**intel.**

# SPDK Vhost Performance Report Release 21.04

**Testing Date:** May 2021

**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*

# *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** | Intel WolfPass **R2224WFTZS**<br><br><br><br>Server board **S2600WFT**<br><br> |
| **Motherboard** | S2600WFT |
| **CPU** | 2 CPU sockets, Intel(R) Xeon(R) Gold 6230N CPU @ 2.30GHz<br><br>Number of cores 20 per socket, number of threads 40 per socket<br>Both sockets populated<br>Microcode: 0x4003003 |
| **Memory** | 12 x 32GB Micron DDR4 36ASF4G72PZ-2G9E2<br><br>Total 384 GBs<br><br>Memory channel population: |
| **Operating System** | Fedora 33 |
| **BIOS** | SE5C620.86B.02.01.0013.121520200651 |

| 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 |

intel.

| Linux kernel version | 5.10.19-200.fc33.x86_64 |
|---|---|
| **SPDK version** | SPDK 21.04 |
| **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<br>CPU Power and Performance Policy = <Performance><br>CPU C-state =  No Limit<br>CPU P-state = Enabled<br>Enhanced Intel® Speedstep® Tech = Enabled<br>Turbo Boost = Enabled<br>Hyper Threading = Enabled |

# Virtual Machine Settings

*Table 3: Guest VM configuration*

| Item | Description |
|---|---|
| CPU | 2vCPU, pass through from physical host server.<br>Explicit core usage enforced using "taskset –a –c" command.<br>QEMU arguments for starting the VM:<br>-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.<br>QEMU arguments:<br>-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 | • Multi queue enabled: scsi_mod.use_blk_mq=1<br>• Spectre-meltdown patches disabled: spectre_v2=off nopti |

# Kernel & BIOS Spectre-Meltdown information

Host server system uses 5.10.19 kernel version which is available from the DNF repository. The default Spectre-Meltdown mitigation patches for this kernel version have been left enabled.

The guest VM systems use 5.10.8 kernel version, which is available from the DNF repository. The default Spectre-Meltdown mitigation patches for this kernel version have been disabled on guest systems by adding the following in their /etc/default/grub file:

spectre_v2=off nopti mitigations=off

# *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 team has leveraged existing SPDK SCSI layer, DPDK Vhost library, QEMU Vhost-scsi and Vhost-user functionality in order to create the high performance SPDK userspace Vhost target.

## SPDK Vhost target architecture

QEMU setups Vhost target via UNIX domain socket. The Vhost target transfers data to/from the guest VM via shared memory. 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 on example of virtio-scsi. It should be noted that there is no QEMU intervention during the I/O submission process. 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.
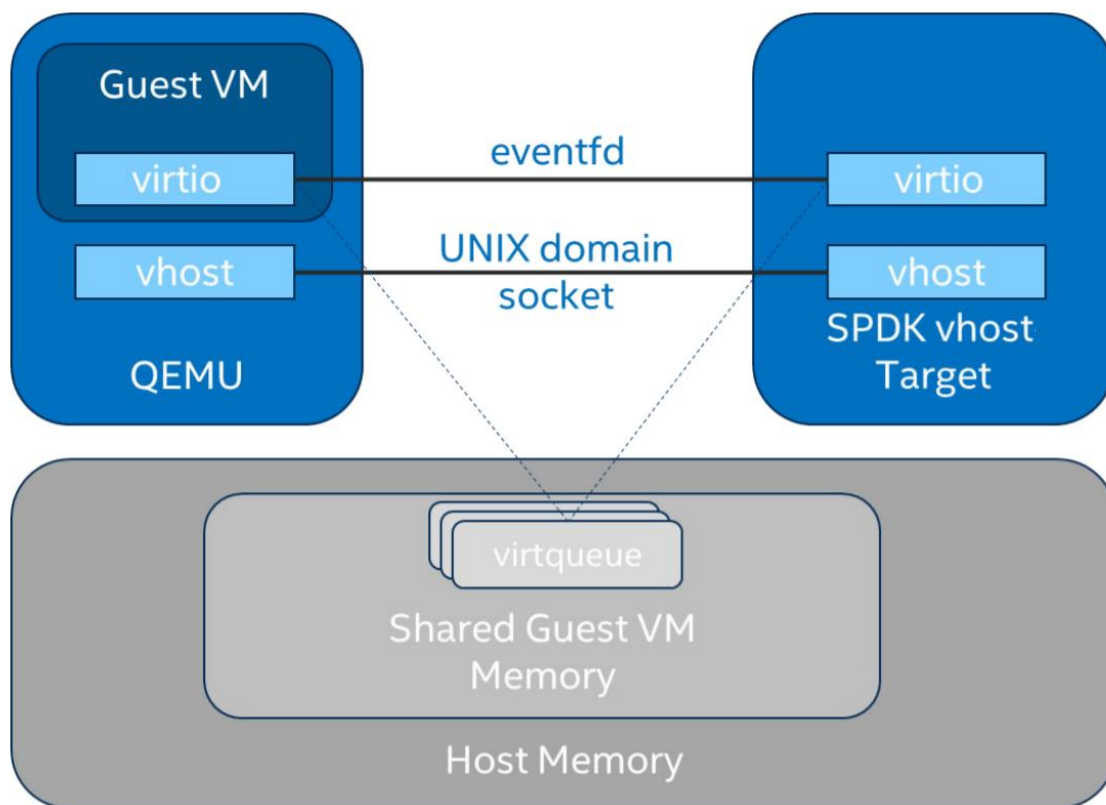


*Figure 1: SPDK Vhost-scsi architecture*

This report shows the performance comparisons between the traditional interrupt-driven kernel Vhost-scsi and the accelerated polled-mode driven 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. The gtod_reduce=1 option was used to disable FIO latency measurements which allowed better IOPS and bandwidth results.

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** | **FIO Version**: fio-3.19<br><br>**VM Configuration**:<br><br>• Common settings are described in the **Virtual Machine Settings** chapter.<br>• Number of VMs: variable (6 VMs per 1 Vhost CPU core, up to 36 VMs max).<br>• 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.<br><br>**SPDK Vhost target configuration:**<br>• Test were run with both the Vhost-scsi and Vhost-blk stacks.<br>• The Vhost-scsi stack was run with Split NVMe bdevs and Logical Volume bdevs.<br>• Vhost-blk stack was run with Logical Volume bdevs.<br>• Tests were ran with 1,2,4,6,8,10 and 12 cores for each stack-bdev combination.<br><br>**Kernel Vhost target configuration:**<br>- N/A |

| FIO configuration | [global]<br>ioengine=libaio<br>direct=1<br>thread=1<br>norandommap=1<br>time_based=1<br>gtod_reduce=1<br>ramp_time=60s<br>runtime=240s<br>numjobs=1<br>bs=4k<br>rw=randrw<br>rwmixread=100 (100% reads), 70 (70% reads, 30% writes), 0 (100% writes)<br>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 | 1.84 |
| | | SCSI / Lvol Bdev | 1.55 |
| | | BLK / Lvol Bdev | 1.55 |
| 2 | 12 | SCSI / Split NVMe Bdev | 3.04 |
| | | SCSI / Lvol Bdev | 2.57 |
| | | BLK / Lvol Bdev | 2.70 |
| 4 | 24 | SCSI / Split NVMe Bdev | 4.92 |
| | | SCSI / Lvol Bdev | 4.13 |
| | | BLK / Lvol Bdev | 4.07 |
| 6 | 36 | SCSI / Split NVMe Bdev | 6.55 |
| | | SCSI / Lvol Bdev | 5.20 |
| | | BLK / Lvol Bdev | 5.58 |
| 8 | 36 | SCSI / Split NVMe Bdev | 7.14 |
| | | SCSI / Lvol Bdev | 6.61 |
| | | BLK / Lvol Bdev | 6.98 |
| 10 | 36 | SCSI / Split NVMe Bdev | 7.18 |
| | | SCSI / Lvol Bdev | 6.76 |
| | | BLK / Lvol Bdev | 7.31 |
| 12 | 36 | SCSI / Split NVMe Bdev | 7.46 |
| | | SCSI / Lvol Bdev | 7.63 |
| | | BLK / Lvol Bdev | 8.35 |



**SPDK Vhost 4KB 100% Random Reads**

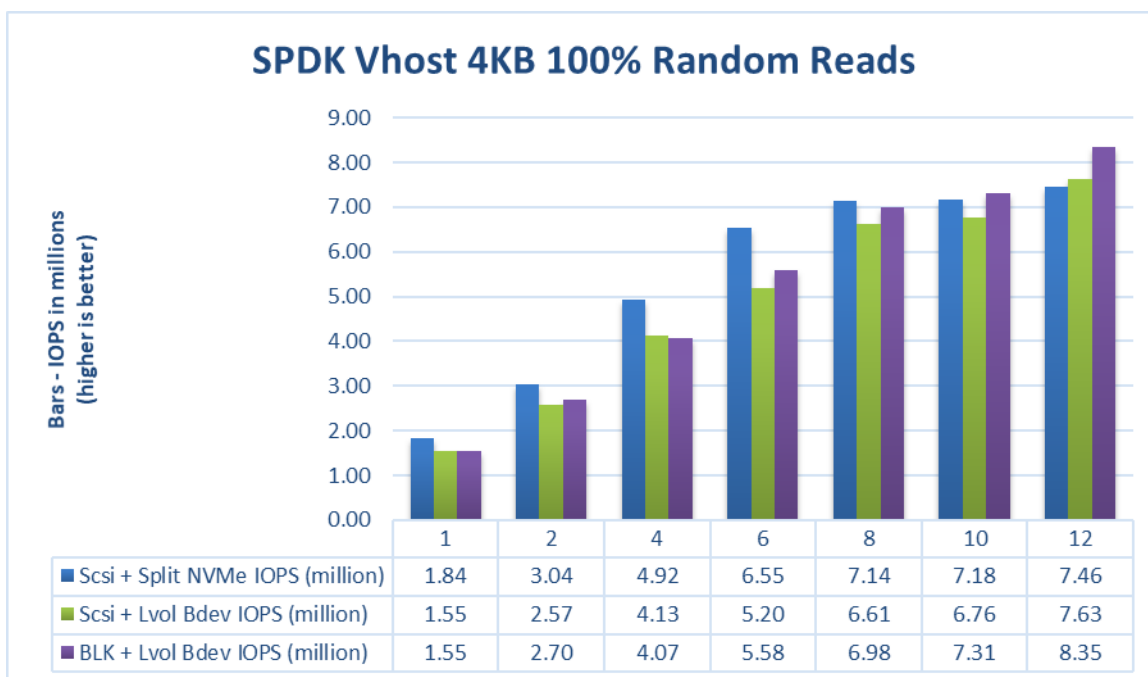| | 1 | 2 | 4 | 6 | 8 | 10 | 12 |
|---|---|---|---|---|---|---|---|
| Scsi + Split NVMe IOPS (million) | 1.84 | 3.04 | 4.92 | 6.55 | 7.14 | 7.18 | 7.46 |
| Scsi + Lvol Bdev IOPS (million) | 1.55 | 2.57 | 4.13 | 5.20 | 6.61 | 6.76 | 7.63 |
| BLK + Lvol Bdev IOPS (million) | 1.55 | 2.70 | 4.07 | 5.58 | 6.98 | 7.31 | 8.35 |

*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*

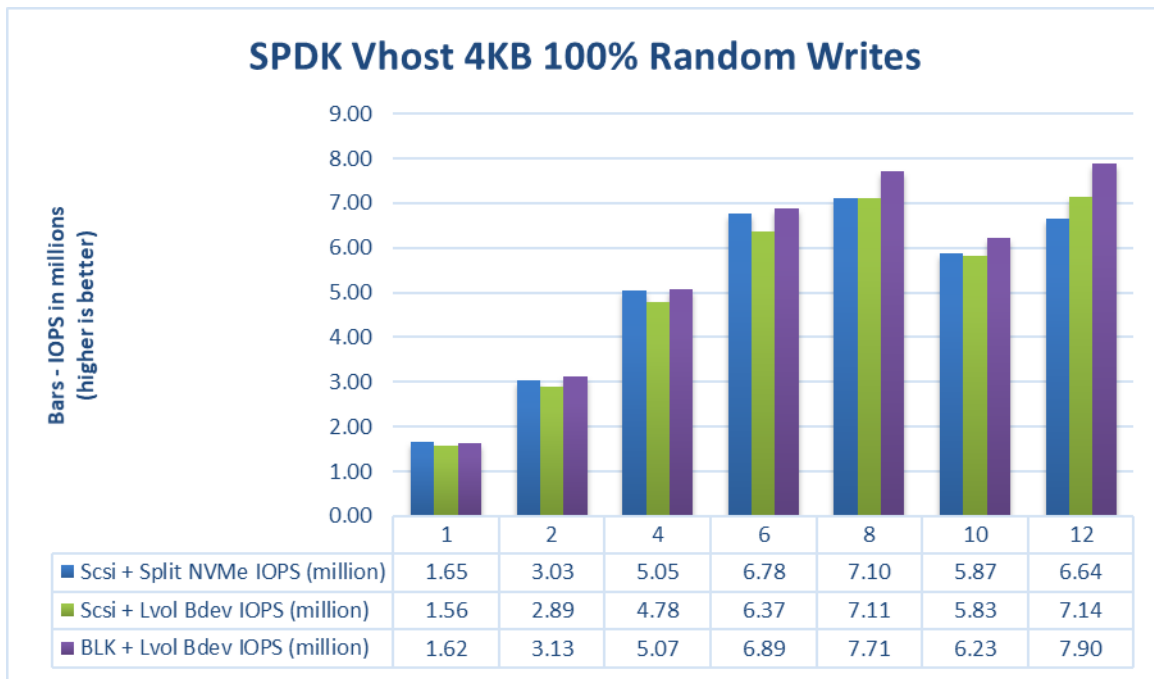| # of CPU cores | # of VMs | Stack / Backend | IOPS (millions) |
|---|---|---|---|
| 1 | 6 | SCSI / Split NVMe Bdev | 1.65 |
| | | SCSI / Lvol Bdev | 1.56 |
| | | BLK / Lvol Bdev | 1.62 |
| 2 | 12 | SCSI / Split NVMe Bdev | 3.03 |
| | | SCSI / Lvol Bdev | 2.89 |
| | | BLK / Lvol Bdev | 3.13 |
| 4 | 24 | SCSI / Split NVMe Bdev | 5.05 |
| | | SCSI / Lvol Bdev | 4.78 |
| | | BLK / Lvol Bdev | 5.07 |
| 6 | 36 | SCSI / Split NVMe Bdev | 6.78 |
| | | SCSI / Lvol Bdev | 6.37 |
| | | BLK / Lvol Bdev | 6.89 |
| 8 | 36 | SCSI / Split NVMe Bdev | 7.10 |
| | | SCSI / Lvol Bdev | 7.11 |
| | | BLK / Lvol Bdev | 7.71 |
| 10 | 36 | SCSI / Split NVMe Bdev | 5.87 |
| | | SCSI / Lvol Bdev | 5.83 |
| | | BLK / Lvol Bdev | 6.23 |
| 12 | 36 | SCSI / Split NVMe Bdev | 6.64 |
| | | SCSI / Lvol Bdev | 7.14 |
| | | BLK / Lvol Bdev | 7.90 |



*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 | 1.64 |
| | | SCSI / Lvol Bdev | 1.51 |
| | | BLK / Lvol Bdev | 1.53 |
| 2 | 12 | SCSI / Split NVMe Bdev | 2.92 |
| | | SCSI / Lvol Bdev | 2.52 |
| | | BLK / Lvol Bdev | 2.69 |
| 4 | 24 | SCSI / Split NVMe Bdev | 5.00 |
| | | SCSI / Lvol Bdev | 4.01 |
| | | BLK / Lvol Bdev | 4.31 |
| 6 | 36 | SCSI / Split NVMe Bdev | 6.19 |
| | | SCSI / Lvol Bdev | 5.22 |
| | | BLK / Lvol Bdev | 5.64 |
| 8 | 36 | SCSI / Split NVMe Bdev | 6.51 |
| | | SCSI / Lvol Bdev | 6.20 |
| | | BLK / Lvol Bdev | 6.59 |
| 10 | 36 | SCSI / Split NVMe Bdev | 6.26 |
| | | SCSI / Lvol Bdev | 6.30 |
| | | BLK / Lvol Bdev | 6.65 |
| 12 | 36 | SCSI / Split NVMe Bdev | 6.53 |
| | | SCSI / Lvol Bdev | 6.42 |
| | | BLK / Lvol Bdev | 7.17 |



SPDK Vhost 4KB 70%/30% Random Read/Write

Bars - IOPS in millions (higher is better)

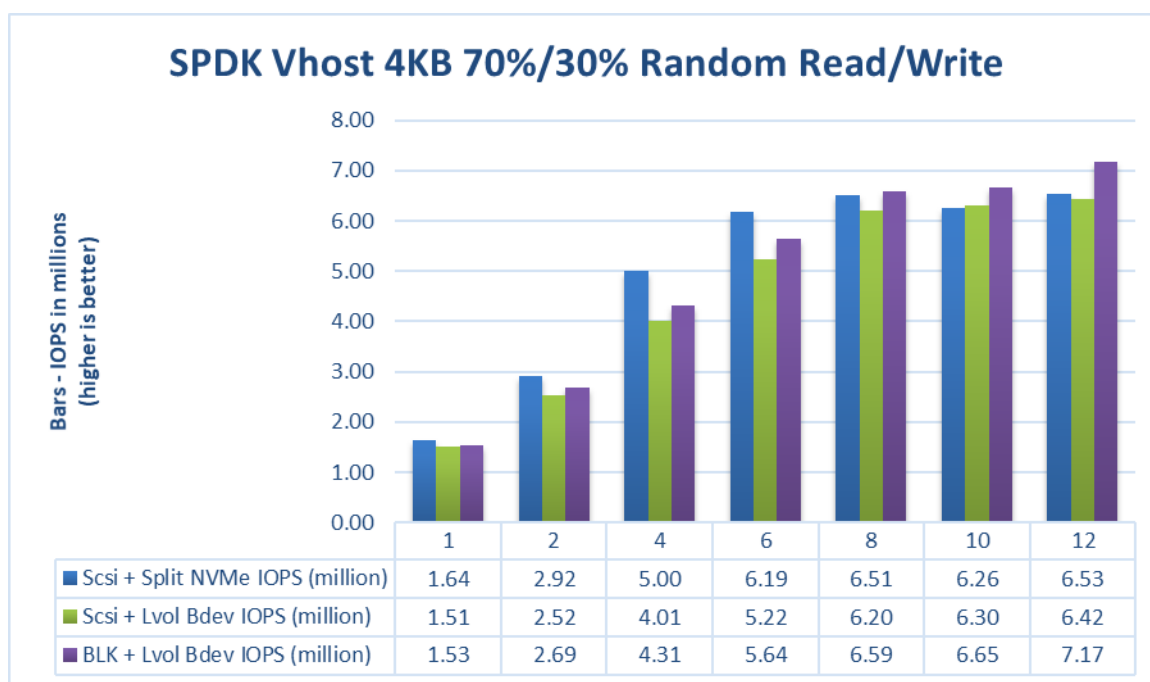| | 1 | 2 | 4 | 6 | 8 | 10 | 12 |
|---|---|---|---|---|---|---|---|
| Scsi + Split NVMe IOPS (million) | 1.64 | 2.92 | 5.00 | 6.19 | 6.51 | 6.26 | 6.53 |
| Scsi + Lvol Bdev IOPS (million) | 1.51 | 2.52 | 4.01 | 5.22 | 6.20 | 6.30 | 6.42 |
| BLK + Lvol Bdev IOPS (million) | 1.53 | 2.69 | 4.31 | 5.64 | 6.59 | 6.65 | 7.17 |

*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 NVMes 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 | 1.84 | 1.55 | -15.89% |
| | 2 | 12 | 3.04 | 2.57 | -15.59% |
| | 4 | 24 | 4.92 | 4.13 | -16.12% |
| | 6 | 36 | 6.55 | 5.20 | -20.69% |
| | 8 | 36 | 7.14 | 6.61 | -7.44% |
| | 10 | 36 | 7.18 | 6.76 | -5.75% |
| | 12 | 36 | 7.46 | 7.63 | 2.23% |
| **4KB 100% Random Write** | 1 | 6 | 1.65 | 1.56 | -5.32% |
| | 2 | 12 | 3.03 | 2.89 | -4.39% |
| | 4 | 24 | 5.05 | 4.78 | -5.35% |
| | 6 | 36 | 6.78 | 6.37 | -6.00% |
| | 8 | 36 | 7.10 | 7.11 | 0.04% |
| | 10 | 36 | 5.87 | 5.83 | -0.71% |
| | 12 | 36 | 6.64 | 7.14 | 7.48% |
| **4KB 70% Random Read 30% Random Write** | 1 | 6 | 1.30 | 1.26 | -3.59% |
| | 2 | 12 | 2.37 | 2.20 | -6.95% |
| | 4 | 24 | 4.27 | 3.80 | -10.92% |
| | 6 | 36 | 5.36 | 4.89 | -8.64% |
| | 8 | 36 | 5.60 | 5.30 | -5.31% |
| | 10 | 36 | 5.41 | 5.23 | -3.31% |
| | 12 | 36 | 5.64 | 5.63 | -0.23% |

# LTO performance impact

Selected test cases were re-run with LTO (Link Time Optimization) enabled for SPDK compilation. This should positively impact overall SPDK performance. The following comparison was done using SPDK Vhost SCSI with Logical Volume bdevs.

*Table 9: LTO performance SPDK Vhost SCSI with Logical Volume bdevs*

| Workload | # of CPU cores | # of VMs | IOPS (millions) LTO Disabled | IOPS (millions) LTO Enabled | LTO Impact (%) |
|---|---|---|---|---|---|
| **4KB 100% Random Read** | 1 | 6 | 1.55 | 1.72 | 11.10% |
| | 2 | 12 | 2.57 | 2.82 | 9.68% |
| | 4 | 24 | 4.13 | 4.63 | 12.28% |
| | 6 | 36 | 5.20 | 5.78 | 11.35% |
| | 8 | 36 | 6.61 | 6.96 | 5.31% |
| | 10 | 36 | 6.76 | 6.84 | 1.06% |
| | 12 | 36 | 7.63 | 7.60 | -0.34% |
| **4KB 100% Random Write** | 1 | 6 | 1.56 | 1.54 | -1.36% |
| | 2 | 12 | 2.89 | 2.86 | -1.10% |
| | 4 | 24 | 4.78 | 4.65 | -2.62% |
| | 6 | 36 | 6.37 | 6.28 | -1.51% |
| | 8 | 36 | 7.11 | 7.03 | -1.11% |
| | 10 | 36 | 5.83 | 5.78 | -0.81% |
| | 12 | 36 | 7.14 | 6.97 | -2.37% |
| **4KB 70% Random Read 30% Random Write** | 1 | 6 | 1.51 | 1.63 | 7.79% |
| | 2 | 12 | 2.52 | 2.73 | 8.03% |
| | 4 | 24 | 4.01 | 4.62 | 15.19% |
| | 6 | 36 | 5.22 | 5.77 | 10.37% |
| | 8 | 36 | 6.20 | 6.40 | 3.09% |
| | 10 | 36 | 6.30 | 6.32 | 0.44% |
| | 12 | 36 | 6.42 | 6.61 | 2.87% |

# 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.

Following results show comparison of running SPDK Vhost-Blk with Packed Ring enabled with fio latency measurements both enabled and disabled. Because other QEMU version was used, base results (Split Ring) were run again to produce a fresh base to compare to.

*Table 10: Packed Ring performance impact on SPDK Vhost BLK controllers. Fio gtod_reduce=disabled*

| 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 | 1.56 | 1.55 | 246.24 | 248.24 | -0.79% | 0.81% |
| | 2 | 12 | 2.73 | 2.77 | 281.06 | 277.43 | 1.36% | -1.29% |
| | 4 | 24 | 4.43 | 4.43 | 346.56 | 342.03 | 0.10% | -1.31% |
| | 6 | 36 | 5.58 | 5.64 | 412.49 | 407.89 | 1.07% | -1.12% |
| | 8 | 36 | 6.83 | 6.87 | 336.42 | 334.51 | 0.50% | -0.57% |
| | 10 | 36 | 7.11 | 7.17 | 331.48 | 326.49 | 0.86% | -1.50% |
| | 12 | 36 | 7.76 | 7.87 | 296.71 | 293.41 | 1.44% | -1.11% |
| **4KB 100% Random Write QD=32** | 1 | 6 | 1.51 | 1.69 | 129.85 | 113.32 | 12.27% | -12.73% |
| | 2 | 12 | 2.87 | 2.94 | 133.74 | 130.86 | 2.20% | -2.15% |
| | 4 | 24 | 4.64 | 4.71 | 162.68 | 161.79 | 1.52% | -0.55% |
| | 6 | 36 | 6.29 | 6.26 | 182.61 | 184.10 | -0.45% | 0.81% |
| | 8 | 36 | 7.11 | 7.15 | 163.28 | 160.56 | 0.49% | -1.66% |
| | 10 | 36 | 5.87 | 7.28 | 199.18 | 159.72 | 24.03% | -19.81% |
| | 12 | 36 | 7.31 | 7.74 | 156.87 | 314.90 | 0.70% | -0.97% |
| **4KB 70% Random Read 30% Random Write QD=64** | 1 | 6 | 1.53 | 1.53 | 247.49 | 250.06 | 0.40% | 1.04% |
| | 2 | 12 | 2.67 | 2.72 | 286.37 | 282.40 | 1.84% | -1.38% |
| | 4 | 24 | 4.47 | 4.49 | 342.24 | 343.03 | 0.51% | 0.23% |
| | 6 | 36 | 5.71 | 5.75 | 402.30 | 400.94 | 0.64% | -0.34% |
| | 8 | 36 | 6.40 | 6.43 | 358.58 | 358.24 | 0.49% | -0.09% |
| | 10 | 36 | 5.97 | 6.12 | 384.85 | 376.29 | 2.54% | -2.23% |
| | 12 | 36 | 6.72 | 6.87 | 340.99 | 335.52 | 2.20% | -1.60% |

# Conclusions

1. For SPDK Vhost SCSI performance with split NVMe bdevs, we measured 1.65 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 1.84 million and 1.64 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. Increasing the number of I/O processing cores further results in non-linear IOPS gains. Peak performance is reached at 12 CPU cores for all workloads.

3. For SPDK Vhost BLK with Logical Volume backend devices, we measured over 1.5 million IOPS on one Vhost core for all workloads. Performance scaled near linearly up to 8 CPU cores. Increasing the number of cores improves performance further, but the gains are not linear.

4. Using Logical Volumes has a noticeable impact on the overall performance. For Vhost tests using 6 or less CPUs (when Vhost is saturated with IO traffic from VMs) performance impact of Logical Volumes is between 10-20%. Further increasing SPDK Vhost CPU cores allow Logical Volumes to perform better and their performance impact is on par with Split NVMe Bdevs (less than 10% difference).

5. LTO compilation option increased SPDK Vhost performance by up to 15% percent in the scaling phase (8 Vhost CPU cores or less) for Random Read and Random Read/Write workloads. Beyond 8 vhost cores, LTO benefit is up to 5% or is negligible. The reason for this behavior is described in point 6. For Random Write workload LTO did not improve performance.

6. 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.

7. Using Packed Ring option instead of default Split Ring mode for SPDK Vhost BLK controllers results in minor performance improvement.

# 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 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 11: Rate Limiting IOPS per VM test case configuration*

| Item | Description |
|---|---|
| **Test case** | Test rate limiting IOPS/VM to 10000 IOPS |
| **Test configuration** | **FIO Version:** fio-3.19<br><br>**VM Configuration**:<br><br>• Common settings are described in the **Virtual Machine Settings** chapter.<br>• Number of VMs: 24 / 48 / 72<br>• 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.<br>    ○ For 24 VMs: 24xNVMe * 1 partition per NVMe = 24 partitions<br>    ○ For 48 VMs: 24xNVMe * 2 partitions per NVMe = 48 partitions<br>    ○ For 72 VMs: 24xNVMe * 3 partitions per NVMe = 72 partitions<br>• Devices on VMs were throttled to run at a maximum of 10k IOPS (read and write)<br><br>**SPDK Vhost target configuration**:<br>• Test were run with both Vhost-scsi and Vhost-blk stacks.<br>• The Vhost-scsi stack was run with Split NVMe bdevs and Logical Volume bdevs.<br>• The Vhost-blk stack was run with Logical Volume bdevs.<br>• Test were run with Vhost using 4 CPU cores (NUMA optimized).<br><br>**Kernel Vhost-scsi configuration:**<br>• Cgroups were used to limit the Vhost process to 4 cores.<br>• NUMA optimization were not explored. |
| **FIO configuration run on each VM** | [global]<br>ioengine=libaio<br>direct=1<br>rw=randrw |

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

# Test Case 2 Results

*Table 12: 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 | 239.96 | 99.39 |
| | SPDK-SCSI | Logical Volume | 239.96 | 99.40 |
| | SPDK-BLK | Logical Volume | 239.97 | 99.40 |
| | Kernel-SCSI | Partitioned NVMe | 169.14 | 141.26 |
| **48 VMs** | SPDK-SCSI | Split NVMe | 479.77 | 99.09 |
| | SPDK-SCSI | Logical Volume | 479.71 | 99.10 |
| | SPDK-BLK | Logical Volume | 479.79 | 99.11 |
| | Kernel-SCSI | Partitioned NVMe | 172.31 | 281.76 |
| **72 VMs** | SPDK-SCSI | Split NVMe | 679.16 | 104.87 |
| | SPDK-SCSI | Logical Volume | 671.36 | 106.10 |
| | SPDK-BLK | Logical Volume | 686.88 | 103.76 |
| | Kernel-SCSI | Partitioned NVMe | 276.60 | 263.51 |



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

*Table 13: 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 | 99.40 |
| | SPDK-SCSI | Logical Volume | 239.99 | 99.38 |
| | SPDK-BLK | Logical Volume | 240.00 | 99.41 |
| | Kernel-SCSI | Partitioned NVMe | 192.39 | 124.80 |
| **48 VMs** | SPDK-SCSI | Split NVMe | 479.98 | 99.24 |
| | SPDK-SCSI | Logical Volume | 479.98 | 99.18 |
| | SPDK-BLK | Logical Volume | 479.99 | 99.26 |
| | Kernel-SCSI | Partitioned NVMe | 175.61 | 274.21 |
| **72 VMs** | SPDK-SCSI | Split NVMe | 719.96 | 99.11 |
| | SPDK-SCSI | Logical Volume | 719.94 | 99.19 |
| | SPDK-BLK | Logical Volume | 719.95 | 99.21 |
| | Kernel-SCSI | Partitioned NVMe | 325.25 | 219.18 |



*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.

3. Average latencies were up to 2.8x times better for Random Read and up to 2.7x times better for Random Write workloads with the SPDK Vhost when compared to Kernel Vhost.

Note: The Kernel-Vhost process was not NUMA-optimized for this scenario.

intel.

# 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 14: Performance per NVMe drive test case configuration*

| Item | Description |
|---|---|
| **Test case** | Test SPDK Vhost target I/O core scaling performance |
| **Test configuration** | **FIO Version:** fio-3.19<br><br>**VM Configuration**:<br><br>• Common settings are described in the **Virtual Machine Settings** chapter.<br>• 2 VMs were tested<br>• Each VM had a single Vhost device which was one of two equal partitions of a single NVMe drive.<br><br>**SPDK Vhost target configuration:**<br>• The SPDK Vhost process was run on a single, physical CPU core.<br>• The Vhost-scsi stack was run with Split NVMe bdevs and Logical Volume bdevs.<br>• The Vhost-blk stack was run with Logical Volume bdevs.<br><br>**Kernel Vhost target configuration**:<br>• The Vhost process was run on a single, physical CPU core using cgroups. |
| **FIO configuration** | [global]<br>ioengine=libaio<br>direct=1<br>rw=randrw<br>rwmixread=100 (100% reads), 70 (70% reads, 30% writes), 0 (100% writes)<br>thread=1<br>norandommap=1<br>time_based=1<br>runtime=240s<br>ramp_time=60s<br>bs=4k<br>iodepth=1 / 8 / 32 / 64<br>numjobs=1 |

# Test Case 3 results

## SPDK Vhost-Scsi

*Table 15: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 | 24644.35 | 80.59 |
| 4k 100% Random Reads | Split NVMe | 8 | 171736.54 | 92.60 |
| 4k 100% Random Reads | Split NVMe | 32 | 444989.30 | 143.54 |
| 4k 100% Random Reads | Split NVMe | 64 | 576174.61 | 221.79 |
| 4k 100% Random Reads | Lvol | 1 | 24670.33 | 80.66 |
| 4k 100% Random Reads | Lvol | 8 | 170383.36 | 93.32 |
| 4k 100% Random Reads | Lvol | 32 | 445830.32 | 143.15 |
| 4k 100% Random Reads | Lvol | 64 | 576085.73 | 222.09 |
| 4k 100% Random Writes | Split NVMe | 1 | 116696.16 | 16.36 |
| 4k 100% Random Writes | Split NVMe | 8 | 413867.52 | 39.00 |
| 4k 100% Random Writes | Split NVMe | 32 | 475329.87 | 134.74 |
| 4k 100% Random Writes | Split NVMe | 64 | 465863.30 | 274.91 |
| 4k 100% Random Writes | Lvol | 1 | 114138.68 | 16.86 |
| 4k 100% Random Writes | Lvol | 8 | 423668.62 | 37.51 |
| 4k 100% Random Writes | Lvol | 32 | 471603.83 | 135.58 |
| 4k 100% Random Writes | Lvol | 64 | 473309.90 | 270.69 |
| 4k 70%/30% Random Read Writes | Split NVMe | 1 | 31672.91 | 62.65 |
| 4k 70%/30% Random Read Writes | Split NVMe | 8 | 169641.70 | 93.46 |
| 4k 70%/30% Random Read Writes | Split NVMe | 32 | 340740.65 | 189.63 |
| 4k 70%/30% Random Read Writes | Split NVMe | 64 | 401866.34 | 320.64 |
| 4k 70%/30% Random Read Writes | Lvol | 1 | 31589.29 | 62.86 |
| 4k 70%/30% Random Read Writes | Lvol | 8 | 160838.33 | 99.22 |
| 4k 70%/30% Random Read Writes | Lvol | 32 | 323809.47 | 197.54 |
| 4k 70%/30% Random Read Writes | Lvol | 64 | 413719.98 | 312.22 |

## SPDK Vhost-Blk

*Table 16: 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 | 24817.79 | 79.91 |
| 4k 100% Random Reads | Lvol | 8 | 172770.59 | 91.97 |
| 4k 100% Random Reads | Lvol | 32 | 450068.44 | 141.71 |
| 4k 100% Random Reads | Lvol | 64 | 584463.26 | 218.98 |
| 4k 100% Random Writes | Lvol | 1 | 107402.19 | 18.28 |
| 4k 100% Random Writes | Lvol | 8 | 411668.40 | 39.40 |
| 4k 100% Random Writes | Lvol | 32 | 473381.59 | 134.99 |
| 4k 100% Random Writes | Lvol | 64 | 470660.92 | 272.56 |
| 4k 70%/30% Random Read Writes | Lvol | 1 | 31670.46 | 62.77 |
| 4k 70%/30% Random Read Writes | Lvol | 8 | 161639.65 | 98.56 |
| 4k 70%/30% Random Read Writes | Lvol | 32 | 343738.67 | 186.13 |
| 4k 70%/30% Random Read Writes | Lvol | 64 | 386573.60 | 332.44 |

## Kernel Vhost-Scsi

*Table 17: 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 | 20753.63 | 91.74 |
| 4k 100% Random Reads | NVMe | 8 | 145636.32 | 109.50 |
| 4k 100% Random Reads | NVMe | 32 | 370743.61 | 172.25 |
| 4k 100% Random Reads | NVMe | 64 | 428334.42 | 298.63 |
| 4k 100% Random Writes | NVMe | 1 | 70074.45 | 27.63 |
| 4k 100% Random Writes | NVMe | 8 | 224220.38 | 71.57 |
| 4k 100% Random Writes | NVMe | 32 | 416000.36 | 154.46 |
| 4k 100% Random Writes | NVMe | 64 | 456087.25 | 280.47 |
| 4k 70%/30% Random Read Writes | NVMe | 1 | 26208.38 | 75.57 |
| 4k 70%/30% Random Read Writes | NVMe | 8 | 144852.53 | 109.94 |
| 4k 70%/30% Random Read Writes | NVMe | 32 | 311034.47 | 205.49 |
| 4k 70%/30% Random Read Writes | NVMe | 64 | 376818.91 | 338.77 |

*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 all 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.

# *List of Tables*

# *List of Figures*

**Notices & Disclaimers**

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

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

§