



Challenges of High-IOPS Enterprise-level NVMeoF-based All Flash Array

From the Viewpoint of Software Vendor

Dr. Weafon Tsao
R&D VP, AccelStor, Inc.

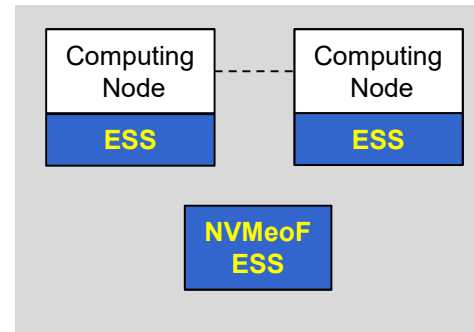
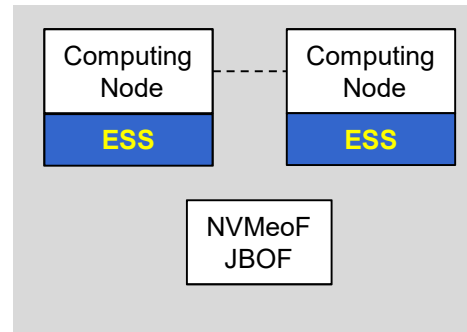
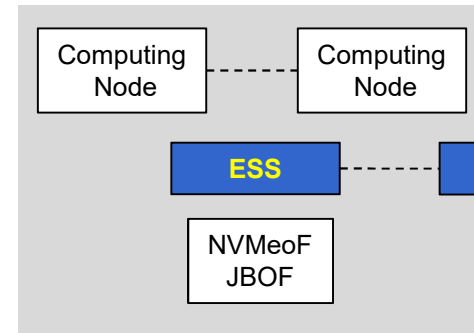
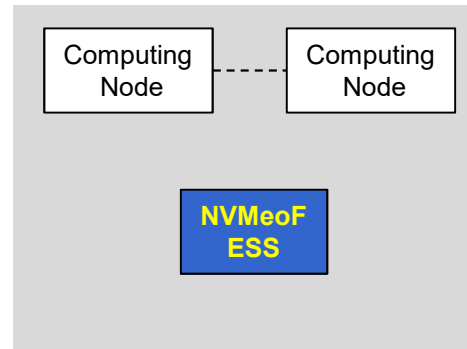
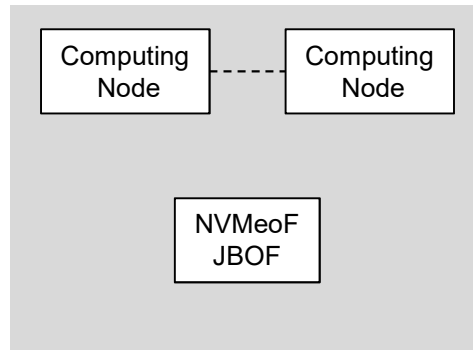


Enterprise-level Storage System (ESS)

- Data Protection (DP)
 - High Availability (HA)
 - Snapshot
 - Replication (Sync/Async)
 - Thin-Provisioning
 - Thin-Clone
 - Deduplication
 - Compression
- Disaster Recovery (DR)
- Data Reduction

NVMeoF-based ESS

Plain Storage System





Why NVMe and NVMf?
Good for Building High-IOPS AFA

Flash Memory Summit 2017
Santa Clara, CA

The photo is designed by fanjianhua / Freepik

Enjoy NVMf for Free!
How you can get a home-made NVMeoF?

Flash Memory Summit 2017
Santa Clara, CA

The photo is designed by mrschmitt / Freepik

Challenges and Solutions

High IOPS with All Features

Flash Memory Summit 2017
Santa Clara, CA

The photo is designed by jigsawstocker / Freepik

Flash Memory Summit 2017
Santa Clara, CA



Why NVMe and NVMeFf?

Good for Building High-IOPS AFA

Flash Memory Summit 2017
Santa Clara, CA

The photo is designed by fanjianhua / Freepik⁵

Dad! Dad!
~ an Interrupt

Flash Memory Summit 2017
Santa Clara, CA

PROFESSOR Robert Kelly's kids invaded a live BBC News interview.

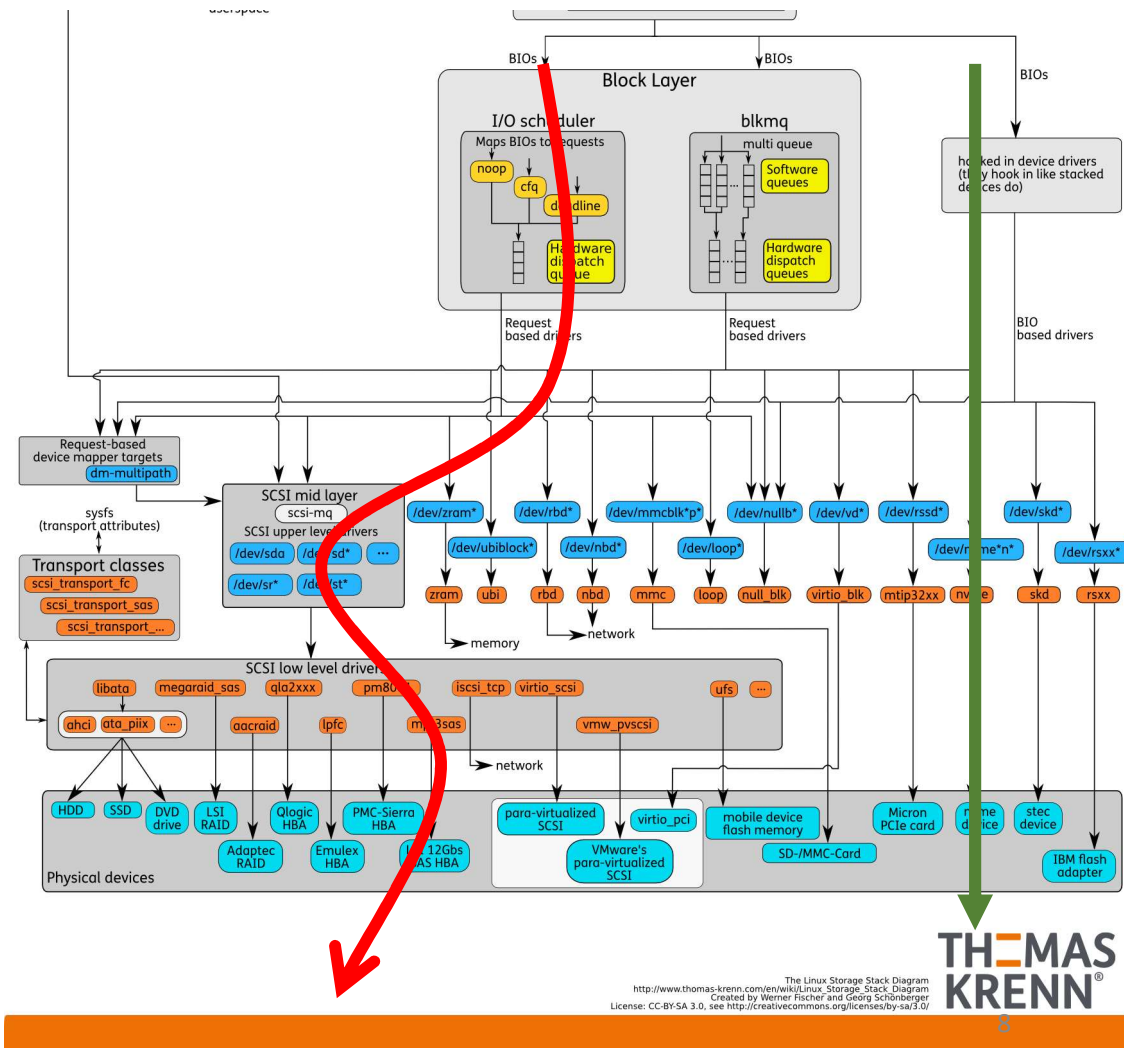


NVMe & NVMeF: Active DMA to Move Data

Legacy & Thick Code Stack



Flash Memory Summit 2017
Santa Clara, CA



The Linux Storage Stack Diagram
Created by Werner Fischer and Georg Schönberger
License: CC-BY-SA 3.0, see <http://creativecommons.org/licenses/by-sa/3.0/>

THOMAS
KRENN®



Enjoy NVMeoF for Free!

How you can get a home-made NVMeoF AFA?

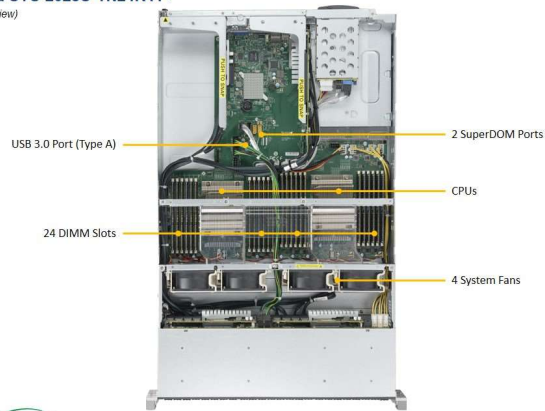
Flash Memory Summit 2017
Santa Clara, CA

The photo is designed by mrsiraphol / Freepik



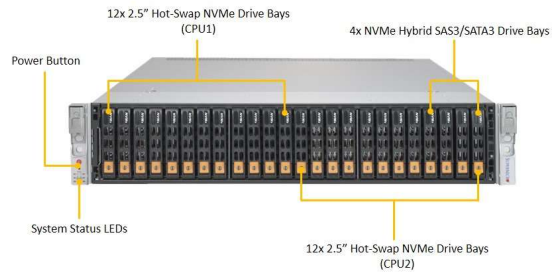
Commodity Hardware

Ultra SYS-2028U-TN24R4T+
(Top View)

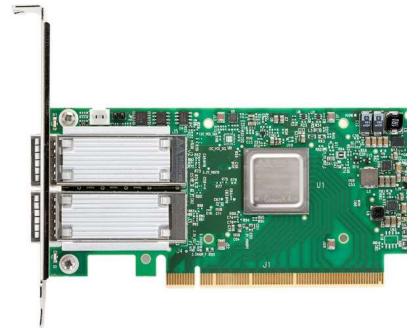


 © Super Micro Computer, Inc. Information in this document is subject to change without notice.

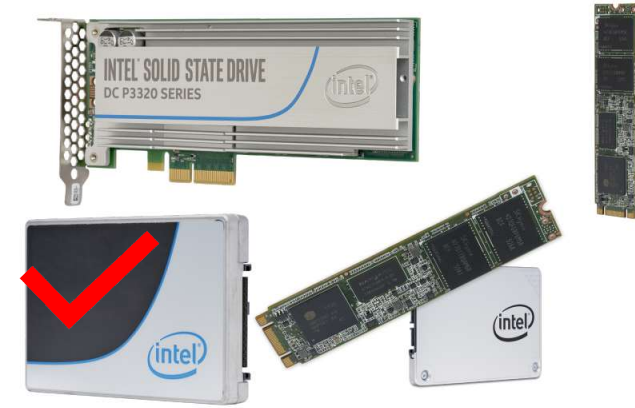
Ultra SYS-2028U-TN24R4T+
(Front View)



 Flash Memory Summit 2017
Santa Clara, CA



RoCE NIC from Mellanox



3 Keys of HW Components
for NVMeoF AFA



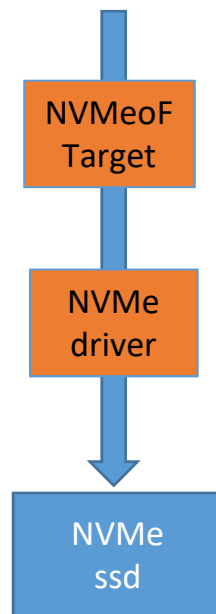
Free Software Resources

- Install Linux Centos 7.2 or above for target-side
- Install Linux Centos 7.2 or above and upgrade kernel to 4.8 for initial-side
- Download SPDK: lock-free software
- Multiple modules in SPDK:
nvme driver, nvmeof target, iscsi target.
- User-layer application
- Replace the role of the Linux driver for nvme

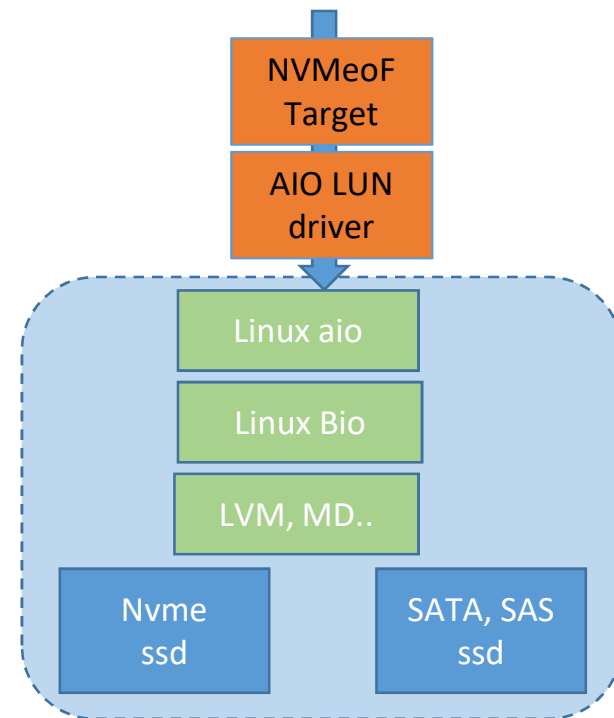
Target Modes: Passthrough vs. Linux AIO

- Passthrough
- Memory
- Linux AIO

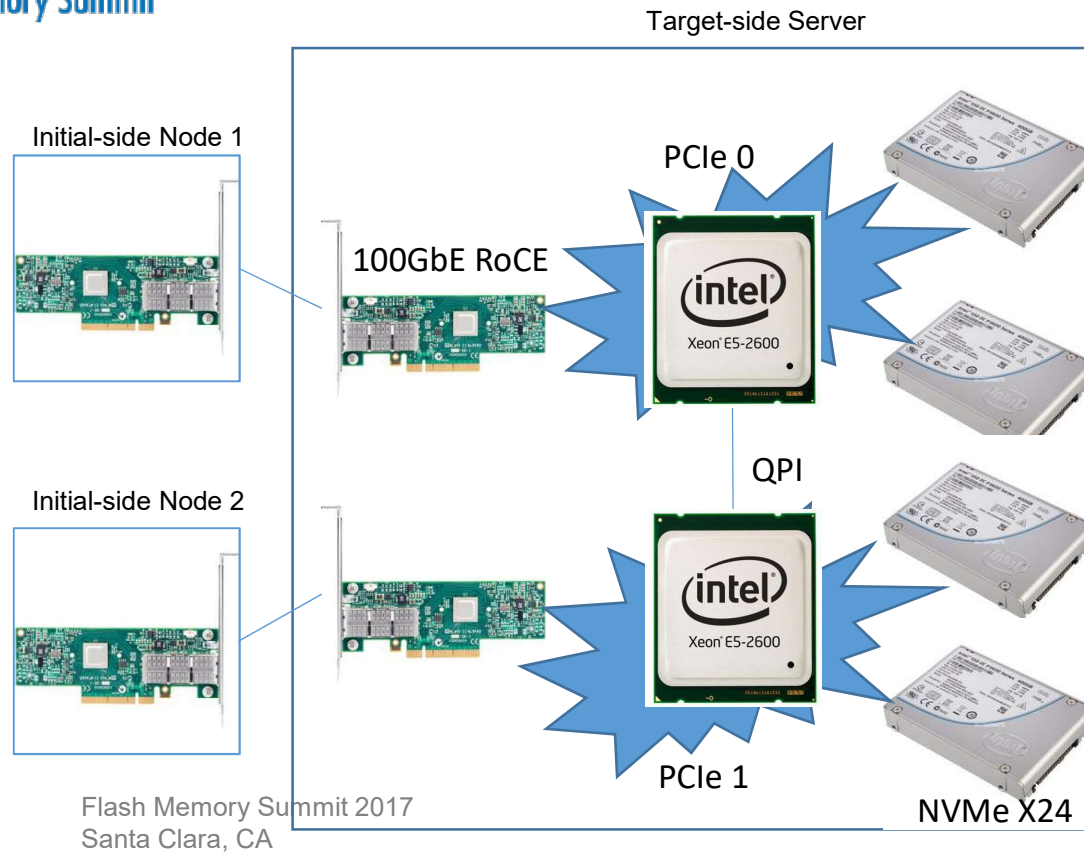
NVMe commands to a physical nvme device



NVMe commands to a virtual nvme device



Test Topology for a NVMeoF AFA



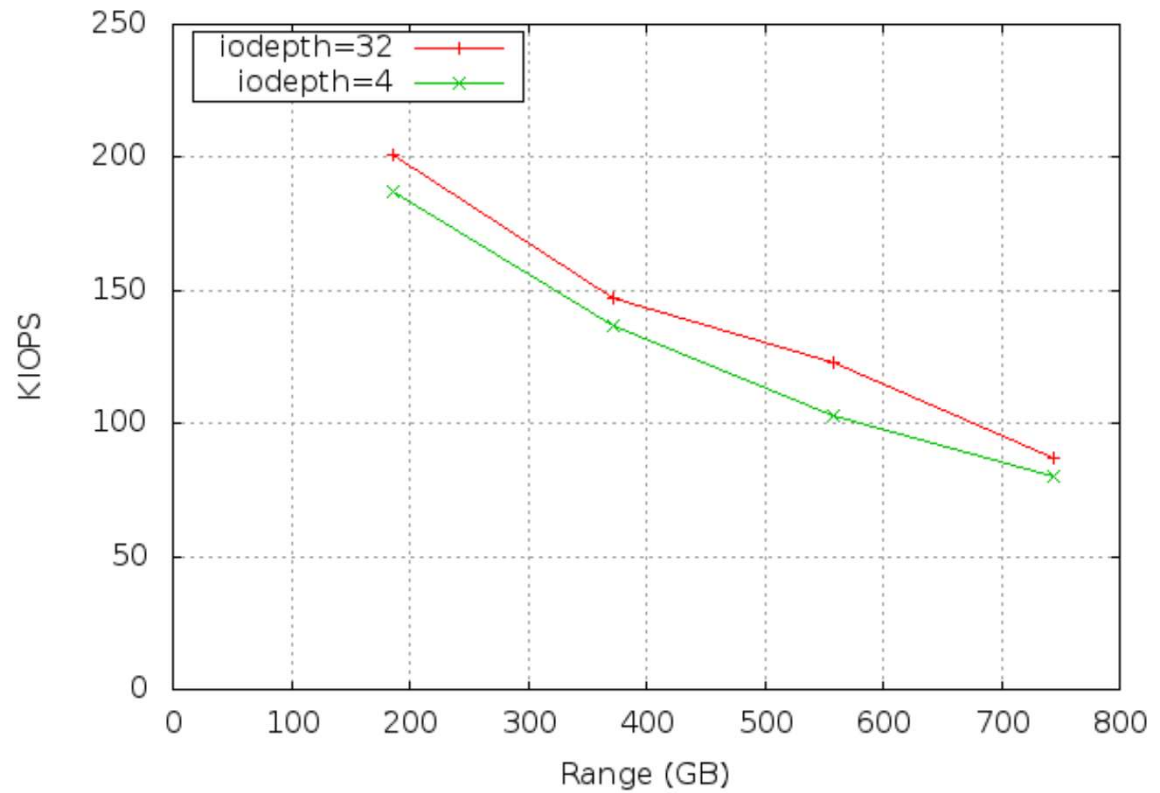
target-side server

- Intel(R) Xeon(R) CPU E5-2667 v3 @ 3.20GHz
- Dc p3600 x 24
- Supermicro
- Mellanox connectX-4

Initial-side node

- Intel(R) Xeon(R) CPU E5-2690 v4 @ 2.60GHz
- Supermicro
- Mellanox connectX-4

Performance on Random WRITE over Different Ranges





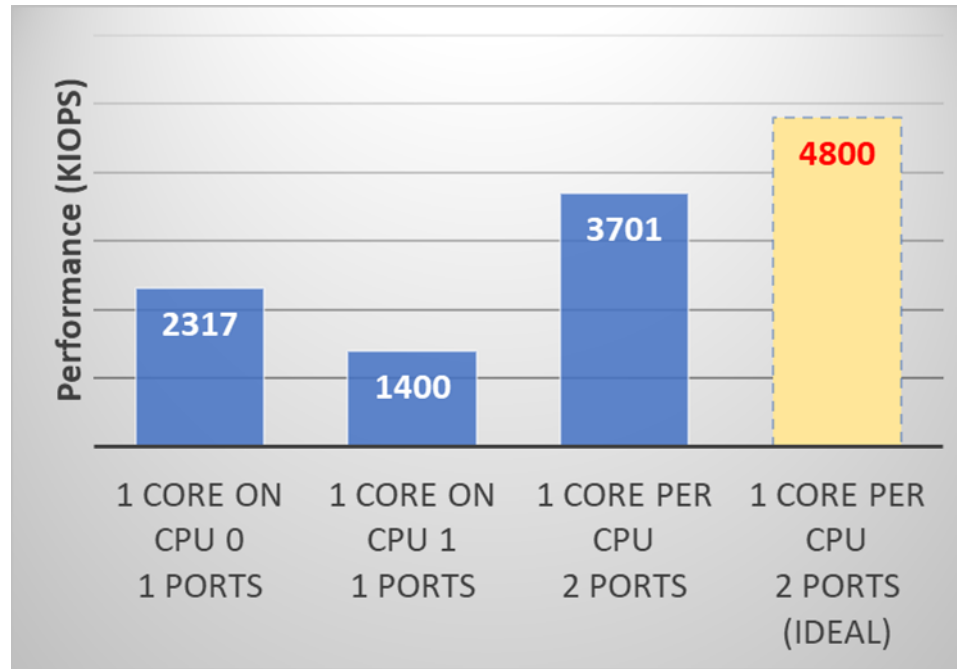
Measured Maximum Random WRITE IOPS

ReactorMask 0x0101
AcceptorCore 0
MaxQueuesPerSession 4
MaxQueueDepth 129

[Subsystem1A]
NQN nqn.2016-06.io.spdk:cnode1
Core 0
Mode Direct
Listen RDMA 50.0.51.1:4420
NVMe 0000:06:00.0

.....
[Subsystem1B]
NQN nqn.2016-06.io.spdk:cnode1B
Core 8
Mode Direct
Listen RDMA 50.0.52.1:4420
NVMe 0000:83:00.0

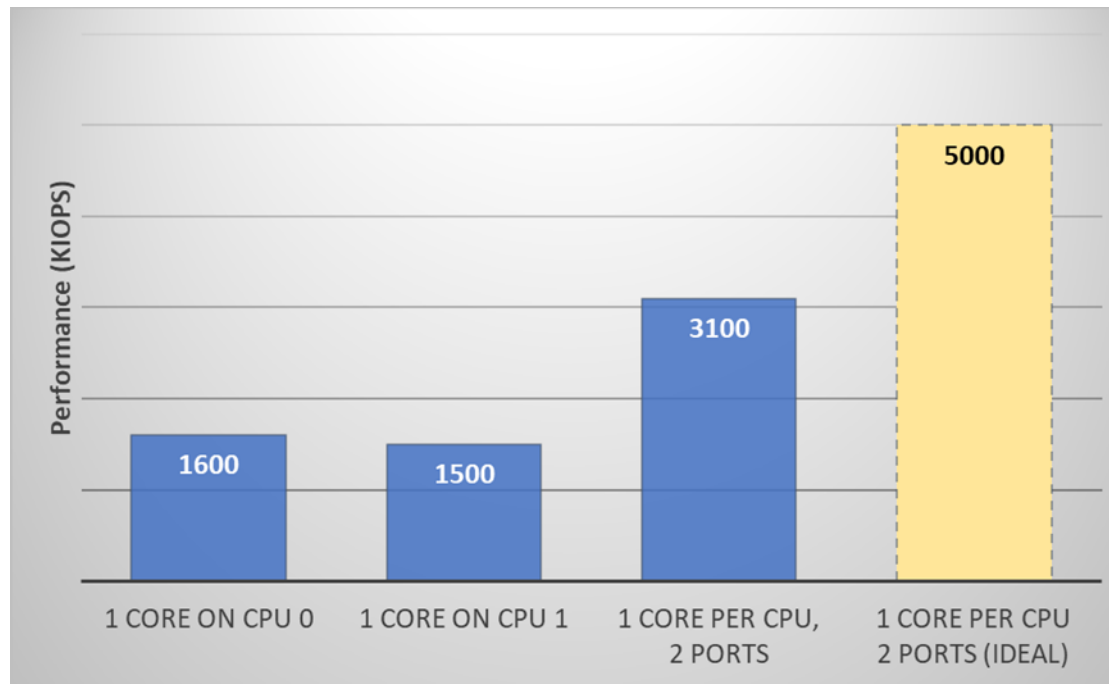
Flash Memory Summit 2017
Santa Clara, CA



Target-side Node X1: Intel(R) Xeon(R) CPU E5-2667 v3 @ 3.20GHz X2
24 nvme drives, 12 for CPU 0 and 12 for CPU 1
Initial-side Node X2: Intel(R) Xeon(R) CPU E5-2690 v4 @ 2.60GHz X2



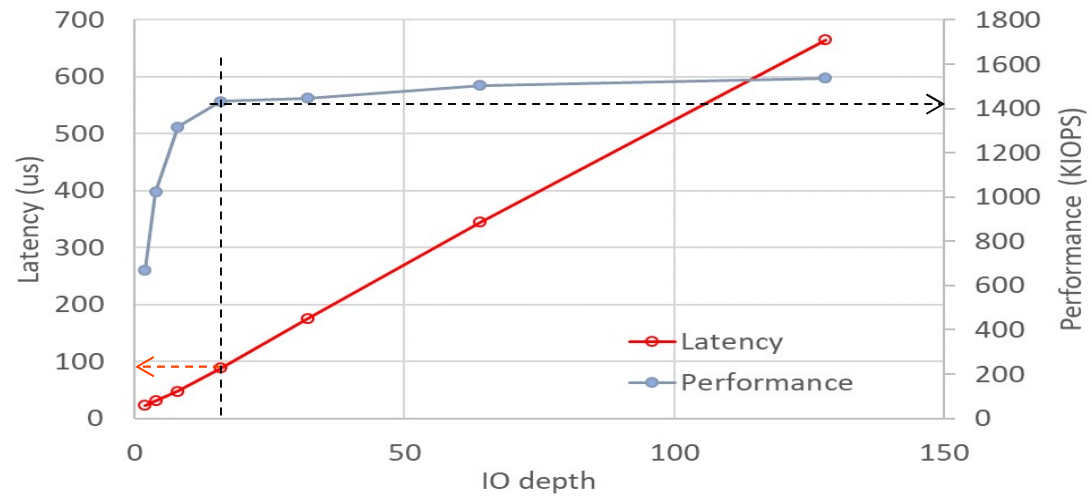
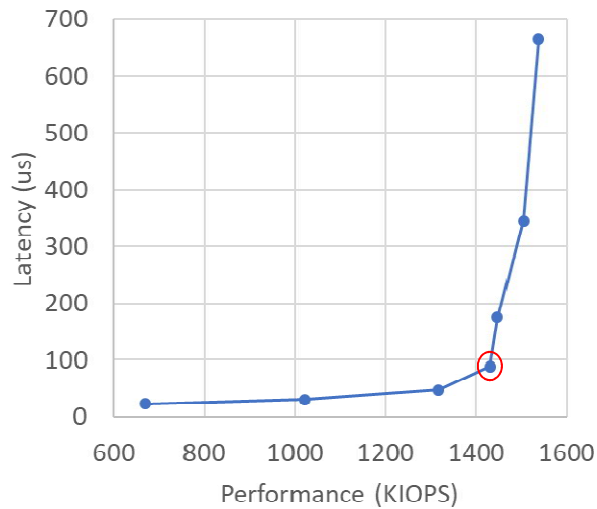
Measured Maximum Random READ IOPS



Flash Memory Summit 2017
Santa Clara, CA

*Target-side Node X1: Intel(R) Xeon(R) CPU E5-2667 v3 @ 3.20GHz X2
24 nvme drives, 12 for CPU 0 and 12 for CPU 1
Initial-side Node X2: Intel(R) Xeon(R) CPU E5-2690 v4 @ 2.60GHz X2*

Measured Latency under Random WRITE over NVMeoF



Target-side Node X1: Intel(R) Xeon(R) CPU E5-2667 v3 @ 3.20GHz X2
 8 nvme drives connecting to CPU 0
 Initial-side Node X2: Intel(R) Xeon(R) CPU E5-2690 v4 @ 2.60GHz X2
 Fio run on the initial host with 8 threads through nvmeof interface.
 IO depth is configured for each fio thread.



Challenges and Solutions

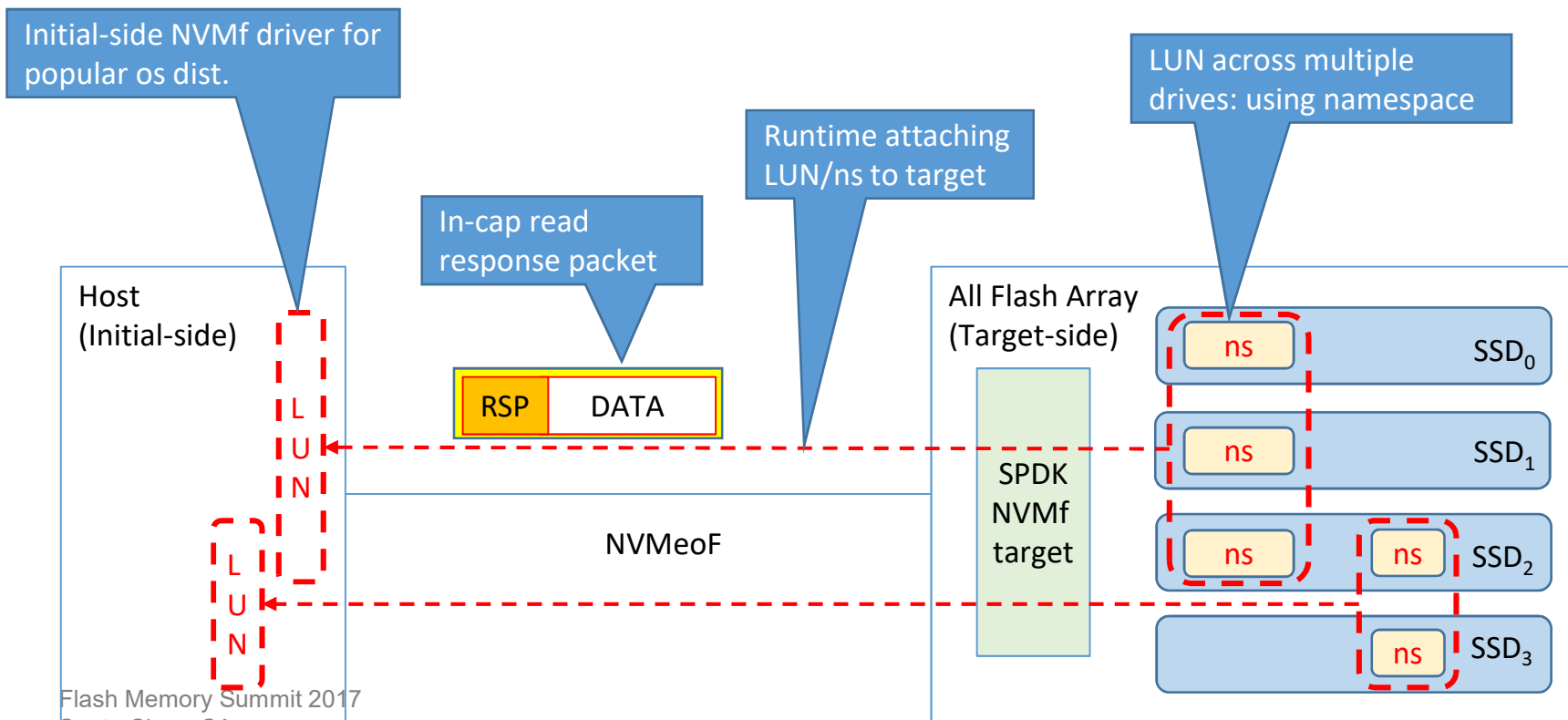
High IOPS with Enterprise Features

Flash Memory Summit 2017
Santa Clara, CA

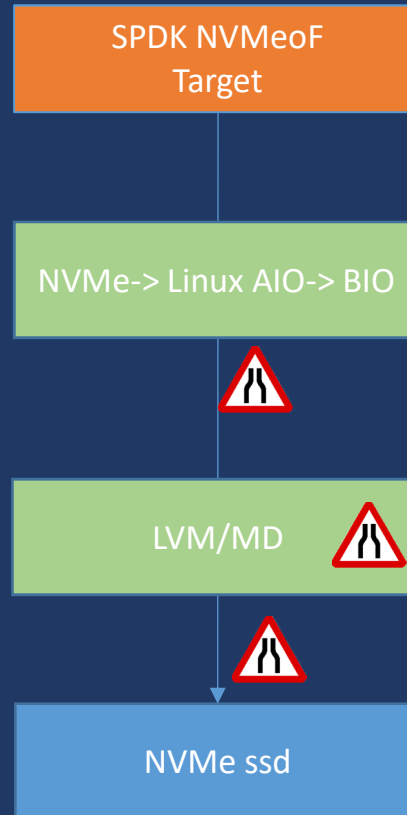
The photo is designed by jigsawstocker / Freepik



Missed Basic Features



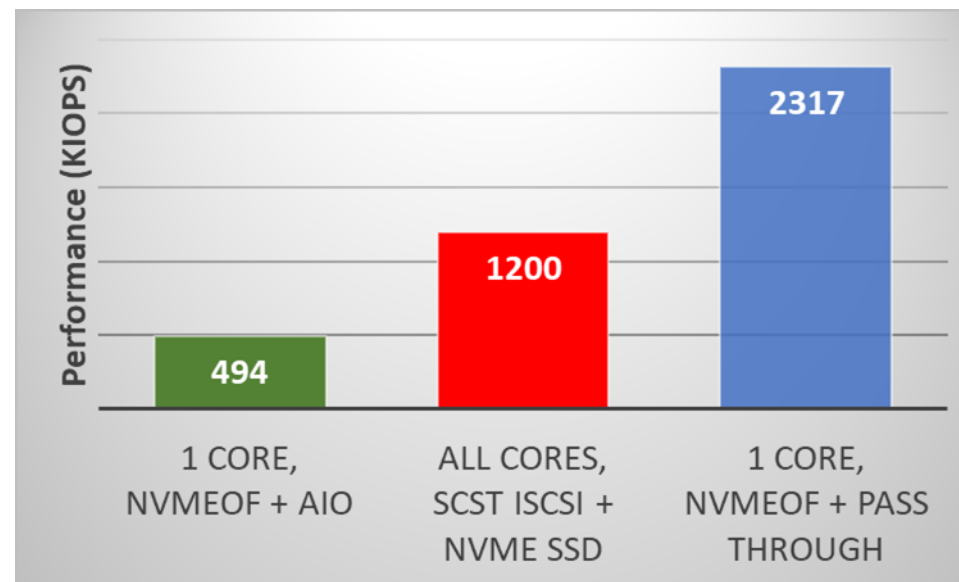
Trade-off between Advanced Features and Performance



- HA
- Scale-up
- Thin-clone
- Thin-provisioning
- Data protection



Performance on 4KB Random WRITE IOPS Under Different Target Modes





IO Amplification (AMP) – Data Protect Mechanism (DPM) on High IOPS SSD

- Write-in-place RAID -> 2R+2W at least
- Solution:
 - RAID-aware device?
 - multi-copy based?
 - Write Buffer:
 - Additional *dual* log devices
 - Additional effort on hw
 - Faster than nvme? Nvdimm
 - RoW RAID

Screenshot of Linux tool: *iostat*
4KB random WRITE:
tps: 4KB transaction per second
User-perceived IOPS: ~ 108K
Actual paid IOPS: ~36K*12, or 432K

```
1 nvme server x 2 nvme server x
avg-cpu:  %user   %nice   %sys
          0.16    0.00    4

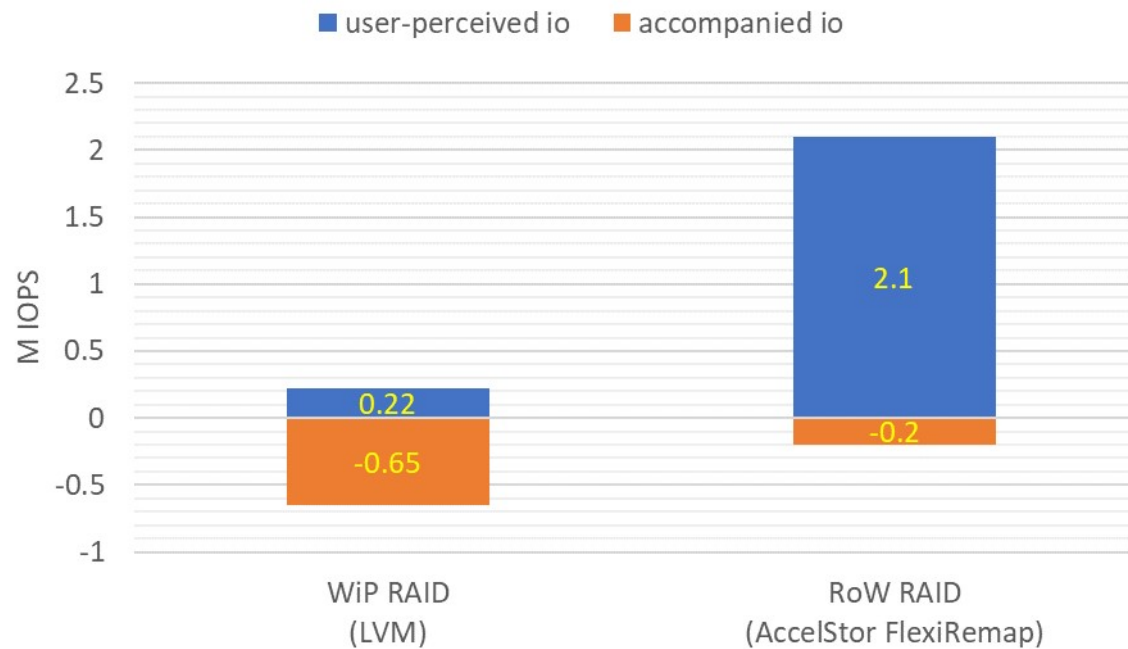
Device:      tps      kB_read/s    kB_wrtn/s    kB_read    kB_wrtn
nvme8n1      36189.00    72400.00     72356.00     72400     72356
nvme10n1     36370.00    72744.00     72736.00     72744     72736
nvme2n1      35723.00    71464.00     71428.00     71464     71428
nvme1n1      35817.00    71660.00     71608.00     71660     71608
nvme9n1      36084.00    72172.00     72164.00     72172     72164
nvme3n1      35759.00    71528.00     71508.00     71528     71508
nvme4n1      36145.00    72300.00     72280.00     72300     72280
nvme11n1     36484.00    72968.00     72968.00     72968     72968
nvme6n1      36441.00    72896.00     72868.00     72896     72868
nvme7n1      35652.00    71340.00     71268.00     71340     71268
nvme5n1      36297.00    72620.00     72568.00     72620     72568
nvme0n1      36335.00    72688.00     72652.00     72688     72652
md0          108376.00      0.00      433504.00      0     433504

^C
[root@weafon7 weafon]#
```

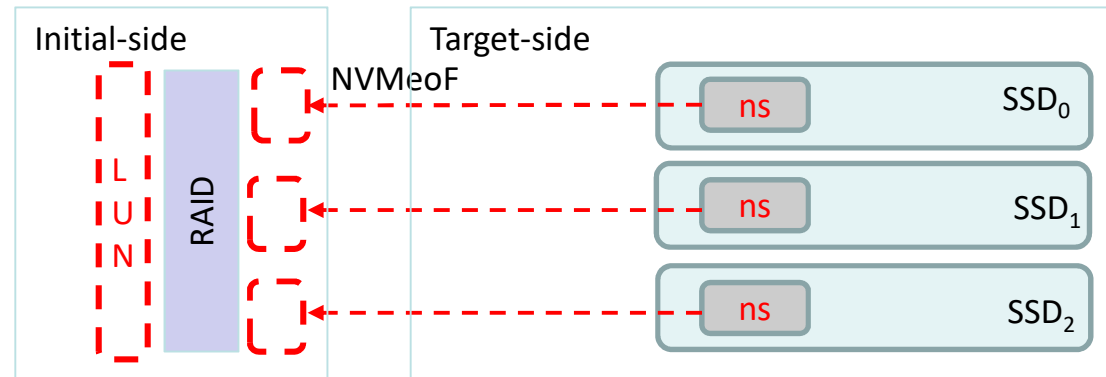
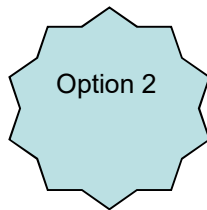
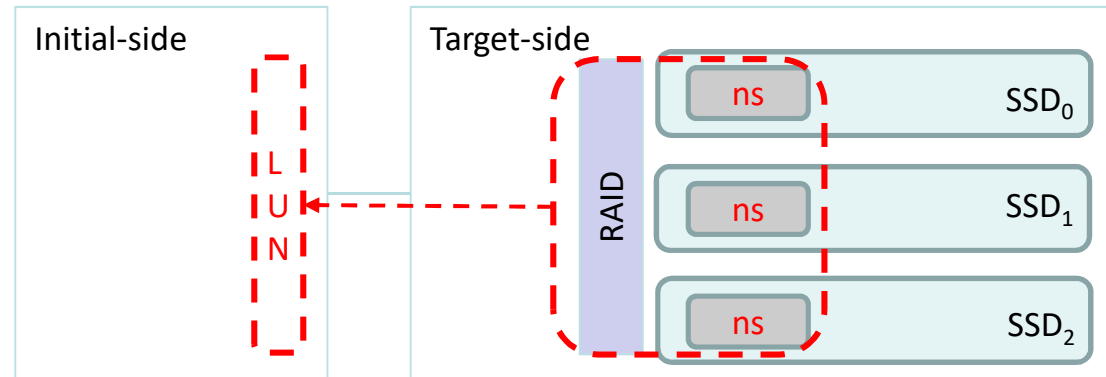
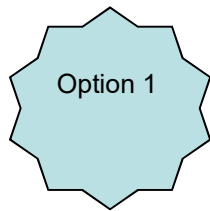
ssh://root@10.144.1.157:22 SSH2 xterm 78x20 20,24 5工作階段 CAP_NUM



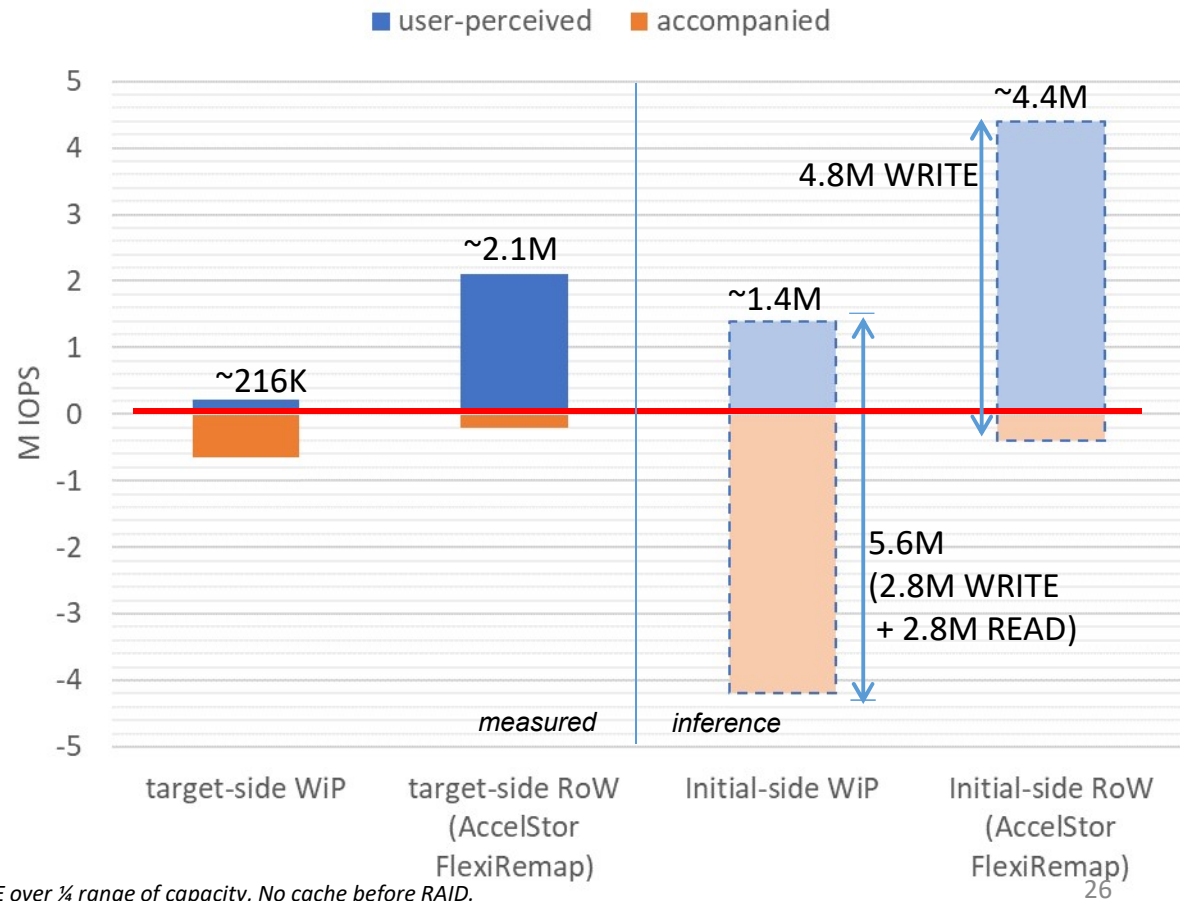
Performance under Different RAID5 over NVMeoF



Where to Deploy DPM?



Benefit of Initial-side RAID

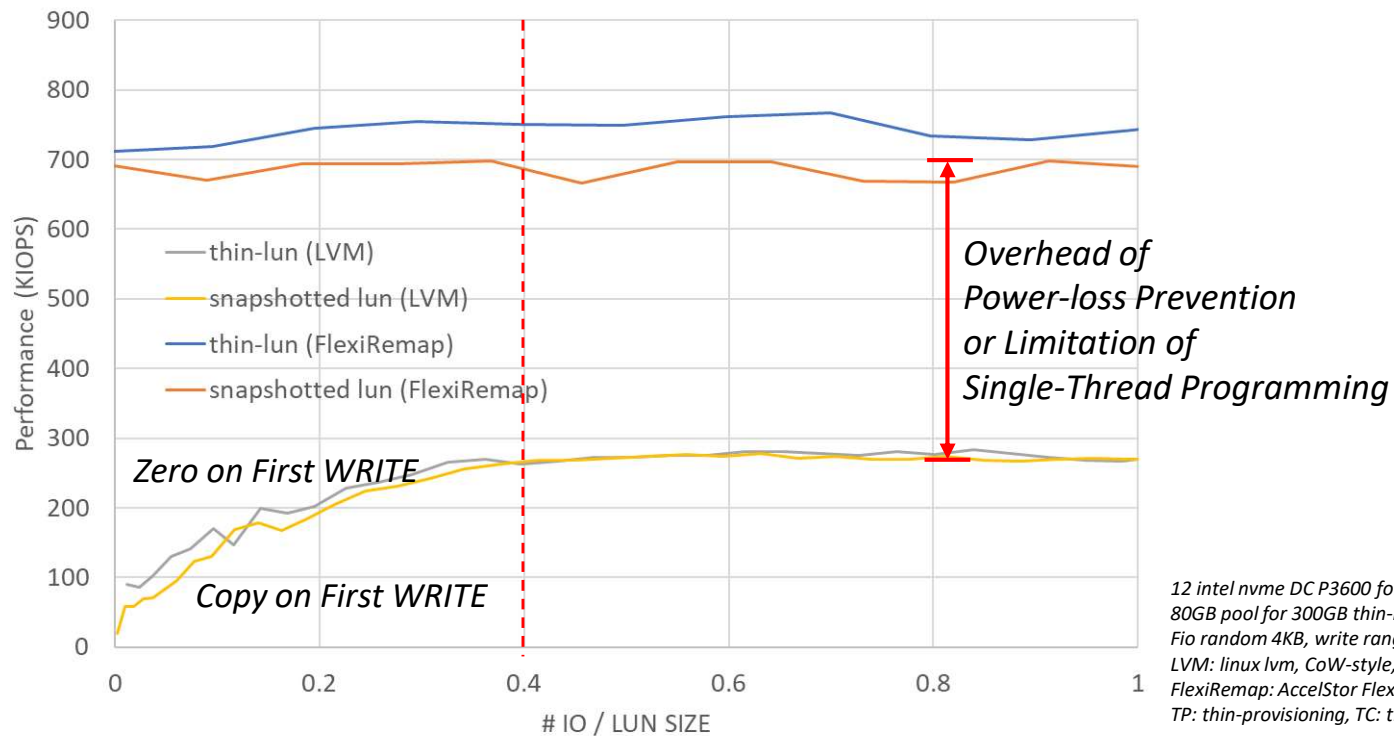


Flash Memory Summit 2017
Santa Clara, CA

4KB Random WRITE over ¼ range of capacity, No cache before RAID.
consisting of 24 NVMe SSDs where 12 SSDs form a raid group and each SSD provide Read 400K iops, write 200K iops.
Fio run on the initial host through nvmeof interface

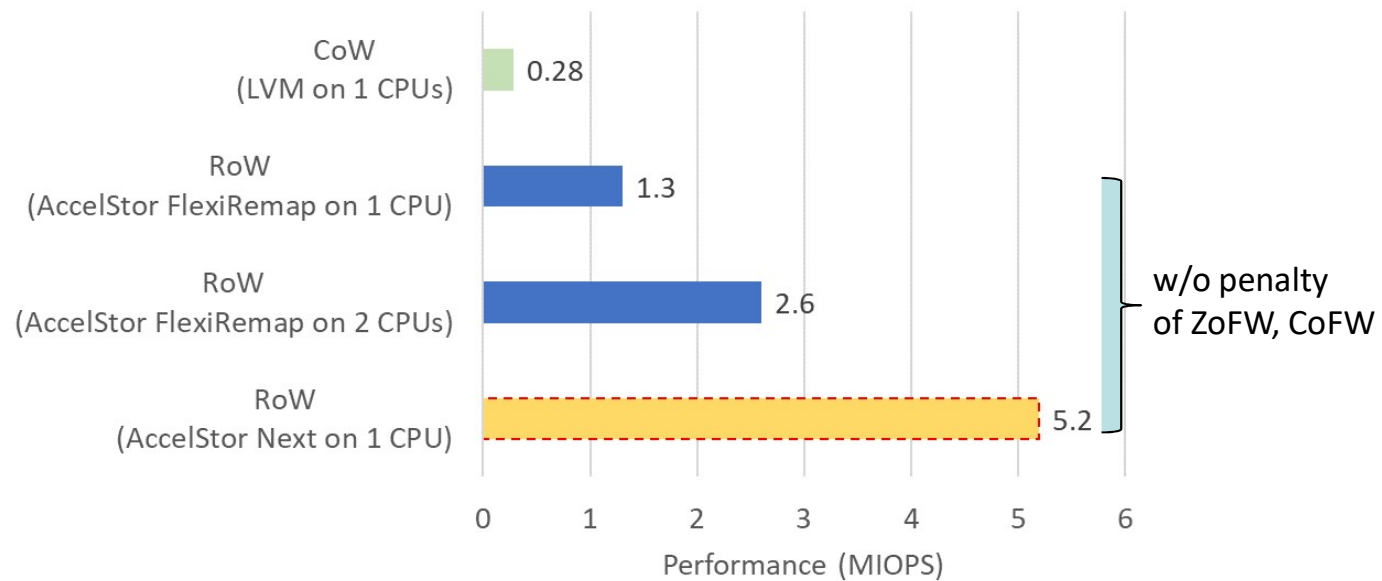


Issues of Small WRITE on Thin-Provisioning/Snapshot/ThinClone





Performance on TP/Snapshot/TC over NVMeoF





Why NVMe and NVMf?
Good for Building High-IOPS AFA

Flash Memory Summit 2017
Santa Clara, CA

The photo is designed by fanjianhua / Freepik

Enjoy NVMf for Free!
How you can get a home-made NVMeoF?

Flash Memory Summit 2017
Santa Clara, CA

The photo is designed by mrschmitt / Freepik

Challenges and Solutions

High IOPS with All Features

Flash Memory Summit 2017
Santa Clara, CA

The photo is designed by jigsawstocker / Freepik

Flash Memory Summit 2017
Santa Clara, CA



Thank You

Weafon.tsao@accelstor.com

AccelStor, Inc.

Booth: 132