

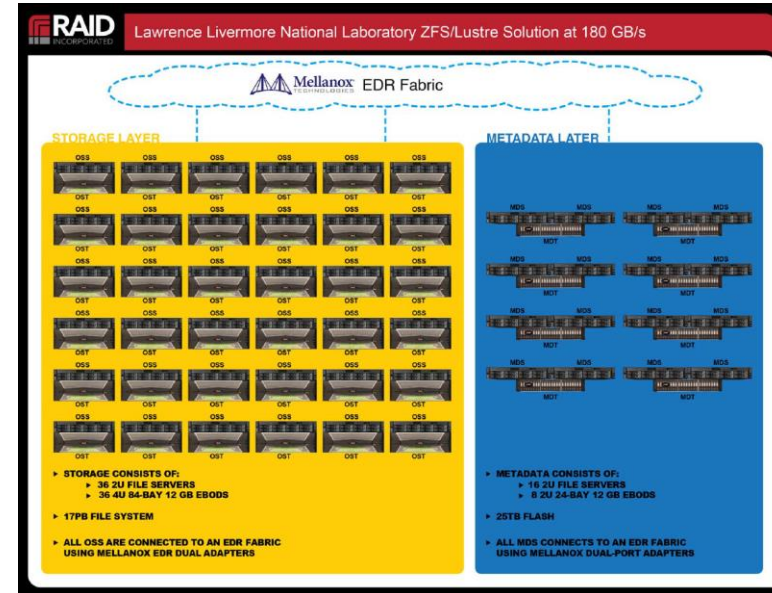
InfiniBand Networked Flash Storage

Superior Performance, Efficiency and Scalability

Motti Beck – Director Enterprise Market Development, Mellanox Technologies

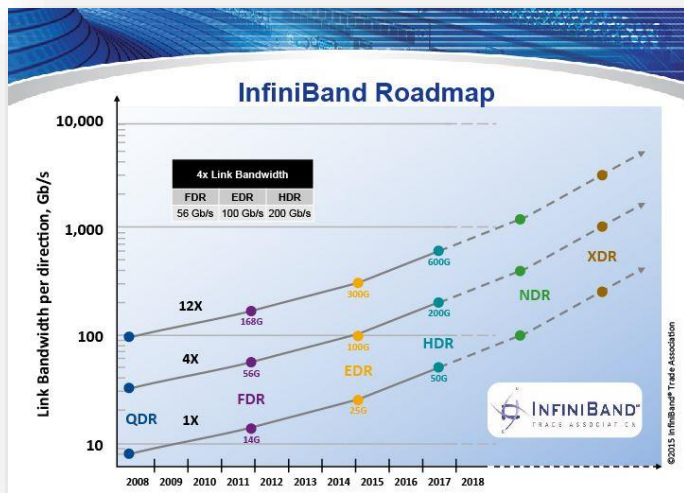
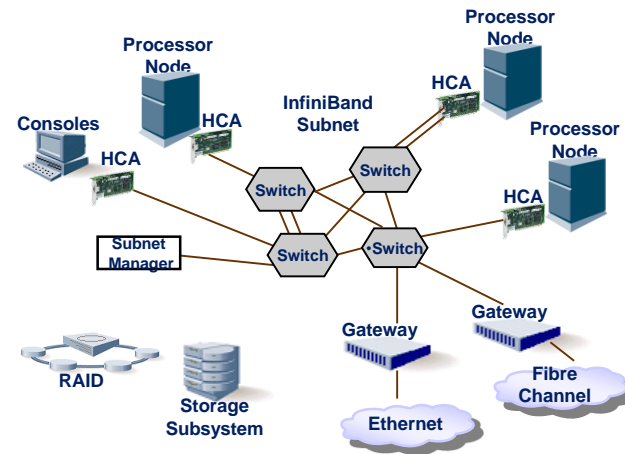
17PB File System with Sustainable 180GB/s

- The parallel file system will run [Lustre 2.8 with ZFS](#) OSDs and multiple metadata servers.
- The Lustre file system contains 36 OSS nodes, with each node capable of 5GB/s of sustained data performance, and 16 metadata servers with **25TB of SSD storage** capacity.
- The solution is anchored by enterprise 4U 84 bay 12Gb SAS JBODs, LSI/Avago 12Gb SAS adapters, **Mellanox EDR IB**, HGST 12Gb Enterprise SAS HDDs, and Intel server technologies.
- The file system incorporates 6 scalable storage units each containing six Lustre OSS and six 4U-84Bay **JBODs with 480 8TB SAS drives**. The solution will be employing ZFS on Linux with raidz2 data parity protection. **Resiliency is provided by multipath and HA failover connectivity**, intended to eliminate single points of failure.



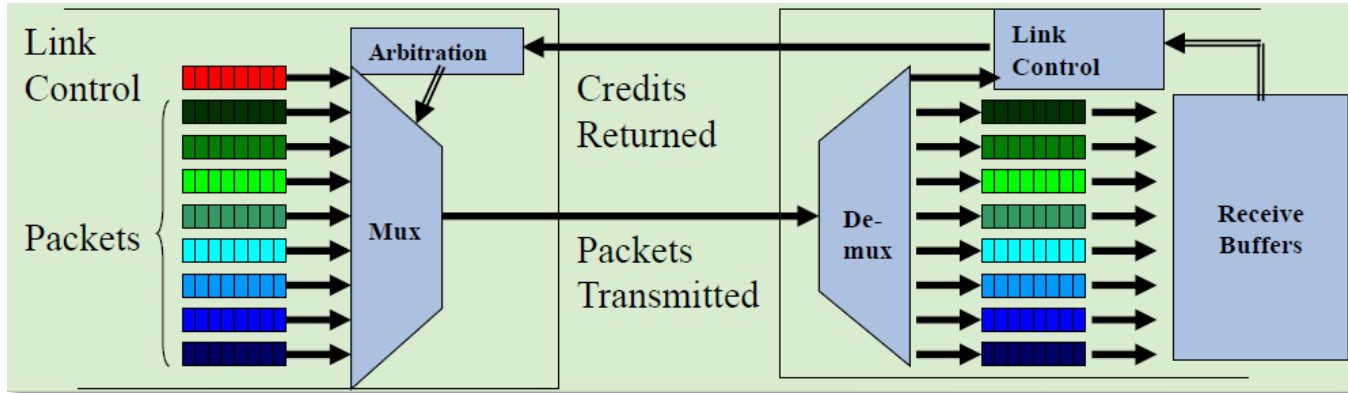
Why InfiniBand in Storage

- Industry standard defined by the InfiniBand Trade Association (IBTA)
- Defines System Area Network architecture
 - Comprehensive specification - from physical to applications
- Simplicity drives higher efficiency and resiliency
 - Reliable, lossless, self-managed fabric
 - Transport offload - Remote Direct Memory Access (RDMA)
 - Centralized fabric management – Subnet Manger (SM)



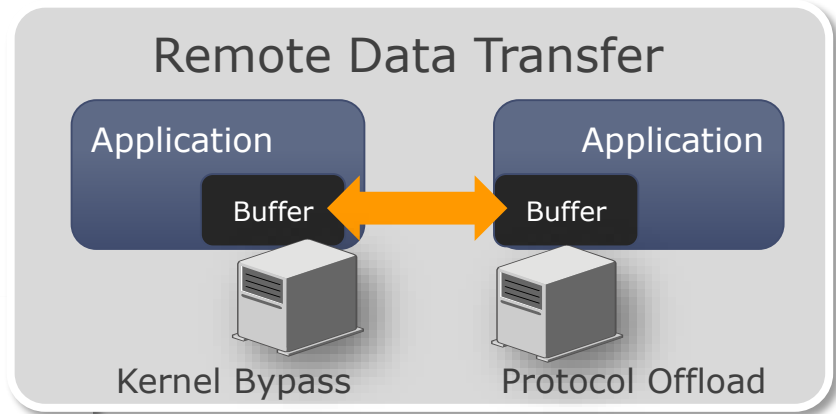
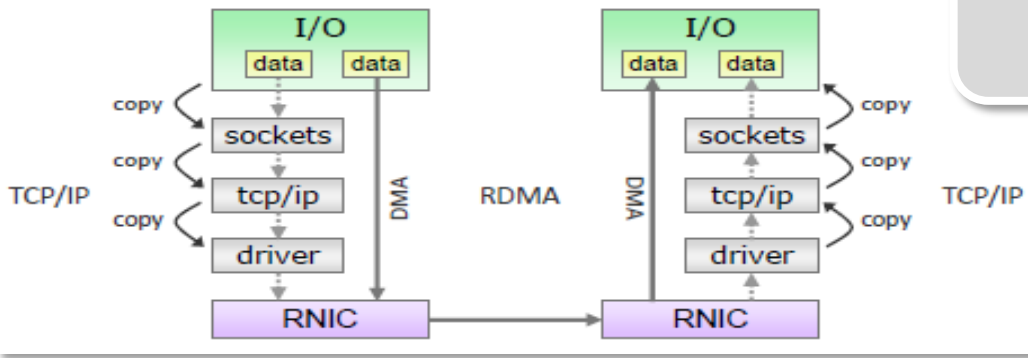
Reliable, Lossless, Self-Managed Fabric

- End to end flow control
 - Credit based flow control for each link
- End to end Congestion management
- Automatic path migration



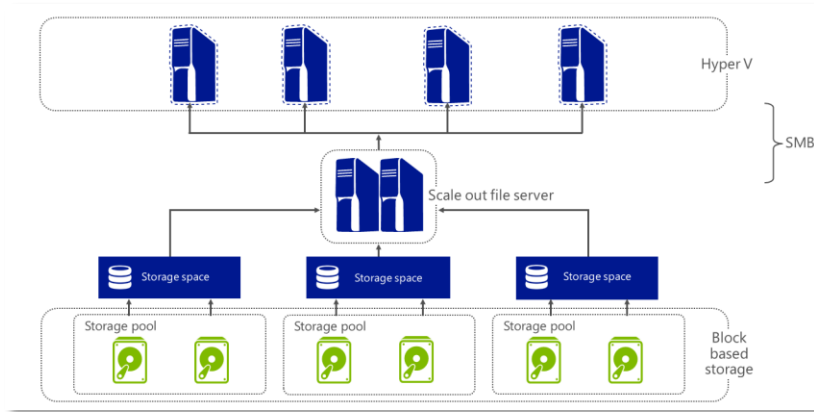
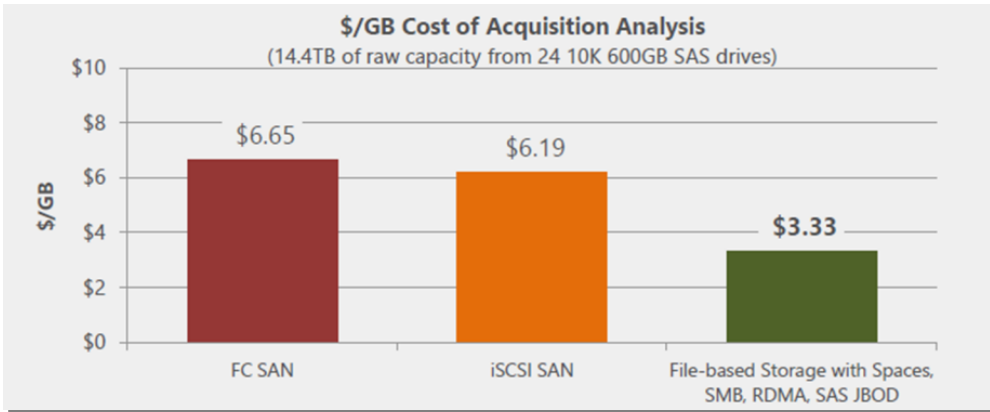
Remote Direct Memory Access RDMA

- Transport offload
- Kernel bypass



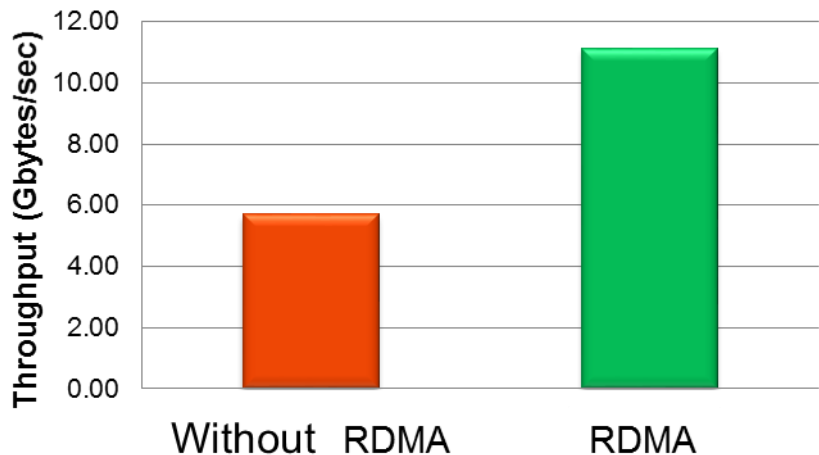
InfiniBand Cuts SAN Cost by 50%

- Delivers SAN-like functionality from the Windows Stack
 - Using SMB Direct (SMB 3.0 over RDMA)
- Utilize inexpensive, industry-standard, commodity hardware
 - Eliminate the cost of proprietary hardware and software from SAN solutions

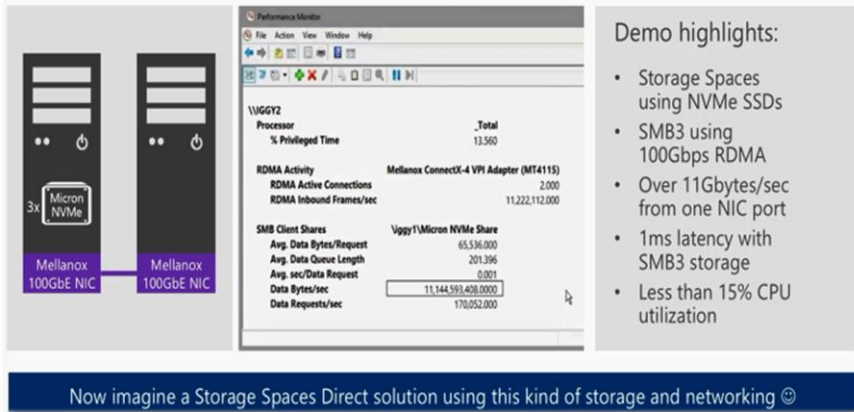


RDMA Delivers Higher Cloud Efficiency

Microsoft Storage Spaces Throughput



Demo Summary: 100GbE and NVMe
A technology demonstration of things to come for Microsoft SDS...



Demo highlights:

- Storage Spaces using NVMe SSDs
- SMB3 using 100Gbps RDMA
- Over 11Gbytes/sec from one NIC port
- 1ms latency with SMB3 storage
- Less than 15% CPU utilization

Now imagine a Storage Spaces Direct solution using this kind of storage and networking ©

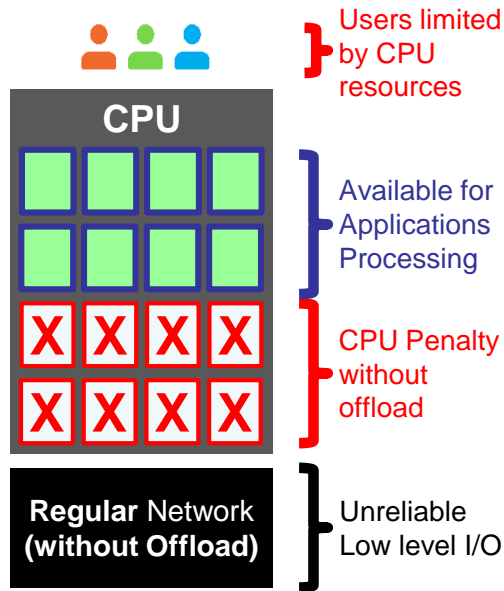
- 2X better performance with RDMA
 - 2X higher bandwidth & 2X better CPU efficiency
- RDMA achieves full Flash storage bandwidth
 - Remote storage without compromises



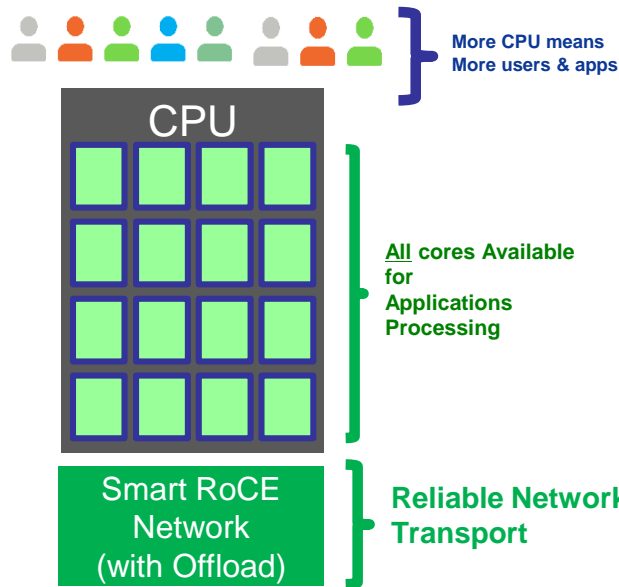
RDMA Frees Up CPU for Application Processing

54Gb/s

90Gb/s



- Half of CPU capacity consumed moving data
 - Even though achieving only half throughput
- Cores unavailable for application processing



- All CPU available to run applications
 - Better efficiency = more users
- Smart network delivers better TCO



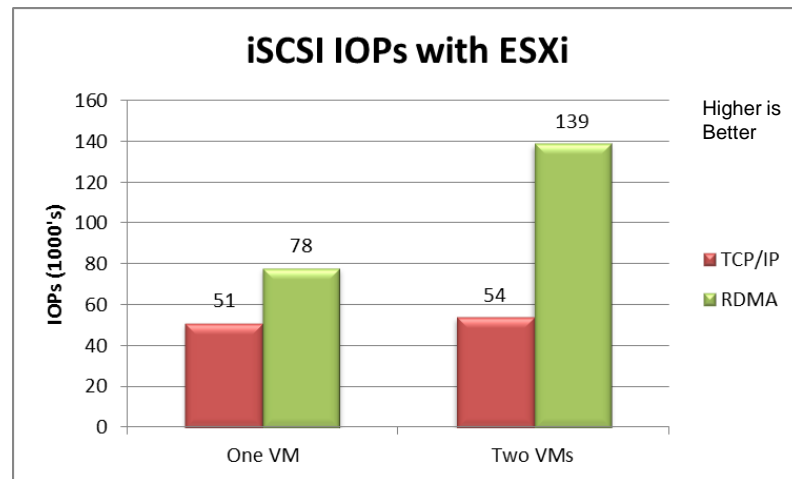
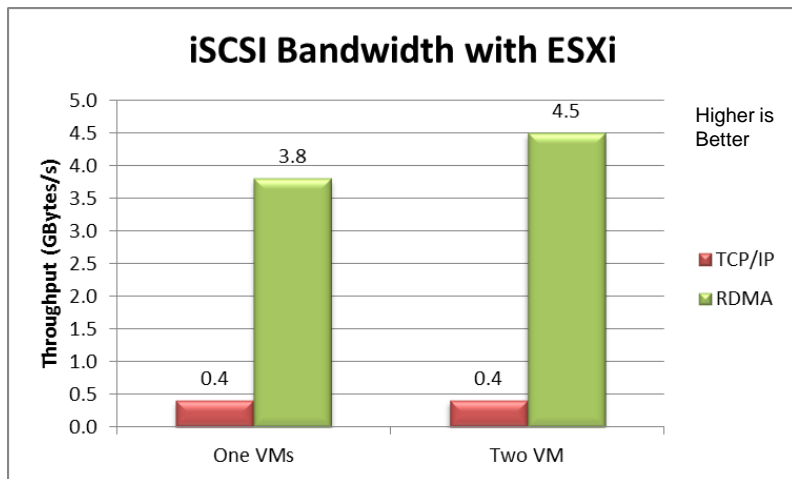
RDMA Removes Storage Bottlenecks in Enterprise

- Stellar Storage Performance with 100Gb/s
- Microsoft labs benchmarked Storage Spaces @ 100Gb/s
 - 4 Node Cluster
 - NVMe Flash & ConnectX-4 100GbE Adapters
- Storage Spaces Direct over RoCE
- Achieved 480Gb/s (60GB/s) throughput
 - Transmit entire content of Wikipedia in 5 seconds
- Storage Spaces Direct will be part of Azure Stack
 - Private cloud package for 2017



<https://blogs.technet.microsoft.com/filecab/2016/05/11/s2dthroughputp5/>

Storage Acceleration Running iSER*



* iSCSI over RDMA

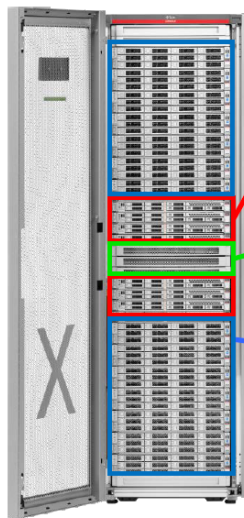
- RDMA Superior Across the Board
 - Throughput & IOP's
 - Efficiency & CPU Utilization
 - Scalability

Test Setup: ESXi 5.0, 2 VMs, 2 LUNS per VM

*iSCSI Extensions for RDMA

10x Bandwidth Performance Advantage vs TCP/IP
2.5x IOPS Performance With iSER Initiator

Exadata X5-2 Product Components



- **Scale-Out Database Servers**
 - **Two 18-core x86 Processors (36 cores)**
 - Oracle Linux 6
 - Oracle Database Enterprise Edition
 - Oracle VM (optional)
 - Oracle Database options (optional)
- **Fastest Internal Fabric**
 - 40 Gb/s InfiniBand
 - Ethernet External Connectivity
- **Scale-Out Intelligent Storage**
 - **High-Capacity Storage Server**
 - **Extreme Flash Storage Server**
 - **Exadata Storage Server Software**



36 cores per server
256 – 768 GB DRAM





Dominant in Storage and Database Platforms

DataDirect
NETWORKS

EMC²
EMC² XtremIO

FUJITSU



ORACLE®

Microsoft



NIMBUS DATA

IBM
xiv tms

SEAGATE

TERADATA

TOSHIBA

violin
MEMORY

WD Western
Digital



Leading in Performance, Efficiency and Scalability