

SPDK Vhost Performance Report

Release 22.01

Testing Date: February 2022
Updated: March 2022

Performed by:

Karol Latecki (karol.latecki@intel.com)

Maciej Wawryk (maciejx.wawryk@intel.com)

Acknowledgments:

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

John Kariuki (john.k.kariuki@intel.com)

Contents

Contents	2
Audience and Purpose.....	3
Test setup	4
Hardware configuration	4
BIOS Settings	5
Virtual Machine Settings.....	5
Introduction to the SPDK Vhost target	6
SPDK Vhost target architecture	6
Test Case 1: SPDK Vhost Core Scaling	8
4KB Random Read Results	10
4KB Random Write Results	11
4KB Random Read-Write Results.....	12
Logical Volumes performance impact	13
Packed Ring performance impact.....	14
Test Case 2: Rate Limiting IOPS per VM.....	16
Test Case 2 Results	18
Conclusions	20
Test Case 3: Performance per NVMe drive	21
Test Case 3 results	22
Conclusions	25
Appendix A - Vfiio-User Performance	27
List of Tables	31
List of Figures	32

Audience and Purpose


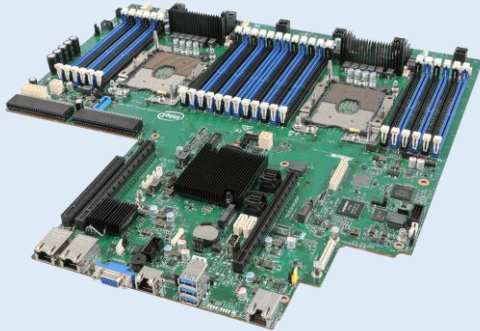
This report is intended for people who are interested in looking at SPDK Vhost scsi and blk stack performance and comparison to its Linux kernel equivalents. It provides performance and efficiency comparisons between SPDK Vhost-scsi and Linux Kernel Vhost-scsi software stacks under various test cases.

The purpose of this report is not to imply a single correct approach, but rather to provide a baseline of well-tested configurations and procedures that produce repeatable and reproducible results. This report can also be viewed as information regarding best known method when performance testing SPDK Vhost-scsi and Vhost-blk stacks.

Test setup

Hardware configuration

Table 1: Hardware setup configuration

Item	Description														
Server Platform	<p>Intel WolfPass R2224WFTZS</p>  <p>Server board S2600WFT</p> 														
Motherboard	S2600WFT														
CPU	<p>2 CPU sockets, Intel(R) Xeon(R) Gold 6230N CPU @ 2.30GHz</p> <p>Number of cores 20 per socket, number of threads 40 per socket</p> <p>Both sockets populated</p> <p>Microcode: 0x4003003</p>														
Memory	<p>12 x 32GB Micron DDR4 36ASF4G72PZ-2G9E2</p> <p>Total 384 GBs</p> <p>Memory channel population:</p> <table> <tr> <th>P1</th><th>P2</th></tr> <tr> <td>CPU1_DIMM_A1</td><td>CPU2_DIMM_A1</td></tr> <tr> <td>CPU1_DIMM_B1</td><td>CPU2_DIMM_B1</td></tr> <tr> <td>CPU1_DIMM_C1</td><td>CPU2_DIMM_C1</td></tr> <tr> <td>CPU1_DIMM_D1</td><td>CPU2_DIMM_D1</td></tr> <tr> <td>CPU1_DIMM_E1</td><td>CPU2_DIMM_E1</td></tr> <tr> <td>CPU1_DIMM_F1</td><td>CPU2_DIMM_F1</td></tr> </table>	P1	P2	CPU1_DIMM_A1	CPU2_DIMM_A1	CPU1_DIMM_B1	CPU2_DIMM_B1	CPU1_DIMM_C1	CPU2_DIMM_C1	CPU1_DIMM_D1	CPU2_DIMM_D1	CPU1_DIMM_E1	CPU2_DIMM_E1	CPU1_DIMM_F1	CPU2_DIMM_F1
P1	P2														
CPU1_DIMM_A1	CPU2_DIMM_A1														
CPU1_DIMM_B1	CPU2_DIMM_B1														
CPU1_DIMM_C1	CPU2_DIMM_C1														
CPU1_DIMM_D1	CPU2_DIMM_D1														
CPU1_DIMM_E1	CPU2_DIMM_E1														
CPU1_DIMM_F1	CPU2_DIMM_F1														
Operating System	Fedora 33														
BIOS	SE5C620.86B.02.01.0013.121520200651														

Linux kernel version	5.11.20-200.fc33.x86_64
SPDK version	SPDK 22.01
Qemu version	QEMU emulator version 5.1.0 (qemu-5.1.0-9.fc33)
Storage	OS: 1x 120GB Intel SSDSC2BB120G4 Storage: 24x Intel® P4610™ 1.6TBs (FW: VDV10170) (6 on CPU NUMA Node 0, 18 on CPU NUMA Node 1)

BIOS Settings

Table 2: Test platform BIOS settings

Item	Description
BIOS	VT-d = Enabled CPU Power and Performance Policy = <Performance> CPU C-state = No Limit CPU P-state = Enabled Enhanced Intel® Speedstep® Tech = Enabled Turbo Boost = Enabled Hyper Threading = Enabled

Virtual Machine Settings

Table 3: Guest VM configuration

Item	Description
CPU	2vCPU, pass through from physical host server. Explicit core usage enforced using “taskset -a -c” command. QEMU arguments for starting the VM: -cpu host -smp 1
Memory	4 GB RAM. Memory is pre-allocated for each VM using Hugepages on host system and used from appropriate NUMA node, to match the CPU which was passed to the VM. QEMU arguments: -m 4096 -object memory-backend-file,id=mem,size=4096M,mem-path=/dev/hugepages,share=on,prealloc=yes,host-nodes=0,policy=bind
Operating System	Fedora 33
Linux kernel version	5.9.16-200.fc33.x86_64
Additional boot options in /etc/default/grub	<ul style="list-style-type: none"> Multi queue enabled: scsi_mod.use_blk_mq=1

Introduction to the SPDK Vhost target

SPDK Vhost is a userspace target designed to extend the performance efficiencies of SPDK into QEMU/KVM virtualization environments. The SPDK Vhost-scsi target presents a broad range of SPDK-managed block devices into virtual machines. SPDK community has leveraged existing SPDK SCSI layer, DPDK Vhost library, QEMU Vhost-scsi and Vhost-blk functionality to create the high performance SPDK userspace Vhost target.

SPDK Vhost target architecture

QEMU sets up the Vhost target via UNIX domain socket. QEMU pre-allocates huge pages for the guest VM to enable DMA by the Vhost target. The guest VM submits I/O directly to the Vhost target via virtqueues in shared memory as shown in Figure 1. The Vhost target transfers data to/from the guest VM via shared memory. The Vhost target then completes I/O to the guest VM via virtqueues in shared memory. There is a completion interrupt sent using eventfd which requires a system call and a guest VM exit. It should be noted that there is no QEMU intervention during the I/O submission process.

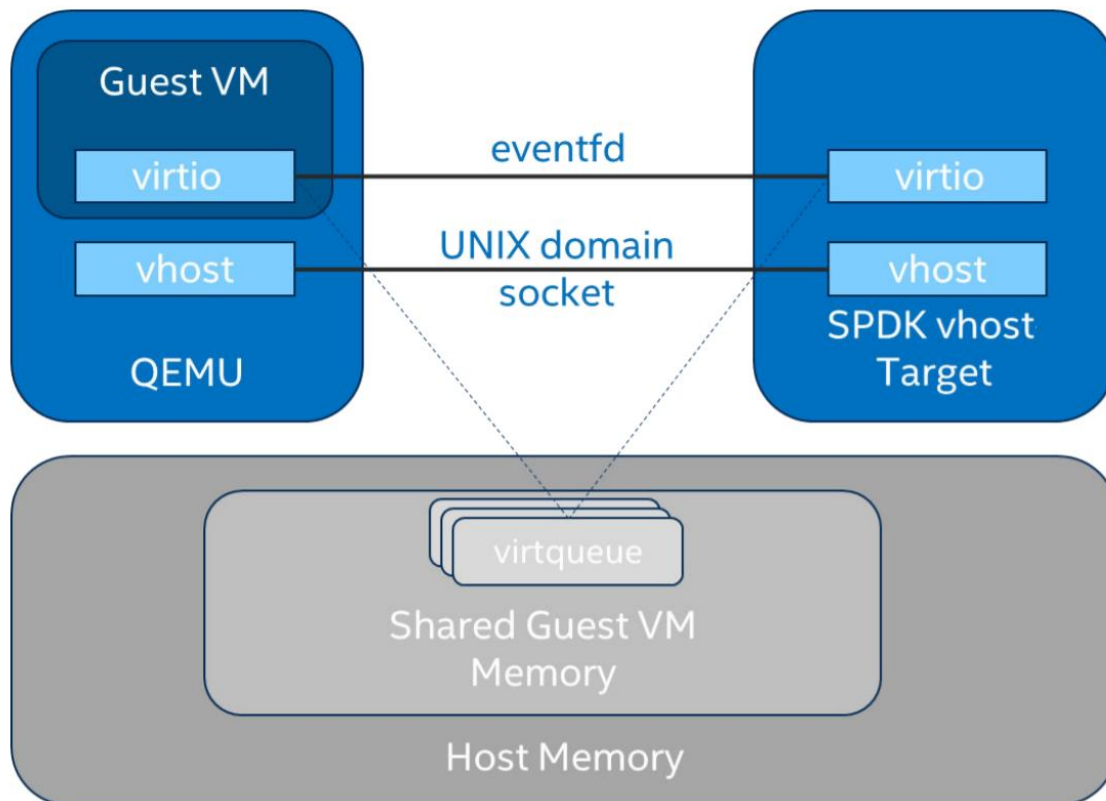


Figure 1: SPDK Vhost-scsi architecture

This report shows the performance comparisons between the traditional interrupt-driven Linux Kernel Vhost-scsi and the accelerated polled-mode SPDK Vhost-scsi under 3 different test cases using local

NVMe storage. Additionally, the SPDK Vhost-blk stack is included in the report for further comparison with the SCSI stack.

Test Case 1: SPDK Vhost Core Scaling

This test case was performed in order to understand aggregate VM performance with SPDK Vhost I/O core scaling. We ran up to 36 virtual machines, each running following FIO workloads:

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

We increased the number of CPU cores used by SPDK Vhost target to process I/O from 1 up to 12 and measured the throughput (in IOPS) and latency. The number of VMs between test runs was not constant and was increased by 6 for each Vhost CPU added, up to a maximum of 36 VMs. VM number was not increased beyond 36 because of the platform capabilities in terms of available CPU cores.

FIO was run in client-server mode. FIO client was run on the host machine and distributed jobs to FIO servers run on each VM. This allowed us to start the FIO jobs across all VMs at the same time.

Results in the table and charts represent aggregate performance (IOPS and average latency) seen across all the VMs. The results are average of 3 runs.

Table 4: SPDK Vhost Core Scaling test configuration

Item	Description
Test case	Test SPDK Vhost target I/O core scaling performance
Test configuration	<p>FIO Version: fio-3.19</p> <p>VM Configuration:</p> <ul style="list-style-type: none"> Common settings are described in the Virtual Machine Settings chapter. Number of VMs: variable (6 VMs per 1 Vhost CPU core, up to 36 VMs max). Each VM has a single Vhost device as a target for the FIO workload. This is achieved by sharing SPDK NVMe bdevs by using either a Split NVMe vbdev or Logical Volume bdev configuration. <p>SPDK Vhost target configuration:</p> <ul style="list-style-type: none"> Test were run with both the Vhost-scsi and Vhost-blk stacks. The Vhost-scsi stack was run with Split NVMe bdevs and Logical Volume bdevs. Vhost-blk stack was run with Logical Volume bdevs. Tests were performed with 1, 2, 4, 6, 8, 10 and 12 Vhost cores for each stack-bdev combination. <p>Kernel Vhost target configuration:</p> <p>- N/A</p>
FIO configuration	<p>[global] ioengine=libaio</p>


```
direct=1
thread=1
norandommap=1
time_based=1
gtod_reduce=0
ramp_time=60s
runtime=240s
numjobs=1
bs=4k
rw=randrw
rwmixread=100 (100% reads), 70 (70% reads, 30% writes), 0 (100%
writes)
iodepth={1, 32, 64}
```

4KB Random Read Results

Table 5: SPDK Vhost core scaling results, 4KB 100% Random Reads IOPS, QD=64

# of CPU cores	# of VMs	Stack / Backend	IOPS (millions)
1	6	SCSI / Split NVMe Bdev	1.90
		SCSI / Lvol Bdev	1.61
		BLK / Lvol Bdev	1.70
2	12	SCSI / Split NVMe Bdev	3.07
		SCSI / Lvol Bdev	2.69
		BLK / Lvol Bdev	2.96
4	24	SCSI / Split NVMe Bdev	5.06
		SCSI / Lvol Bdev	4.20
		BLK / Lvol Bdev	4.52
6	36	SCSI / Split NVMe Bdev	6.62
		SCSI / Lvol Bdev	5.45
		BLK / Lvol Bdev	5.89
8	36	SCSI / Split NVMe Bdev	7.03
		SCSI / Lvol Bdev	6.55
		BLK / Lvol Bdev	7.03
10	36	SCSI / Split NVMe Bdev	6.93
		SCSI / Lvol Bdev	6.47
		BLK / Lvol Bdev	7.15
12	36	SCSI / Split NVMe Bdev	6.74
		SCSI / Lvol Bdev	6.94
		BLK / Lvol Bdev	7.55

Figure 2: Comparison of performance between various SPDK Vhost stack-bdev combinations for 4KB Random Read QD=64 workload

4KB Random Write Results

Table 6: SPDK Vhost core scaling results, 4KB 100% Random Write IOPS, QD=32

# of CPU cores	# of VMs	Stack / Backend	IOPS (millions)
1	6	SCSI / Split NVMe Bdev	1.70
		SCSI / Lvol Bdev	1.58
		BLK / Lvol Bdev	1.65
2	12	SCSI / Split NVMe Bdev	3.03
		SCSI / Lvol Bdev	2.77
		BLK / Lvol Bdev	3.01
4	24	SCSI / Split NVMe Bdev	4.99
		SCSI / Lvol Bdev	4.44
		BLK / Lvol Bdev	4.92
6	36	SCSI / Split NVMe Bdev	5.88
		SCSI / Lvol Bdev	5.58
		BLK / Lvol Bdev	5.99
8	36	SCSI / Split NVMe Bdev	6.65
		SCSI / Lvol Bdev	6.52
		BLK / Lvol Bdev	7.03
10	36	SCSI / Split NVMe Bdev	6.23
		SCSI / Lvol Bdev	6.31
		BLK / Lvol Bdev	6.61
12	36	SCSI / Split NVMe Bdev	6.41
		SCSI / Lvol Bdev	6.46
		BLK / Lvol Bdev	7.23

Figure 3: Comparison of performance between various SPDK Vhost stack-bdev combinations for 4KB Random Write QD=32 workload

4KB Random Read-Write Results

Table 7: SPDK Vhost core scaling results, 4KB Random 70% Read 30% Write IOPS, QD=64

# of CPU cores	# of VMs	Stack / Backend	IOPS (millions)
1	6	SCSI / Split NVMe Bdev	1.74
		SCSI / Lvol Bdev	1.54
		BLK / Lvol Bdev	1.61
2	12	SCSI / Split NVMe Bdev	2.93
		SCSI / Lvol Bdev	2.56
		BLK / Lvol Bdev	2.83
4	24	SCSI / Split NVMe Bdev	4.72
		SCSI / Lvol Bdev	4.20
		BLK / Lvol Bdev	4.59
6	36	SCSI / Split NVMe Bdev	6.16
		SCSI / Lvol Bdev	5.45
		BLK / Lvol Bdev	5.88
8	36	SCSI / Split NVMe Bdev	6.39
		SCSI / Lvol Bdev	6.02
		BLK / Lvol Bdev	6.57
10	36	SCSI / Split NVMe Bdev	5.79
		SCSI / Lvol Bdev	5.85
		BLK / Lvol Bdev	6.38
12	36	SCSI / Split NVMe Bdev	5.91
		SCSI / Lvol Bdev	6.13
		BLK / Lvol Bdev	6.63

Figure 4: Comparison of performance between various SPDK Vhost stack-bdev combinations for 4KB Random 70% Read 30% Write QD=64 workload

Logical Volumes performance impact

The SPDK Vhost SCSI tests were run using two bdev backends – Split NVMe and Logical Volumes. Both “Split NVMe Bdevs” and “Logical Volume Bdevs” allow to logically partition NVMe SSDs, the latter being more flexible in configuration. Here we measure the overhead of extra flexibility afforded by Logical Volumes.

Table 8: Logical Volumes performance impact for SPDK Vhost SCSI

Workload	# of CPU cores	# of VMs	Vhost SCSI + Split NVMe IOPS (millions)	Vhost SCSI + Lvol IOPS (millions)	Lvol Impact (%)
4KB 100% Random Read	1	6	1.90	1.61	-15.06%
	2	12	3.07	2.69	-12.62%
	4	24	5.06	4.20	-17.01%
	6	36	6.62	5.45	-17.71%
	8	36	7.03	6.55	-6.83%
	10	36	6.93	6.47	-6.53%
	12	36	6.74	6.94	2.94%
4KB 100% Random Write	1	6	1.70	1.58	-6.83%
	2	12	3.03	2.77	-8.87%
	4	24	4.99	4.44	-11.04%
	6	36	5.88	5.58	-5.10%
	8	36	6.65	6.52	-1.98%
	10	36	6.23	6.31	1.31%
	12	36	6.41	6.46	0.75%
4KB 70% Random Read 30% Random Write	1	6	1.74	1.54	-11.67%
	2	12	2.93	2.56	-12.38%
	4	24	4.72	4.20	-10.99%
	6	36	6.16	5.45	-11.41%
	8	36	6.39	6.02	-5.71%
	10	36	5.79	5.85	1.00%
	12	36	5.91	6.13	3.64%

Packed Ring performance impact

Selected test cases were re-run to show benefits of using Packed Rings as an option when configuring SPDK Vhost BLK controllers. For this, an optional parameter “—packed_ring” must be used when creating a SPDK Vhost BLK controller. Packed Ring feature requires QEMU 4.2.0 or later.

Table 9: Packed Ring performance impact on SPDK Vhost BLK controllers.

Workload	# of CPU cores	# of VMs	IOPS (millions) Split Ring	IOPS (millions) Packed Ring	Avg. Latency (usec) Split Ring	Avg. Latency (usec) Packed Ring	Packed Ring IOPS impact (%)	Packed Ring Avg. Latency impact (%)
4KB 100% Random Read QD=64	1	6	1.70	1.70	225.95	224.98	0.39%	-0.43%
	2	12	2.96	3.06	258.95	250.94	3.25%	-3.09%
	4	24	4.52	4.56	337.35	338.58	1.07%	0.36%
	6	36	5.89	6.04	388.98	382.18	2.52%	-1.75%
	8	36	7.03	7.29	324.95	317.65	3.67%	-2.25%
	10	36	7.15	7.34	325.02	312.99	2.59%	-3.70%
	12	36	7.55	7.70	305.13	297.57	2.04%	-2.48%
4KB 100% Random Write QD=32	1	6	1.65	1.54	117.52	128.27	-6.28%	9.14%
	2	12	3.01	3.06	128.71	124.45	1.89%	-3.31%
	4	24	4.92	4.87	153.90	156.89	-1.15%	1.94%
	6	36	5.99	5.99	193.27	193.17	0.02%	-0.05%
	8	36	7.03	7.00	163.14	163.95	-0.46%	0.50%
	10	36	6.61	6.84	173.47	168.90	3.47%	-2.64%
	12	36	7.23	7.22	159.37	158.07	-0.23%	-0.82%
4KB 70% Random Read 30% Random Write QD=64	1	6	1.61	1.60	239.86	243.80	-0.68%	1.65%
	2	12	2.83	2.92	270.65	263.78	3.04%	-2.54%
	4	24	4.59	4.61	334.80	334.69	0.50%	-0.03%
	6	36	5.88	5.92	392.16	388.22	0.82%	-1.01%
	8	36	6.57	6.64	348.40	346.41	1.03%	-0.57%
	10	36	6.38	6.36	366.32	361.49	-0.26%	-1.32%
	12	36	6.63	6.64	346.82	345.17	0.18%	-0.47%

Conclusions

1. For SPDK Vhost SCSI performance with split NVMe bdevs, we measured 1.70 million IOPS on one Vhost core for the 4KB 100% Random Write workload. The single Vhost core IOPS for 4 KB Random Read and 4KB Random 70/30 Read/Write were 1.90 million and 1.74 million IOPS respectively. For all workloads, the IOPS scaled near linearly with addition of I/O processing cores up to 6 CPU cores. Peak performance was achieved at 8 CPU cores for all workloads. Further increasing the number of cores does not result in performance improvement or it is not significant.
2. For SPDK Vhost SCSI with Logical Volume backend devices, we measured about 1.5 million IOPS on one Vhost core for all 3 workloads. Performance scaled near linearly with addition of I/O processing cores up to 8 CPU cores for Random Read and up to 6 CPU cores for other workloads. Increasing the number of I/O processing cores further results in non-linear IOPS gains. Peak performance is reached at 12 CPU cores for Random Read and Random Read/Write workloads and at 8 CPU cores for Random Write workload.
3. For SPDK Vhost BLK with Logical Volume backend devices, we measured about 1.6 million IOPS on one Vhost core for all workloads. Performance scaled near linearly up to 8 CPU cores for Random Read and up to 6 CPU cores for other workloads. Increasing the number of cores improves performance further, but the gains are not linear.
4. Using Logical Volumes has an impact of 5-17% lower IOPS when Vhost saturates the I/O cores at 6 or less CPUs. Further increasing SPDK Vhost CPU cores allow Logical Volumes to perform better and their performance is on par with Split NVMe Bdevs (less than 10% difference).
5. For some workloads there is a slight performance drop when Vhost is run with 10 or 12 CPU cores. The platform has 80 CPU threads available, and 36 VMs use 72 CPU threads. Therefore, when 10 or 12 are used for the Vhost process there is not enough left to accommodate all the VMs. Some of the VMs share CPU threads, thus becoming less efficient.
6. Using Packed Ring option instead of default Split Ring mode for SPDK Vhost BLK controllers results in minor performance improvement.
7. Using vfio-user along with vfio-user-pci devices attached to Qemu guest VMs results in 20-60% lower performance than using SPDK vhost-user-scsi-pci. Delta between SPDK Vhost SCSI and SPDK vfio-user decreases with addition of CPU cores to target process.
8. For 4k Random Read and 4k Random Write workload SPDK vfio-user performance trend is rising with addition of CPU cores to target process. This suggests that IOPS throughput at some point might reach the same values as for SPDK Vhost stack. However, we were unable to verify this, as there were no free CPU cores left in the system. Increasing the number of CPU cores for Target would have to be done at expense of number of running Qemu guest VMs, which would lower the overall performance.

Test Case 2: Rate Limiting IOPS per VM

This test case was geared towards understanding how many VMs can be supported at a pre-defined Quality of Service of IOPS per Vhost device. Both read and write IOPS were rate limited for each Vhost device on each of the VMs and then VM density was compared between SPDK & the Linux Kernel. 10K IOPS per VM were chosen as the rate limiter using Linux cgroups2.

Each individual VM was running FIO with the following workloads:

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

The results in tables are average of 3 runs.

Table 10: Rate Limiting IOPS per VM test case configuration

Item	Description
Test case	Test rate limiting IOPS/VM to 10000 IOPS
Test configuration	<p>FIO Version: fio-3.19</p> <p>VM Configuration:</p> <ul style="list-style-type: none"> Common settings are described in the Virtual Machine Settings chapter. Number of VMs: 24, 48, and 72 Each VM has a single Vhost device which is one of equal partitions of NVMe drive. Total number of partitions depends on run test case. <ul style="list-style-type: none"> For 24 VMs: 24xNVMe * 1 partition per NVMe = 24 partitions For 48 VMs: 24xNVMe * 2 partitions per NVMe = 48 partitions For 72 VMs: 24xNVMe * 3 partitions per NVMe = 72 partitions Devices on VMs were throttled to run at a maximum of 10k IOPS (read or write) <p>SPDK Vhost target configuration:</p> <ul style="list-style-type: none"> Test were run with both Vhost-scsi and Vhost-blk stacks. The Vhost-scsi stack was run with Split NVMe bdevs and Logical Volume bdevs. The Vhost-blk stack was run with Logical Volume bdevs. Test were run with the Vhost target using 4 CPU cores (NUMA optimized). <p>Kernel Vhost-scsi configuration:</p> <ul style="list-style-type: none"> Cgroups were used to limit the Vhost process to 4 cores. NUMA optimizations were not explored.
FIO configuration run on each VM	<pre>[global] ioengine=libaio direct=1 rw=randrw</pre>

	<pre>rwmixread=100 (100% reads), 0 (100% writes) thread=1 norandommap=1 time_based=1 runtime=300s ramp_time=10s bs=4k iodepth=1 numjobs=1</pre>
--	---

Test Case 2 Results

Table 11: 4KB 100% Random Reads QD=1 rate limiting test results

# of VMs	Stack	Backend bdev	IOPS (k)	Avg Lat. (usec)
24 VMs	SPDK-SCSI	Split NVMe	239.99	99.46
	SPDK-SCSI	Logical Volume	239.95	99.48
	SPDK-BLK	Logical Volume	239.98	99.46
	Kernel-SCSI	Partitioned NVMe	212.88	113.12
48 VMs	SPDK-SCSI	Split NVMe	479.83	99.13
	SPDK-SCSI	Logical Volume	479.82	99.13
	SPDK-BLK	Logical Volume	479.89	99.15
	Kernel-SCSI	Partitioned NVMe	208.37	230.40
72 VMs	SPDK-SCSI	Split NVMe	683.71	104.32
	SPDK-SCSI	Logical Volume	676.17	105.48
	SPDK-BLK	Logical Volume	690.80	103.21
	Kernel-SCSI	Partitioned NVMe	310.58	231.56

Figure 5: 4KB 100% Random Reads IOPS, QD=1, throttling = 10k IOPS

Table 12: 4KB 100% Random Writes QD=1 rate limiting test results

# of VMs	Stack	Backend bdev	IOPS (k)	Avg Lat. (usec)
24 VMs	SPDK-SCSI	Split NVMe	239.98	99.48
	SPDK-SCSI	Logical Volume	240.00	99.46
	SPDK-BLK	Logical Volume	239.99	99.49
	Kernel-SCSI	Partitioned NVMe	239.96	99.35
48 VMs	SPDK-SCSI	Split NVMe	479.99	99.25
	SPDK-SCSI	Logical Volume	479.98	99.25
	SPDK-BLK	Logical Volume	479.98	99.33
	Kernel-SCSI	Partitioned NVMe	253.79	189.43
72 VMs	SPDK-SCSI	Split NVMe	719.93	99.09
	SPDK-SCSI	Logical Volume	719.93	99.08
	SPDK-BLK	Logical Volume	719.96	99.11
	Kernel-SCSI	Partitioned NVMe	334.66	213.46

Figure 6: 4KB 100% Random Writes IOPS, QD=1, throttling = 10k IOPS

Conclusions

1. Using just 4 I/O processing cores, the SPDK vhost served 10,000 IOPS/VM to up to 48 VMs for 4 KB random read workload and 72 VMs for the 4 KB random write workload.
2. The Kernel Vhost was not able to serve IO at 10K IOPS/VM with just 4 I/O processing cores for the 4KB random read workload. For the 4K Rand write the Kernel Vhost target served 10K IOPS/VM to up to 24 VMs
3. Average latencies were up to 2.3x times better for Random Read and up to 2.1x times better for Random Write workloads with the SPDK Vhost when compared to Kernel Vhost.

Test Case 3: Performance per NVMe drive

This test case was performed in order to understand performance and efficiency of the Vhost scsi and blk process using SPDK vs. Linux Kernel with a single NVMe drive on 2 VMs. Each VM had a single Vhost device which is one of two equal partitions of an NVMe drive. Results in the table represent performance (IOPS, avg. latency & CPU %) seen from the VM. The VM was running FIO with the following workloads:

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

The results in tables are average of 3 runs.

Table 13: Performance per NVMe drive test case configuration

Item	Description
Test case	Test SPDK Vhost target I/O core scaling performance
Test configuration	<p>FIO Version: fio-3.19</p> <p>VM Configuration:</p> <ul style="list-style-type: none"> Common settings are described in the Virtual Machine Settings chapter. 2 VMs were tested Each VM had a single Vhost device which was one of two equal partitions of a single NVMe drive. <p>SPDK Vhost target configuration:</p> <ul style="list-style-type: none"> The SPDK Vhost process was run on a single, physical CPU core. The Vhost-scsi stack was run with Split NVMe bdevs and Logical Volume bdevs. The Vhost-blk stack was run with Logical Volume bdevs. <p>Kernel Vhost target configuration:</p> <ul style="list-style-type: none"> The Vhost process was run on a single, physical CPU core using cgroups.
FIO configuration	<pre>[global] ioengine=libaio direct=1 rw=randrw rwmixread=100 (100% reads), 70 (70% reads, 30% writes), 0 (100% writes) thread=1 norandommap=1 time_based=1 runtime=240s ramp_time=60s bs=4k iodepth=1 / 8 / 32 / 64 numjobs=1</pre>

Test Case 3 results

SPDK Vhost-Scsi

Table 14: Performance per NVMe drive IOPS and latency results, SPDK SCSI stack

Access pattern	Backend	QD	Throughput (IOPS)	Avg. latency (usec)
4k 100% Random Reads	Split NVMe	1	24940.71	79.84
4k 100% Random Reads	Split NVMe	8	171601.30	93.03
4k 100% Random Reads	Split NVMe	32	443223.44	144.05
4k 100% Random Reads	Split NVMe	64	566021.45	225.65
4k 100% Random Reads	Lvol	1	24917.84	79.92
4k 100% Random Reads	Lvol	8	171080.44	93.05
4k 100% Random Reads	Lvol	32	444044.27	143.74
4k 100% Random Reads	Lvol	64	569534.07	224.19
4k 100% Random Writes	Split NVMe	1	113847.01	17.41
4k 100% Random Writes	Split NVMe	8	336923.32	49.28
4k 100% Random Writes	Split NVMe	32	458691.25	139.76
4k 100% Random Writes	Split NVMe	64	466782.14	274.04
4k 100% Random Writes	Lvol	1	113099.42	17.63
4k 100% Random Writes	Lvol	8	336316.82	49.40
4k 100% Random Writes	Lvol	32	451963.46	142.42
4k 100% Random Writes	Lvol	64	457748.12	278.65
4k 70%/30% Random Read Writes	Split NVMe	1	31510.92	63.17
4k 70%/30% Random Read Writes	Split NVMe	8	168155.90	95.17
4k 70%/30% Random Read Writes	Split NVMe	32	331638.40	192.65
4k 70%/30% Random Read Writes	Split NVMe	64	377127.78	338.46
4k 70%/30% Random Read Writes	Lvol	1	31529.15	63.12
4k 70%/30% Random Read Writes	Lvol	8	167909.44	96.14
4k 70%/30% Random Read Writes	Lvol	32	320177.18	202.05
4k 70%/30% Random Read Writes	Lvol	64	409097.92	315.05

SPDK Vhost-Blk

Table 15: Performance per NVMe drive IOPS and latency results, SPDK BLK stack

Access pattern	Backend	QD	Throughput (IOPS)	Avg. latency (usec)
4k 100% Random Reads	Lvol	1	25139.99	79.21
4k 100% Random Reads	Lvol	8	173452.74	92.10
4k 100% Random Reads	Lvol	32	448502.13	142.43
4k 100% Random Reads	Lvol	64	575938.33	221.88
4k 100% Random Writes	Lvol	1	120640.92	15.98
4k 100% Random Writes	Lvol	8	339365.27	48.95
4k 100% Random Writes	Lvol	32	449096.60	142.16
4k 100% Random Writes	Lvol	64	457994.11	279.89
4k 70%/30% Random Read Writes	Lvol	1	31647.84	63.01
4k 70%/30% Random Read Writes	Lvol	8	181488.48	87.49
4k 70%/30% Random Read Writes	Lvol	32	340726.11	188.52
4k 70%/30% Random Read Writes	Lvol	64	374466.14	339.86

Kernel Vhost-Scsi

Table 16: Performance per NVMe drive IOPS and latency results, Kernel Vhost-Scsi

Access pattern	Backend	QD	Throughput (IOPS)	Avg. latency (usec)
4k 100% Random Reads	NVMe	1	20721.08	95.89
4k 100% Random Reads	NVMe	8	146840.57	108.23
4k 100% Random Reads	NVMe	32	368601.43	173.10
4k 100% Random Reads	NVMe	64	440877.00	291.53
4k 100% Random Writes	NVMe	1	81459.67	24.06
4k 100% Random Writes	NVMe	8	227738.93	68.61
4k 100% Random Writes	NVMe	32	369800.71	173.80
4k 100% Random Writes	NVMe	64	449037.82	284.09
4k 70%/30% Random Read Writes	NVMe	1	31647.84	63.01
4k 70%/30% Random Read Writes	NVMe	8	181488.48	87.49
4k 70%/30% Random Read Writes	NVMe	32	340726.11	188.52
4k 70%/30% Random Read Writes	NVMe	64	374466.14	339.86

Figure 7: 4KB 100% Random Reads IOPS and latency

Figure 8: 4KB 100% Random Writes IOPS and latency

Figure 9: 4KB 70%/30% Random Read/Write IOPS and latency

Conclusions

1. SPDK Vhost-scsi with NVMe Split bdevs has lower latency and higher throughput than Kernel Vhost-scsi in majority of workload / queue depth combinations.

Summary

This report compared performance results while running Vhost-scsi using traditional interrupt-driven kernel Vhost-scsi against the accelerated polled-mode driven SPDK implementation. Various local ephemeral configurations were demonstrated, including rate limiting IOPS, performance per VM, and maximum performance from an underlying system when comparing kernel vs. SPDK Vhost-scsi target implementations.

In addition, performance impacts of using SPDK Logical Volume Bdevs and the SPDK Vhost-blk stack were presented.

This report provided information regarding methodologies and practices while benchmarking Vhost-scsi and Vhost-blk using both SPDK and 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 - Vfiio-User Performance

Vfio-user is a framework that allows implementing PCI devices in userspace. Clients, such as Qemu, talk to the server over a UNIX socket using vfio-user protocol. For more information see: [libvfio-user](#).

SPDK has added a vfio-user transport to the nvme target. This allows presenting a virtual NVMe PCIe device to a client, such as a Qemu VM. This is analogous to how the SPDK vhost target presents virtual virtio PCIe devices to a Qemu VM.

The performance comparison between SPDK Vhost stacks versus vfio-user transport layer was done using [oracle/qemu fork](#) (vfio-user-v0.93 branch), which supports running Qemu guest VMs with “vfio-user-pci” devices.

In SPDK v22.01, the vfio-user transport triggers an interrupt (vfio_irq_trigger) for every NVMe completion – it has no ability to coalesce interrupts. This is the primary reason why the vfio-user results collected for v22.01 are significantly lower than vhost-blk. The vhost-blk target is able to coalesce interrupts, resulting in much higher performance. Eliminating unnecessary interrupts reduces interrupt overhead in the guest VM, but also reduces system call overhead in the SPDK target application. The SPDK v22.05 release will use the NVMe queue pair’s cq_head to identify opportunities to coalesce interrupts – preliminary testing shows that the performance puts NVMe vfio-user on par with vhost-blk.

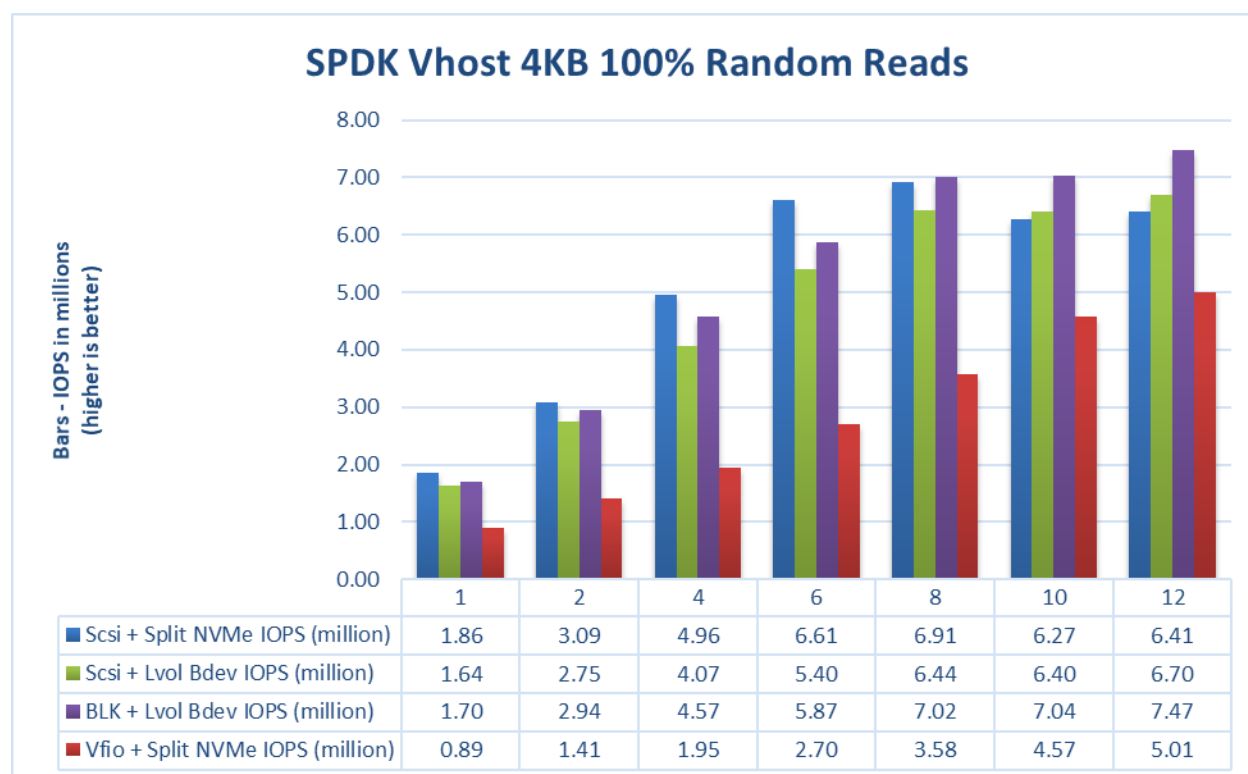


Figure 10: Comparison of performance between various SPDK Vhost stack-bdev combinations and Vfiio-User transport for 4KB Random Read QD=64 workload

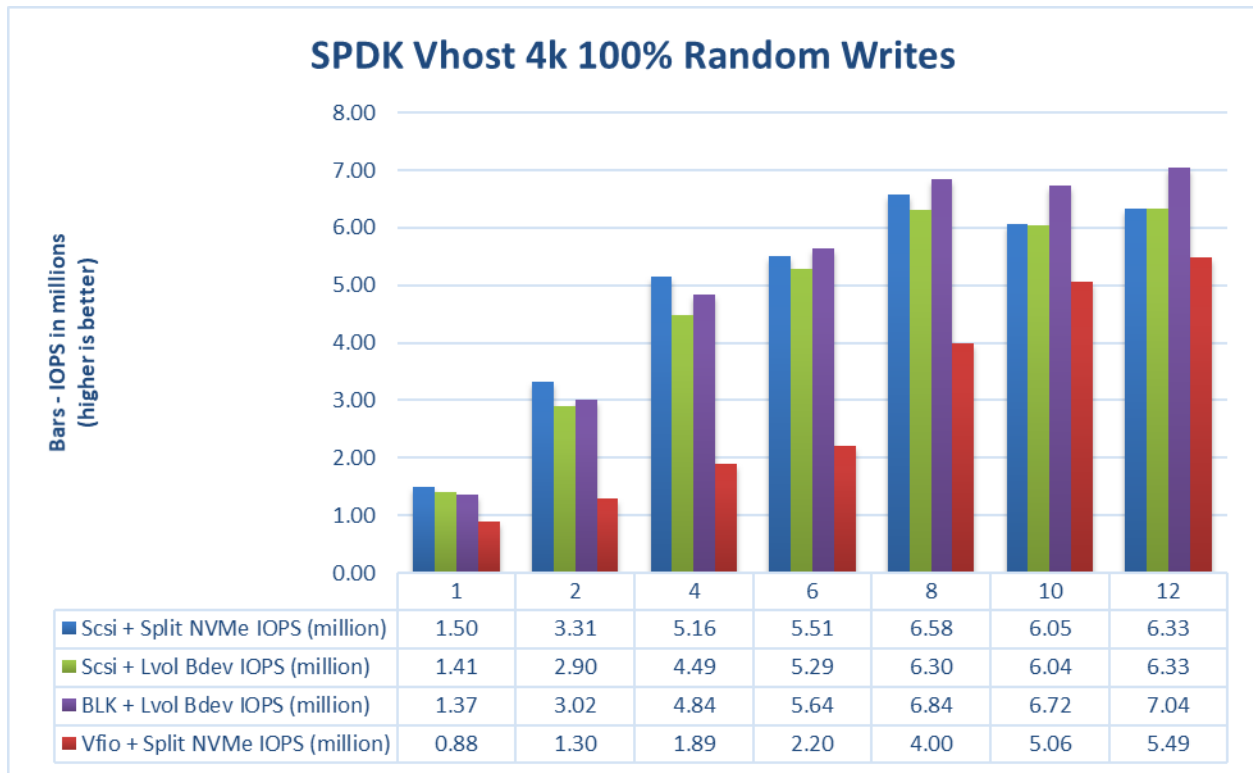


Figure 11: Comparison of performance between various SPDK Vhost stack-bdev combinations and Vfiio-User transport for 4KB Random Write QD=64 workload

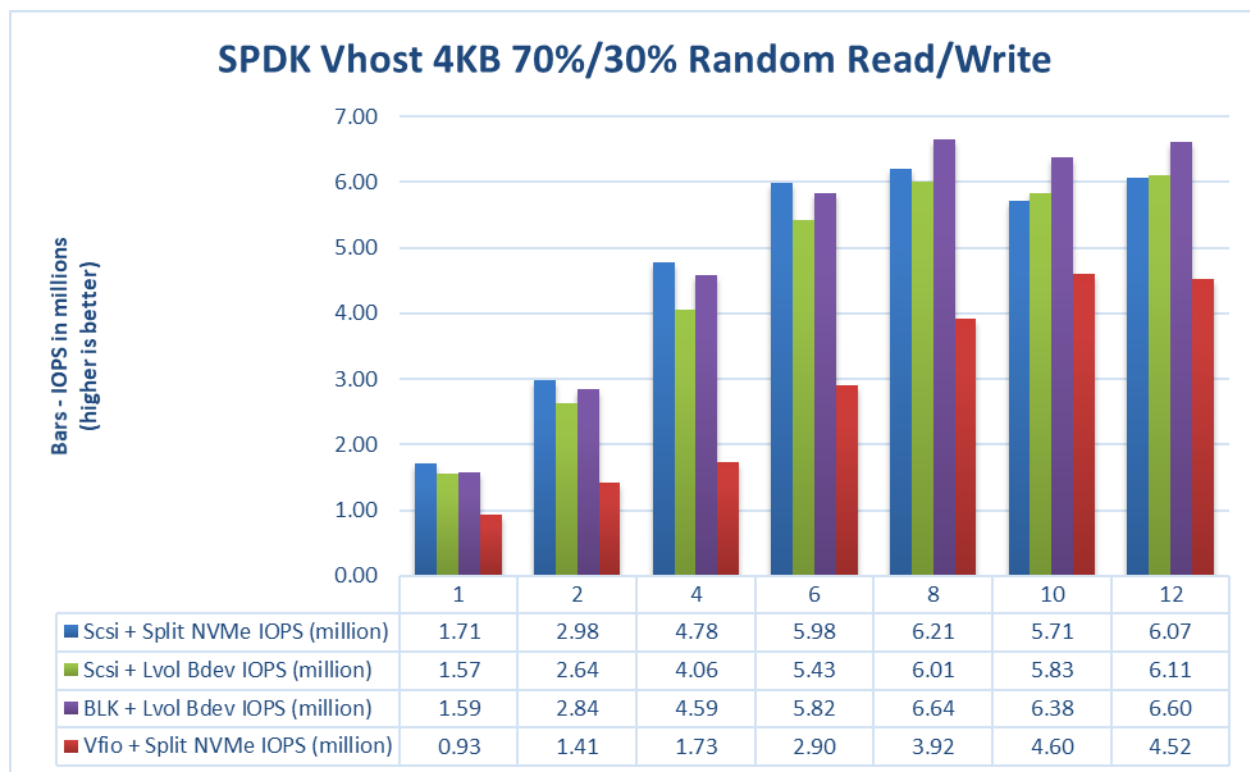


Figure 12: Comparison of performance between various SPDK Vhost stack-bdev combinations and Vfiio-User transport for 4KB Random 70% Read / 30% Write QD=64 workload

Table 17: SPDK Vhost-SCSI vs SPDK Vfiio-user transport layer performance

Workload	# of CPU cores	# of VMs	IOPS (millions) SCSI w. Split NVMe	IOPS (millions) vfiio-user w. Split NVMe	Avg. Latency (usec) SCSI w. Split NVMe	Avg. Latency (usec) vfiio-user w. Split NVMe	Vfiio-user vs. SCSI IOPS (%)	Vfiio-user vs. SCSI Avg. Latency impact (%)
4KB 100% Random Read QD=64	1	6	1.86	0.89	206.54	429.16	-51.79%	107.79%
	2	12	3.09	1.41	248.06	593.00	-54.30%	139.06%
	4	24	4.96	1.95	310.43	774.34	-60.66%	149.44%
	6	36	6.61	2.70	349.21	857.51	-59.16%	145.56%
	8	36	6.91	3.58	332.92	664.65	-48.17%	99.65%
	10	36	6.27	4.57	365.99	505.05	-27.06%	38.00%
	12	36	6.41	5.01	358.37	466.14	-21.78%	30.07%
4KB 100% Random Write QD=64	1	6	1.50	0.88	264.40	128.27	-41.27%	64.79%
	2	12	3.31	1.30	232.34	435.71	-60.70%	154.95%
	4	24	5.16	1.89	295.56	592.36	-63.36%	177.32%
	6	36	5.51	2.20	417.34	819.62	-60.02%	150.24%
	8	36	6.58	4.00	349.51	1044.35	-39.19%	69.38%
	10	36	6.05	5.06	380.13	591.99	-16.45%	19.40%
	12	36	6.33	5.49	362.84	453.86	-13.33%	16.28%
4KB 70% Random Read 30% Random Write QD=64	1	6	1.71	0.93	221.30	413.06	-45.75%	86.65%
	2	12	2.98	1.41	259.80	548.11	-52.62%	110.97%
	4	24	4.78	1.73	321.51	889.44	-63.82%	176.64%
	6	36	5.98	2.90	383.89	803.16	-51.57%	109.21%
	8	36	6.21	3.92	369.70	586.11	-36.79%	58.54%
	10	36	5.71	4.60	403.00	499.23	-19.46%	23.88%
	12	36	6.07	4.52	378.22	508.51	-25.53%	34.45%

List of Tables

Table 1: Hardware setup configuration	4
Table 2: Test platform BIOS settings	5
Table 3: Guest VM configuration	5
Table 4: SPDK Vhost Core Scaling test configuration	8
Table 5: SPDK Vhost core scaling results, 4KB 100% Random Reads IOPS, QD=64	10
Table 6: SPDK Vhost core scaling results, 4KB 100% Random Write IOPS, QD=32	11
Table 7: SPDK Vhost core scaling results, 4KB Random 70% Read 30% Write IOPS, QD=64	12
Table 8: Logical Volumes performance impact for SPDK Vhost SCSI	13
Table 9: Packed Ring performance impact on SPDK Vhost BLK controllers.	14
Table 10: Rate Limiting IOPS per VM test case configuration	16
Table 11: 4KB 100% Random Reads QD=1 rate limiting test results.....	18
Table 12: 4KB 100% Random Writes QD=1 rate limiting test results	19
Table 13: Performance per NVMe drive test case configuration	21
Table 14: Performance per NVMe drive IOPS and latency results, SPDK SCSI stack	22
Table 15: Performance per NVMe drive IOPS and latency results, SPDK BLK stack	23
Table 16: Performance per NVMe drive IOPS and latency results, Kernel Vhost-Scsi	23
Table 17: SPDK Vhost-SCSI vs SPDK Vfiio-user transport layer performance.....	30

List of Figures

Figure 1: SPDK Vhost-scsi architecture	6
Figure 2: Comparison of performance between various SPDK Vhost stack-bdev combinations for 4KB Random Read QD=64 workload	10
Figure 3: Comparison of performance between various SPDK Vhost stack-bdev combinations for 4KB Random Write QD=32 workload	11
Figure 4: Comparison of performance between various SPDK Vhost stack-bdev combinations for 4KB Random 70% Read 30% Write QD=64 workload	12
Figure 5: 4KB 100% Random Reads IOPS, QD=1, throttling = 10k IOPS	18
Figure 6: 4KB 100% Random Writes IOPS, QD=1, throttling = 10k IOPS	19
Figure 7: 4KB 100% Random Reads IOPS and latency	24
Figure 8: 4KB 100% Random Writes IOPS and latency	24
Figure 9: 4KB 70%/30% Random Read/Write IOPS and latency	25
Figure 10: Comparison of performance between various SPDK Vhost stack-bdev combinations and Vfiio-User transport for 4KB Random Read QD=64 workload	27
Figure 11: Comparison of performance between various SPDK Vhost stack-bdev combinations and Vfiio-User transport for 4KB Random Write QD=64 workload	28
Figure 12: Comparison of performance between various SPDK Vhost stack-bdev combinations and Vfiio-User transport for 4KB Random 70% Read / 30% Write QD=64 workload	29

Notices & Disclaimers

Performance varies by use, configuration and other factors. Learn more at [www.Intel.com/PerformanceIndex](https://www.intel.com/PerformanceIndex).

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.

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.