



Flash Memory Summit

A Kubernetes based Platform for IOT

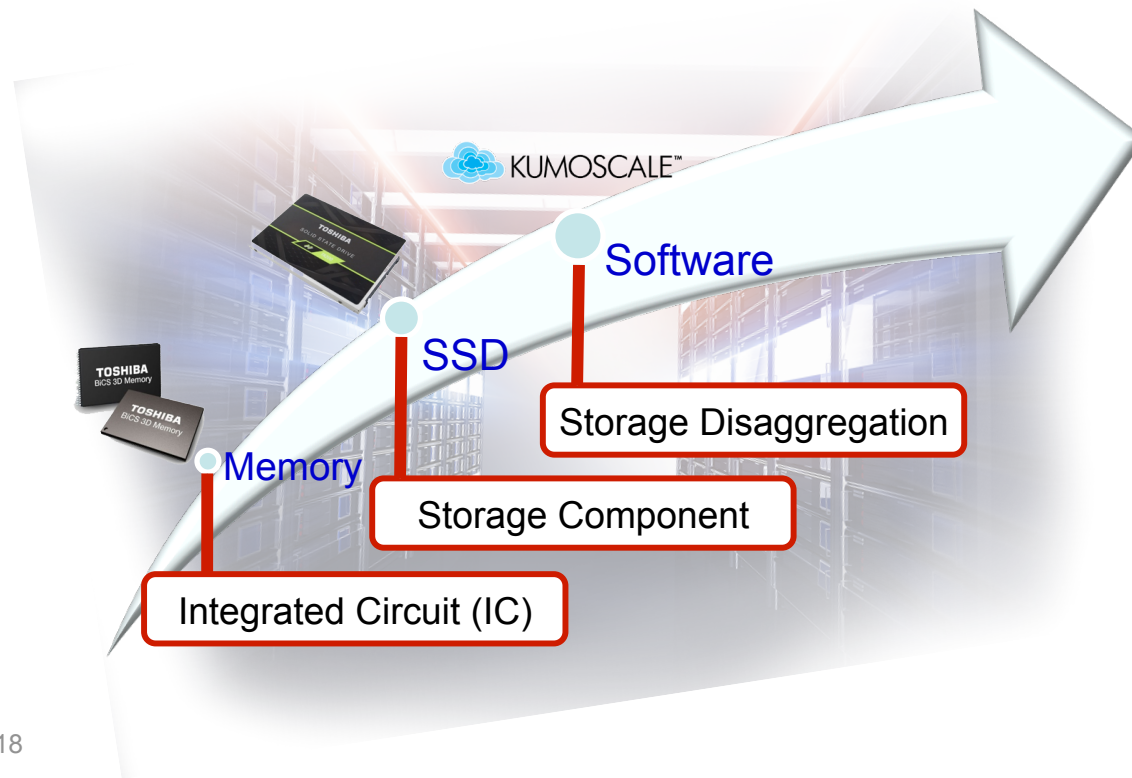
TOSHIBA



KUMOSCALE™



About Toshiba Memory Corporation



Flash Memory Summit 2018
Santa Clara, CA

TOSHIBA



Flash Memory Summit

Challenges with a Scalable Architecture for IOT

DAS



direct-attached
flash

PROBLEM 1



stranded flash & IOPS

PROBLEM 2



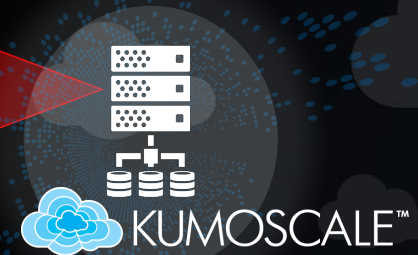
stranded compute

PROBLEM 3



lost operational
agility & revenue

KumoScale™ enables customers to create an NVMe-oF™ storage node using standard x86 compute platform.



shared accelerated
storage

SOLUTION

Flash Memory Summit 2018
Santa Clara, CA

TOSHIBA
Leading Innovation >>>

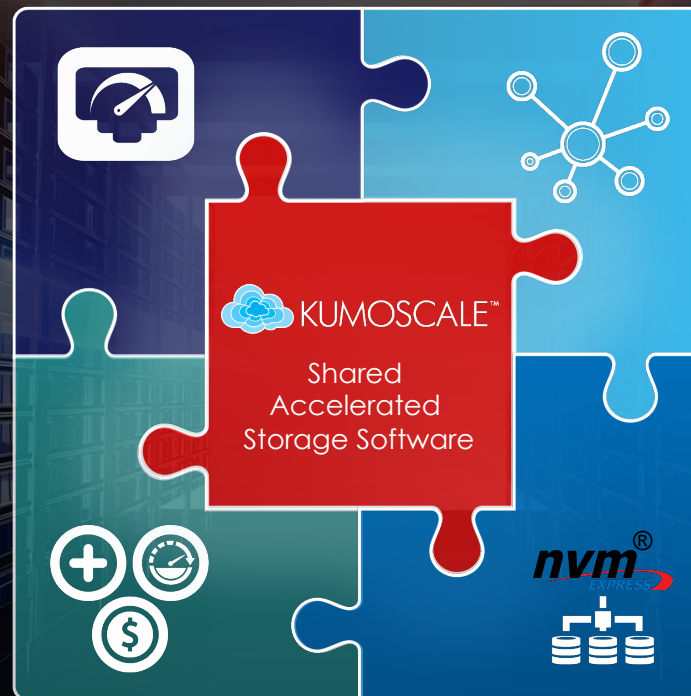


Flash Memory Summit

Shared Accelerated Storage for Cloud Now Possible

High-bandwidth,
low-latency
networks

Mature
orchestration
frameworks



Bigger, faster &
more cost-
effective
NVMe™ SSDs

NVMe-oF™
protocol

Flash Memory Summit 2018
Santa Clara, CA

maximum data center efficiency

TOSHIBA
Leading Innovation >>>



Flash Memory Summit

Storage Disaggregation Benchmark Tests

Objective

- Establish NoSQL (MongoDB™) YCSB Benchmark for KumoScale versus DAS versus Competitive NVMe SSDs on
- Bare metal DB hosts
- Containerized Database hosts (docker : openshift)

Goals

- Bare Metal
- KumoScale is within <10% latency adder over DAS NVMe
- Greater performance

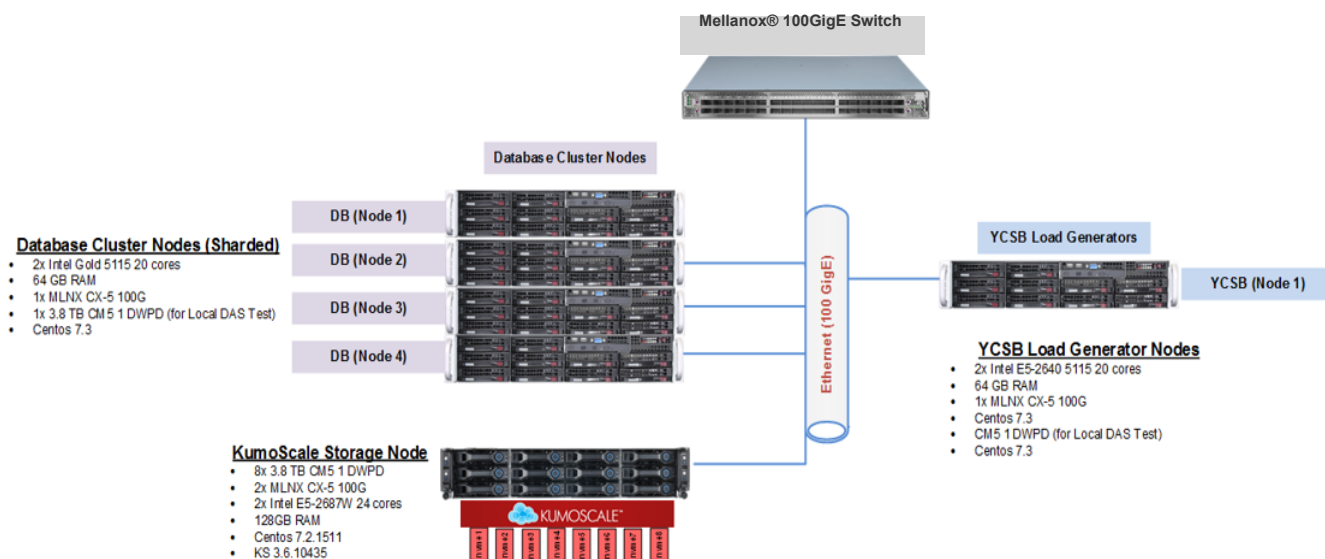
Containerized

- Demonstrate Failover of Compute Node without service disruption
- Demonstrate Failover of Storage Node without service disruption
- Container Storage carved out of 1 KS Node/2 subsystems or 2 KS nodes



Flash Memory Summit

Benchmark Test Setup – YCSB: MongoDB™



Test Profile

- YCSB ver: 0.12
- MongoDB v3.6.5
- 300M-500M records of 1KB each
- Key-size = 100-bytes
- 3 Shards
- Workload A (Mixed) 50% read + 50% update

Metrics

- Avg. Read Latency
- Avg. Update Latency
- 99th %tile Read Latency
- 99th %tile Update Latency
- Throughput Ops/s

Results

- YCSB metrics for:
- Mongo on KumoScale vs
 - Mongo on Local NVMe
 - Mongo on Local SATA SSD

Flash Memory Summit 2018
Santa Clara, CA

TOSHIBA
Leading Innovation >>>



Clustered MongoDB YCSB Results: Test 1 - 50r:50w Mixed workloadA

**Mongo and YCSB Not Tuned for Performance*

Metrics	Local NVMe (CM5 3.8T)	KumoScale NVMe-oF (3T Abstracted NS)
Configuration	3 MongoDB Shards : 1 MongoS router : 1 YCSB Engine	
Load Phase: Insert	300 Million records : 1K record size : 200 parallel Threads	
• Average Insert Latency (ms) <small>(lower is good)</small>	2.95	3.10
• 99 th percentile Insert latency (ms)	5.15	5.61
• Insert Speed (operations/sec) <small>(higher is good)</small>	67,605	64,277
Run Phase: Read & Update		
• Avg Read Latency (ms)	2.74	2.72
• Avg Update Latency (ms)	2.83	2.83
• 99 th percentile Read latency (ms)	5.25	5.18
• 99 th percentile Update latency (ms)	5.36	5.28
• Read & Update speed (operations/sec) <small>(higher is good)</small>	71,768	71,868



Flash Memory Summit

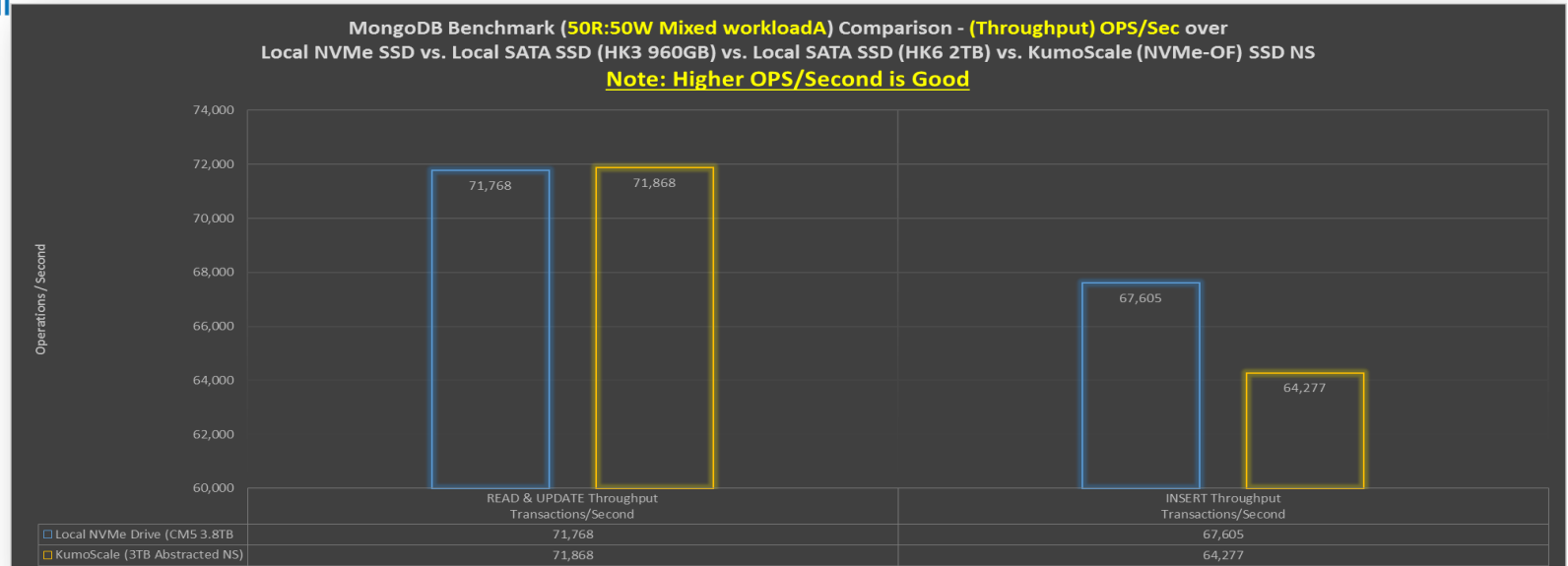
MongoDB™ Benchmark (50R:50W Mixed workloadA) Comparison – (INSERT, READ, UPDATE) AVG Latency



The above chart illustrates the AVG Latency Operations over the given Local Target(NVMe) vs. KumoScale (NVMe-OF) using Mellanox 100GigE Switches. KumoScale Latency shows overall better performance in comparison with others.



MongoDB™ Benchmark (50R:50W Mixed workloadA) Comparison – (INSERT, READ, UPDATE) OPS/Second → Throughput

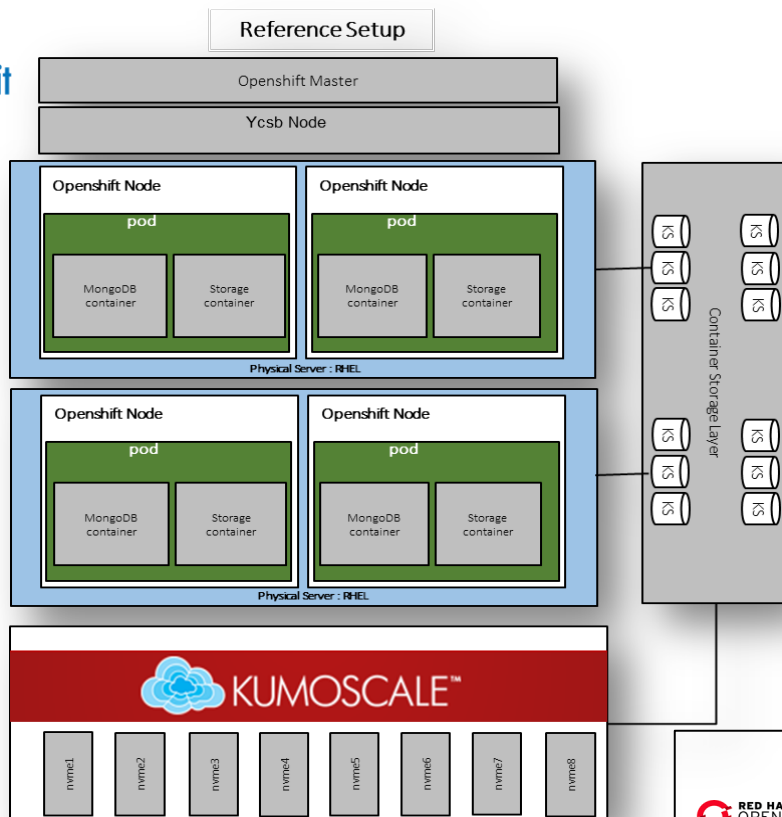


The above chart illustrates the Operations / Second over the given Local Target(NVMe) vs. KumoScale (NVMe-OF) using Mellanox 100GigE Switches. KumoScale Throughput shows overall better performance in comparison with others.



Flash Memory Summit

Test 2: MongoDB™ on RED HAT® OpenShift® Platform



OpenShift Node

- One or more nodes per physical server.
- One or more containers run in one node.

Problem 1: Red Hat Gluster Storage (RHGS)

- Container Storage carved out of 1 KS Node/2 subsystems or 2 KS nodes.
- RHGS dynamically manages pool of storage (KS volumes) to hold data for stateful containers
- RHGS provides replication of data across volumes.
- Demonstrate Failover of Storage Node without service disruption

Problem 2 & Problem 3: OpenShift Master

- Manages scheduling and orchestration
- Container replication and restart to defined policy.
- Demonstrate Failover of Compute Node without service disruption
- Load balancing
- Scale Out

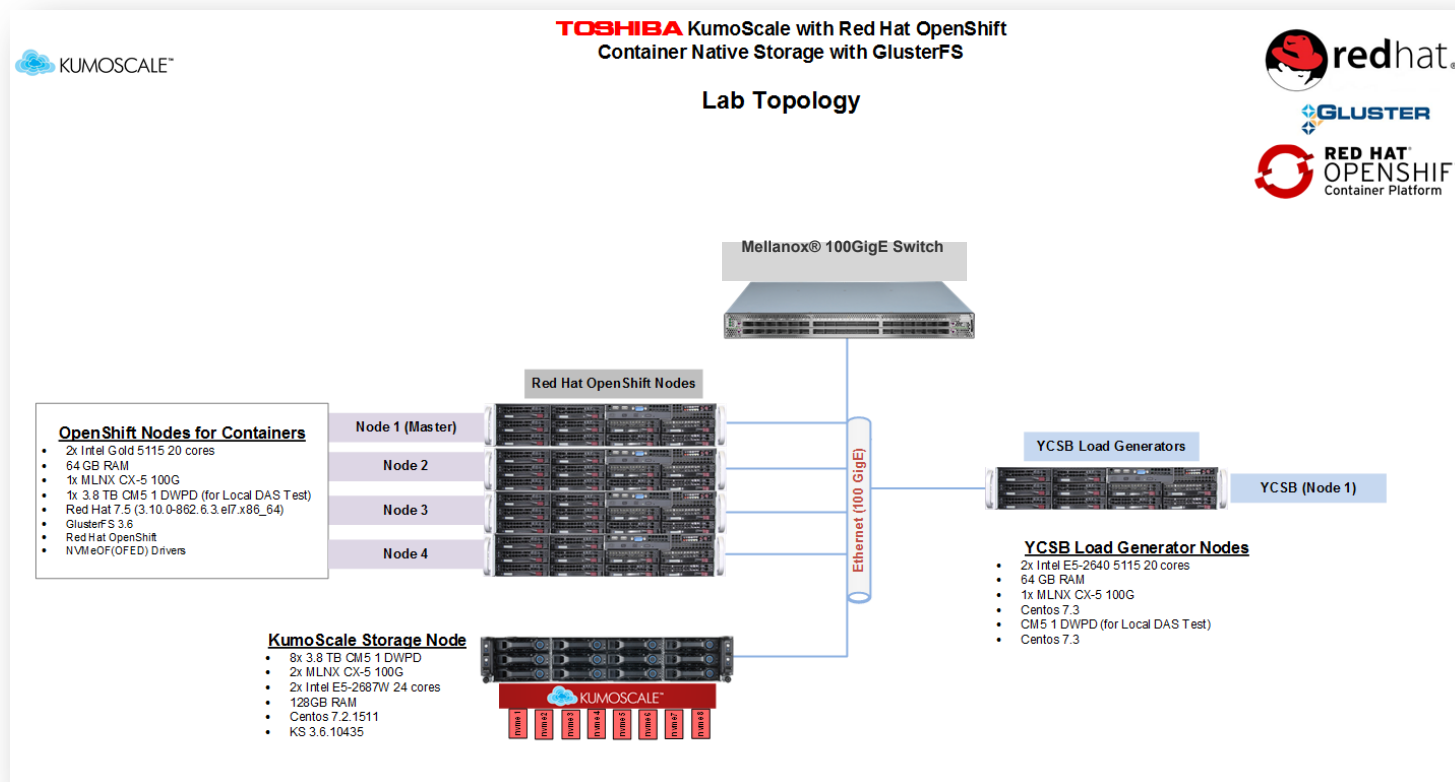
Why OpenShift? – Integrated Container Platform



Flash Memory Summit 2018
Santa Clara, CA

TOSHIBA
Leading Innovation >>>

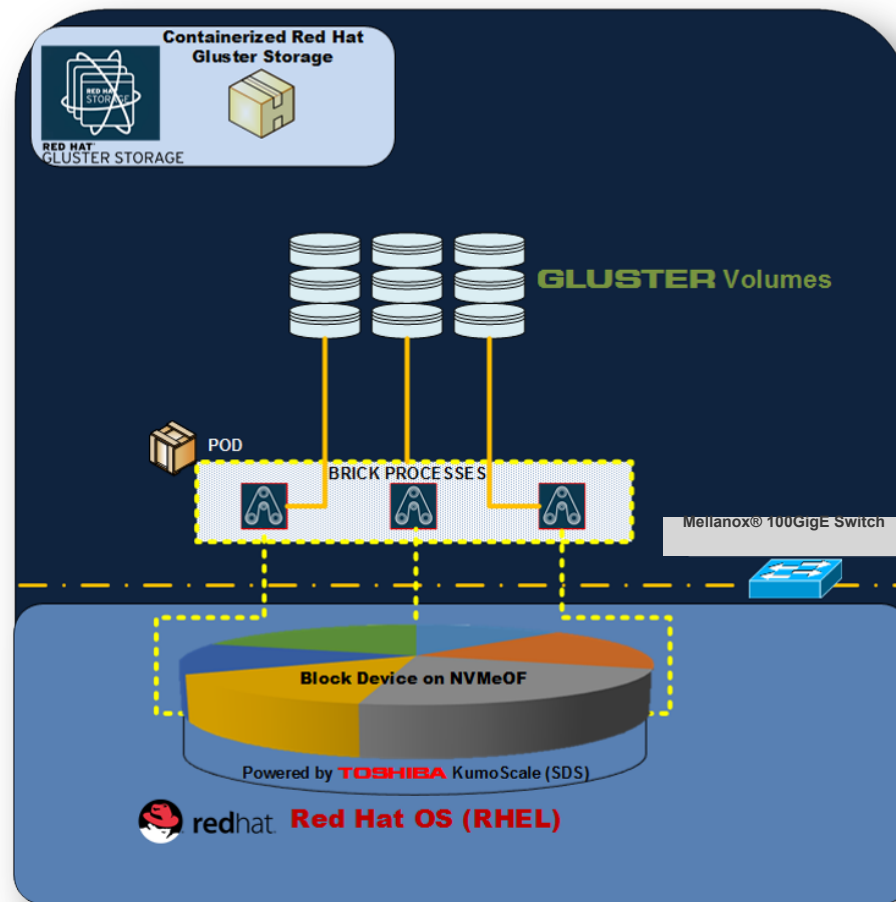
Storage Disaggregation with Container Native Storage



Operating System: RED HAT® Enterprise Linux 7.5
Container Platform: RED HAT® OpenShift®
File System: GlusterFS®



Storage Disaggregation with Container Native Storage with GlusterFS

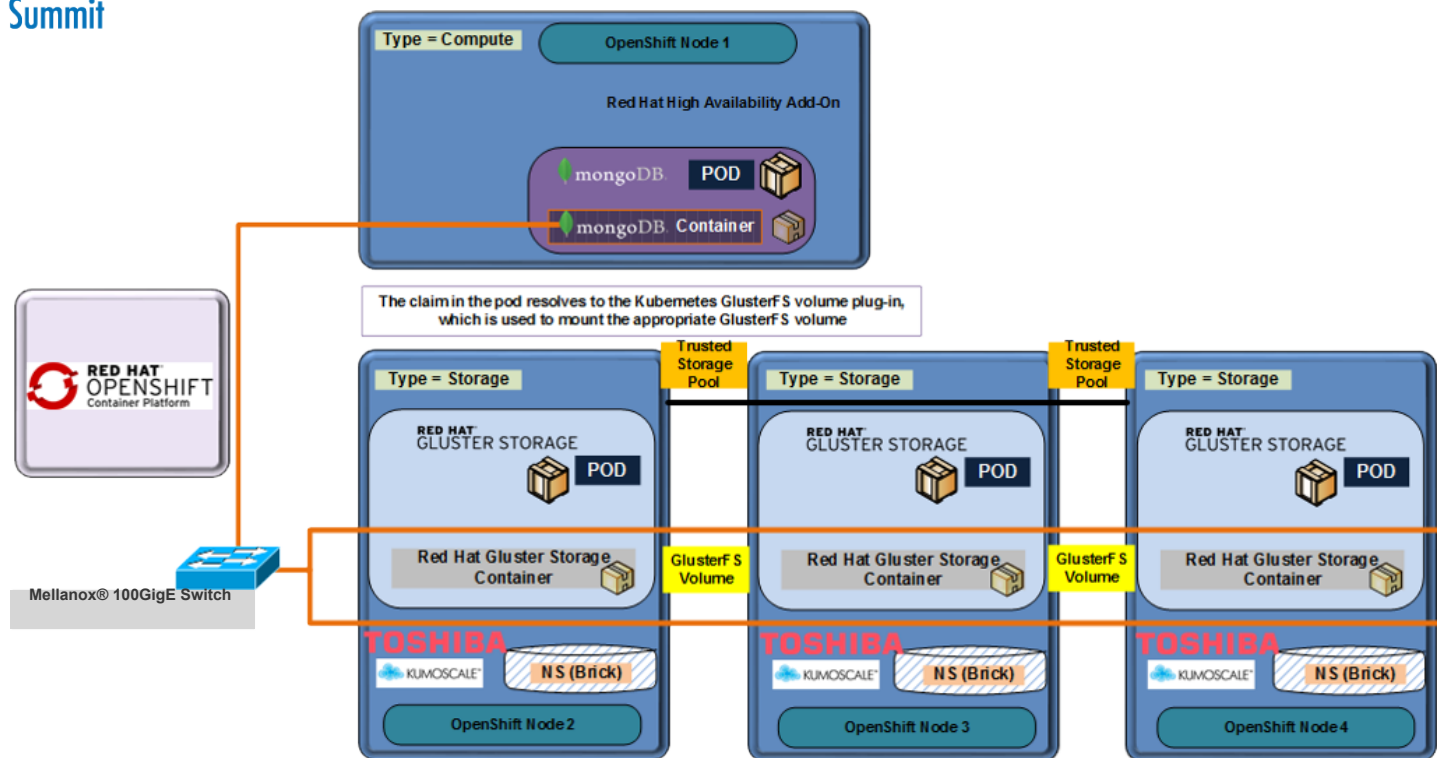


TOSHIBA
Leading Innovation >>>



Flash Memory Summit

Storage Disaggregation with Container Native Storage with GlusterFS



Flash Memory Summit 2018
Santa Clara, CA

TOSHIBA
Leading Innovation >>>



Flash Memory Summit

Backup Slides

For reference only



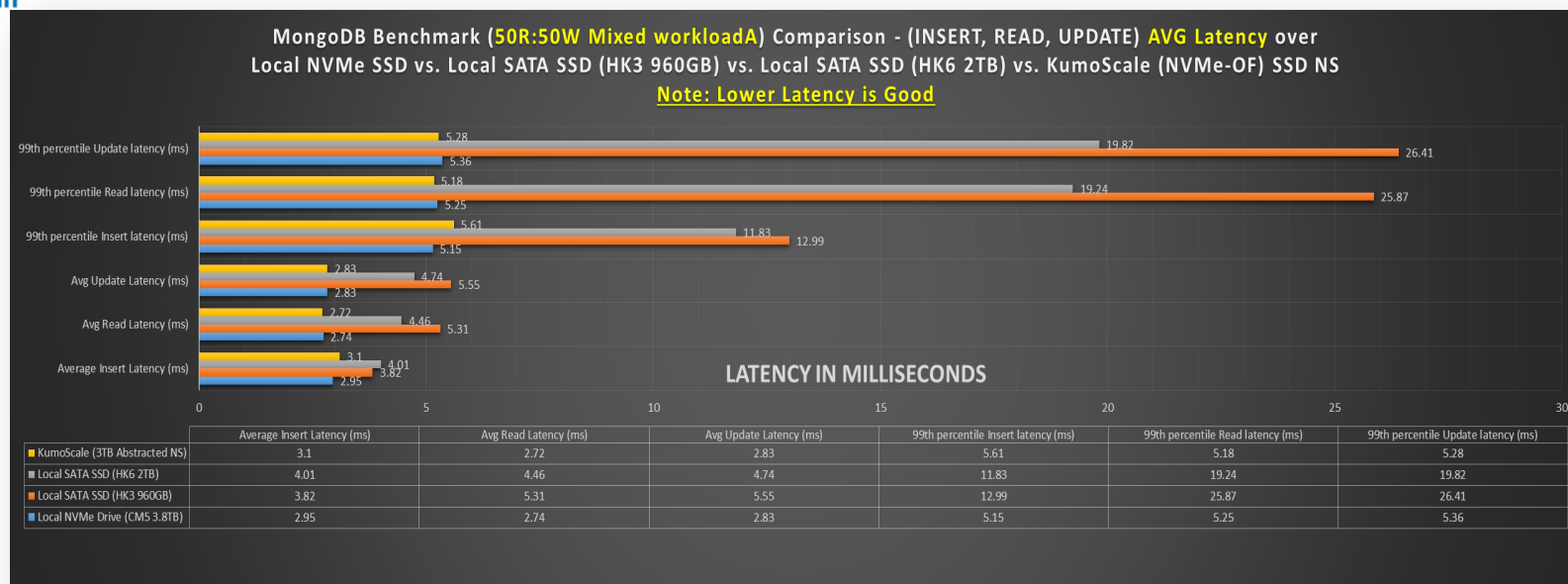
Clustered MongoDB YCSB Results: Test 1 - 50r:50w Mixed workloadA

Metrics	Local NVMe (CM5 3.8T)	KumoScale NVMe-oF (3T Abstracted NS)	Local SATA SSD (HK3 960GB)	Local SATA SSD (HK6 2TB)
Configuration	3 MongoDB Shards : 1 MongoS router : 1 YCSB Engine			
Load Phase: Insert	300 Million records : 1K record size : 200 parallel Threads			
• Average Insert Latency (ms) <small>(lower is good)</small>	2.95	3.10	3.82	4.01
• 99th percentile Insert latency (ms)	5.15	5.61	12.99	11.83
• Insert Speed (operations/sec) <small>(higher is good)</small>	67,605	64,277	52,104	53,354
Run Phase: Read & Update				
• Avg Read Latency (ms)	2.74	2.72	5.31	4.46
• Avg Update Latency (ms)	2.83	2.83	5.55	4.74
• 99th percentile Read latency (ms)	5.25	5.18	25.87	19.24
• 99th percentile Update latency (ms)	5.36	5.28	26.41	19.82
• Read & Update speed (operations/sec) <small>(higher is good)</small>	71,768	71,868	36,751	43,404



Flash Memory Summit

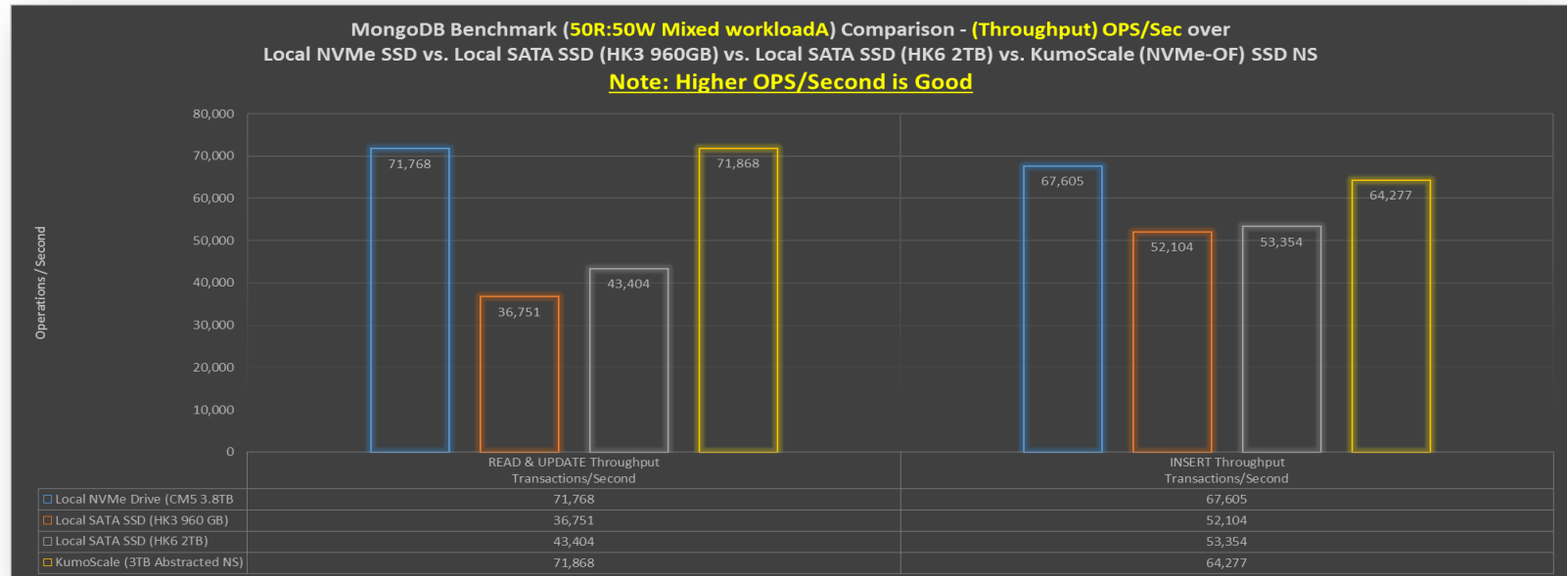
MongoDB™ Benchmark (50R:50W Mixed workloadA) Comparison – (INSERT, READ, UPDATE) AVG Latency



The above chart illustrates the AVG Latency Operations over the given Local Target(s) vs. KumoScale (NVMe-OF) using Mellanox 100GigE Switches. KumoScale Latency shows overall better performance in comparison with others.



MongoDB™ Benchmark (50R:50W Mixed workloadA) Comparison – (INSERT, READ, UPDATE) OPS/Second → Throughput



The above chart illustrates the Operations / Second over the given Local Target(s) vs. KumoScale (NVMe-OF) using Mellanox 100GigE Switches. KumoScale Throughput shows overall better performance in comparison with others.



Flash Memory Summit

Q & A

Flash Memory Summit 2018
Santa Clara, CA

TOSHIBA
Leading Innovation >>>