



Flash Memory Summit



Erasure Code Offload for Distributed Software Defined Storage

Dror Goldenberg
VP Software Architecture
Mellanox Technologies



Software Defined Storage – Why?



Scale out

- More capacity
- More performance

• MORE PERFORMANCE

Cheaper

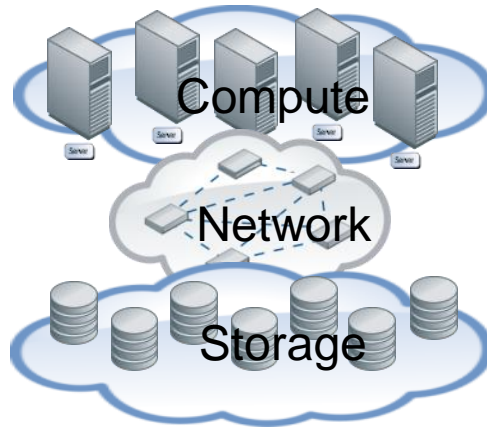
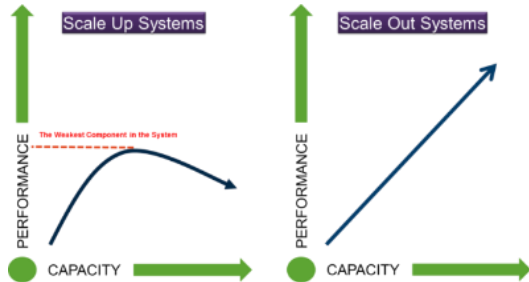
- Resource utilization
- Software only

• LOW COST

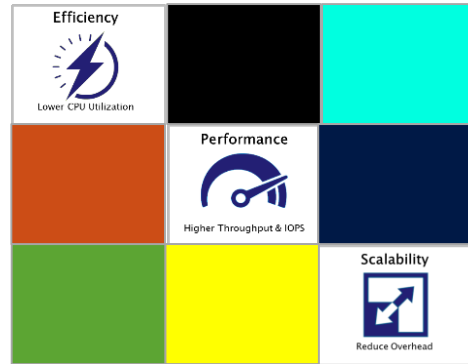
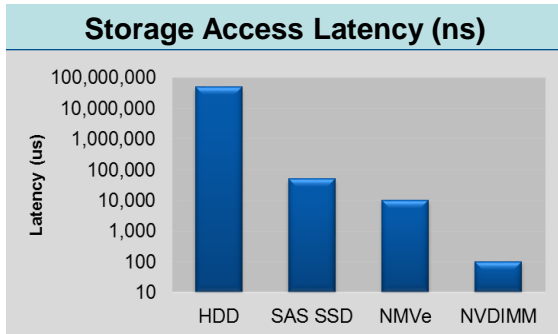
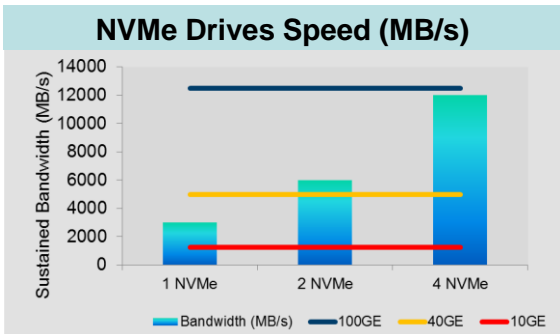
Flexible

- Multiple architectures
- Upgradeable

• ADAPTIVE



Software Defined Storage: Scaling and Performing



Network BW
Efficient Scale Out

Low latency
Efficient Storage Disaggregation

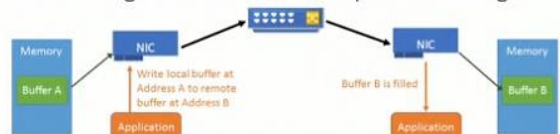
Offloading
Efficiency, Lower TCO

Azure Cloud Storage: Erasure Coding at Scale

“To make storage cheaper we use lots more network!
How do we make Azure Storage scale? RoCE (RDMA over Converged Ethernet) enabled at 40GbE for Windows Azure Storage, achieving **massive COGS savings**”

Keynote
Albert Greenberg, Microsoft
SDN Azure Infrastructure

RDMA – High Performance Transport for Storage



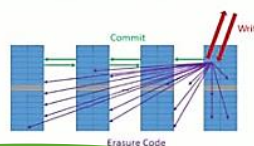
- Remote DMA primitives (e.g. Read address, Write address) implemented on-NIC
 - Zero Copy (NIC handles all transfers via DMA)
 - Zero CPU Utilization at 40Gbps** (NIC handles all packetization)
 - <2µs E2E latency
- RoCE enables Infiniband RDMA transport over IP/Ethernet network (all L3)
- Enabled at 40GbE for Windows Azure Storage, achieving massive COGS savings by eliminating many CPUs in the rack

All the logic is in the host:

Software Defined Storage now scales with the Software Defined Network

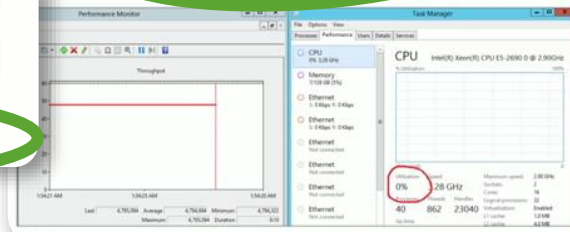
Storage is Software Defined, Too

- We want to make storage clusters scale cheaply on commodity servers
- Erasure Coding provides durability of 3-copy writes with small (<1.5x) overhead by distributing coded blocks over many servers
- Lots of network I/O for each storage I/O



To make storage cheaper, we use lots more network!

Just so we're clear...
40Gbps of I/O with 0% CPU





Ensuring Data Availability



Standalone



Cluster



Hot swap



RAID 0



RAID 1



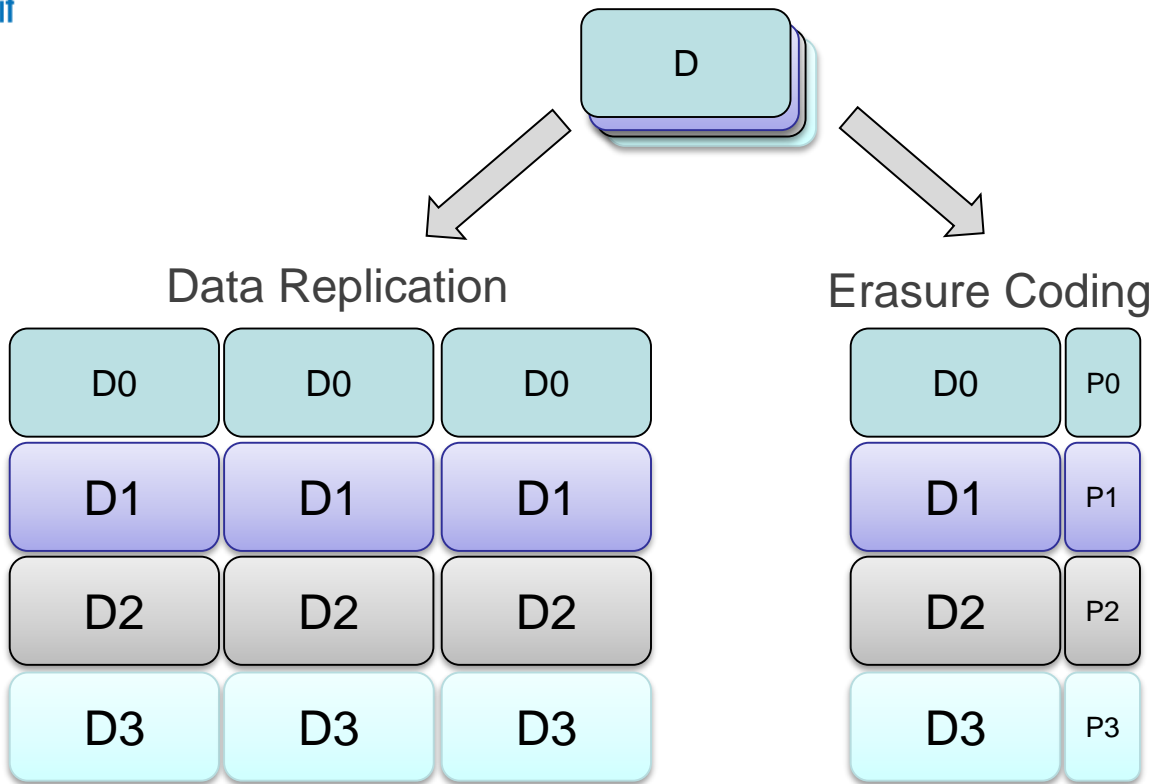
RAID 5



RAID 0+1



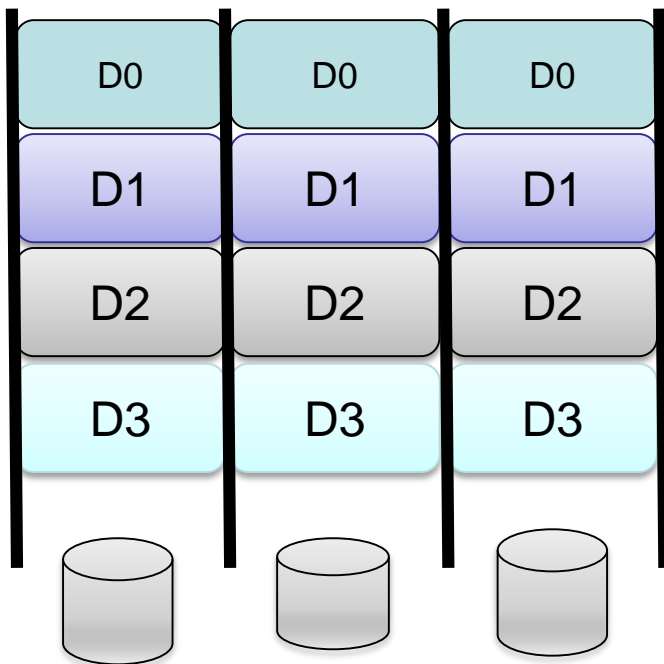
Ensuring Data Availability



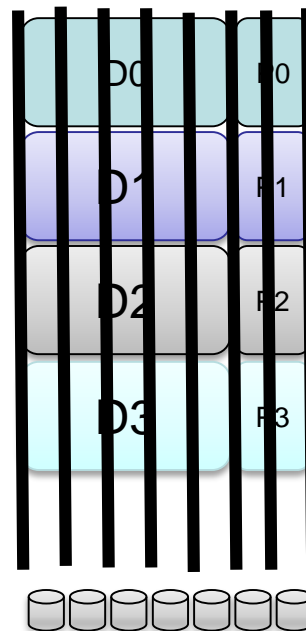


Ensuring Data Availability

Data Replication



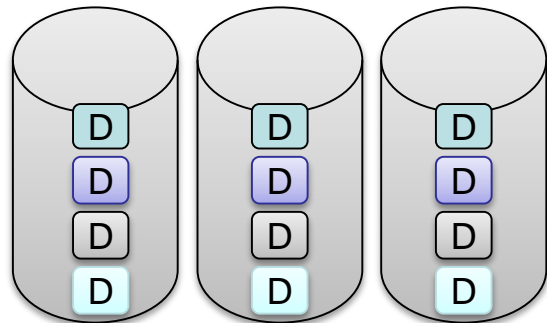
Erasure Coding





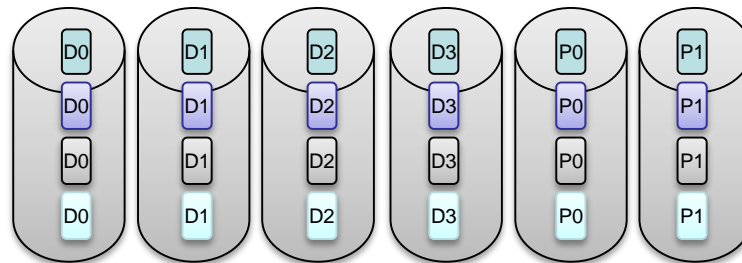
Ensuring Data Availability - Summary

Data Replication



- Capacity: 3x data (typical)
- Resilient to 2 failures

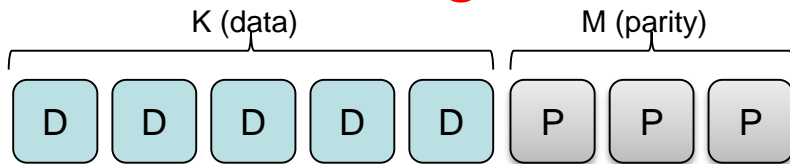
Erasure Coding



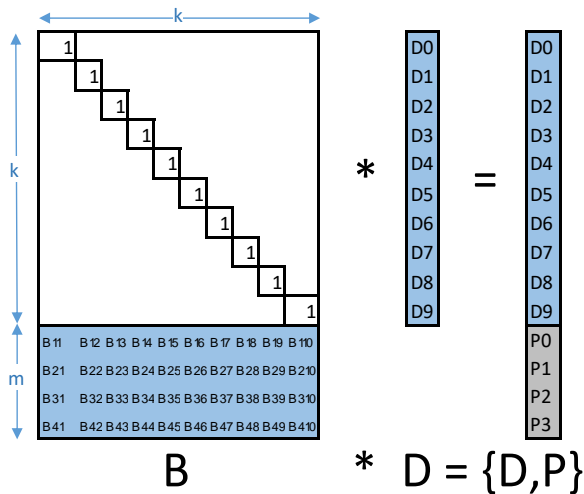
- Capacity: 1.4x data (typical)
- Better failure resilience
 - e.g. 10+4: 1.4x capacity, 4 failures
- CPU & network hungry
- Longer & data intensive rebuild



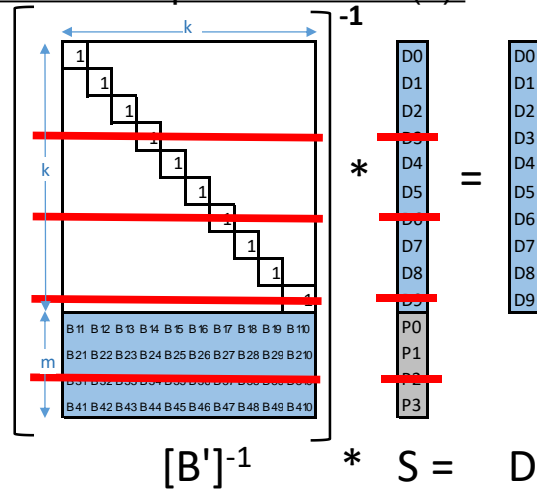
Erasure Coding Calculation



Calculating parity – Encoding
Matrix multiplication in GF(8):



Failure recovery (rebuild) – Decoding
Matrix multiplication in GF(8):





Offload APIs - 101

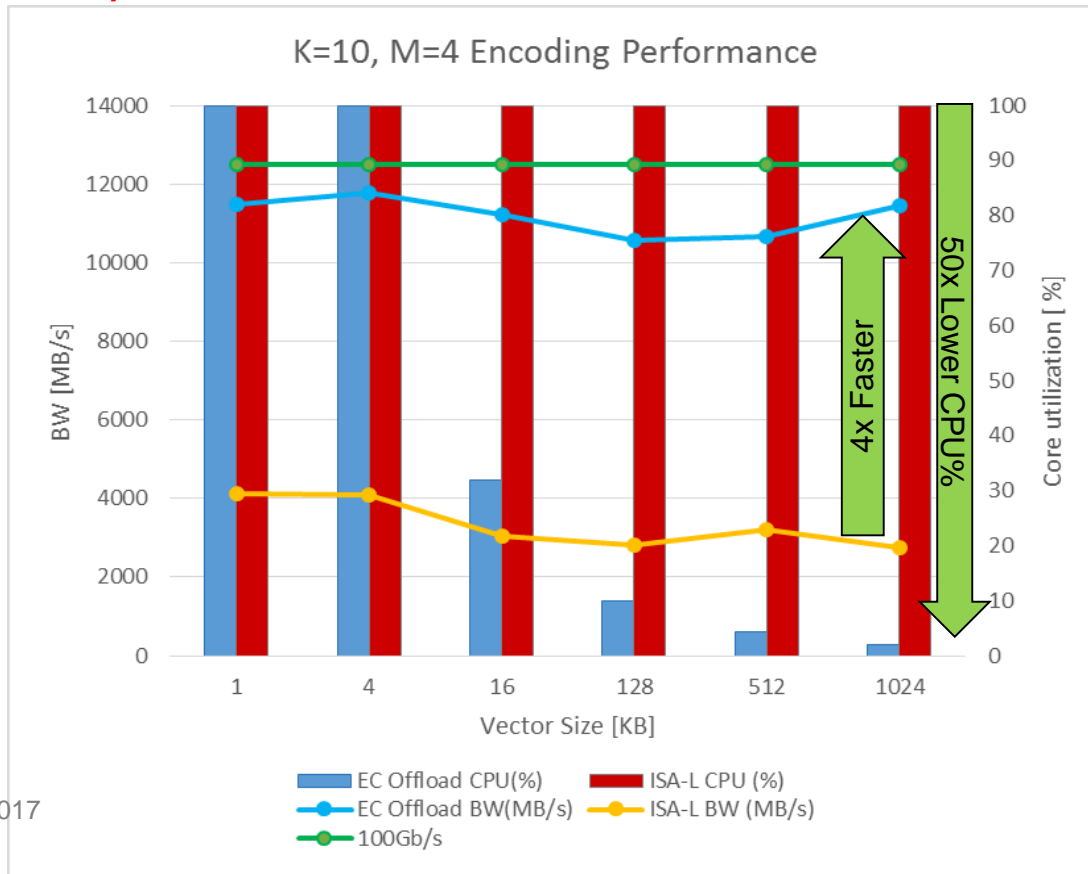
- Erasure Coding **Onload**
 - Computation all done on CPU (CPU at 100%)
 - Cache/TLB pollution
 - Example: ISA-L

- Erasure Coding **Offload**
 - Computation all done in accelerator (CPU at 0%)
 - Cache/TLB unaffected
 - Example: ec_offload APIs

Efficient data movement - critical element to enable distributed scale out erasure coding
RDMA - asynchronous offload networking API



Encoding Performance (Single Core) – x86 ISA-L vs Offload



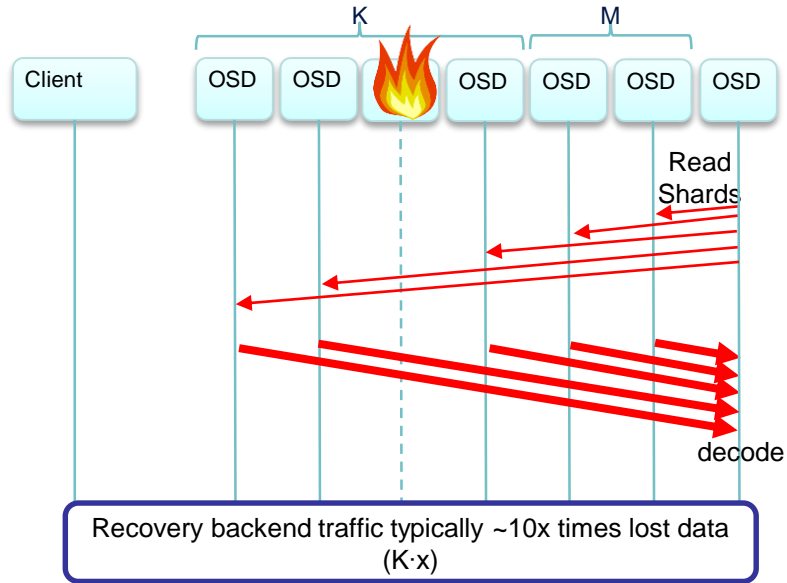


Encoding Performance (Single Core) – ARM vs Offload





Network Traffic – Rebuild (Erasure Coding)



- Example - Time to rebuild (10+4)
 - Net networking time to move data
 - 20TB system @40GE 14.4hrs
 - 200TB system @40GE 144.4hrs
- Similar flows for scrubbing

Data recovery for erasure coding can be 10x more network intensive
Efficient data movement (RDMA) and efficient network are critical elements



Summary

Flash Memory Summit



- Cloud infrastructure requires efficient and scalable storage
- Software Defined Storage drives scale out storage
 - Performance, capacity, flexibility, TCO
- Erasure Codes enable data availability at lower capacity
 - Tradeoff: CPU & Network intensive
- Erasure Codes offload offers
 - Better performance (per I/O, efficient rebuild)
 - Lower cost
- Network efficiency (bandwidth, latency, offload) – an important enabler
- Library available today
 - Integration to scale out storage systems underway (Reed Solomon, LRC)
 - Can be used for local (host based) erasure coding, e.g. RAID5 codes



Flash Memory Summit



Thank You!