



meets *nvm* EXPRESS[®] : All-Flash Ceph!

Brent Compton, Red Hat Inc.

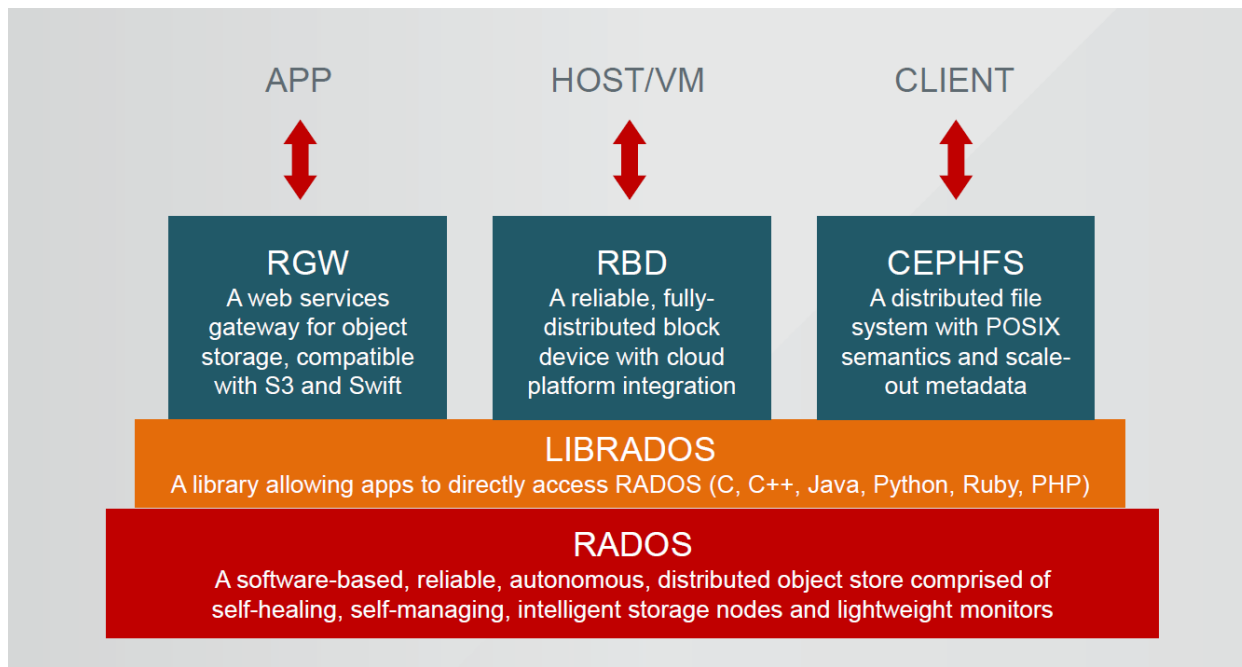
Gunna Marripudi, Samsung Semiconductor Inc.



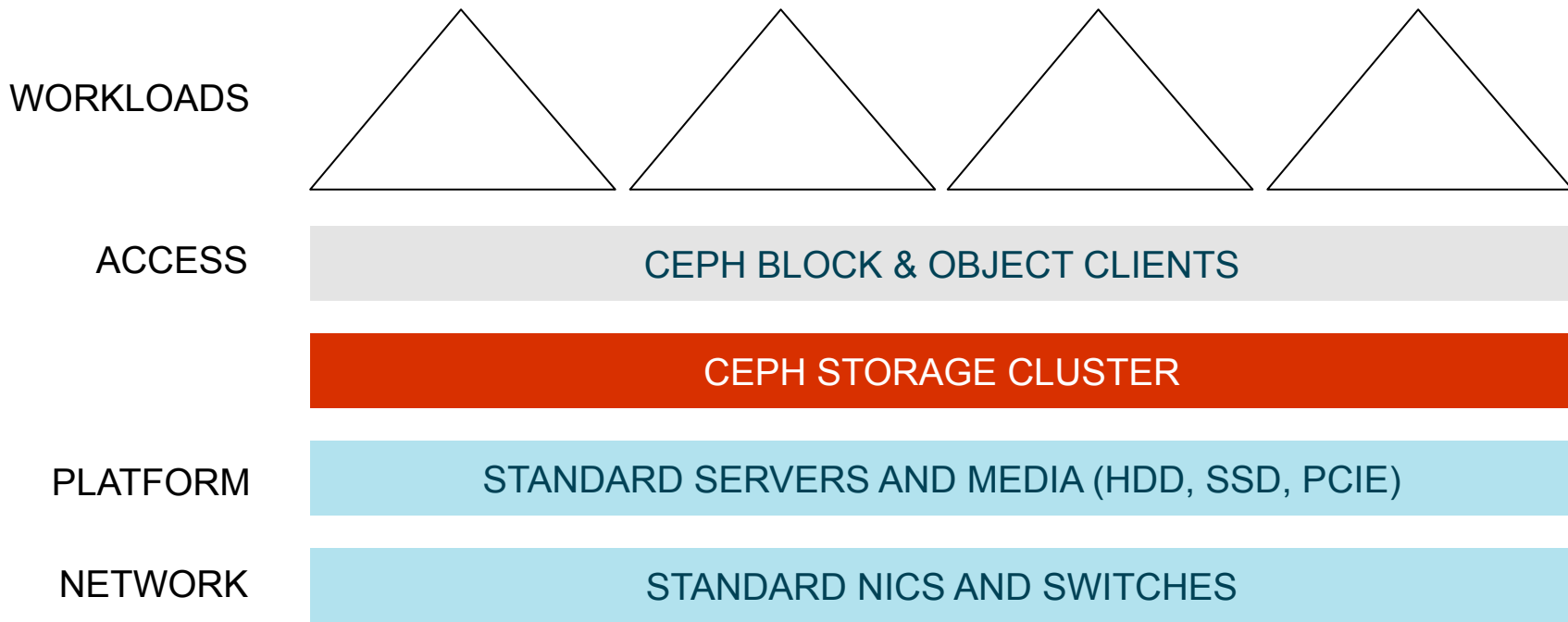
Legal Disclaimer

- This presentation is intended to provide information concerning SSD and memory industry. We do our best to make sure that information presented is accurate and fully up-to-date. However, the presentation may be subject to technical inaccuracies, information that is not up-to-date or typographical errors. As a consequence, Samsung does not in any way guarantee the accuracy or completeness of information provided on this presentation.
- The information in this presentation or accompanying oral statements may include forward-looking statements. These forward-looking statements include all matters that are not historical facts, statements regarding the Samsung Electronics' intentions, beliefs or current expectations concerning, among other things, market prospects, growth, strategies, and the industry in which Samsung operates. By their nature, forward-looking statements involve risks and uncertainties, because they relate to events and depend on circumstances that may or may not occur in the future. Samsung cautions you that forward looking statements are not guarantees of future performance and that the actual developments of Samsung, the market, or industry in which Samsung operates may differ materially from those made or suggested by the forward-looking statements contained in this presentation or in the accompanying oral statements. In addition, even if the information contained herein or the oral statements are shown to be accurate, those developments may not be indicative developments in future periods.

Ceph Architecture Overview



Ceph Cluster Building Blocks



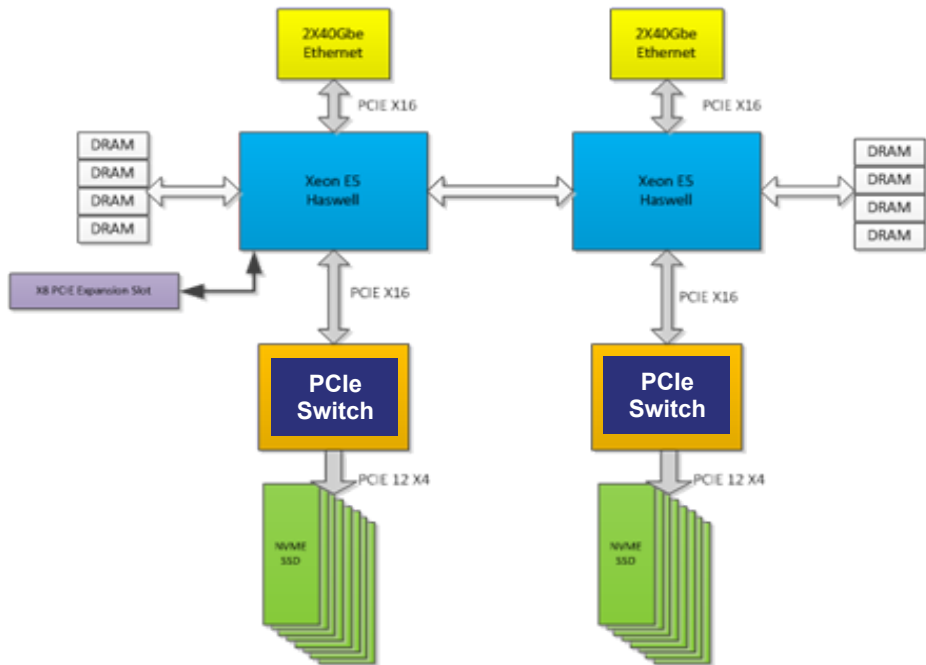
Ceph Cluster Possibilities

	OpenStack Starter 100TB	S 500TB	M 1PB	L 2PB
IOPS OPTIMIZED				
THROUGHPUT OPTIMIZED				
COST-CAPACITY OPTIMIZED				

Flash in Ceph Clusters

	OpenStack Starter 100TB	S 500TB	M 1PB	L 2PB
IOPS OPTIMIZED	SSDs			
THROUGHPUT OPTIMIZED				
COST-CAPACITY OPTIMIZED				

NVMe™ Reference Design*

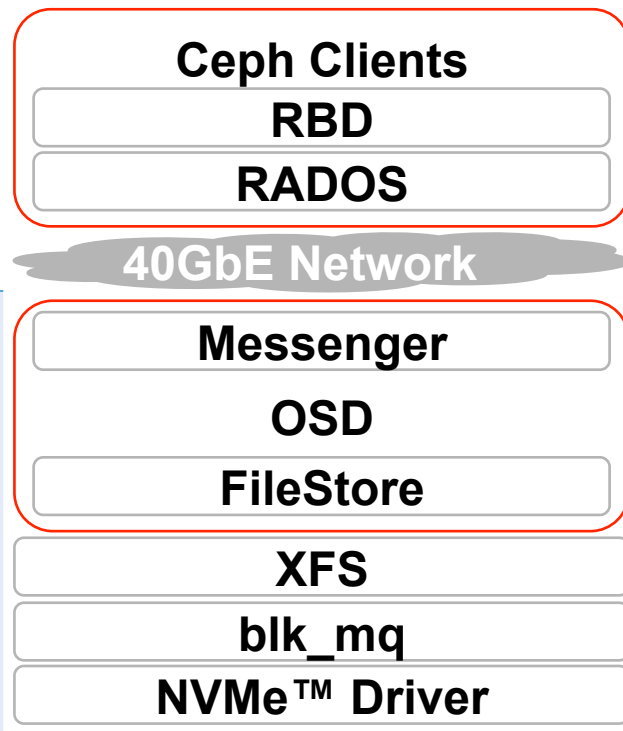
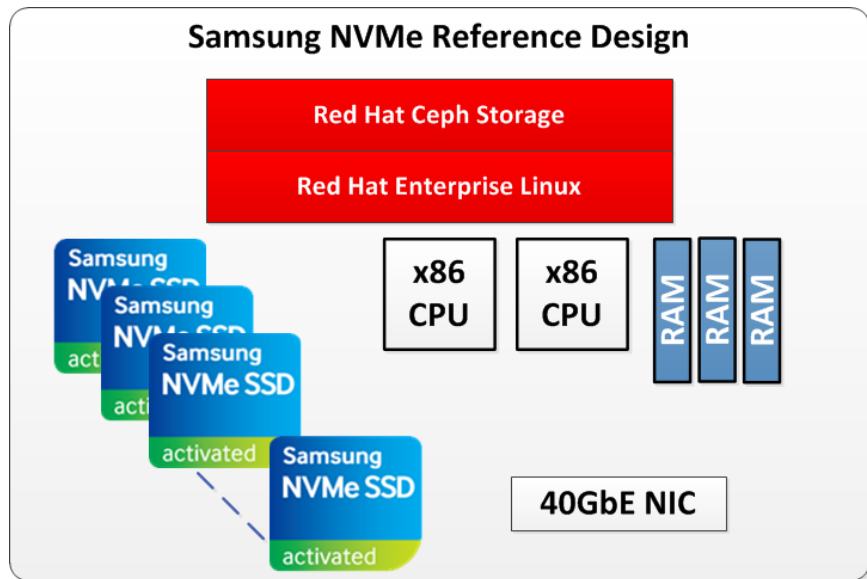


**EIA 310-D 19"
2U Form Factor**

24x 2.5" NVMe™ SSDs (U.2)

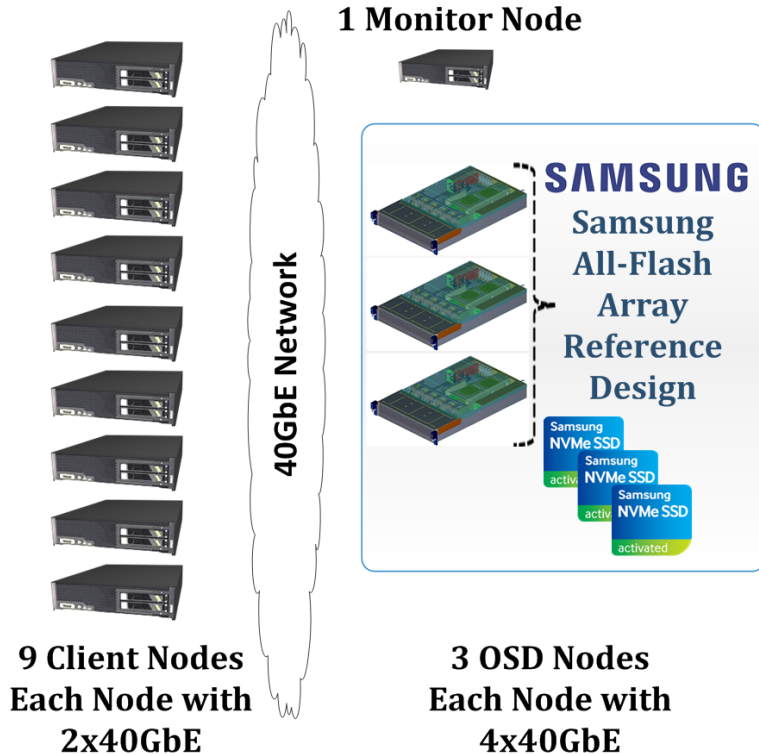
2x PCIe Gen3 x16 Slots

High-Performance Ceph over NVMe™



Ceph Reference Design over NVMe™

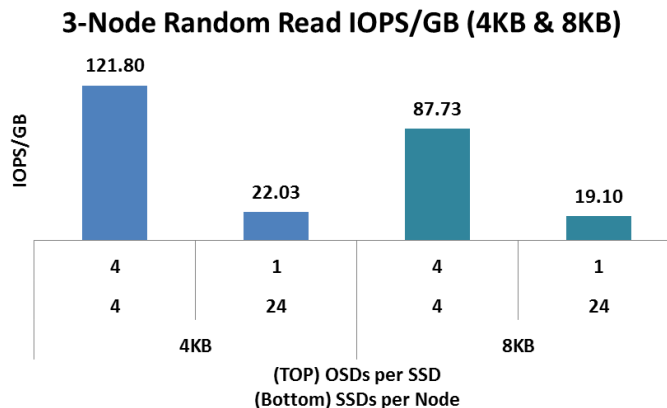
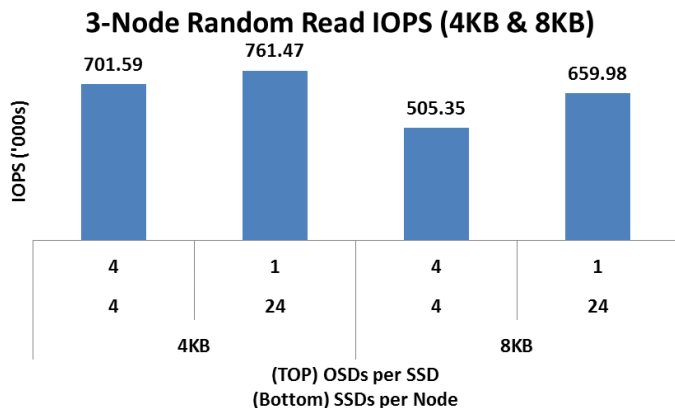
- Test configuration
 - Ceph Jewel w/jemalloc
 - 2x replication
 - Default debug values set
 - OSD nodes with
 - 2x Xeon® E5-2699 v3 CPUs
 - 24x Samsung PM953 2.5” 960GB SSDs
 - CBT test framework
 - RHEL 7.2



IOPS Optimized Configuration

- Random Read

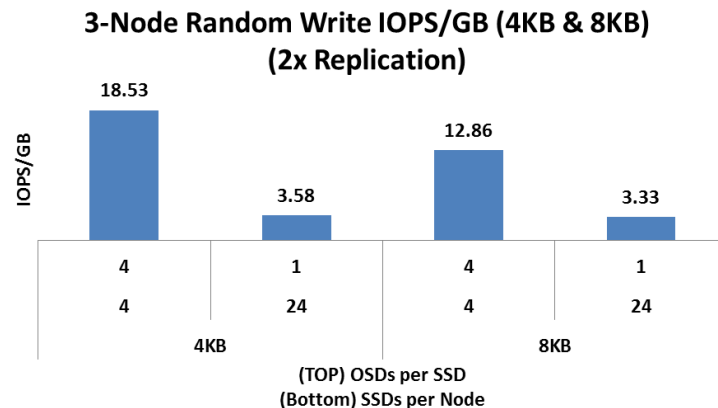
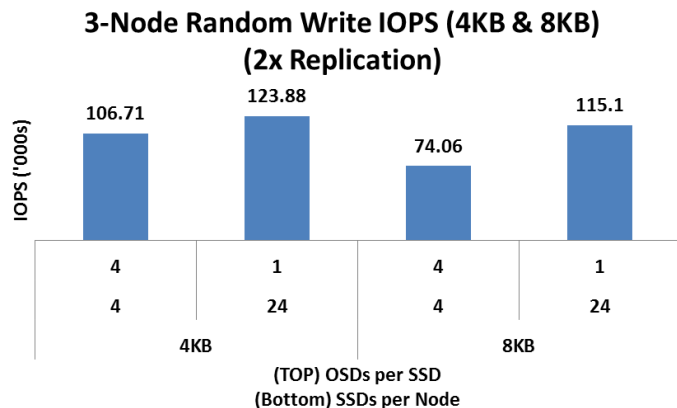
- >700K IOPS with 12 PM953 SSDs in the cluster for 4KB IOs
- >500K IOPS with 12 PM953 SSDs in the cluster for 8KB IOs



IOPS Optimized Configuration

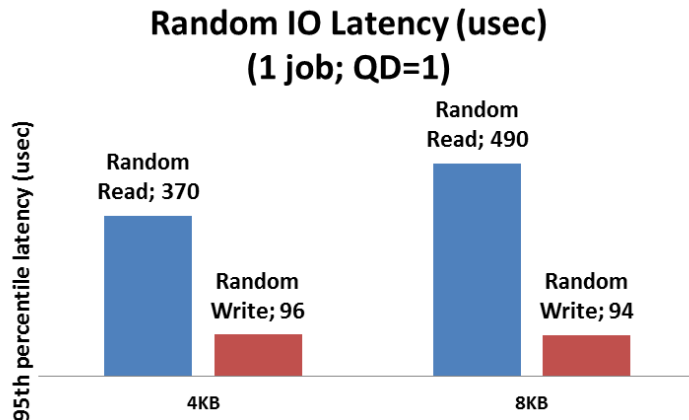
- Random Write

- >100K IOPS with 12 PM953 SSDs in the cluster for 4KB IOs
- >70K IOPS with 12 PM953 SSDs in the cluster for 8KB IOs



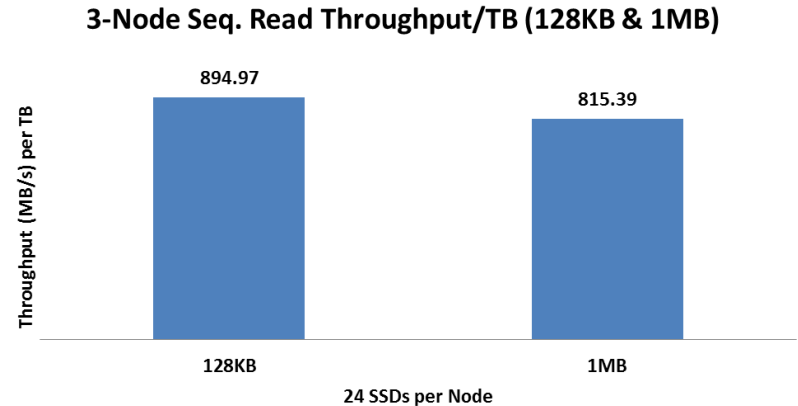
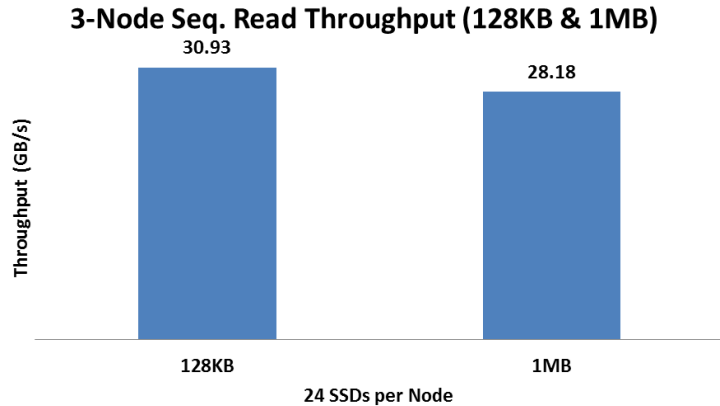
IOPS Optimized Configuration

- Random IO Latency (1 job, QD=1)
 - 95th percentile <500usec for Random Read of 4KB and 8KB
 - 95th percentile <100usec for Random Write of 4KB and 8KB



Throughput Optimized Configuration

- Sequential Read
 - >28GB/s throughput with 72 PM953 SSDs in the cluster
 - >800MB/s throughput per TB cluster capacity

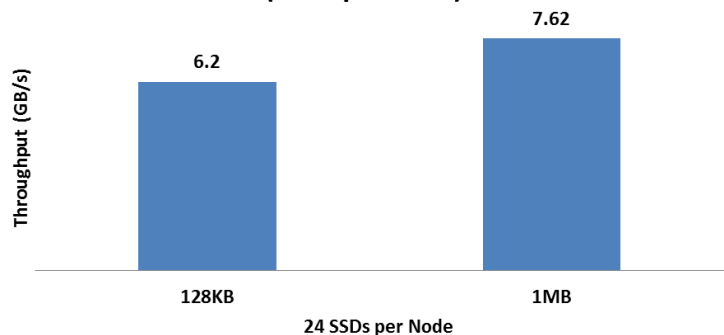


Throughput Optimized Configuration

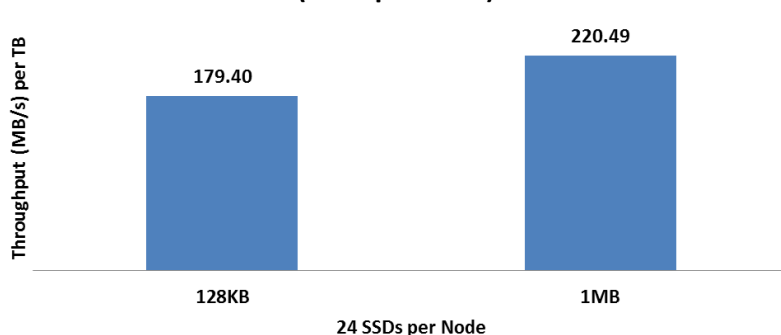
- Sequential Write

- >6GB/s throughput with 72 PM953 SSDs in the cluster
- >179MB/s throughput per TB cluster capacity

3-Node Seq. Write Throughput (128KB & 1MB)
(2x Replication)



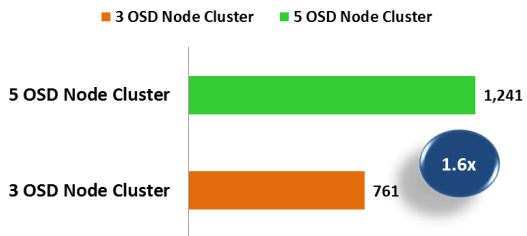
3-Node Seq. Write Throughput/TB (128KB & 1MB)
(2x Replication)



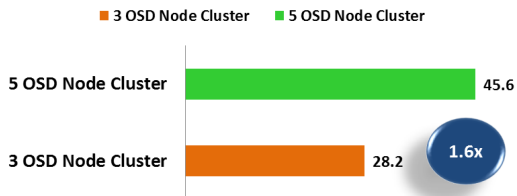
Ceph Scalability: 3 Nodes to 5 Nodes

- When cluster scales 1.6x, performance scales 1.6x (write: 1.3x)

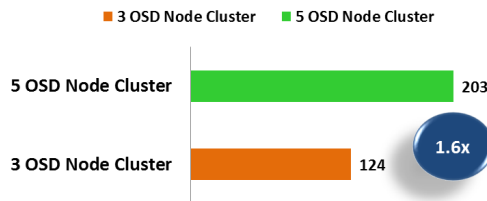
Random Read K IOPS (4KB)



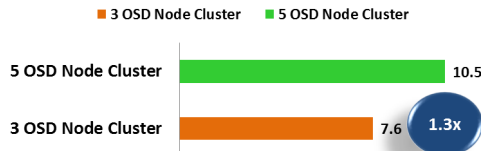
Seq. Read Bandwidth in GB/s (1MB)



Random Write K IOPS (4KB)
(Replication Factor = 2)

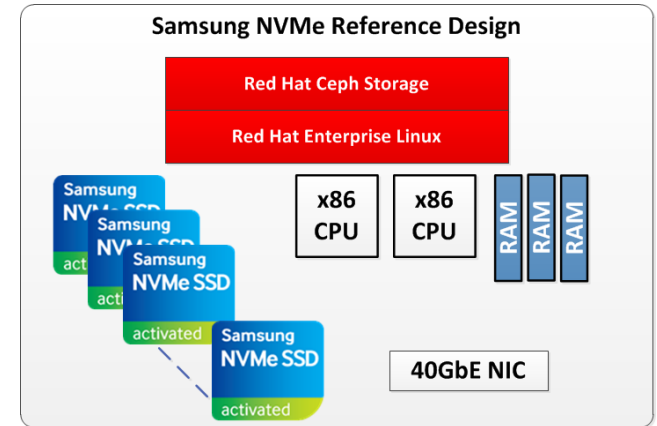


Seq. Write Bandwidth in GB/s (1MB)
(Replication Factor = 2)



Ceph Reference Design for NVMe™

- Flexible and scalable platform for IOPS and Throughput optimized configuration
- 5 Node Ceph cluster using Samsung NVMe™ reference design
 - <1msec latency for QD=1
 - 1.2M IOPS
 - 45GB/s throughput (72Gb/s per node)
- Ready for end-user deployments!



References

- <http://www.samsung.com/semiconductor/support/tools-utilities/All-Flash-Array-Reference-Design/>
- <http://www.samsung.com/semiconductor/support/tools-utilities/All-Flash-Array-Reference-Design/downloads/High-Performance-Red-Hat-Ceph-Storage-Using-Samsung-NVMe-SSDs-WP-20160712.pdf>

Thanks

Test Configuration

- System tuning, Ceph configuration and CBT test methodology as detailed in
 - <http://www.samsung.com/semiconductor/support/tools-utilities/All-Flash-Array-Reference-Design/downloads/High-Performance-Red-Hat-Ceph-Storage-Using-Samsung-NVMe-SSDs-WP-20160622.pdf>