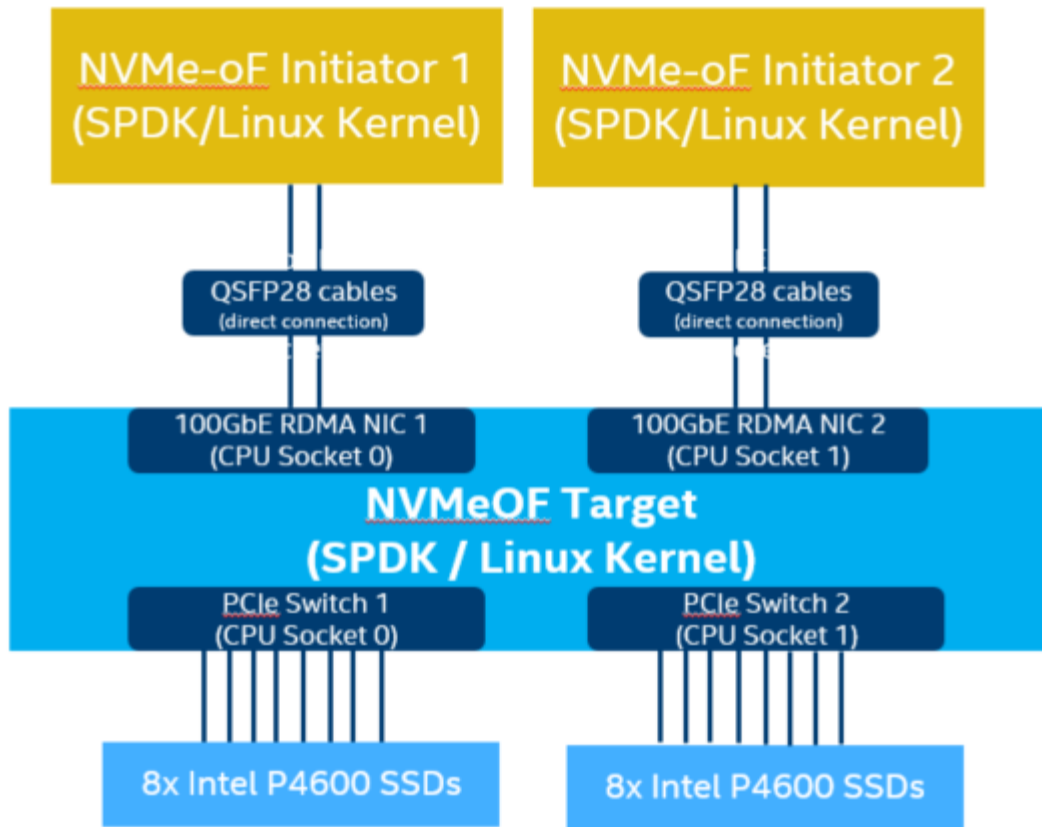


Figure 1: High-Level NVMe-oF performance testing setup





# 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 1, 2, 3, 4, 5 and 6 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 the test configuration. The SPDK NVMe-oF Target was configured using JSON-RPC; the table contains a sequence of commands used by `spdk/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.

Item	Description
Test Case	SPDK NVMe-oF Target I/O core scaling
SPDK NVMe-oF Target configuration	<p>The commands below were executed with <code>spdk/scripts/rpc.py</code> script.</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 construct_nvme_bdev -t PCIe -b Nvme11 -a 0000:b8:00.0                     </pre>



```
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

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)

nvmf_subsystem_create nqn.2018-09.io.spdk:cnode1 -s SPDK001 -a -m 8
nvmf_subsystem_create nqn.2018-09.io.spdk:cnode2 -s SPDK002 -a -m 8
nvmf_subsystem_create nqn.2018-09.io.spdk:cnode3 -s SPDK003 -a -m 8
nvmf_subsystem_create nqn.2018-09.io.spdk:cnode4 -s SPDK004 -a -m 8
nvmf_subsystem_create nqn.2018-09.io.spdk:cnode5 -s SPDK005 -a -m 8
nvmf_subsystem_create nqn.2018-09.io.spdk:cnode6 -s SPDK006 -a -m 8
nvmf_subsystem_create nqn.2018-09.io.spdk:cnode7 -s SPDK007 -a -m 8
nvmf_subsystem_create nqn.2018-09.io.spdk:cnode8 -s SPDK008 -a -m 8
nvmf_subsystem_create nqn.2018-09.io.spdk:cnode9 -s SPDK009 -a -m 8
nvmf_subsystem_create nqn.2018-09.io.spdk:cnode10 -s SPDK0010 -a -m 8
nvmf_subsystem_create nqn.2018-09.io.spdk:cnode11 -s SPDK0011 -a -m 8
nvmf_subsystem_create nqn.2018-09.io.spdk:cnode12 -s SPDK0012 -a -m 8
nvmf_subsystem_create nqn.2018-09.io.spdk:cnode13 -s SPDK0013 -a -m 8
nvmf_subsystem_create nqn.2018-09.io.spdk:cnode14 -s SPDK0014 -a -m 8
nvmf_subsystem_create nqn.2018-09.io.spdk:cnode15 -s SPDK0015 -a -m 8
nvmf_subsystem_create nqn.2018-09.io.spdk:cnode16 -s SPDK0016 -a -m 8

nvmf_subsystem_add_ns nqn.2018-09.io.spdk:cnode1 Nvme0n1
nvmf_subsystem_add_ns nqn.2018-09.io.spdk:cnode2 Nvme1n1
nvmf_subsystem_add_ns nqn.2018-09.io.spdk:cnode3 Nvme2n1
nvmf_subsystem_add_ns nqn.2018-09.io.spdk:cnode4 Nvme3n1
nvmf_subsystem_add_ns nqn.2018-09.io.spdk:cnode5 Nvme4n1
nvmf_subsystem_add_ns nqn.2018-09.io.spdk:cnode6 Nvme5n1
nvmf_subsystem_add_ns nqn.2018-09.io.spdk:cnode7 Nvme6n1
nvmf_subsystem_add_ns nqn.2018-09.io.spdk:cnode8 Nvme7n1
nvmf_subsystem_add_ns nqn.2018-09.io.spdk:cnode9 Nvme8n1
nvmf_subsystem_add_ns nqn.2018-09.io.spdk:cnode10 Nvme9n1
nvmf_subsystem_add_ns nqn.2018-09.io.spdk:cnode11 Nvme10n1
nvmf_subsystem_add_ns nqn.2018-09.io.spdk:cnode12 Nvme11n1
nvmf_subsystem_add_ns nqn.2018-09.io.spdk:cnode13 Nvme12n1
nvmf_subsystem_add_ns nqn.2018-09.io.spdk:cnode14 Nvme13n1
nvmf_subsystem_add_ns nqn.2018-09.io.spdk:cnode15 Nvme14n1
nvmf_subsystem_add_ns nqn.2018-09.io.spdk:cnode16 Nvme15n1

nvmf_subsystem_add_listener nqn.2018-09.io.spdk:cnode1 -t rdma -f ipv4 -s 4420 -a 20.0.0.1
nvmf_subsystem_add_listener nqn.2018-09.io.spdk:cnode2 -t rdma -f ipv4 -s 4420 -a 20.0.0.1
nvmf_subsystem_add_listener nqn.2018-09.io.spdk:cnode3 -t rdma -f ipv4 -s 4420 -a 20.0.0.1
nvmf_subsystem_add_listener nqn.2018-09.io.spdk:cnode4 -t rdma -f ipv4 -s 4420 -a 20.0.0.1
```



	<pre> nvme_subsystem_add_listener nqn.2018-09.io.spdk:cnode5 -t rdma -f ipv4 -s 4420 -a 20.0.1.1 nvme_subsystem_add_listener nqn.2018-09.io.spdk:cnode6 -t rdma -f ipv4 -s 4420 -a 20.0.1.1 nvme_subsystem_add_listener nqn.2018-09.io.spdk:cnode7 -t rdma -f ipv4 -s 4420 -a 20.0.1.1 nvme_subsystem_add_listener nqn.2018-09.io.spdk:cnode8 -t rdma -f ipv4 -s 4420 -a 20.0.1.1 nvme_subsystem_add_listener nqn.2018-09.io.spdk:cnode9 -t rdma -f ipv4 -s 4420 -a 10.0.0.1 nvme_subsystem_add_listener nqn.2018-09.io.spdk:cnode10 -t rdma -f ipv4 -s 4420 -a 10.0.0.1 nvme_subsystem_add_listener nqn.2018-09.io.spdk:cnode11 -t rdma -f ipv4 -s 4420 -a 10.0.0.1 nvme_subsystem_add_listener nqn.2018-09.io.spdk:cnode12 -t rdma -f ipv4 -s 4420 -a 10.0.0.1 nvme_subsystem_add_listener nqn.2018-09.io.spdk:cnode13 -t rdma -f ipv4 -s 4420 -a 10.0.1.1 nvme_subsystem_add_listener nqn.2018-09.io.spdk:cnode14 -t rdma -f ipv4 -s 4420 -a 10.0.1.1 nvme_subsystem_add_listener nqn.2018-09.io.spdk:cnode15 -t rdma -f ipv4 -s 4420 -a 10.0.1.1 nvme_subsystem_add_listener nqn.2018-09.io.spdk:cnode16 -t rdma -f ipv4 -s 4420 -a 10.0.1.1 </pre>
<p><b>SPDK NVMe-oF Initiator - FIO plugin configuration</b></p>	<pre> <b>BDEV.conf</b> [Nvme] TransportId "trtype:RDMA adrfam:IPv4 traddr:20.0.0.1 trsvcid:4420 subnqn:nqn.2018-09.io.spdk:cnode1" Nvme0 TransportId "trtype:RDMA adrfam:IPv4 traddr:20.0.0.1 trsvcid:4420 subnqn:nqn.2018-09.io.spdk:cnode2" Nvme1 TransportId "trtype:RDMA adrfam:IPv4 traddr:20.0.0.1 trsvcid:4420 subnqn:nqn.2018-09.io.spdk:cnode3" Nvme2 TransportId "trtype:RDMA adrfam:IPv4 traddr:20.0.0.1 trsvcid:4420 subnqn:nqn.2018-09.io.spdk:cnode4" Nvme3 TransportId "trtype:RDMA adrfam:IPv4 traddr:20.0.1.1 trsvcid:4420 subnqn:nqn.2018-09.io.spdk:cnode5" Nvme4 TransportId "trtype:RDMA adrfam:IPv4 traddr:20.0.1.1 trsvcid:4420 subnqn:nqn.2018-09.io.spdk:cnode6" Nvme5 TransportId "trtype:RDMA adrfam:IPv4 traddr:20.0.1.1 trsvcid:4420 subnqn:nqn.2018-09.io.spdk:cnode7" Nvme6 TransportId "trtype:RDMA adrfam:IPv4 traddr:20.0.1.1 trsvcid:4420 subnqn:nqn.2018-09.io.spdk:cnode8" Nvme7  <b>FIO.conf</b> [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, 8, 16, 32} time_based=1 ramp_time=60 runtime=300  [filename0] filename=Nvme0n1 [filename1] filename=Nvme1n1 [filename2] filename=Nvme2n1 [filename3] filename=Nvme3n1 [filename4] filename=Nvme4n1 [filename5] filename=Nvme5n1 [filename6] filename=Nvme6n1 [filename7] filename=Nvme7n1 </pre>

## 4KB Random Read Results

Test Result: 4KB 100% Random Read IOPS QD=64

# of Cores	Bandwidth (MBps)	Throughput (IOPS k)	Avg. Latency (usec)
1 core	4136.13	1058.8	966.8
2 cores	9629.23	2465.1	415.2
3 cores	15862.90	4060.9	255.5
4 cores	21323.86	5458.9	187.4
5 cores	21438.12	5488.2	186.5
6 cores	21448.79	5490.9	186.3

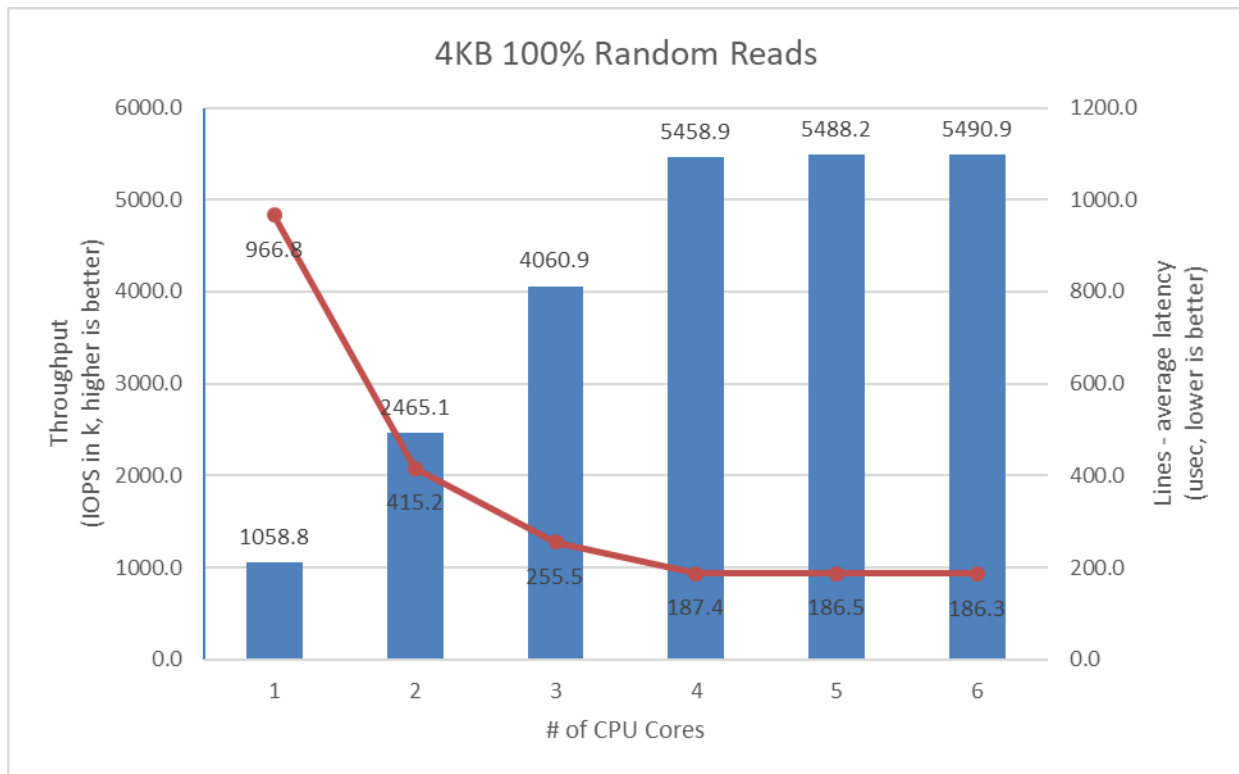


Figure 2: SPDK NVMe-oF RDMA Target I/O core scaling: IOPS vs. Latency while running 4KB 100% Random Read workload at QD=64



## 4KB Random Write Results

Test Result: 4KB 100% Random Writes IOPS QD=32

# of Cores	Bandwidth (MBps)	Throughput (IOPS k)	Avg. Latency (usec)
1 core	5845.37	1496.4	340.1
2 cores	12719.51	3256.2	156.0
3 cores	18924.27	4844.6	104.2
4 cores	21857.14	5595.4	91.3
5 cores	22433.72	5743.0	88.5
6 cores	22545.45	5771.6	88.3

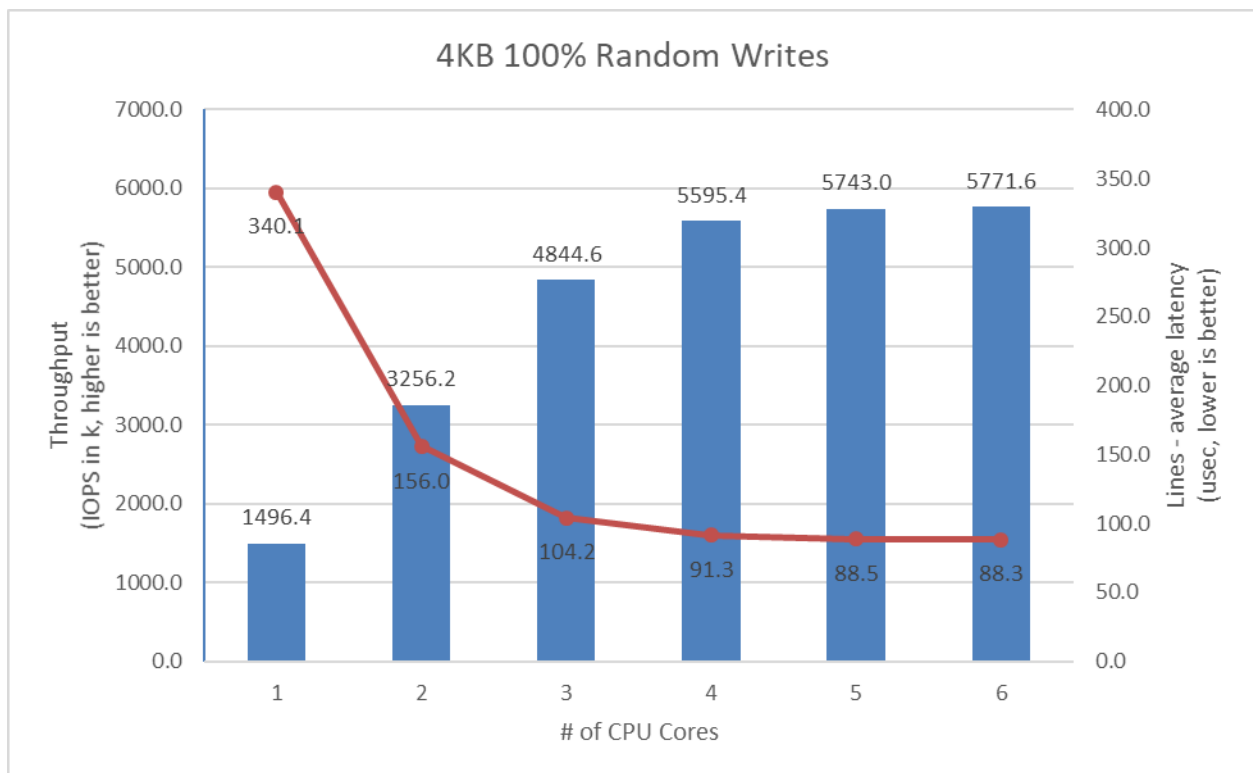
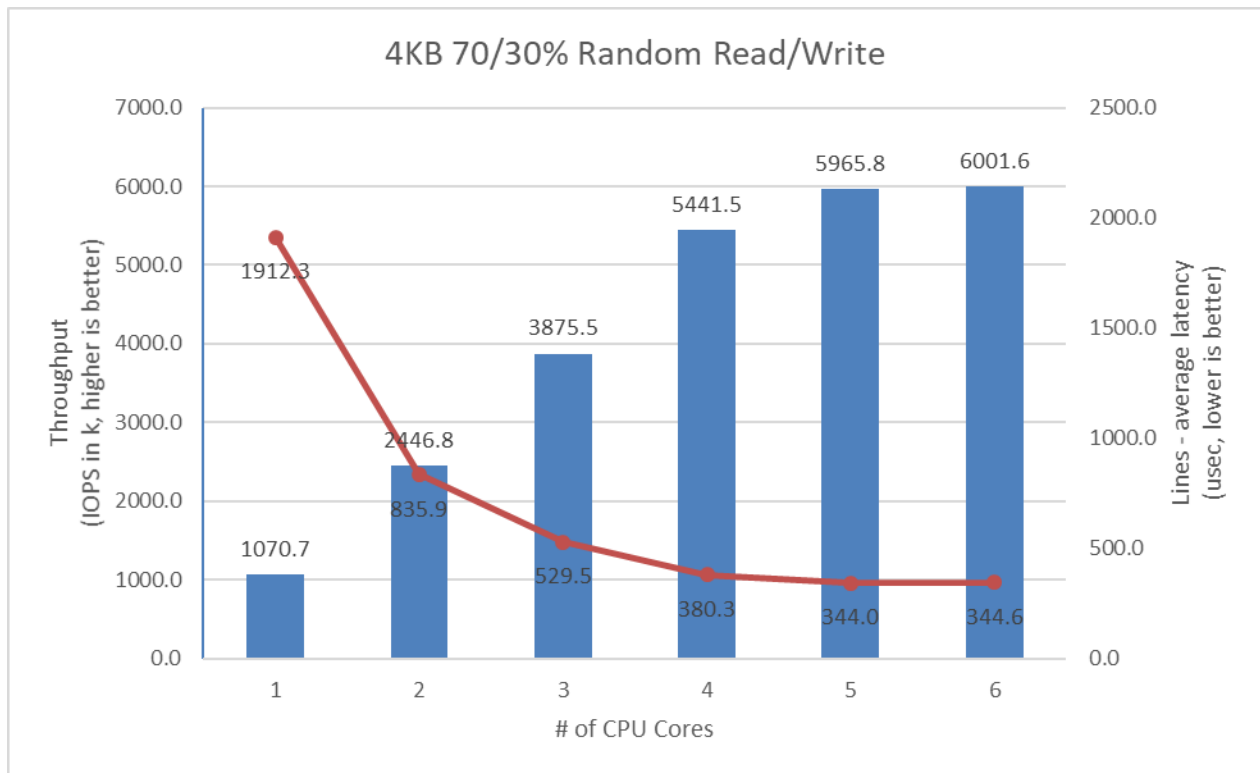


Figure 3: SPDK NVMe-oF RDMA Target I/O core scaling: IOPS vs. Latency while running 4KB 100% Random Write workload at QD=32

## 4KB Random Read-Write Results

Test Result: 4KB 70% Read 30% Write IOPS, QD=128

# of Cores	Bandwidth (MBps)	Throughput (IOPS k)	Avg. Latency (usec)
1 core	4182.31	1070.7	1912.3
2 cores	9557.81	2446.8	835.9
3 cores	15138.60	3875.5	529.5
4 cores	21255.81	5441.5	380.3
5 cores	23304.02	5965.8	344.0
6 cores	23443.90	6001.6	344.6



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



## LTO Impact

Test cases presented in previous sub-chapters were run with LTO option enabled. LTO is “Link-Time Optimization” and is a SPDK compilation option aimed at improving SPDK efficiency. Below table presents impact LTO has on the performance test results.

Workload	# of Cores	Bandwidth (MBps)			IOPS (thousands)			Avg. Latency (usec)		
		Non-LTO	LTO	LTO Diff. (%)	Non-LTO	LTO	LTO Diff. (%)	Non-LTO	LTO	LTO Diff. (%)
4KB 100% Random Read, QD=64	1	4136.13	4104.81	-0.76%	1058.85	1050.83	-0.76%	966.81	974.19	0.76%
	2	9629.23	9763.84	1.40%	2465.08	2499.54	1.40%	415.21	409.36	-1.41%
	3	15862.90	15956.74	0.59%	4060.90	4084.92	0.59%	255.51	252.67	-1.11%
	4	21323.86	21697.67	1.75%	5458.91	5554.60	1.75%	187.39	184.15	-1.73%
	5	21438.12	20943.86	-2.31%	5488.15	5361.63	-2.31%	186.53	191.16	2.48%
	6	21448.79	23412.01	9.15%	5490.89	5993.47	9.15%	186.33	171.37	-8.03%
4KB 100% Random Write, QD=32	1	5606.04	5559.52	-0.83%	1435.14	1423.23	-0.83%	711.21	717.18	0.84%
	2	12608.85	12630.12	0.17%	3227.86	3233.31	0.17%	314.57	314.06	-0.16%
	3	19823.19	19977.09	0.78%	5074.73	5114.13	0.78%	200.00	200.85	0.42%
	4	22119.64	22199.71	0.36%	5662.62	5683.12	0.36%	179.58	179.12	-0.26%
	5	21941.56	22633.01	3.15%	5617.04	5794.05	3.15%	181.62	176.05	-3.07%
	6	22065.39	22794.61	3.30%	5648.74	5835.42	3.30%	180.98	175.02	-3.30%
4KB 70%/30% Random Read/Write, QD=64	1	4182.31	4201.42	0.46%	1070.66	1075.56	0.46%	1912.31	1903.61	-0.45%
	2	9557.81	9446.39	-1.17%	2446.79	2418.27	-1.17%	835.93	845.81	1.18%
	3	15138.60	15499.71	2.39%	3875.47	3967.92	2.39%	529.46	523.00	-1.22%
	4	21255.81	21706.87	2.12%	5441.48	5556.95	2.12%	380.35	373.95	-1.68%
	5	23304.02	22272.14	-4.43%	5965.82	5701.66	-4.43%	344.05	359.88	4.60%
	6	23443.90	25284.17	7.85%	6001.63	6472.74	7.85%	344.58	317.14	-7.96%



## Conclusions

1. For 4KB 100% Random Reads workload, the 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 minimal performance gain as the network is almost saturated.
2. For 4KB 100% Random Write workload, throughput scales up and latency decreases linearly with the scaling of I/O cores until the network is saturated at just 4 CPU cores. The SSDs were not preconditioned prior to executing the test to ensure the 4KB random write workload saturated the network; preconditioning the SSDs limits the workload performance to the SSDs steady state performance.
3. For 4K Random 70/30 Reads/Writes workload, performance scales up and latency decreases close to linearly with the scaling of SPDK NVMe-OF Target I/O cores up to 4 CPU cores. Beyond that there is a negligible performance increase when increasing the number of cores up to 5 and 6. This workload is close to reaching network saturation at about 5M IOPS.





## 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.

### Test Result: 128KB 100% Sequential Reads QD=8

# of Cores	Bandwidth (MBps)	Throughput (IOPS k)	Avg. Latency (usec)
1 core	24688.13	197.5	648.1
2 cores	24762.93	198.1	646.1
3 cores	24787.76	198.3	645.4
4 cores	24736.94	197.9	646.8

### Test Result: 128KB 100% Sequential Writes QD=1

# of Cores	Bandwidth (MBps)	Throughput (IOPS k)	Avg. Latency (usec)
1 core	23037.78	184.3	86.6
2 cores	23279.03	186.2	85.7
3 cores	23292.70	186.3	85.7
4 cores	23285.87	186.3	85.7

### Test Result: 128KB 70% Reads 30% Writes QD=8

# of Cores	Bandwidth (MBps)	Throughput (IOPS k)	Avg. Latency (usec)
1 core	25680.21	205.4	626.3
2 cores	25793.45	206.3	624.7
3 cores	25841.67	206.7	624.8
4 cores	25924.10	207.4	623.4



## **Conclusions**

1. For large sequential I/Os, a single CPU core 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.



## 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.

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> [Nvme] TransportId "trtype:RDMA adrfam:IPv4 traddr:20.0.0.1 trsvcid:4420 subnqn:nqn.2018-09.io.spdk:cnode1" Nvme0 TransportId "trtype:RDMA adrfam:IPv4 traddr:20.0.0.1 trsvcid:4420 subnqn:nqn.2018-09.io.spdk:cnode2" Nvme1 TransportId "trtype:RDMA adrfam:IPv4 traddr:20.0.0.1 trsvcid:4420 subnqn:nqn.2018-09.io.spdk:cnode3" Nvme2 TransportId "trtype:RDMA adrfam:IPv4 traddr:20.0.0.1 trsvcid:4420 subnqn:nqn.2018-09.io.spdk:cnode4" Nvme3 TransportId "trtype:RDMA adrfam:IPv4 traddr:20.0.1.1 trsvcid:4420 subnqn:nqn.2018-09.io.spdk:cnode5" Nvme4 TransportId "trtype:RDMA adrfam:IPv4 traddr:20.0.1.1 trsvcid:4420 subnqn:nqn.2018-09.io.spdk:cnode6" Nvme5 TransportId "trtype:RDMA adrfam:IPv4 traddr:20.0.1.1 trsvcid:4420 subnqn:nqn.2018-09.io.spdk:cnode7" Nvme6 TransportId "trtype:RDMA adrfam:IPv4 traddr:20.0.1.1 trsvcid:4420 subnqn:nqn.2018-09.io.spdk:cnode8" Nvme7</p> <p><b>FIO.conf</b> [global] ioengine=/tmp/spdk/examples/bdev/fio_plugin/fio_plugin spdk_conf=/tmp/spdk/bdev.conf</p>

	<pre> thread=1 group_reporting=1 direct=1  norandommap=1 rw=randrw 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>	<pre> <b>BDEV.conf</b> [Nvme] TransportId "trtype:RDMA adrfam:IPv4 traddr:10.0.0.1 trsvcid:4420 subnqn:nqn.2018-09.io.spdk:cnode9" Nvme0 TransportId "trtype:RDMA adrfam:IPv4 traddr:10.0.0.1 trsvcid:4420 subnqn:nqn.2018-09.io.spdk:cnode10" Nvme1 TransportId "trtype:RDMA adrfam:IPv4 traddr:10.0.0.1 trsvcid:4420 subnqn:nqn.2018-09.io.spdk:cnode11" Nvme2 TransportId "trtype:RDMA adrfam:IPv4 traddr:10.0.0.1 trsvcid:4420 subnqn:nqn.2018-09.io.spdk:cnode12" Nvme3 TransportId "trtype:RDMA adrfam:IPv4 traddr:10.0.1.1 trsvcid:4421 subnqn:nqn.2018-09.io.spdk:cnode13" Nvme4 TransportId "trtype:RDMA adrfam:IPv4 traddr:10.0.1.1 trsvcid:4421 subnqn:nqn.2018-09.io.spdk:cnode14" Nvme5 TransportId "trtype:RDMA adrfam:IPv4 traddr:10.0.1.1 trsvcid:4421 subnqn:nqn.2018-09.io.spdk:cnode15" Nvme6 TransportId "trtype:RDMA adrfam:IPv4 traddr:10.0.1.1 trsvcid:4421 subnqn:nqn.2018-09.io.spdk:cnode16" Nvme7  <b>FIO.conf</b> Similar as Initiator 1. The only differences are were in the filename section which are shown below.  <b>FIO.conf filename section for 1 CPU test run</b> N/A (only Initiator 1 is used)  <b>FIO.conf filename section for 2 CPU test run</b> [filename0] filename=Nvme0n1 filename=Nvme1n1 filename=Nvme2n1 filename=Nvme3n1 </pre>



**FIO.conf filename section for 3 CPU test run**

```
[filename0]  
filename=Nvme0n1  
filename=Nvme1n1  
filename=Nvme2n1  
filename=Nvme3n1
```

**FIO.conf filename section for 4 CPU test run**

```
[filename0]  
filename=Nvme0n1  
filename=Nvme1n1  
filename=Nvme2n1  
filename=Nvme3n1  
[filename1]  
filename=Nvme4n1  
filename=Nvme5n1  
filename=Nvme6n1  
filename=Nvme7n1
```

## 4KB Random Read Results

Test Result: 4KB 100% Random Read, QD=64, SPDK Target 6 CPU cores

# of Initiator CPU Cores	Bandwidth (MBps)	Throughput (IOPS k)	Avg. Latency (usec)
1 core	5818.90	1489.6	279.4
2 cores	15420.15	3947.6	224.1
3 cores	19302.99	4941.6	198.1
4 cores	22493.35	5758.3	177.5
5 cores	20071.76	5138.4	199.5
6 cores	22709.26	5813.6	176.9
7 cores	22118.47	5662.3	181.3
8 cores	22824.53	5843.1	175.7

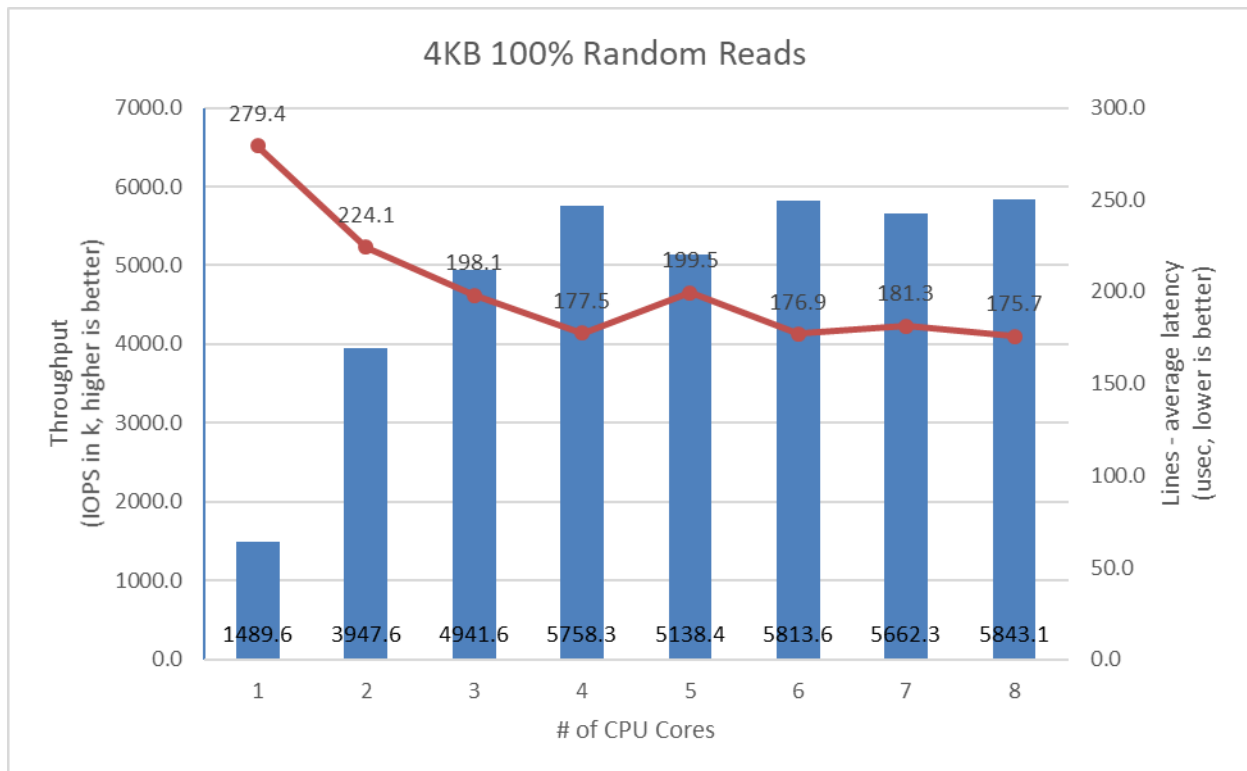


Figure 5: SPDK NVMe-oF RDMA Initiator I/O core scaling: IOPS vs. Latency while running 4KB 100% Random Read workload at QD=64



## 4KB Random Write Results

**Note:** The SSDs were pre-conditioned just once with 2x Sequential Write workload before running all 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).

Test Result: 4KB 100% Random Write, QD=32, SPDK Target 6 CPU Cores

# of Initiator CPU Cores	Bandwidth (MBps)	Throughput (IOPS k)	Avg. Latency (usec)
1 core	5875.95	1504.2	135.7
2 cores	13937.19	3567.9	114.4
3 cores	18260.68	4674.7	100.5
4 cores	22941.55	5873.0	84.1
5 cores	22911.96	5865.5	85.0
6 cores	22380.28	5729.3	88.4
7 cores	22933.87	5871.1	86.4
8 cores	22643.68	5796.8	87.7

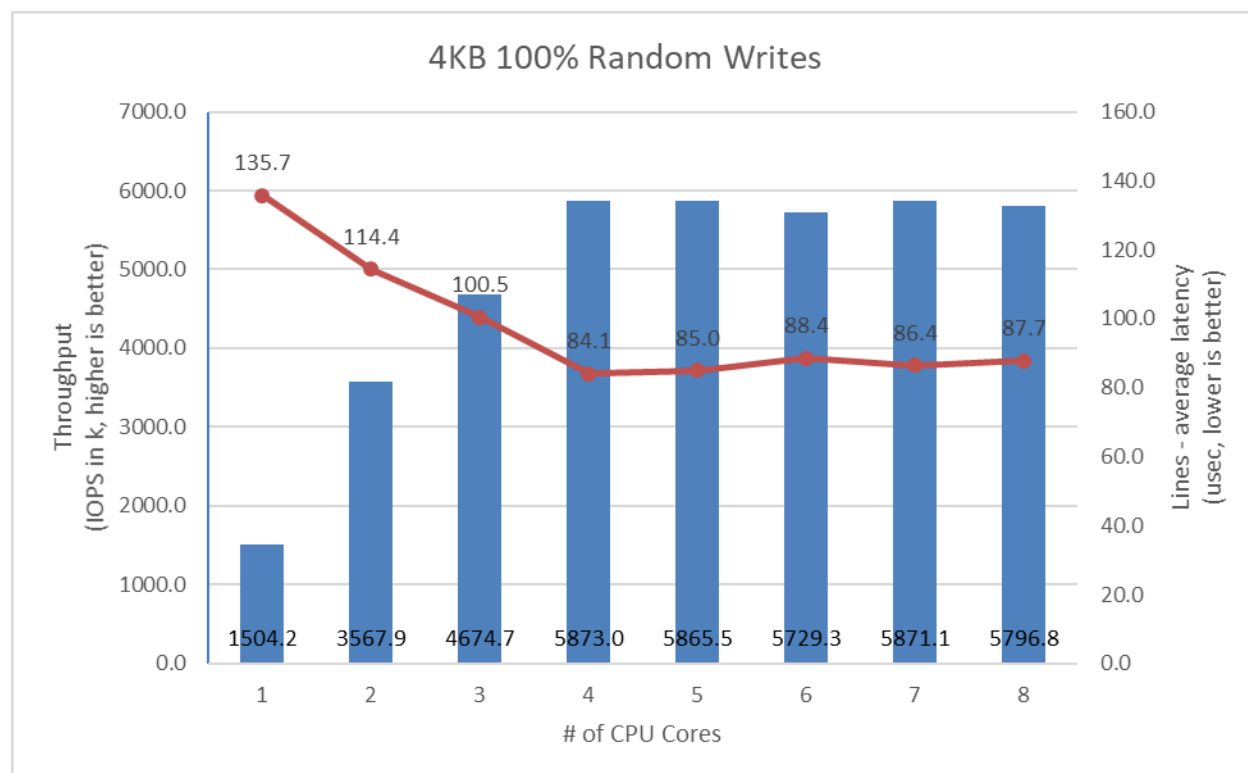


Figure 6: SPDK NVMe-oF RDMA Initiator I/O core scaling: IOPS vs. Latency while running 4KB 100% Random Write workload at QD=32



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

Test Result: 4KB 70% Random Read 30% Random Write QD=128, SPDK Target 6 CPU Cores

# of Initiator CPU Cores	Bandwidth (MBps)	Throughput (IOPS k)	Avg. Latency (usec)
1 core	4844.89	1240.3	660.8
2 cores	10848.77	2777.3	605.1
3 cores	16474.93	4217.6	476.5
4 cores	22241.05	5693.7	352.4
5 cores	22408.20	5736.5	355.0
6 cores	20977.05	5370.1	380.7
7 cores	20963.83	5366.7	381.9
8 cores	22104.74	5658.8	363.4

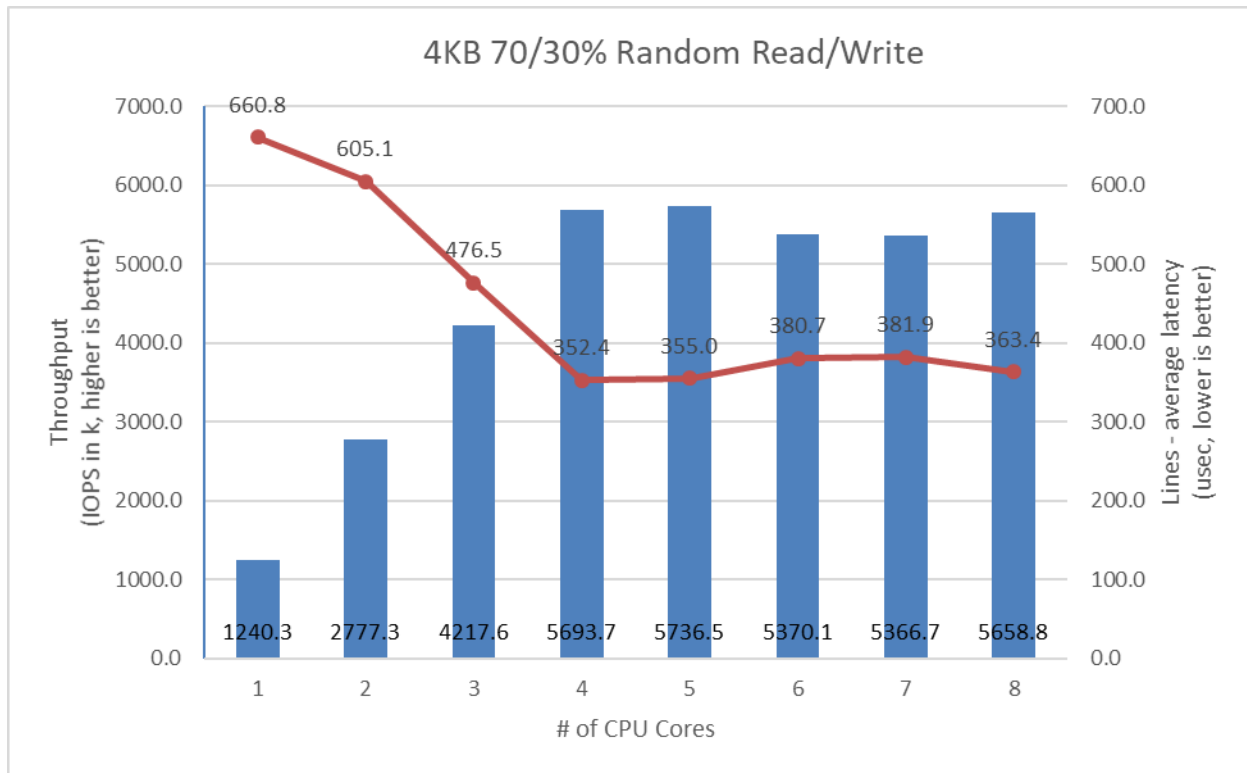
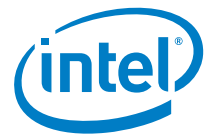


Figure 7: SPDK NVMe-oF RDMA Initiator I/O core scaling: IOPS vs. Latency while running 4KB Random 70% Read 30% Write workload at QD=128





## Conclusions

1. Scaling is linear or close to linear in case of all tested workloads.
2. Peak performance of 5.7M IOPS was reached at 4 Initiator CPU cores for all workloads, reaching network saturation.

## 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.

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>nvmf-cli</code> tool.</p> <pre> {   "ports": [     {       "addr": {         "adrfam": "ipv4",         "traddr": "20.0.0.1",         "trsvcid": "4420",         "trtype": "rdma"       },       "portid": 1,       "referrals": [],     }   ] }           </pre>



	<pre> "subsystems": [   "nqn.2018-09.io.spdk:cnode1" ] }, "hosts": [], "subsystems": [   {     "allowed_hosts": [],     "attr": {       "allow_any_host": "1",       "version": "1.3"     },     "namespaces": [       {         "device": {           "path": "/dev/nullb0",           "uuid": "621e25d2-8334-4c1a-8532-b6454390b8f9"         },         "enable": 1,         "nsid": 1       }     ],     "nqn": "nqn.2018-09.io.spdk:cnode1"   } ] } </pre>
<b>FIO configuration</b>	
<p><b>SPDK NVMe-oF Initiator FIO plugin configuration</b></p>	<p><b>BDEV.conf</b> [Nvme] TransportId "trtype:RDMA adrfam:IPv4 traddr:20.0.0.1 trsvcid:4420 subnqn:nqn.2018-09.io.spdk:cnode1" Nvme0</p> <p><b>FIO.conf</b> [global] ioengine=/tmp/spdk/examples/bdev/fio_plugin/fio_plugin spdk_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=Nvme0n1</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>

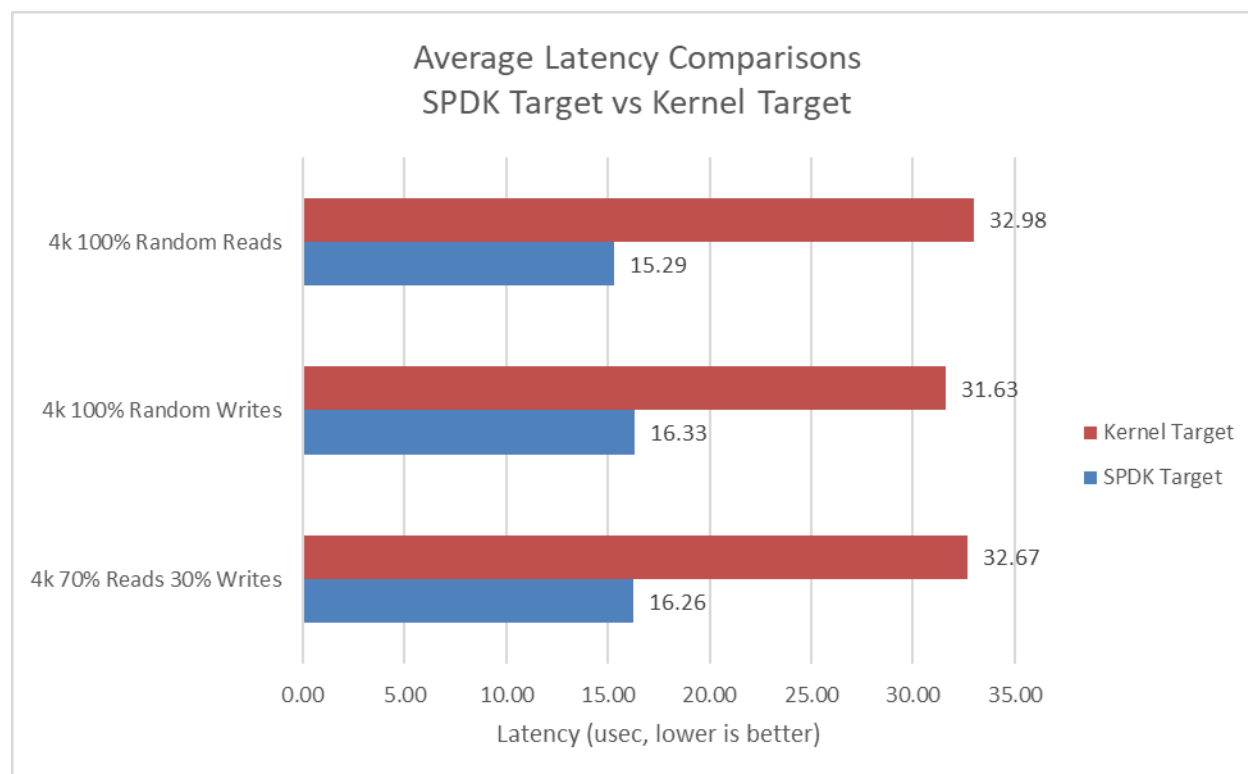


```
FIO.conf  
[global]  
ioengine=libaio  
thread=1  
group_reporting=1  
direct=1  
  
norandommap=1  
rw=randrw  
rwmixread={100, 70, 0}  
bs=4k  
iodepth=1  
time_based=1  
ramp_time=60  
runtime=300  
  
[filename0]  
filename=/dev/nvme0n1
```



## 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:** Average I/O Latency comparisons at QD=1 between SPDK and Linux Kernel NVMe-oF RDMA Target for various workloads using the Linux Kernel Initiator

### SPDK NVMe-oF RDMA Target QD=1 Latency for a 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>	15.29	63860	17.2	20.9	37.5	139.6
<b>4KB 100% Random Write</b>	16.33	59601	17.9	21.1	40.7	144.4
<b>4KB 70/30% Random Read/Write</b>	16.26	59889	17.0	20.0	39.1	146.9

### Linux Kernel NVMe-oF RDMA Target QD=1 Latency for a 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>	32.98	29548	32.9	38.0	49.7	99.8
<b>4KB 100% Random Write</b>	31.63	30746	32.0	36.9	64.5	161.5
<b>4KB 70/30% Random Read/Write</b>	32.67	29770	32.8	38.1	46.8	98.2



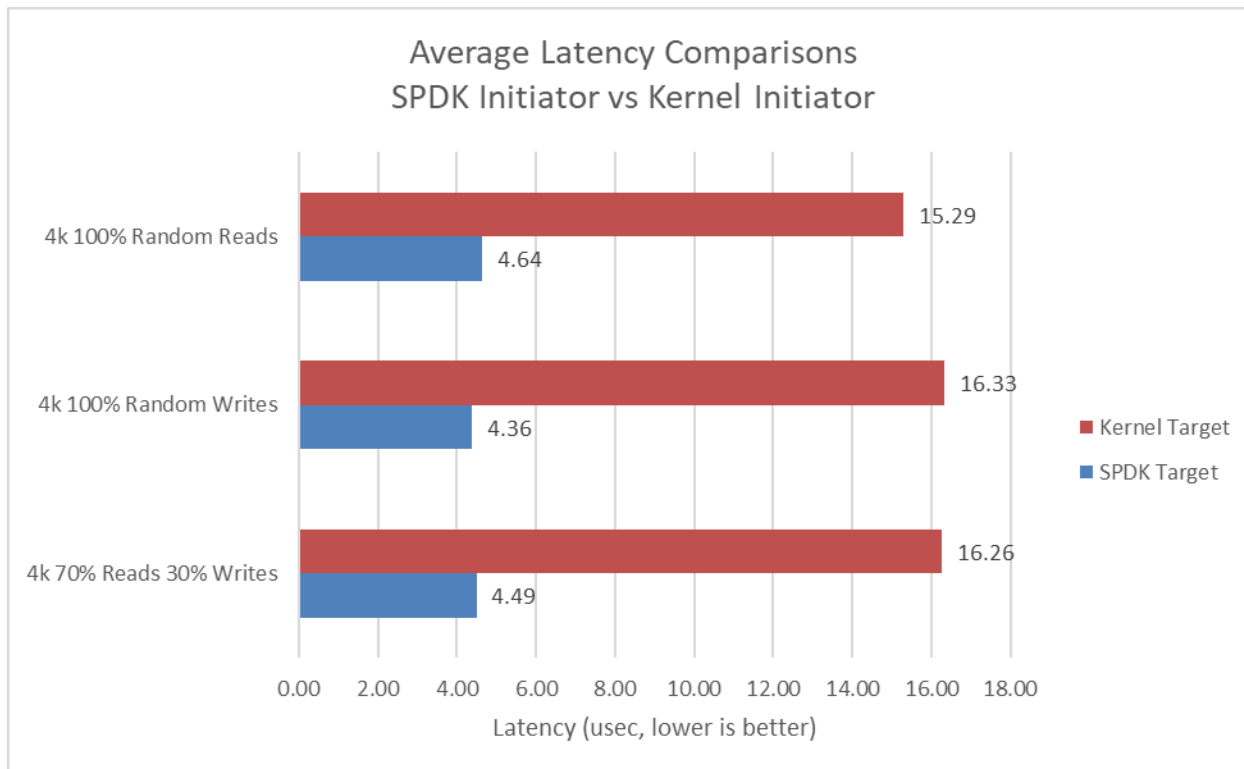
## **Conclusions**

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 18 usec vs. the Linux Kernel NVMe-oF target used in Fedora 30 5.4.14 setup. This is entirely software overhead, therefore, using the SPDK NVMe-oF target reduces the NVMe-oF software overhead by approximately 50% vs. the Linux Kernel NVMe-oF target.



## 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:** Average I/O Latency Comparisons at QD=1 between SPDK and Linux Kernel NVMe-oF RDMA Initiator for various workloads against SPDK NVMe-OF Target

### SPDK NVMe-oF RDMA Initiator QD=1 Latency for a 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>	4.64	208149	4.8	14.3	20.2	45.3
<b>4KB 100% Random Write</b>	4.36	220680	4.5	8.2	19.7	45.8
<b>4KB 70/30% Random Read/Write</b>	4.49	214302	4.6	17.6	20.2	46.5

### Linux Kernel NVMe-oF RDMA Initiator QD=1 Latency for a 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>	15.29	63860	17.2	20.9	37.5	139.6
<b>4KB 100% Random Write</b>	16.33	59601	17.9	21.1	40.7	144.4
<b>4KB 70/30% Random Read/Write</b>	16.26	59889	17.0	20.0	39.1	146.9



## **Conclusions**

- 1.** The SPDK NVMe-oF initiator reduces the NVMe-oF software overhead by up to 3.8x times vs. the Linux Kernel NVMe-oF Initiator for the RDMA transport.





## 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

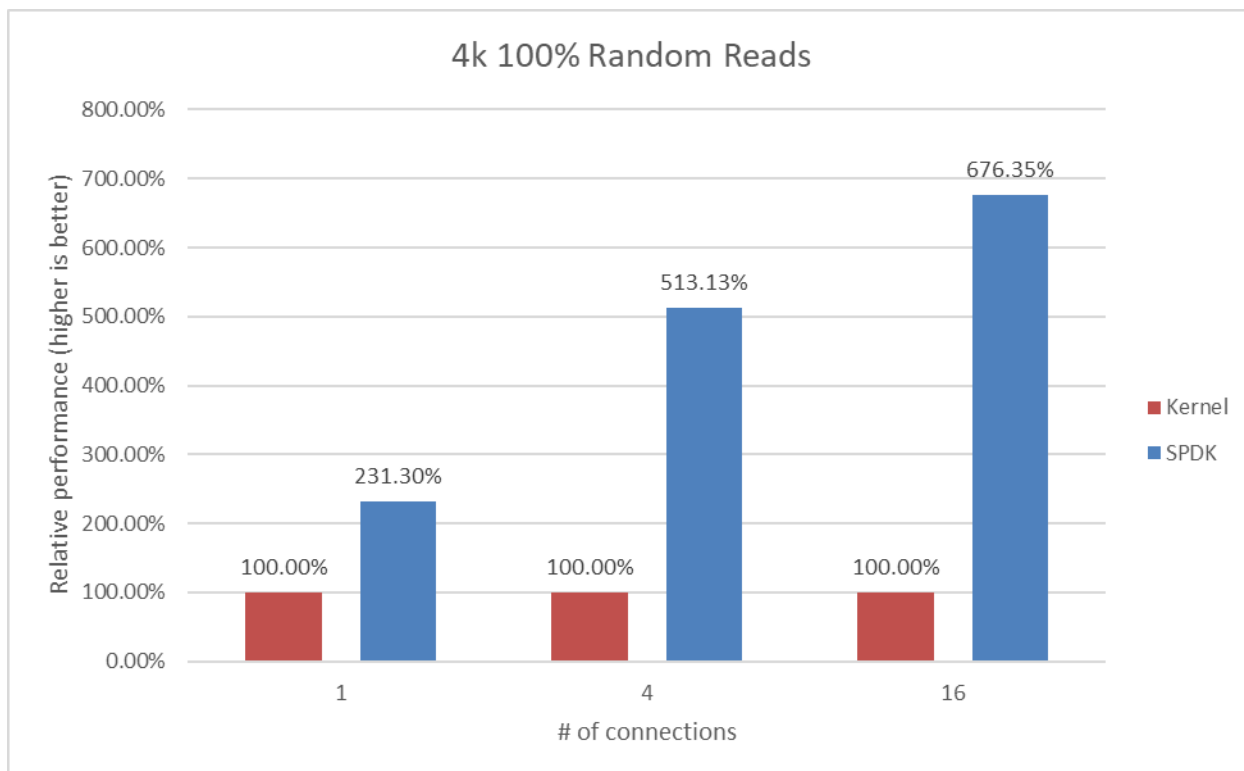
Item	Description
<b>Test Case</b>	NVMe-oF RDMA Target performance with increasing # of connections
<b>SPDK NVMe-oF Target configuration</b>	Same as in Test Case #1, using 4 CPU cores.
<b>Kernel NVMe-oF Target configuration</b>	Target configuration file loaded using nvmet-cli tool. For detail configuration file contents please see Appendix 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 performance 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 10:** Relative Performance Comparison of Linux Kernel vs. SPDK NVMe-oF RDMA Target for 4KB 100% Random Read QD=64, using Kernel Initiators

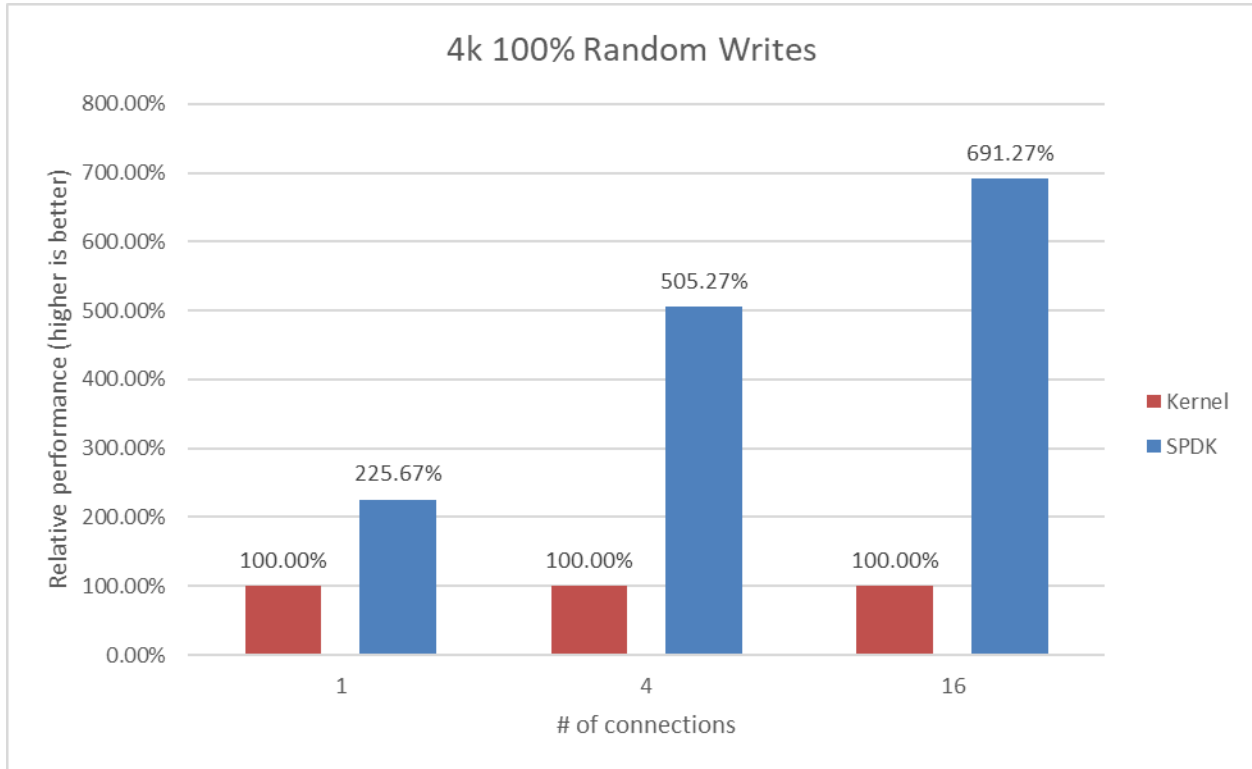
### Linux Kernel NVMe-oF RDMA Target: 4KB 100% Random Reads, QD=64

Connections per subsystem	Bandwidth (MBps)	IOPS (k)	Avg. Latency (usec)	# CPU Cores
<b>1</b>	22396.58	5733.5	356.8	11.2
<b>4</b>	15846.66	4056.7	507.4	19.8
<b>16</b>	13209.90	3381.7	602.0	26.8

### 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>	18066.32	4625.0	442.6	4.0
<b>4</b>	16188.24	4144.2	492.9	4.0
<b>16</b>	13613.13	3485.0	584.0	4.0

## 4KB Random Write results



**Figure 11:** Relative Performance Comparison of Linux Kernel vs. SPDK NVMe-oF RDMA Target for 4KB 100% Random Write QD=64, using Kernel Initiators

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

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

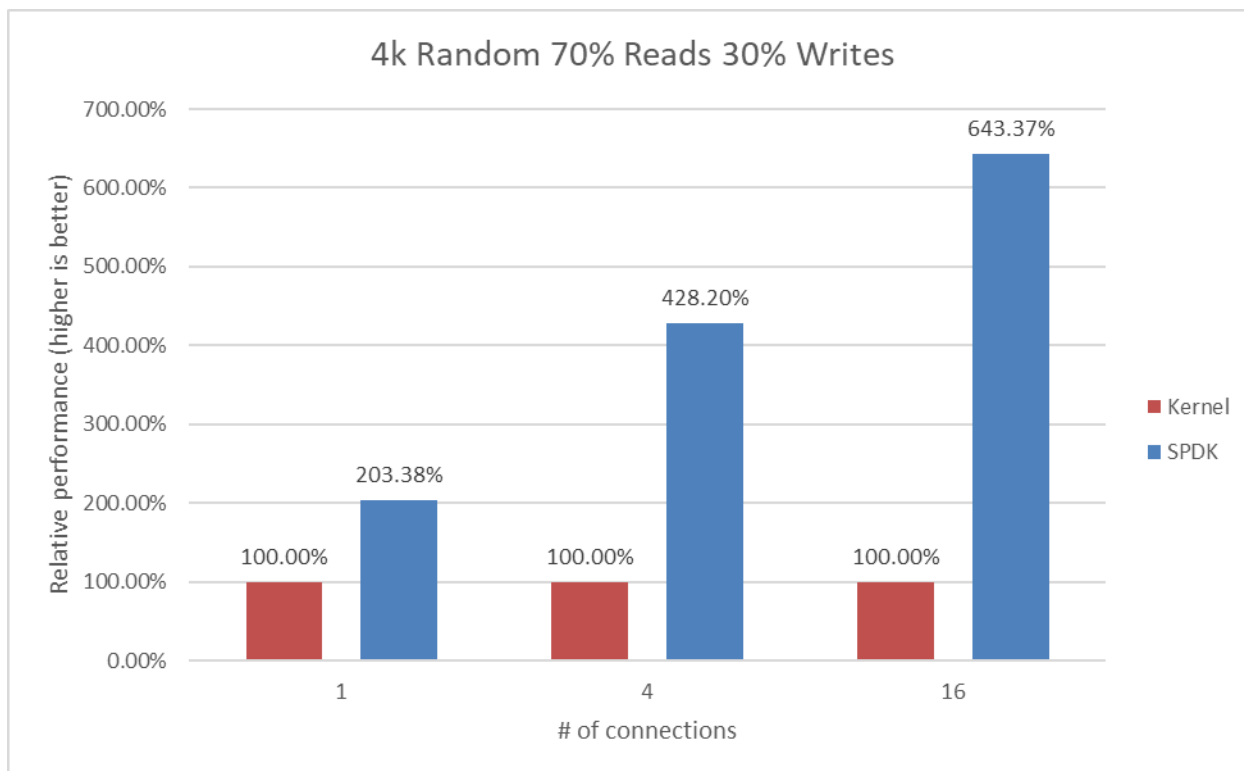
Connections per subsystem	Bandwidth (MBps)	IOPS (k)	Avg. Latency (usec)	# CPU Cores
<b>1</b>	22396.58	5733.5	356.8	11.2
<b>4</b>	15846.66	4056.7	507.4	19.8
<b>16</b>	13209.90	3381.7	602.0	26.8

### SPDK NVMe-oF RDMA Target: 4KB 100% Random Writes, QD=32

Connections per subsystem	Bandwidth (MBps)	IOPS (k)	Avg. Latency (usec)	# CPU Cores
<b>1</b>	18066.32	4625.0	442.6	4.0
<b>4</b>	16188.24	4144.2	492.9	4.0
<b>16</b>	13613.13	3485.0	584.0	4.0



## 4KB Random Read-Write Results



**Figure 12:** Relative Performance Comparison of Linux Kernel vs. SPDK NVMe-oF RDMA Target for 4KB Random 70% Read 30% Write QD=32, using Kernel Initiators

### Linux Kernel NVMe-oF RDMA Target: 4KB 70% Random Read 30% Random Write, QD=32

Connections per subsystem	Bandwidth (MBps)	IOPS (k)	Avg. Latency (usec)	# CPU Cores
<b>1</b>	17374.28	4447.8	381.8	10.4
<b>4</b>	17260.07	4418.6	398.9	21.6
<b>16</b>	14251.12	3648.3	500.3	31.0

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

Connections per subsystem	Bandwidth (MBps)	IOPS (k)	Avg. Latency (usec)	# CPU Cores
<b>1</b>	13576.94	3475.7	524.0	4.0
<b>4</b>	13698.54	3506.8	538.0	4.0
<b>16</b>	11838.78	3030.7	639.8	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 6.9x 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 cores. In case of SPDK NVMe-oF Initiator core scaling it was observed that throughput scales almost linearly, but the latency remains rather stable which might suggest that the Target side might be able to process slightly more IO.

It was also observed that the SPDK NVMe-oF Target average latency is up to 17 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 up to 4 times lower than Kernel initiator.

The SPDK NVMe-oF Target performed up to 6.5 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.

## Appendix A

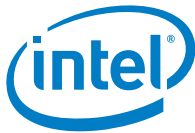
Example Linux Kernel NVMe-oF RDMA Target configuration for Test Case 4.

```
{
  "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"
      ]
    }
  ]
}
```





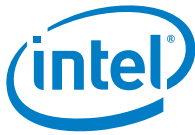
```
]
},
{
  "addr": {
    "adrfam": "ipv4",
    "traddr": "20.0.1.1",
    "trsvcid": "4424",
    "trtype": "rdma"
  },
  "portid": 5,
  "referrals": [],
  "subsystems": [
    "nqn.2018-09.io.spdk:cnode5"
  ]
},
{
  "addr": {
    "adrfam": "ipv4",
    "traddr": "20.0.1.1",
    "trsvcid": "4425",
    "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",
```

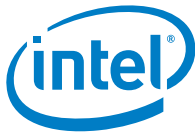


```
    "uuid": "b53be81d-6f5c-4768-b3bd-203614d8cf20"
  },
  "enable": 1,
  "nsid": 1
}
],
"nqn": "nqn.2018-09.io.spdk:cnode1"
},
{
  "allowed_hosts": [],
  "attr": {
    "allow_any_host": "1",
    "version": "1.3"
  },
  "namespaces": [
    {
      "device": {
        "path": "/dev/nvme1n1",
        "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"
},
{
  "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

Software and workloads used in performance tests may have been optimized for performance only on Intel microprocessors.

Performance tests, such as SYSmark and MobileMark, are measured using specific computer systems, components, software, operations and functions. Any change to any of those factors may cause the results to vary. You should consult other information and performance tests to assist you in fully evaluating your contemplated purchases, including the performance of that product when combined with other products. For more complete information visit [www.intel.com/benchmarks](http://www.intel.com/benchmarks).

Performance results are based on testing as of dates shown in configurations and may not reflect all publicly available updates. See backup for configuration details. No product or component can be absolutely secure.

Your costs and results may vary.

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.

§