ceph

# Accelerating Ceph Performance with High Speed Networks and Protocols

Qingchun Song

Sr. Director of Market Development-APJ & China

Mellanox® TECHNOLOGIES

# Mellanox Overview

## Company Headquarters

- Yokneam, Israel
- Sunnyvale, California
- Worldwide Offices

~2,900

**Employees worldwide**
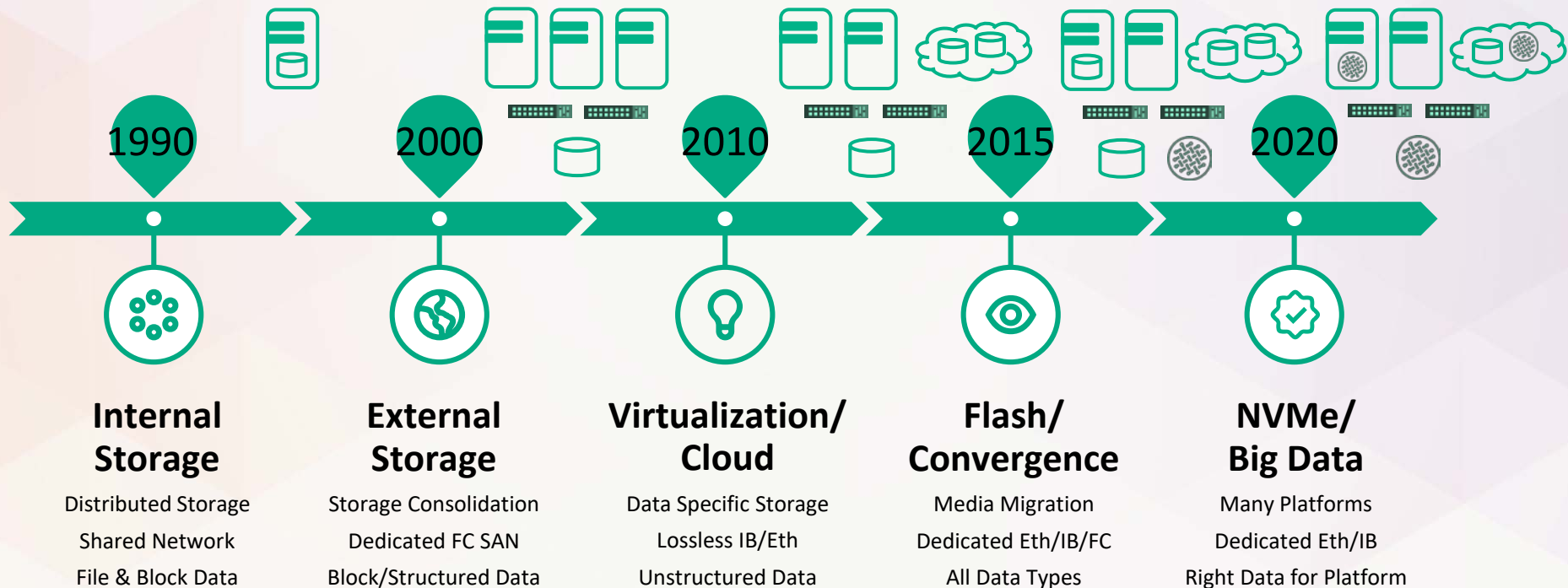
**NASDAQ®**

Ticker: MLNX

# Leadership in Storage Platforms
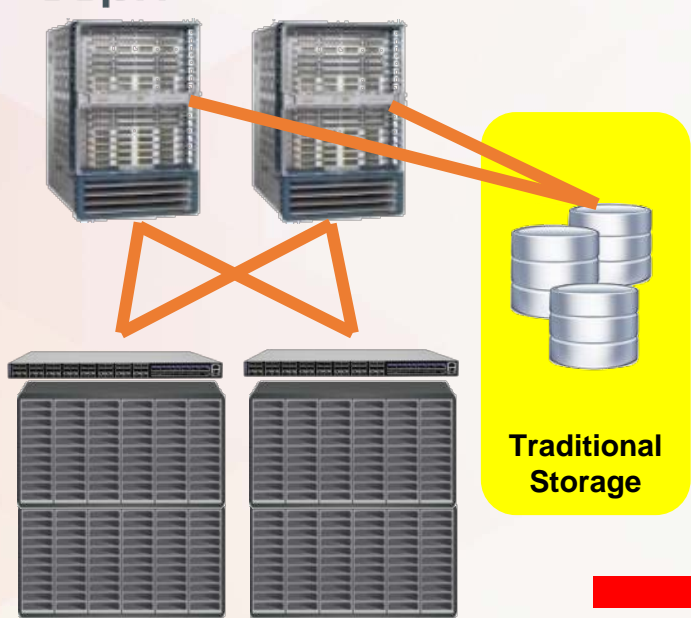
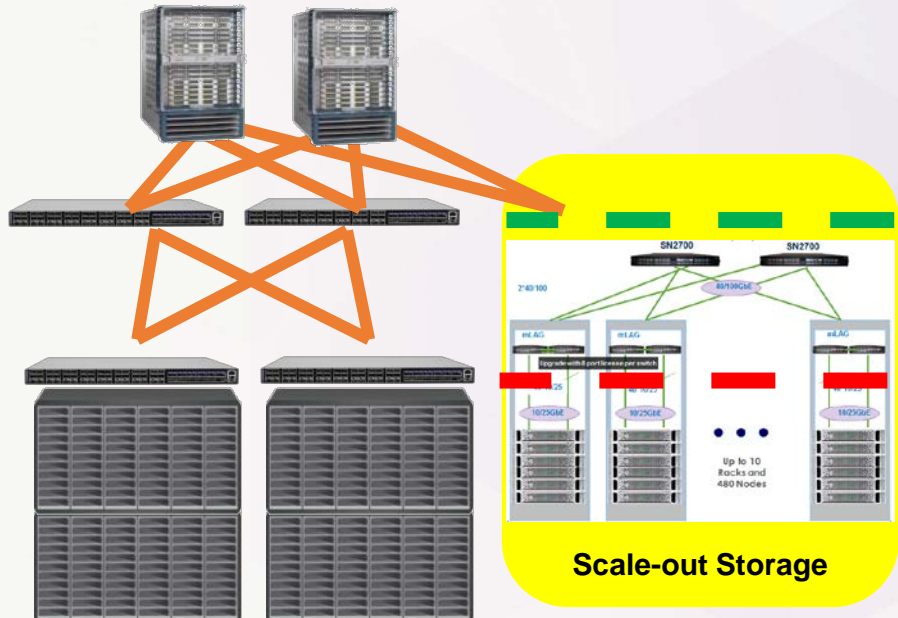Delivering the Highest Data Center Return on Investment

# Storage & Connectivity Evolution

| 1990 | 2000 | 2010 | 2015 | 2020 |
|------|------|------|------|------|
| **Internal Storage** | **External Storage** | **Virtualization/ Cloud** | **Flash/ Convergence** | **NVMe/ Big Data** |
| Distributed Storage | Storage Consolidation | Data Specific Storage | Media Migration | Many Platforms |
| Shared Network | Dedicated FC SAN | Lossless IB/Eth | Dedicated Eth/IB/FC | Dedicated Eth/IB |
| File & Block Data | Block/Structured Data | Unstructured Data | All Data Types | Right Data for Platform |

# Where to Draw the Line?



**Traditional Storage**

**Scale-out Storage**

Legacy DC – FC SAN

Modern DC – Ethernet Storage Fabric

# CePH Work Flow



## Object Store Daemon (OSD) Read and Write Flow

**Read**

**1** Client app issues read request, RADOS sends request to primary OSD

**2** Primary OSD reads data from local disk and completes read request

**Write**

**1** Client app writes data, RADOS sends data to primary OSD

**2** Primary OSD identifies replica OSDs and sends data, writes data to local disk

**3** Replica OSDs write data to local disk, signal completion to primary

**4** Primary OSD signals completion to client app

Compute Node: App, RBD, RADOS

OSD / Disk — Server

# Storage or Data Bottlenec... Bandwidth

## Ceph Cluster Overview

- **Ceph Clients**
  - Block/Object/File system storage
  - User space or kernel driver
- **Peer to Peer via Ethernet**
  - Direct access to storage
  - No centralized metadata = no bottlenecks
- **Ceph Storage Nodes**
  - Data distributed and replicated across nodes
  - No single point of failure
  - Scale capacity and performance with additional nodes

**Client Servers**

| Application Guest OS | Application Guest OS | Application Guest OS | Application Guest OS | Application Guest OS | Application Guest OS |
|---|---|---|---|---|---|
| KVM | | KVM | | KVM | |
| RBD Object File | | RBD Object File | | RBD Object File | |

Spine ---------- Spine

Leaf -------------------------------- Leaf

| OSD | OSD | OSD | OSD | OSD | OSD | OSD | OSD | OSD | OSD | OSD | OSD |
|---|---|---|---|---|---|---|---|---|---|---|---|

MON SSD SSD ... SSD SSD ... MON SSD SSD ... SSD SSD

**Storage Servers**

OSD read:
- Client(App <-> RBD <-> RADOS) <-> NIC <-> Leaf <-> Spine <-> Leaf <-> NIC <->OSD <-> NVMe

OSD write:
- Client(App <-> RBD <-> RADOS) <-> NIC <-> Leaf <-> Spine <-> Leaf <-> NIC <->OSD <-> NVMe <-> OSD <-> NIC <-> Leaf <-> Spine <-> Leaf <-> NIC <->OSD <-> NVMe

# Ceph Bandwidth Performance Improvement

- Aggregate performance of 4 Ceph servers
  - 25GbE has 92% more bandwidth than 10GbE
  - 25GbE has 86% more IOPS than 10GbE
- Internet search results seem to recommend one 10GbE NIC for each ~15 HDDs in an OSD
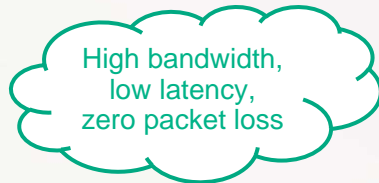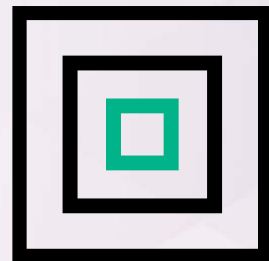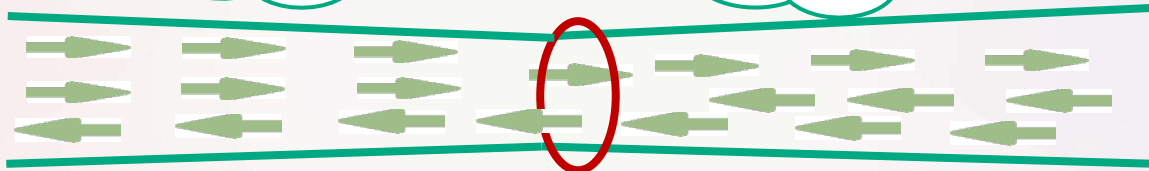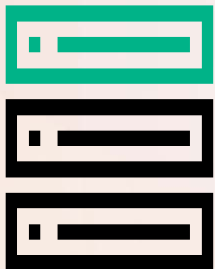  - Mirantis, Red Hat, Supermicro, etc.

# Storage or Data Bottleneck: L[  ]y

**Predictable Performance, Deterministic & Secure Fabrics**

| Servers | Fabrics | Storage |
|---|---|---|

Block, file, and object storage

Faster, more predictable performance

Simplified security and management

High bandwidth, low latency, zero packet loss

- Higher processing capability
- High-density virtualization

- Move to All-flash
- Faster protocols – NVMe-oF

**Data Center modernization requires Future Proof, faster, lossless Ethernet Storage Fabrics**

Ceph中国社区

IT大咖说
知识共享平台

# RDMA Is The Key For Storage Latency?



**Efficient Data Movement (RDMA)**

Kernel Bypass     Protocol Offload

Throughput (Gbytes/sec)

adapter based transport

# RDMA Enables Efficient Data Movement

**100GbE With CPU Onload**

**100 GbE With Network Offload**



**CPU Onload Penalties**
- *Half the Throughput*
- *Twice the Latency*
- *Higher CPU Consumption*

2X Better Bandwidth

Half the Latency

33% Lower CPU

See the demo: https://www.youtube.com/watch?v=u8ZYhUjSUoI

- Without RDMA
  - 5.7 GB/s throughput
  - 20-26% CPU utilization
  - 4 cores 100% consumed by moving data

- With Hardware RDMA
  - 11.1 GB/s throughput at half the latency
  - 13-14% CPU utilization
  - More CPU power for applications, better ROI

- Conservative Results: 44%~60% more IOPS

- RDMA offers significant benefits to Ceph performance for small block size (4KB) IOPS.
  - 2 OSDs with 4 clients, RDMA allowed 44% more IOPS.
  - 4 OSDs and 4 clients, RDMA allowed 60% more IOPS.

- Best Results: 3x Higher IOPS
- RDMA's biggest benefit for Ceph block storage
  - High IOPS workloads
  - Small block sizes (<32KB)
- Enable > 10GB/s from single node
- Enable < 10usec latency under load





Ceph Read IOPS: TCP vs. RDMA

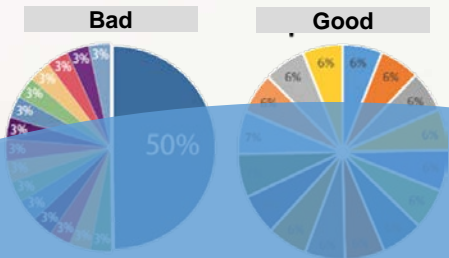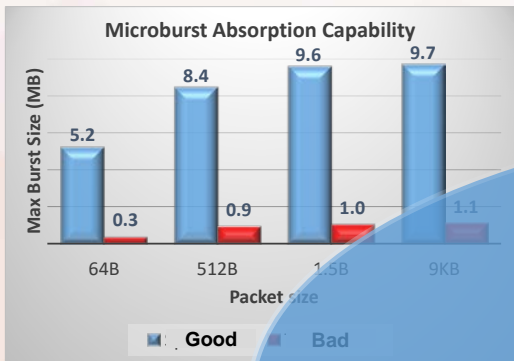# RDMA: Mitigates Meltdown Mess, Stops Spectre Security Slowdown

# CePH RDMA Status

- CePH RDMA working group
  - Mellanox
  - Xsky
  - Samsung
  - SanDisk
  - RedHat
- The latest stable CePH RDMA version
  - https://github.com/Mellanox/ceph/tree/luminous-12.1.0-rdma
- Bring Up Ceph RDMA - Developer's Guide
  - https://community.mellanox.com/docs/DOC-2721
- RDMA/RoCE Configuration Guide
  - https://community.mellanox.com/docs/DOC-2283

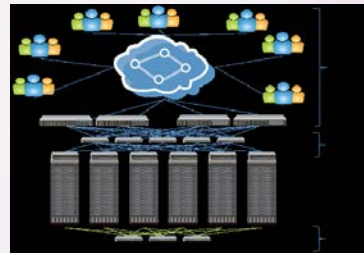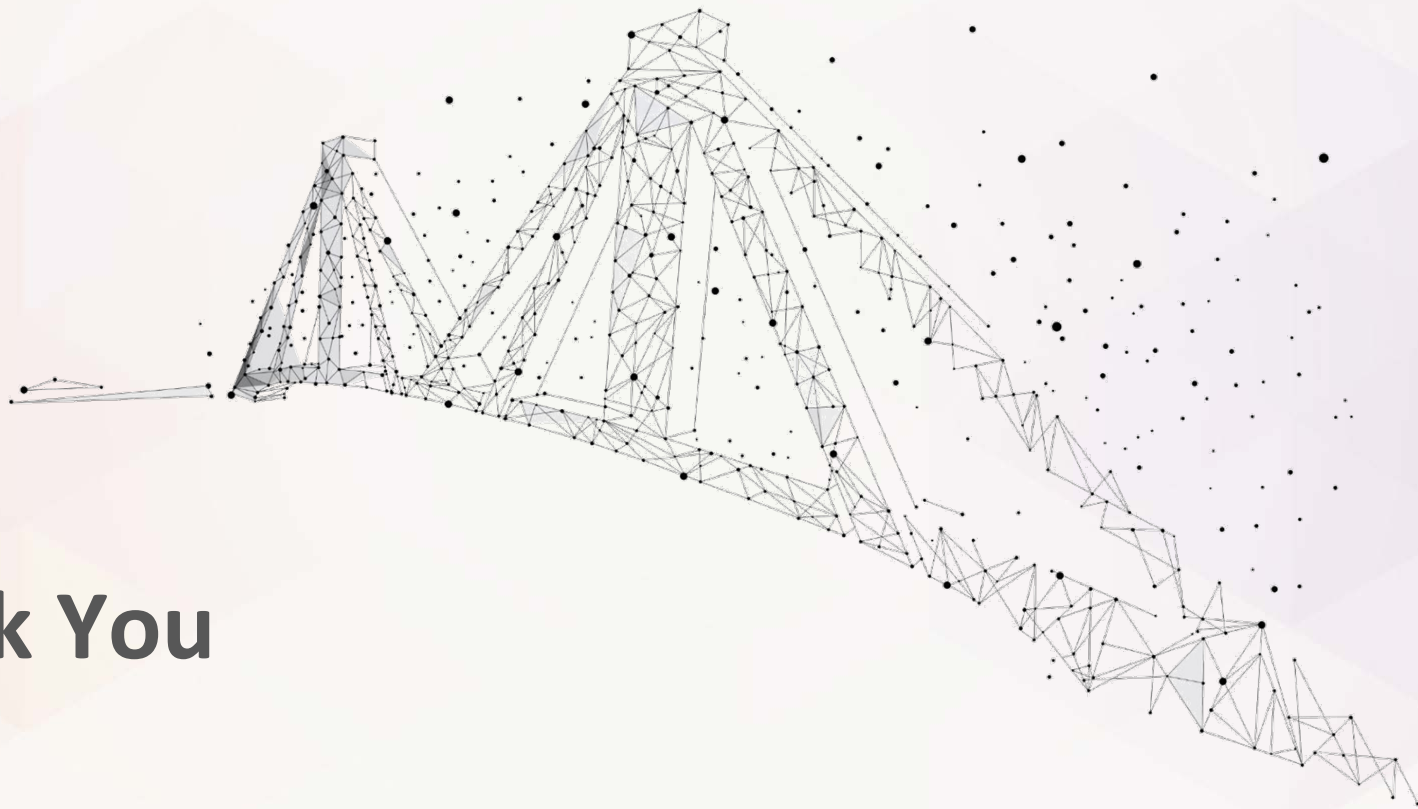# **Summary**

- Ceph Benefits from Faster Network
  - 10GbE is not enough!
- RDMA further optimizes Ceph performance
- Reduce the impact from Meltdown/Spectre fixes
- ESF(Ethernet Storage Fabric) is trend

**Thank You**