



Flash Memory Summit

INVT-202A-1: Handling the Network Requirements of High-Speed NVMe SSDs

Manoj Wadekar, Storage Engineer, Facebook
Rob Davis, VP Storage Technology, Mellanox



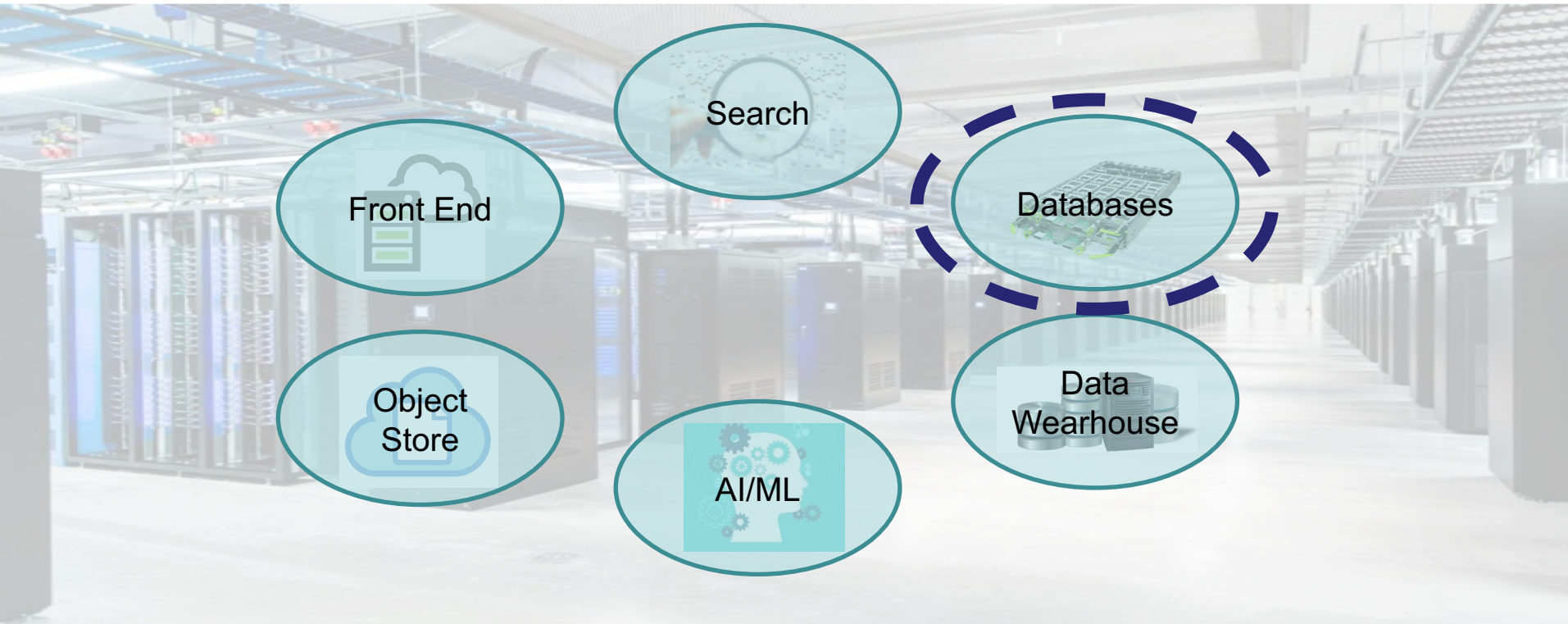
Flash Memory Summit 2019
Santa Clara, CA





Flash Memory Summit

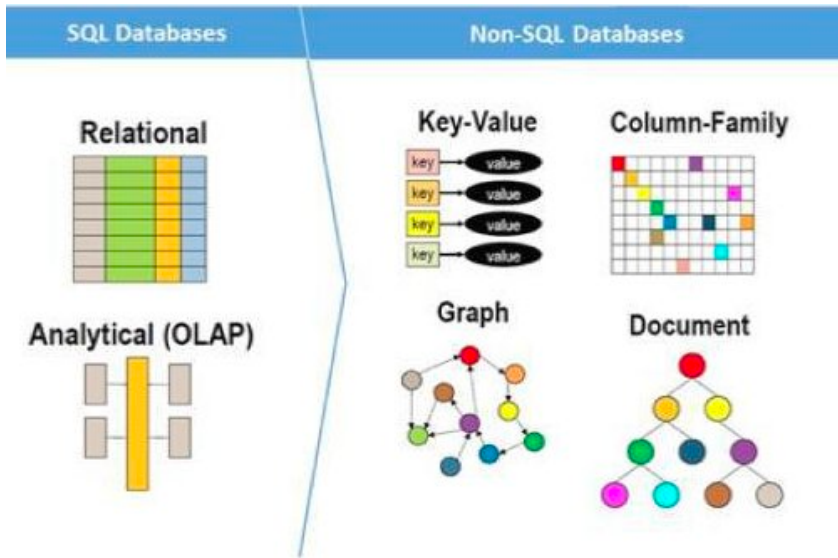
Typical Hyperscale Server Infrastructure





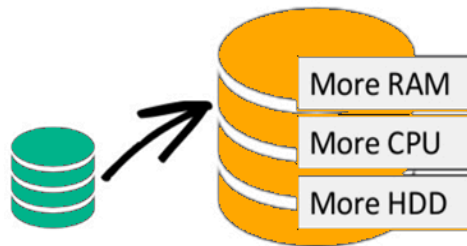
Databases for Hyperscale

Flash Memory Summit



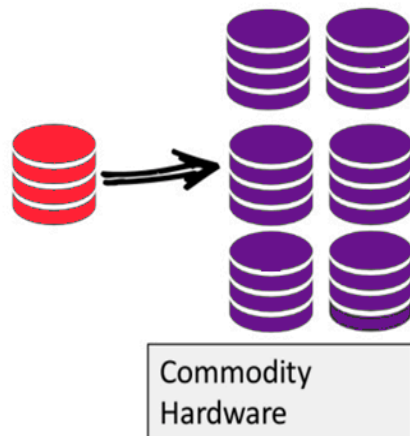
Open Source Software

Scale-Up (*vertical* scaling):



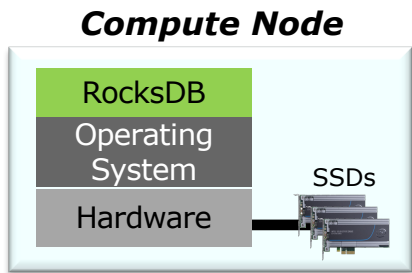
Open Hardware: E.g. OCP

Scale-Out (*horizontal* scaling):





Disaggregated Storage Architecture

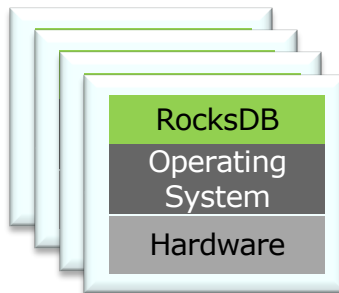


Disaggregated



- **Local attached storage**
- **Static binding**
- **Stranded capacity, IOPS**
- Inefficient, increased TCO

Compute Node

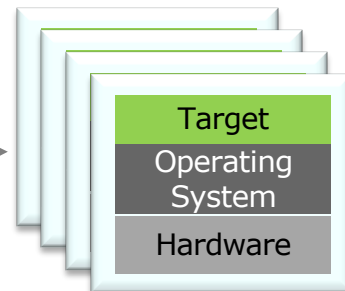


- **Logical disaggregation**
- Consumes physical or logical block devices
- **Dynamic binding** based on workload requirements
- Efficient, improved TCO

iSCSI, NVMe-oF, etc.



Head Node



PCIe



JBOF



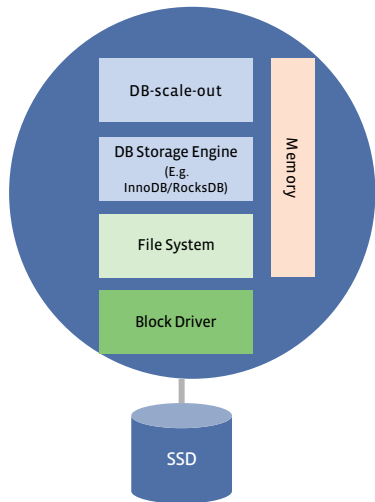
- **Physical disaggregation**
- **Static binding**
- **Shared resources**
- Target can expose physical or logical devices



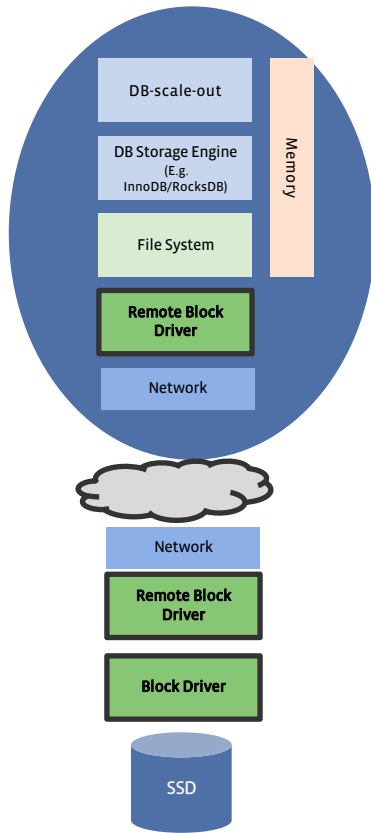
Storage Disaggregation

Flash Memory Summit

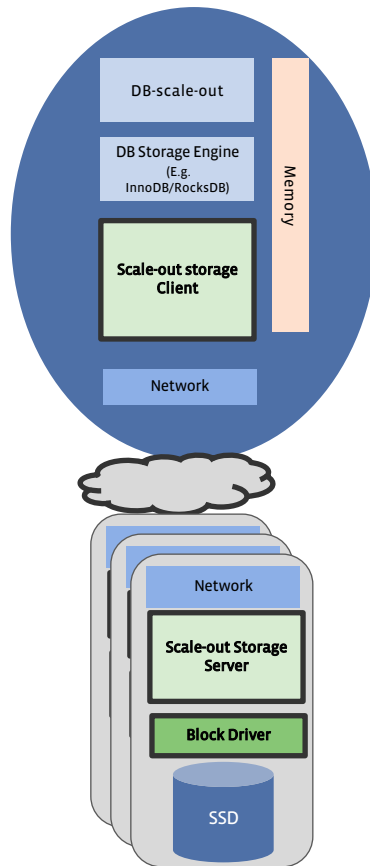
Local Storage



Remote Block



Remote storage Service



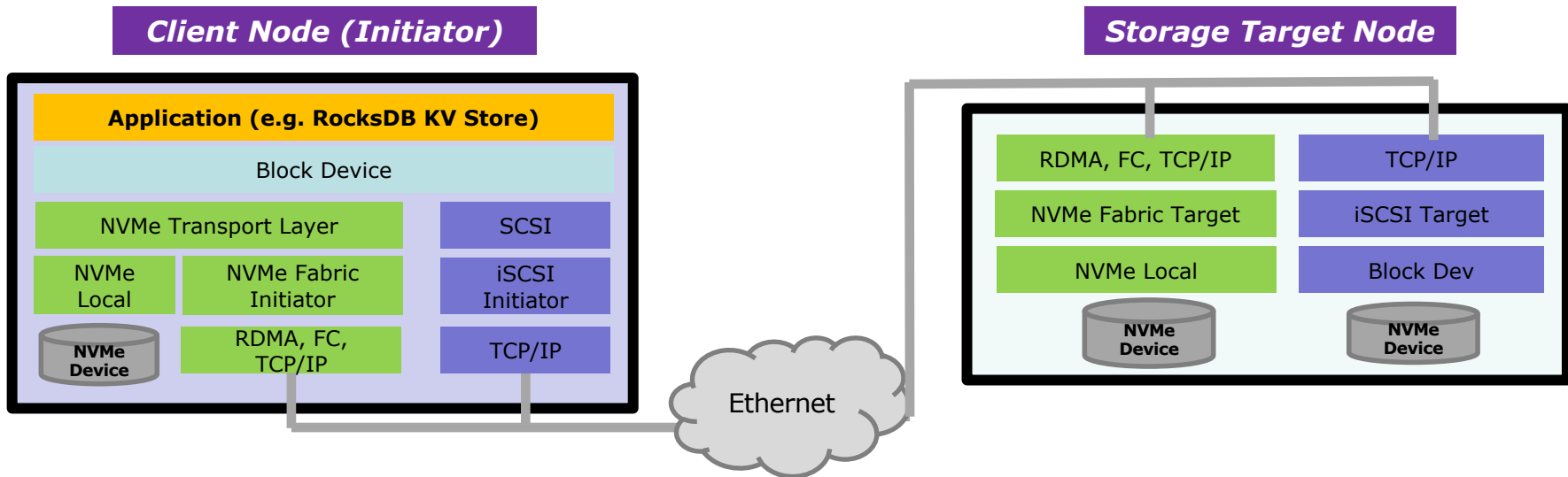


Remote Block Storage



STORAGE

Flash Memory Summit

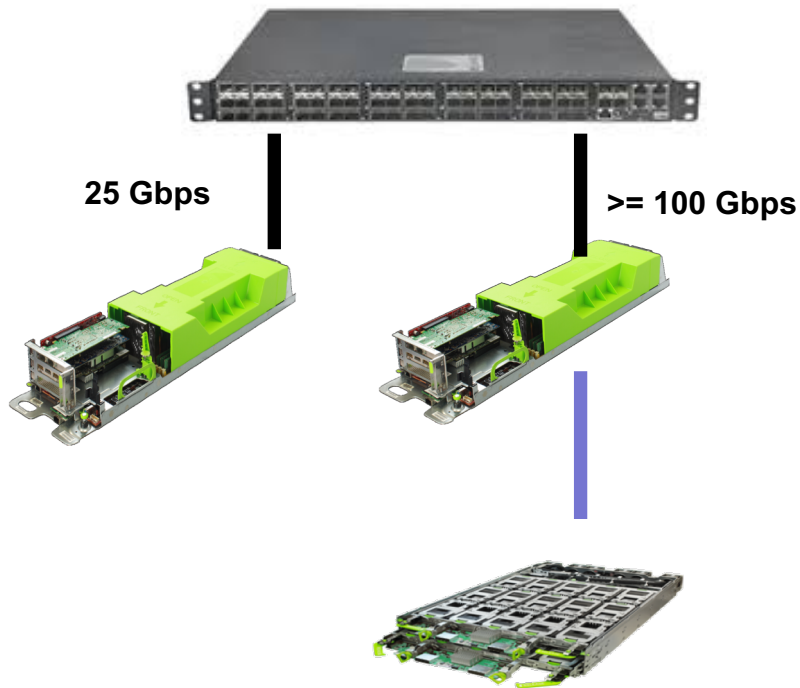


- **Enables sharing** of NVMe flash storage over network
- Can use **traditional** block protocols (e.g. iSCSI) or **NVMe optimized** protocols (e.g., NVMe/TCP)
- **NVMe over Fabrics** – supports multiple transports, extends **NVMe efficiency over network**
 - Poll and interrupt mode architecture
 - Kernel and user mode implementations



Flash Memory Summit

Remote Block Storage: OCP HW

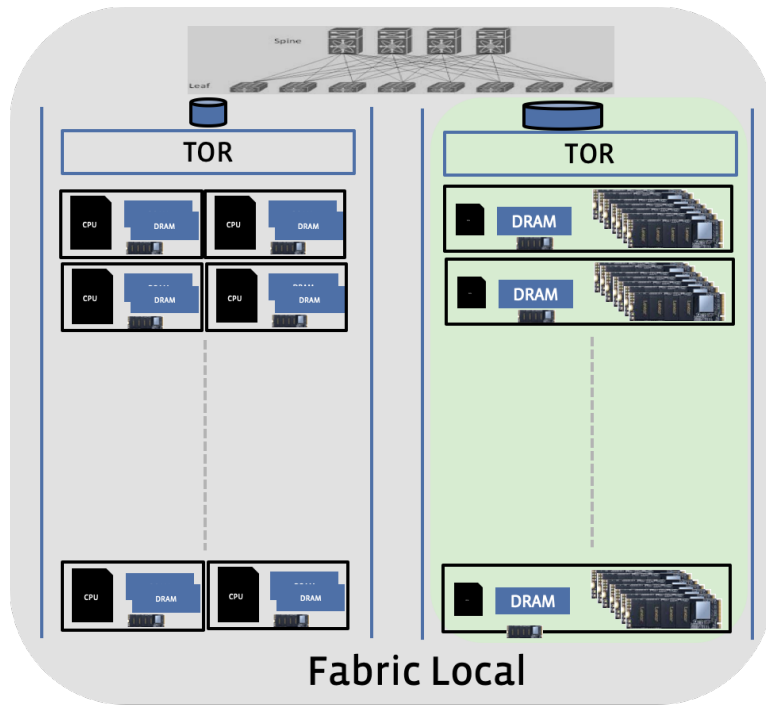
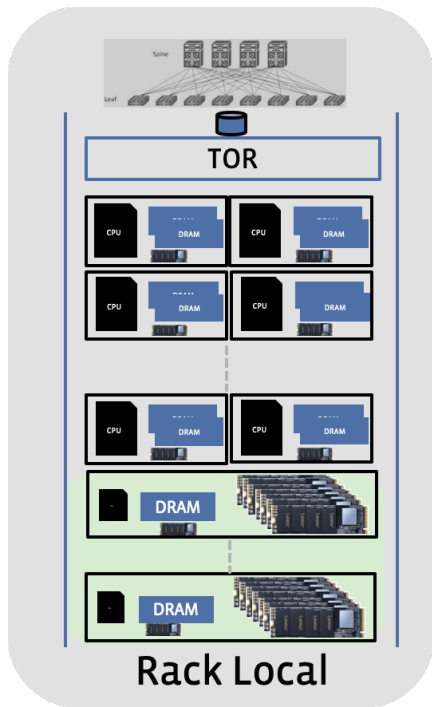


- [FB Lightning](#) supports 30 M.2 NVMe SSDs
- Storage can be accessed over Ethernet using [Tioga Pass server](#)



Network: Speed

Flash Memory Summit



- Disaggregation
- Multi-tenant access
- High bandwidth applications
- ..

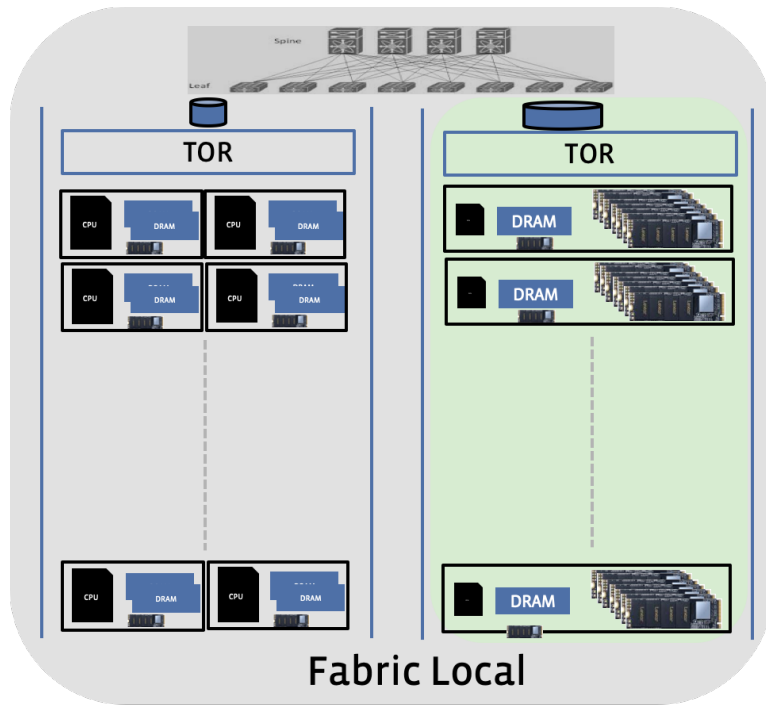
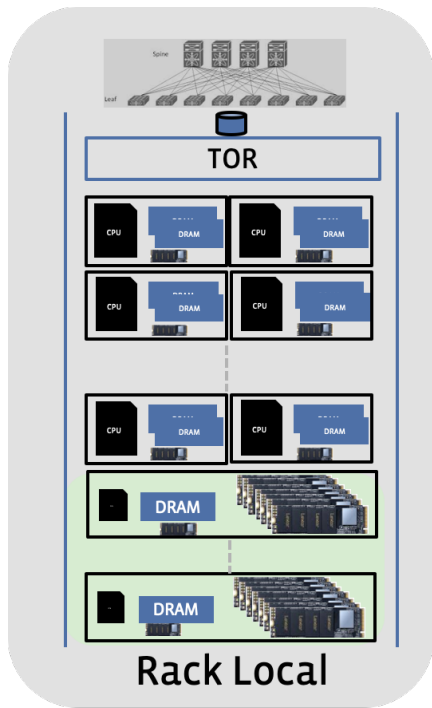
100s Gbps

Tbps



Network: Latency

Flash Memory Summit



- Link Latency
- Stack Latency
- Number of hops
- Queuing
- Congestion
- ..

100s Gbps

Tbps



Flash Memory Summit

Summary – hyperscale storage

- Storage Disaggregation important for hyperscale efficiency
- It depends on high performance network
- Is Ethernet Ready?

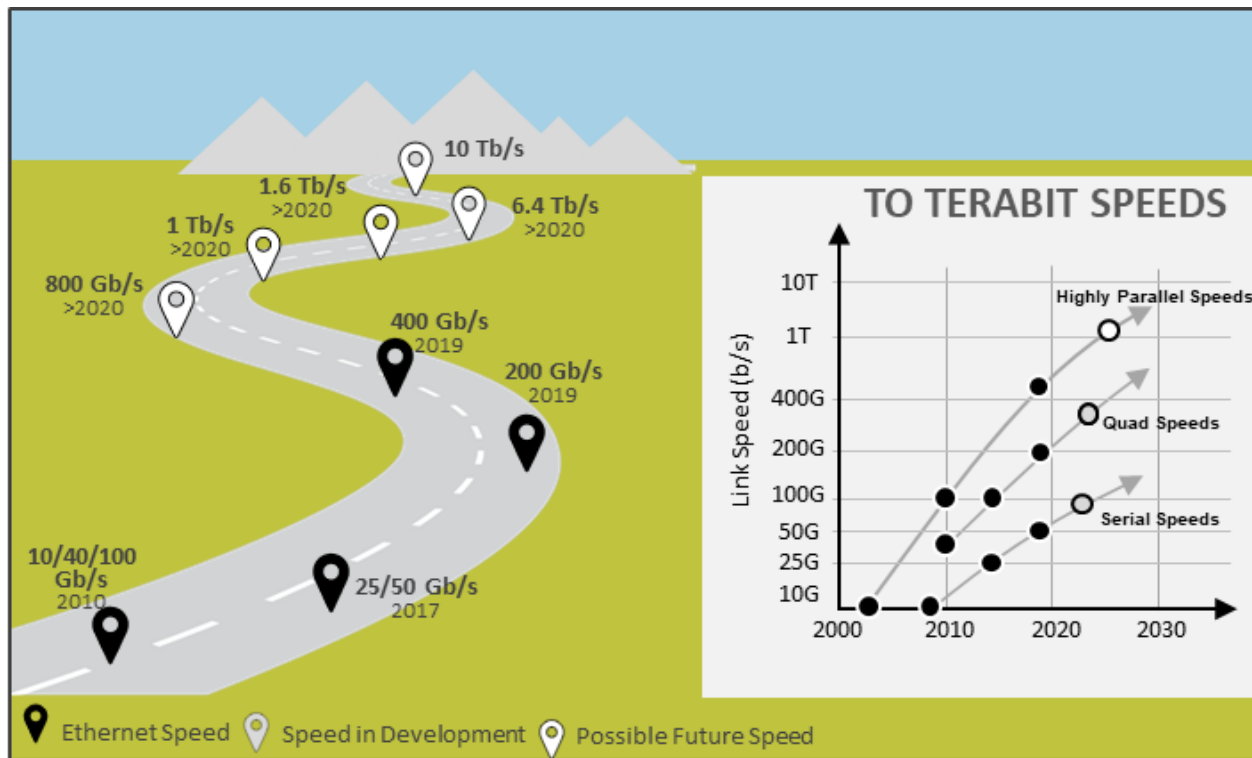


Flash Memory Summit

Ethernet is Ready!

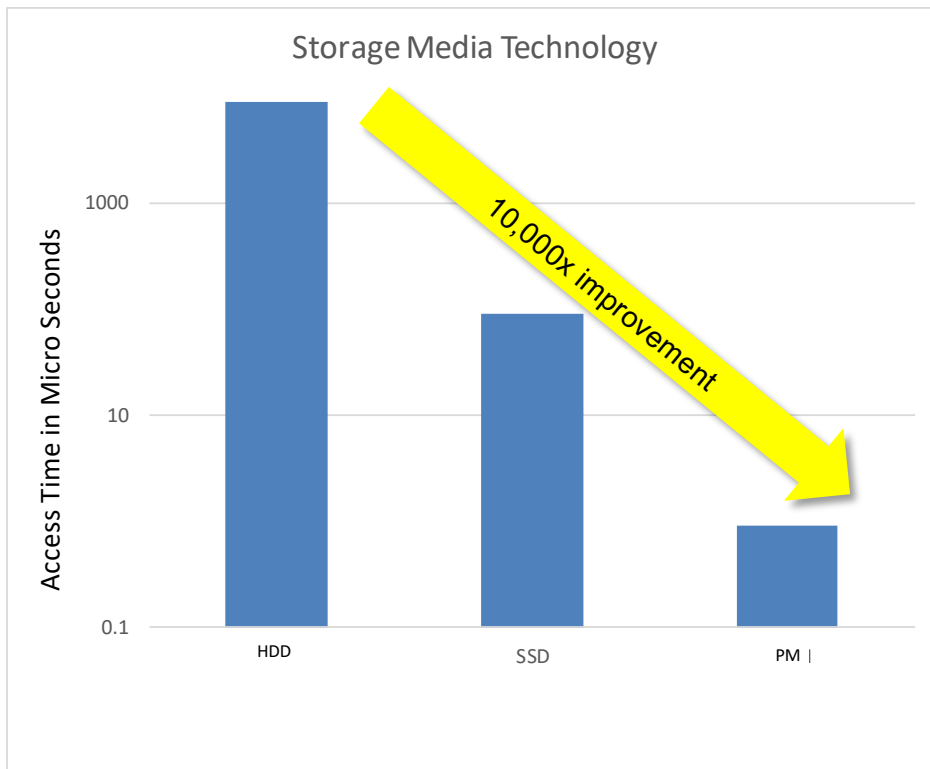


Ethernet Technology Roadmap





We Need Network Speed for Flash!





High Performance Networking is Here Today

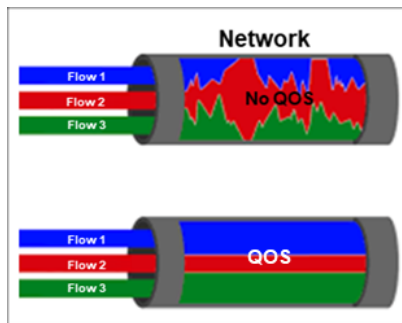


End-to-End 25, 40, 50, 100, 200GbE and soon 400Gb



But High Bandwidth is Only Part of the Solution

- We also Need Low Latency
- Effective Performant Network Congestion Control
- QOS
- Security





Low Latency

Network Congestion Control

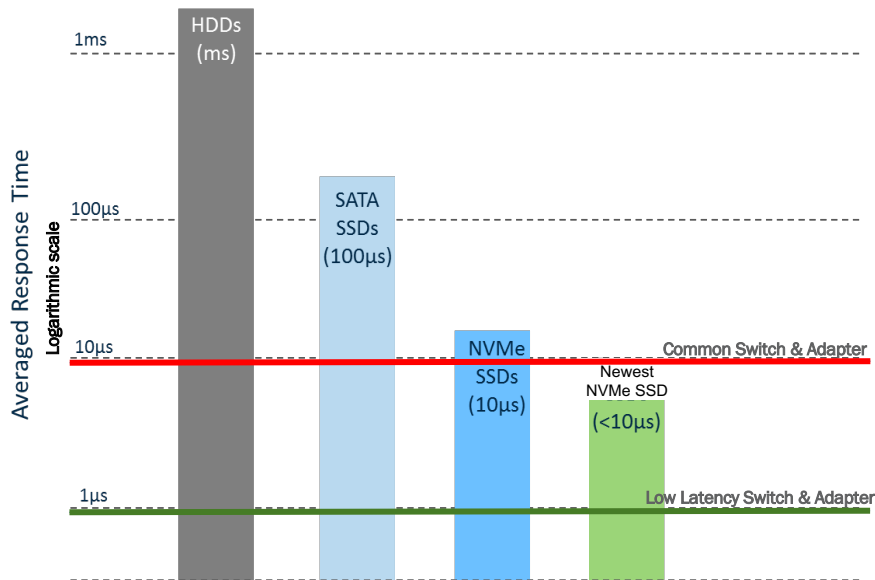
QOS

Security

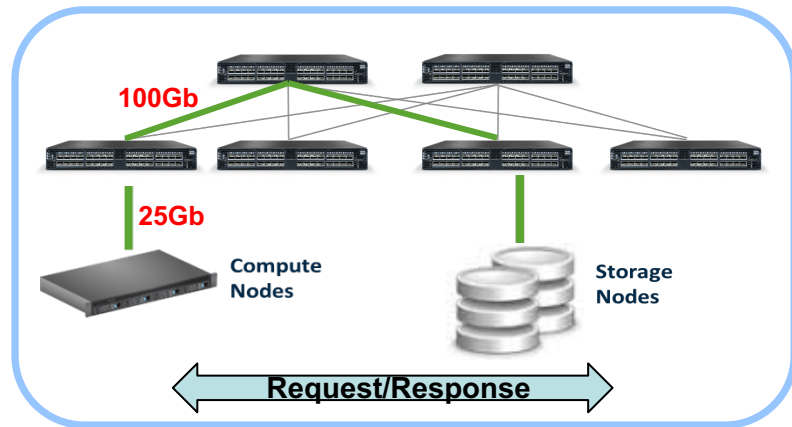
At Scale



Importance of Latency with Flash Storage

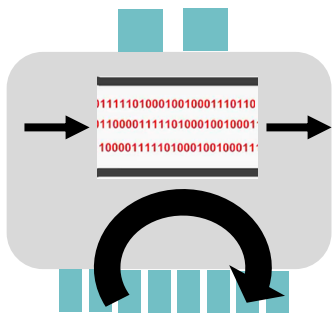


Network hops multiply latency

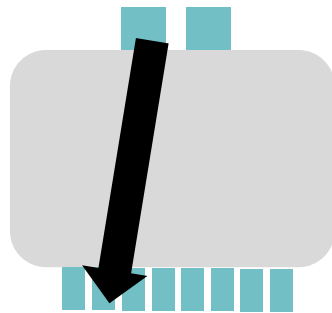




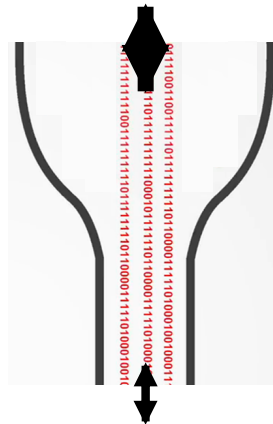
Intelligent Cut-Through Reduces or Eliminates Store & Forward Latency



Downlink to Downlink
Full Cut-Through



Uplink to Downlink
Full Cut-Through



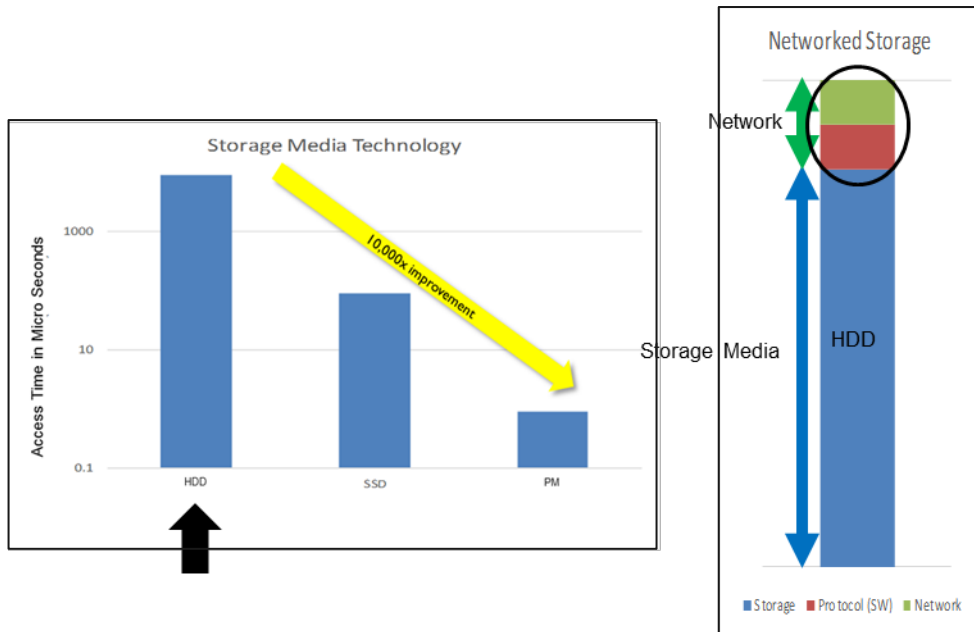
Downlink to Uplink
Smart Store and Forward

Intelligent cut-through also reduces network congestion

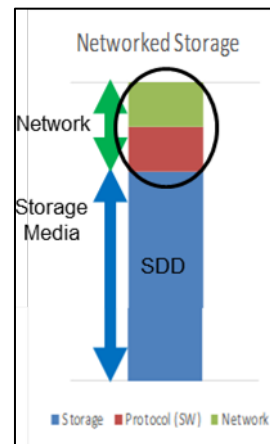
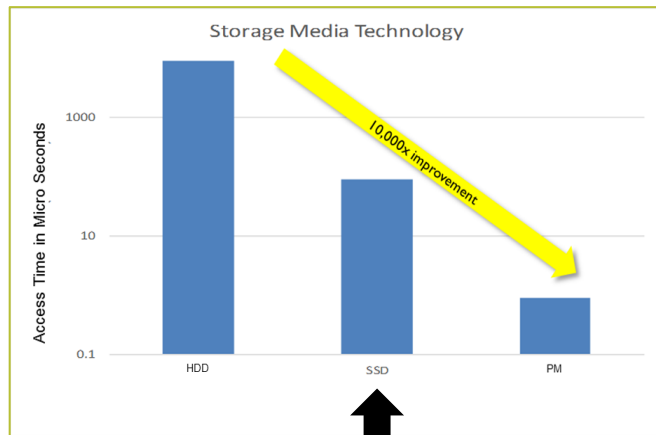




Latency is Not Only About the Hardware

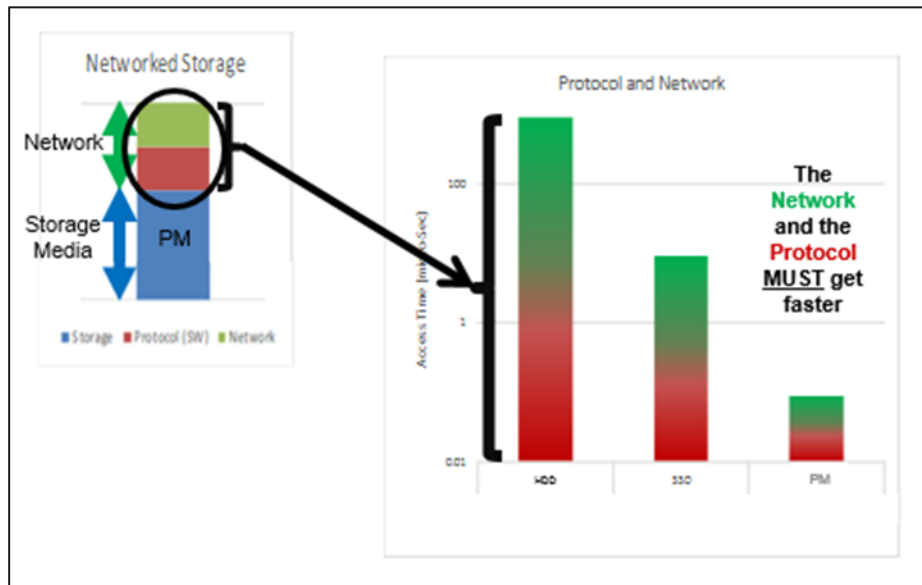
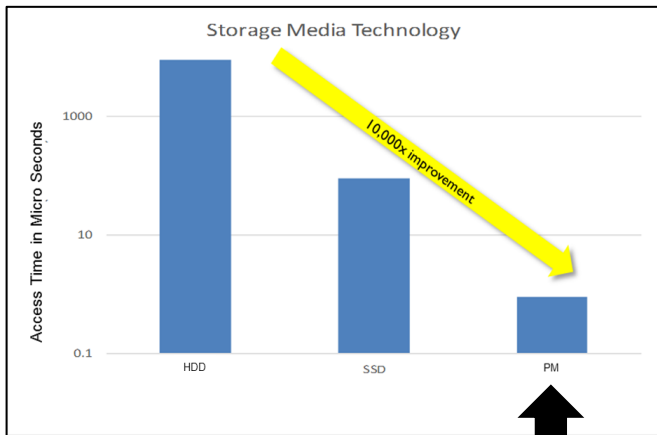


Latency is Not Only About the Hardware





Latency is Not Only About the Hardware





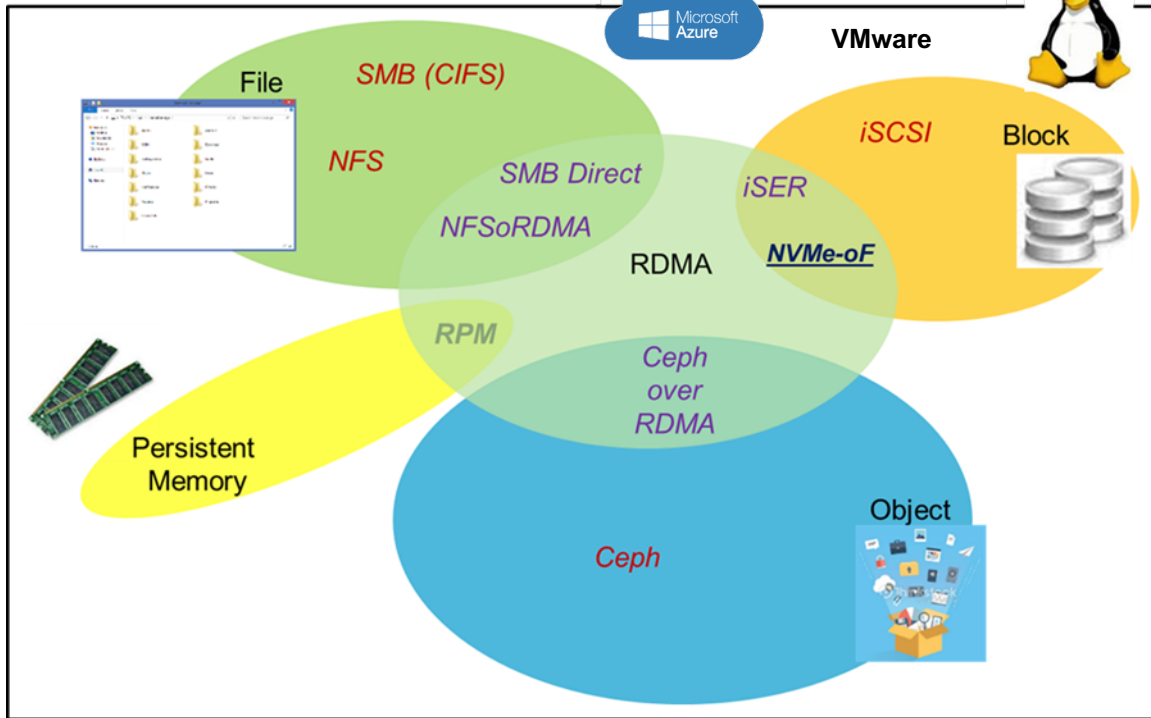
RDMA protocol

Windows

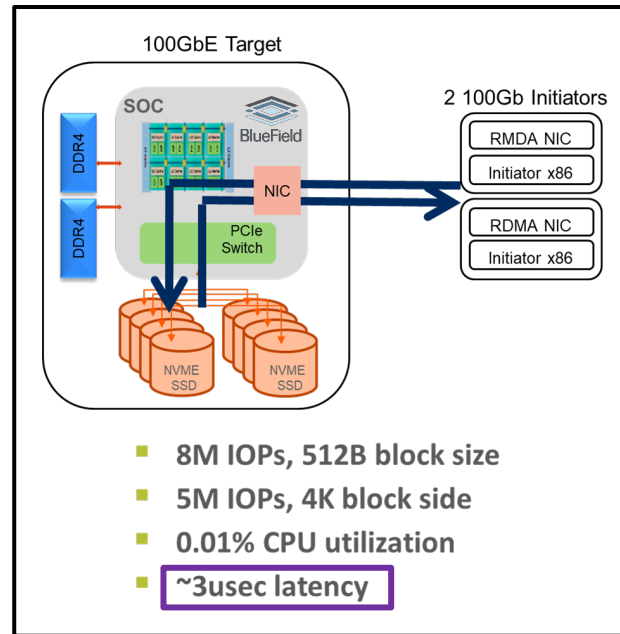


Alibaba Group
阿里巴巴集团

VMware



NVMe-oF over RoCE





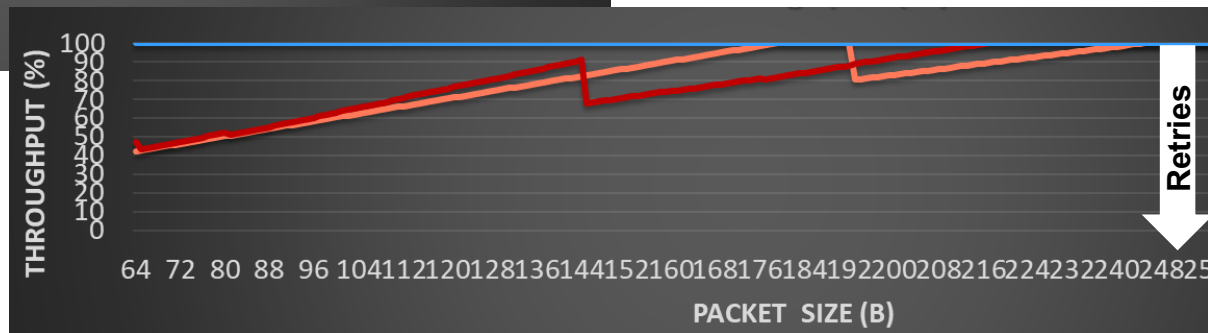
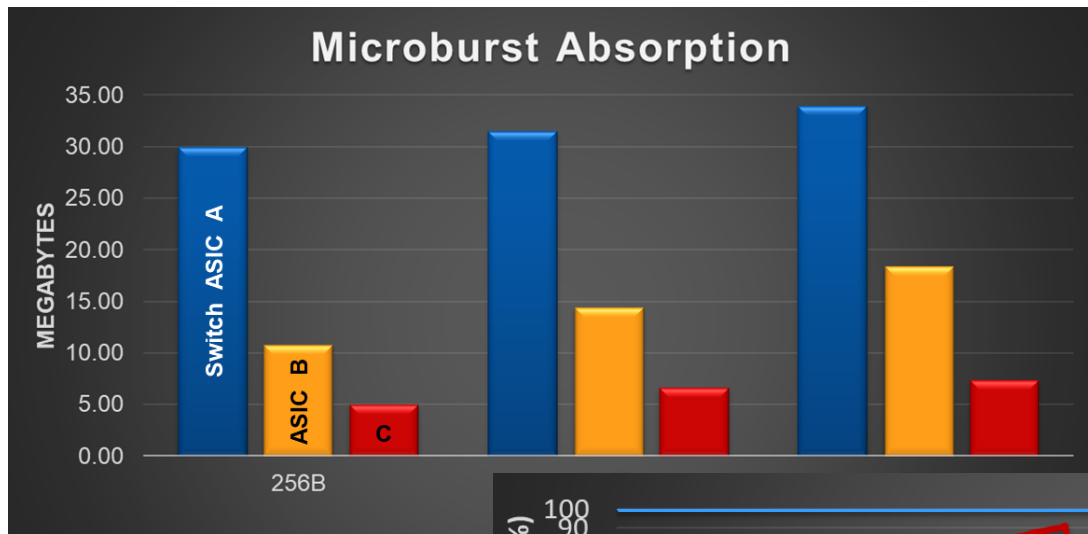
Flash Memory Summit

Effective Performant Network Congestion Control



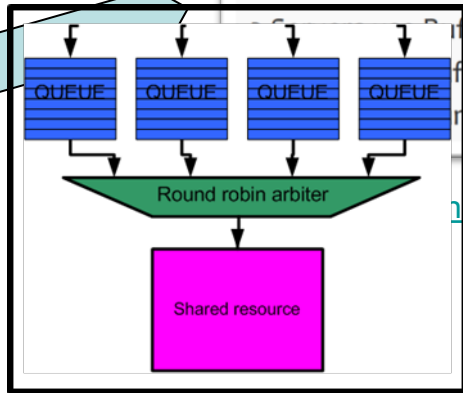
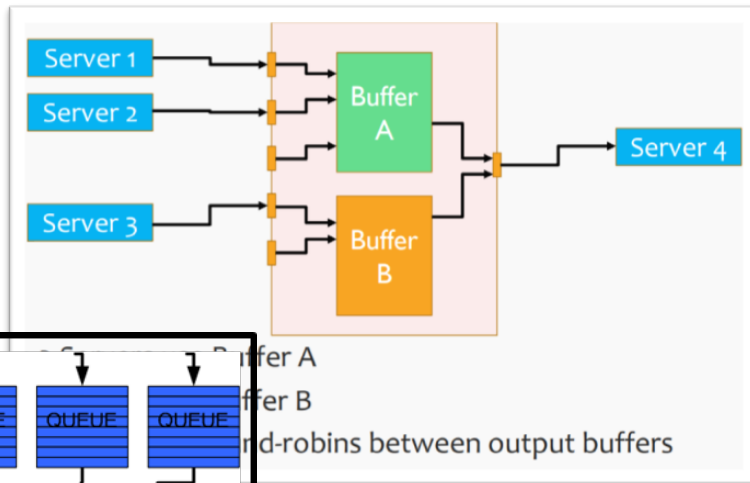
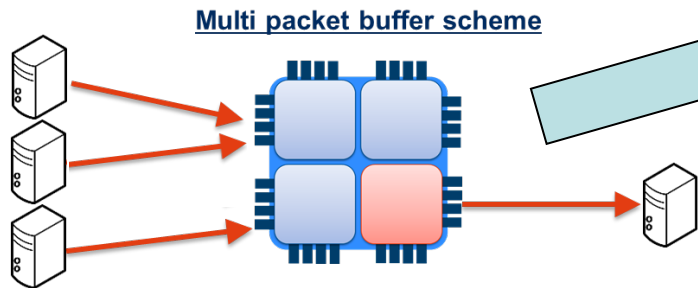
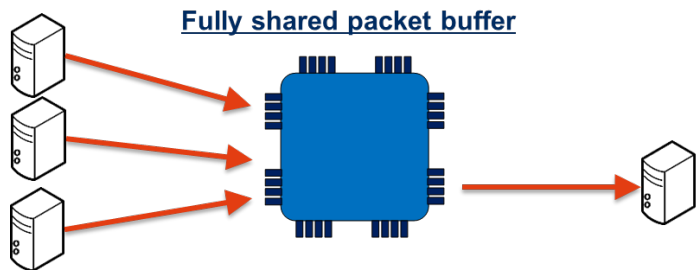


Switch Buffer Size and Congestion





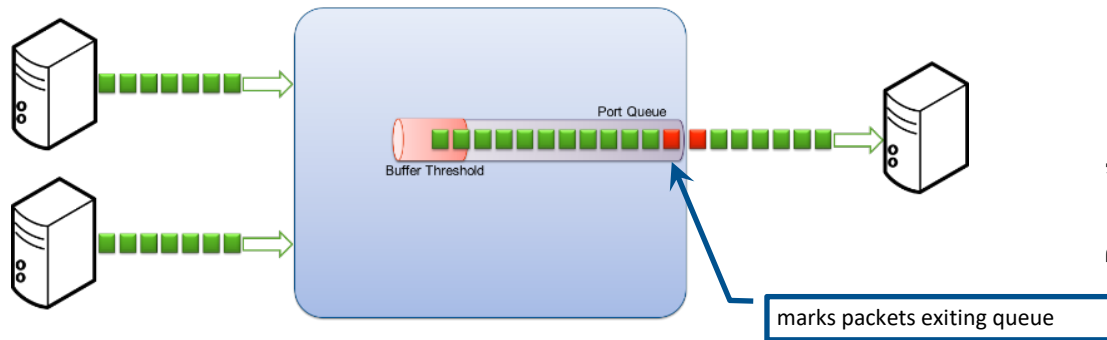
Fairness in Switch Architecture



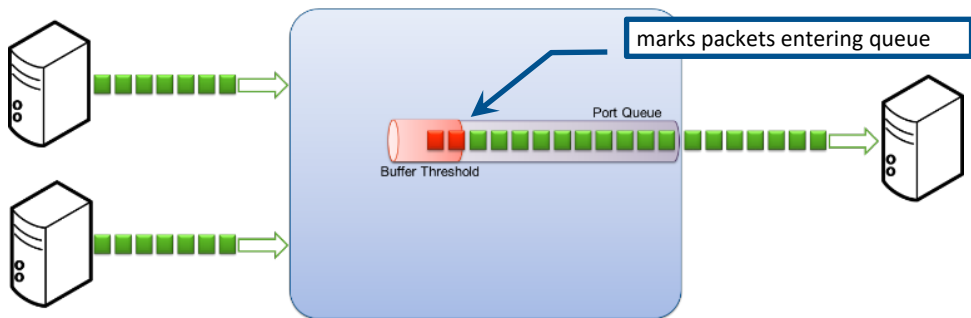
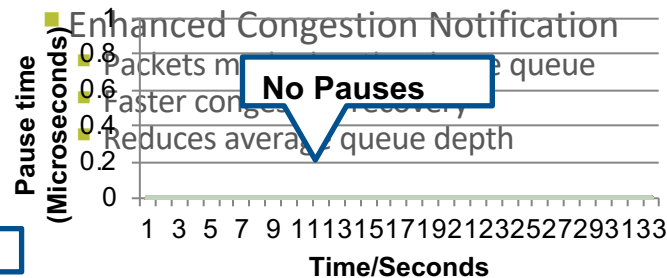
<https://www.flash-memory-summit.com/2018/talks/LPC%20DC-TCP%20Eval.pdf>



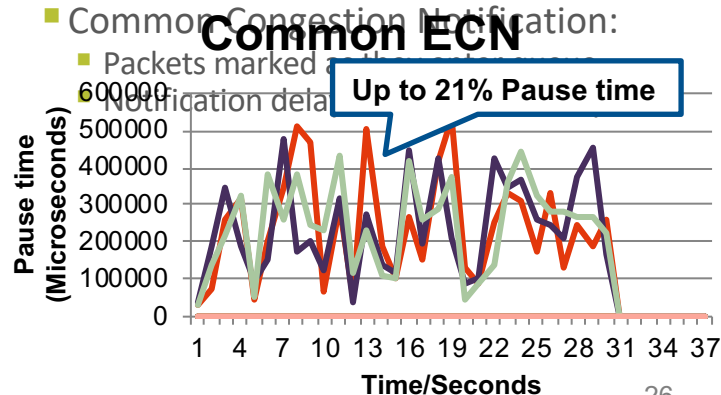
All ECN is Not Equal



Enhanced ECN

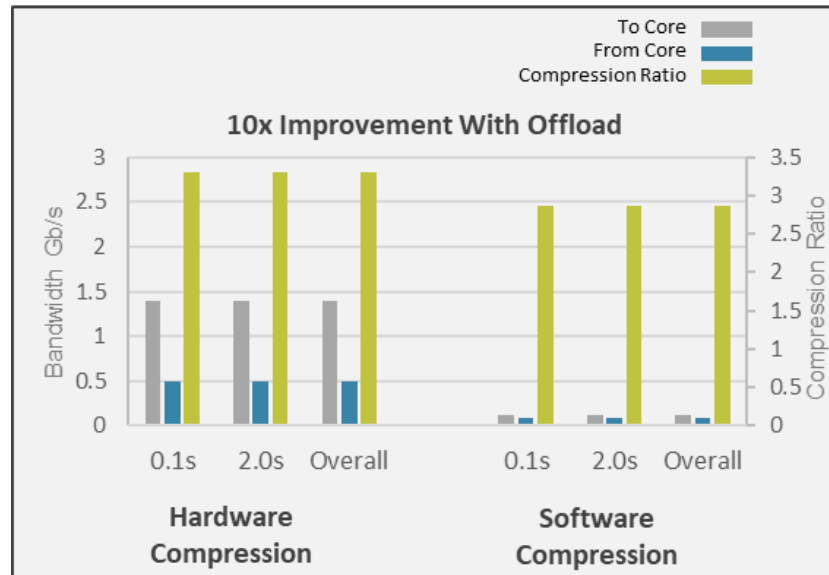
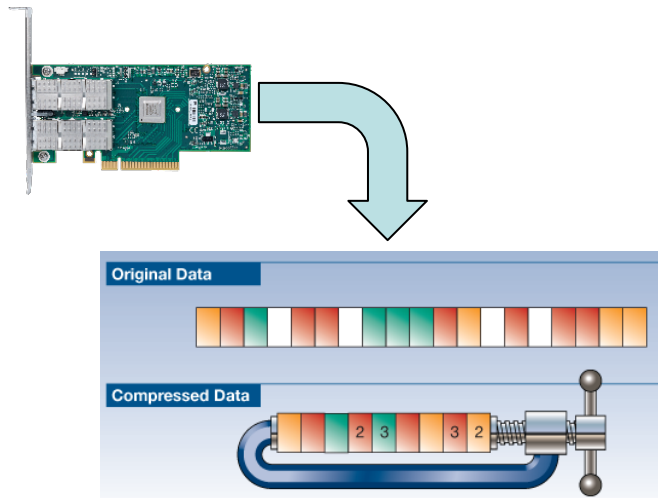


Common ECN





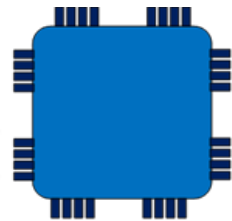
Reduce the Data Before Sending



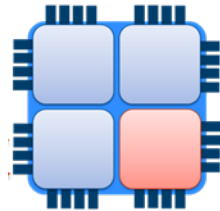


Quality of Service (QOS)

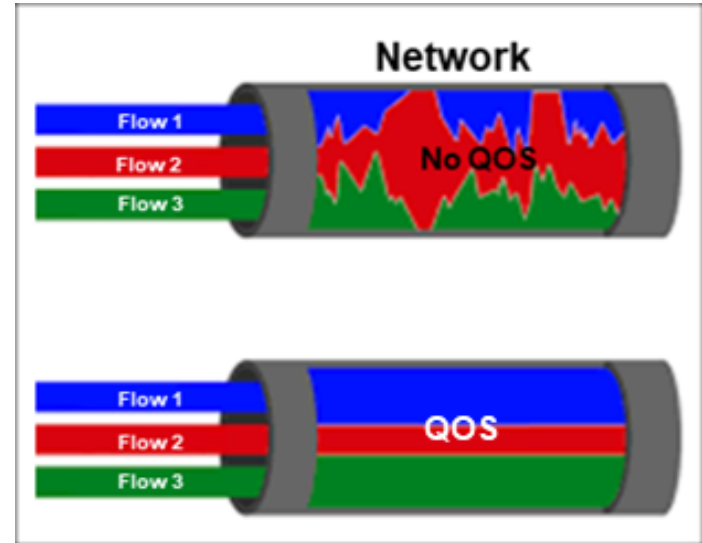
- A fully shared buffer architecture is best for implementing effective QOS
- The QOS algorithm must adapt the bandwidth allocation to the incoming priorities at wire speed



Fully Shared
buffer



Multi-buffer





Adaptive Flow Prioritization

- Egress flow prioritization
- Short flows get benefits
- Reduce flow completion time

Standard Queueing



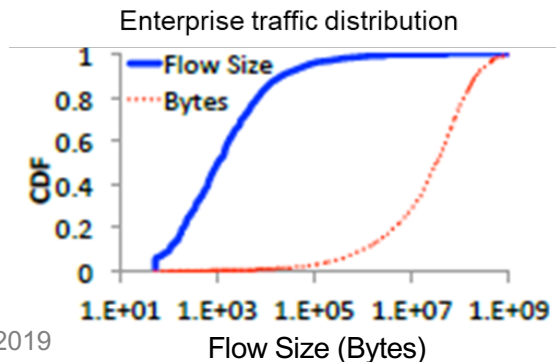
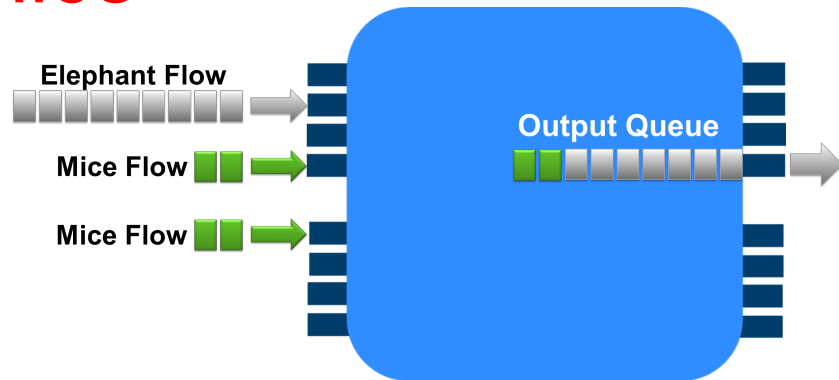
Adaptive Flow Prioritization





Elephants and Mice

- Majority of flows in the datacenter are small – Mice Flows
- Majority of packets belong to a few large flows – Elephant Flows





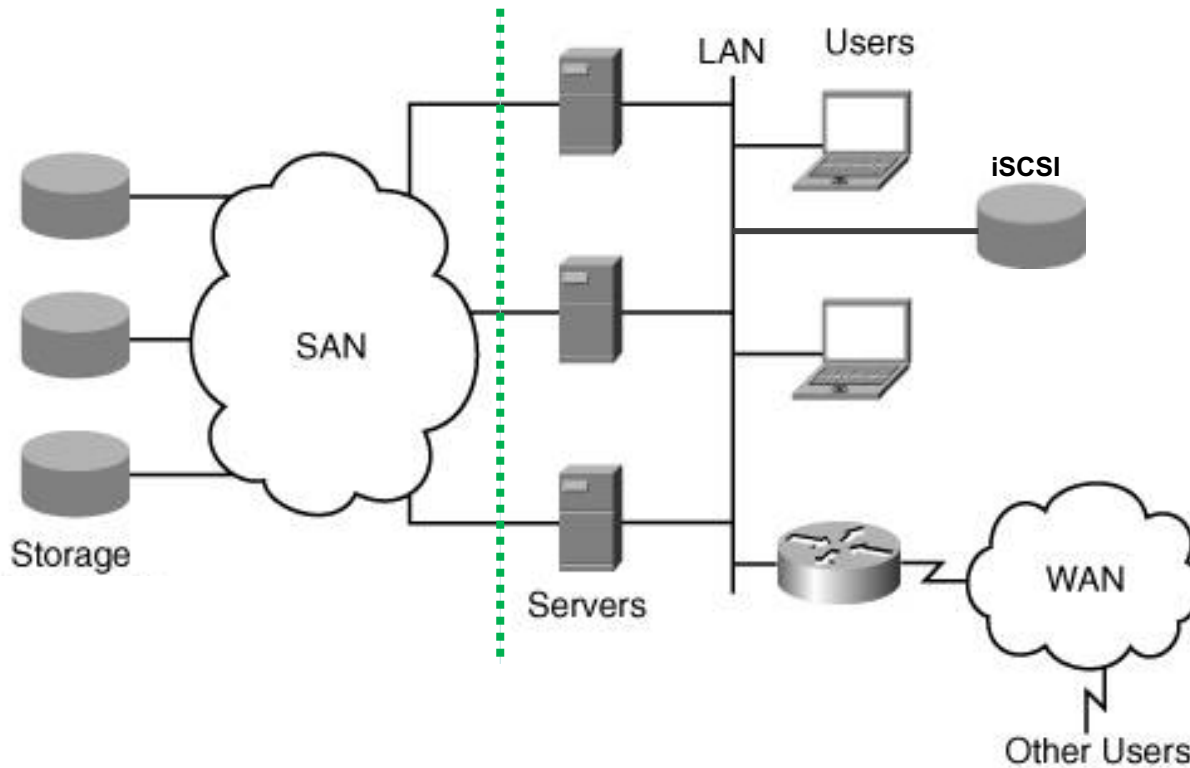
Flash Memory Summit

Security



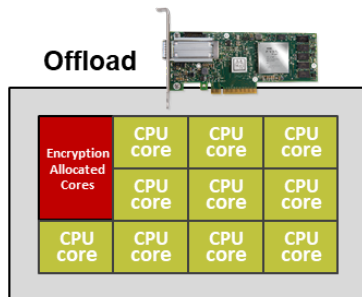


Isolated Fibre Channel SAN

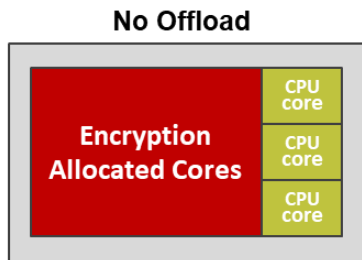




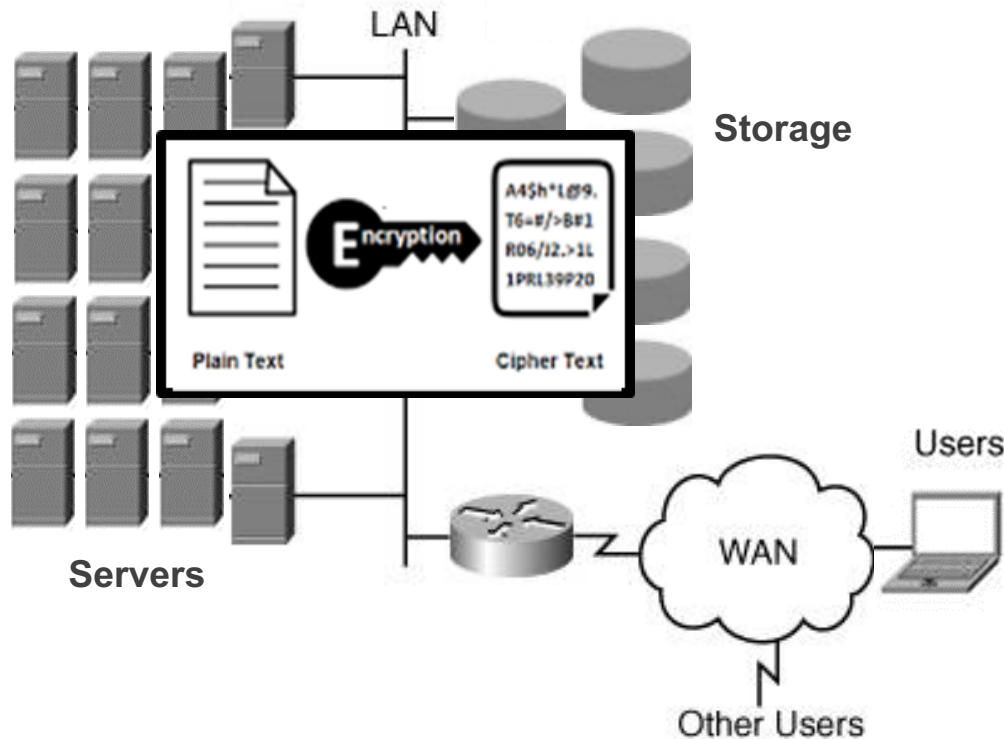
No Longer Isolated Storage



Only 16% CPU Overhead

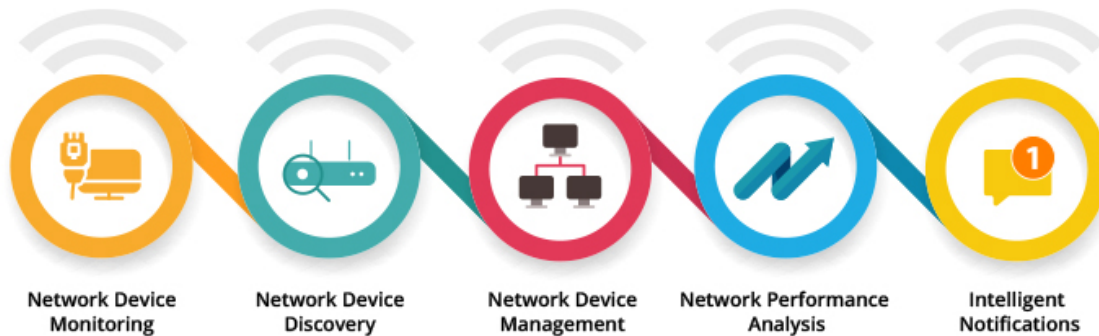


75% CPU Overhead

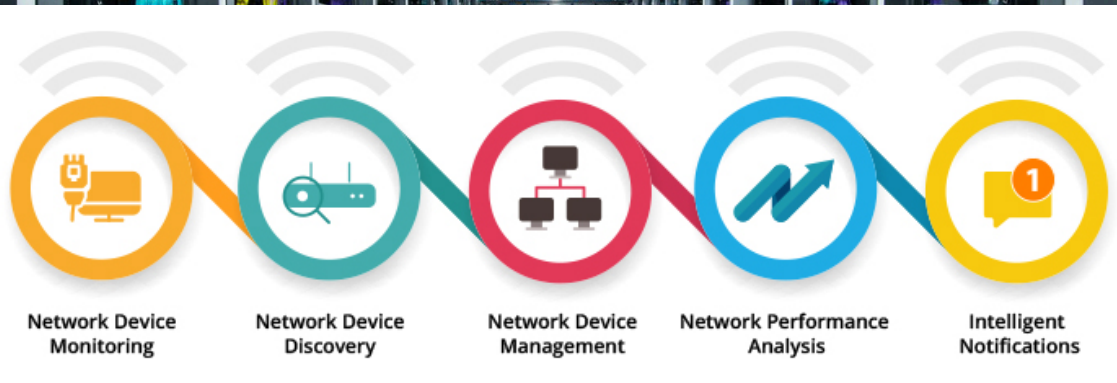




Management



Management



Network Device
Monitoring

Network Device
Discovery

Network Device
Management

Network Performance
Analysis

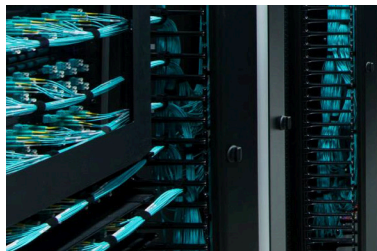
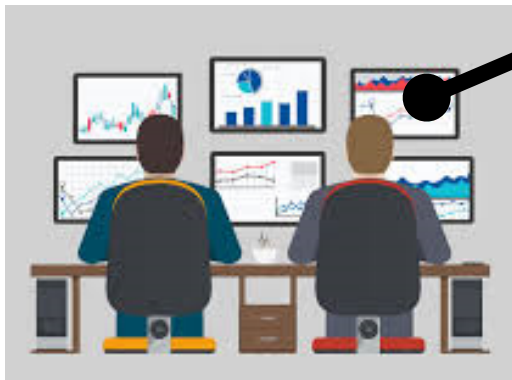
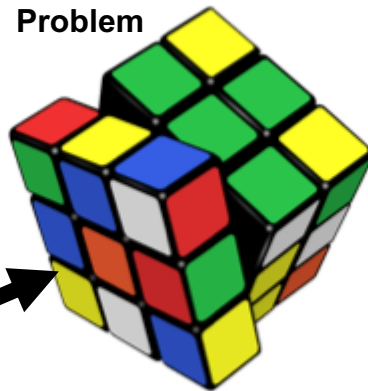
Intelligent
Notifications

At Scale



Normal Operation is Not the Issue

Problem



At Scale





Flash Memory Summit

The Answer is Automated Telemetry Capture and Analysis



**Successful analysis of the
telemetry data...**



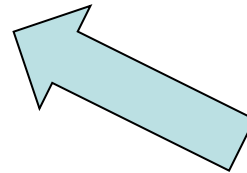
...leads to root cause



What is Telemetry?

The Important Questions

- ✓ WHO is being impacted
- ✓ WHEN it happened
- ✓ WHAT is causing the problem
- ✓ WHERE is the problem
- ✓ WHY it is happening



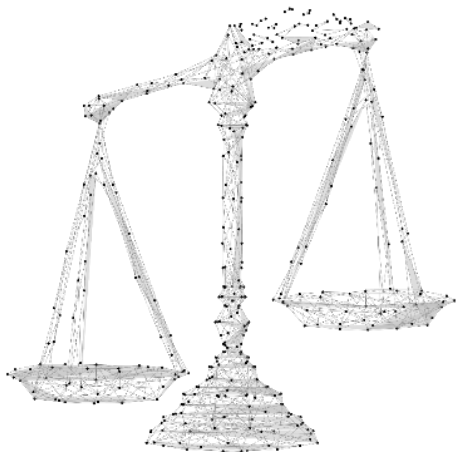
Telemetry is an automated communications process by which measurements and other data are collected at remote or inaccessible points and transmitted to receiving equipment for monitoring and analysis.



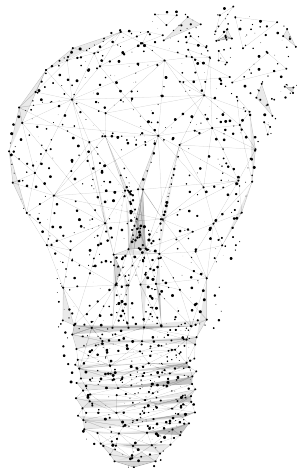


Telemetry Tell “What Just Happen?”

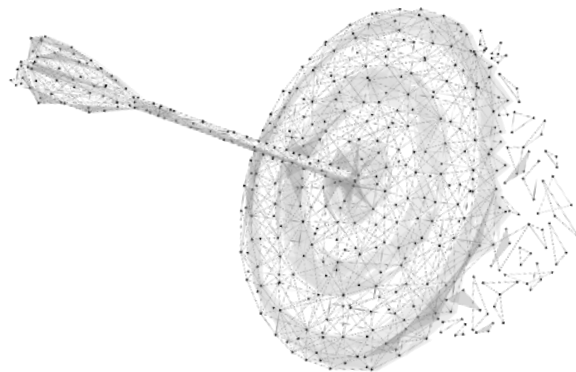
Faster Time to
Innocence

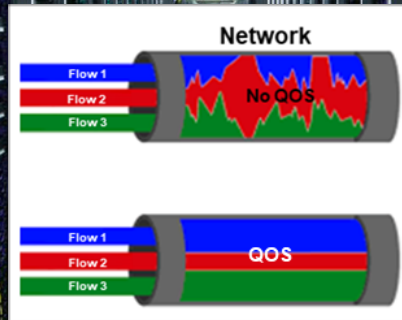
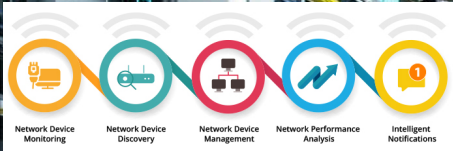


Faster Time To
Resolution



Get more out of the
Network





Ethernet is READY



Flash Memory Summit

Questions?



Flash Memory Summit

Thank You!