



September 23-26, 2019  
Santa Clara, CA

# Security, Integrity and Choices for NVMe over Fabrics

**Nishant Lodha**  
**Marvell**



# Agenda

- NVMe-oF®, the choices and the confusion
- Use Cases by Fabric
- Securing NVMe-oF
- Key Takeaways



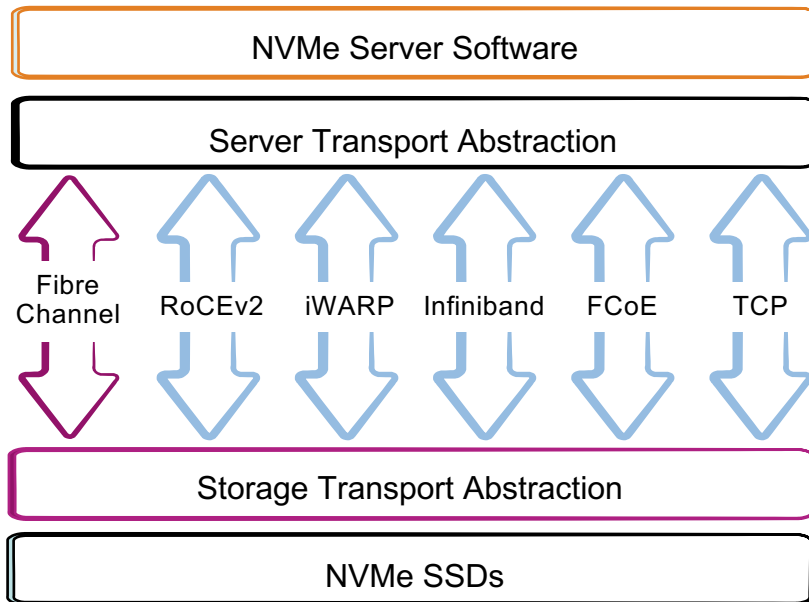
**NVMe-oF**

# Scaling our NVMe Requires a (Real) Network

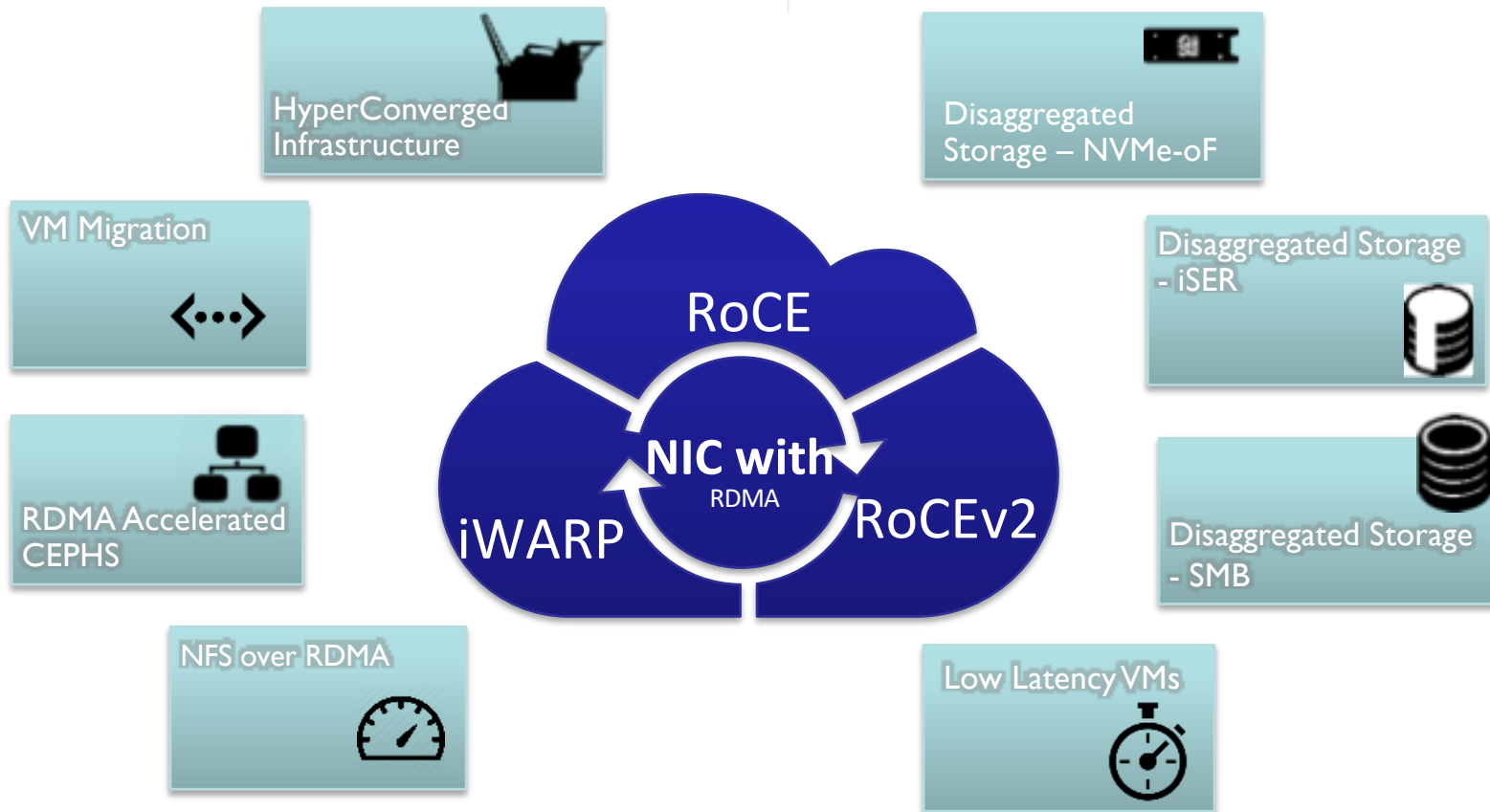
September 23-26, 2019  
Santa Clara, CA

SDC<sup>19</sup>

- Many options, plenty of confusion
- Fibre Channel is the transport for the vast majority of today's all flash arrays  
FC-NVMe Standardized in Mid-2017
- RoCEv2, iWARP and InfiniBand are RDMA-based but not compatible with each other  
NVMe-oF RDMA Standardized in 2016
- FCoE fabric is an option
- NVMe/TCP – is here! Standardized in NOV2018 🎉



# RDMA Use Cases by Application



# NVMe-oF™ RDMA – potential challenges

## Infrastructure and Skillset change?

**Not Automatic**

**Not Precise**

**Not for everyone**

**Congestion**



Keeping the  
network **'lossless'**

**RDMA/OEFD**  
expertise

**Skillset Requirements**



**RNIC Upgrade  
Required**

**RDMA Camps**

**Creates Islands**

**Backward Compatibility**



# Relationship Status: Microsoft and RoCE



Elden Christensen  
@EldenCluster

Follow

After endless support calls with customers struggling with the configuration complexity of RoCE, we have updated our RDMA network recommendations:

[docs.microsoft.com/en-us/windows-...](https://docs.microsoft.com/en-us/windows-...)



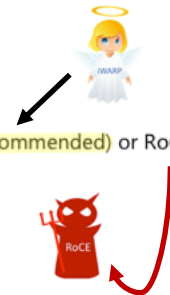
## Storage Spaces Direct hardware requirements

04/11/2018 • 3 minutes to read • Contributors all

### Networking

Recommended (for high performance, at scale, or deployments of 4+ nodes)

- NICs that are remote-direct memory access (RDMA) capable, **iWARP (recommended)** or RoCE
- Two or more NICs for redundancy and performance
- **25 Gbps network interface or higher**



See the Microsoft Blog – comparing the RDMA types

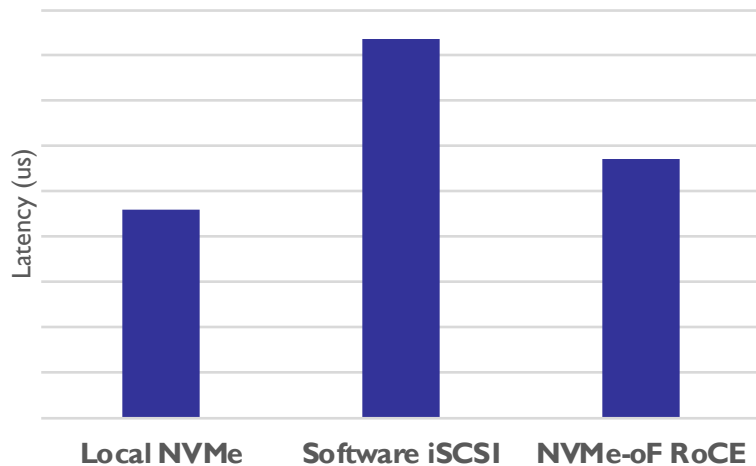
<https://blogs.technet.microsoft.com/filecab/2017/09/21/storage-spaces-direct-with-cavium-fastling-41000/>

# NVMe Transport Performance Comparisons

**iSCSI adds 82% more latency, Delivers fewer IOPS**

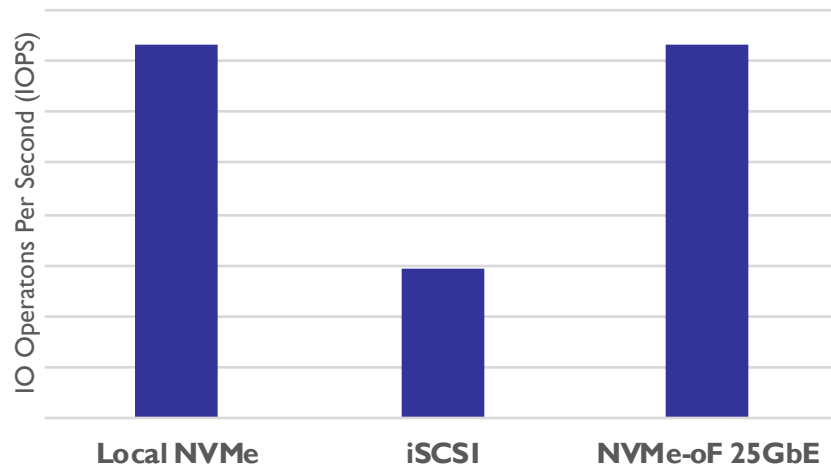
## NVMe-oF Latency Comparisons

4KB Random Reads Single Thread and IO Depth



## NVMe-oF IOPS Comparisons

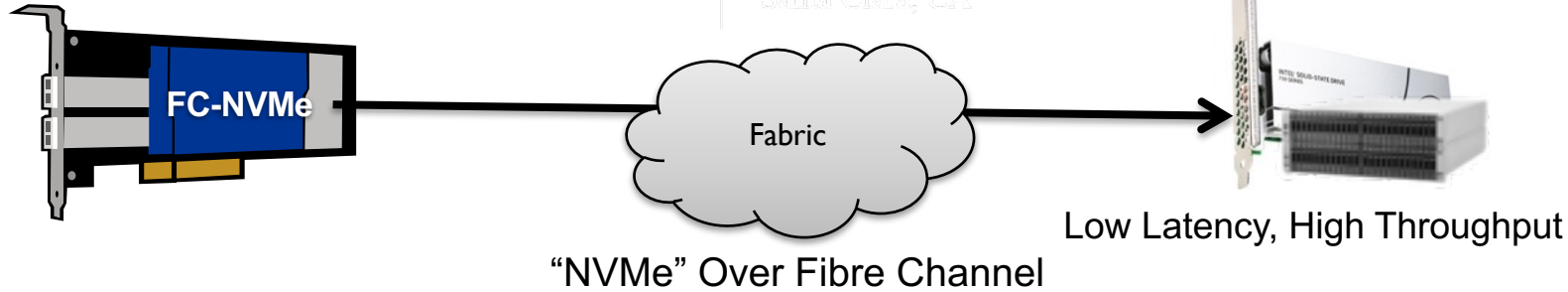
32KB Random Reads 8 Threads and 32 IO Depth





# FC-NVMe!

Storage Developer Conference, 2019  
Santa Clara, CA



Transport NVMe Natively over Fibre Channel

Low Latency

Reliable, Secure,  
Available

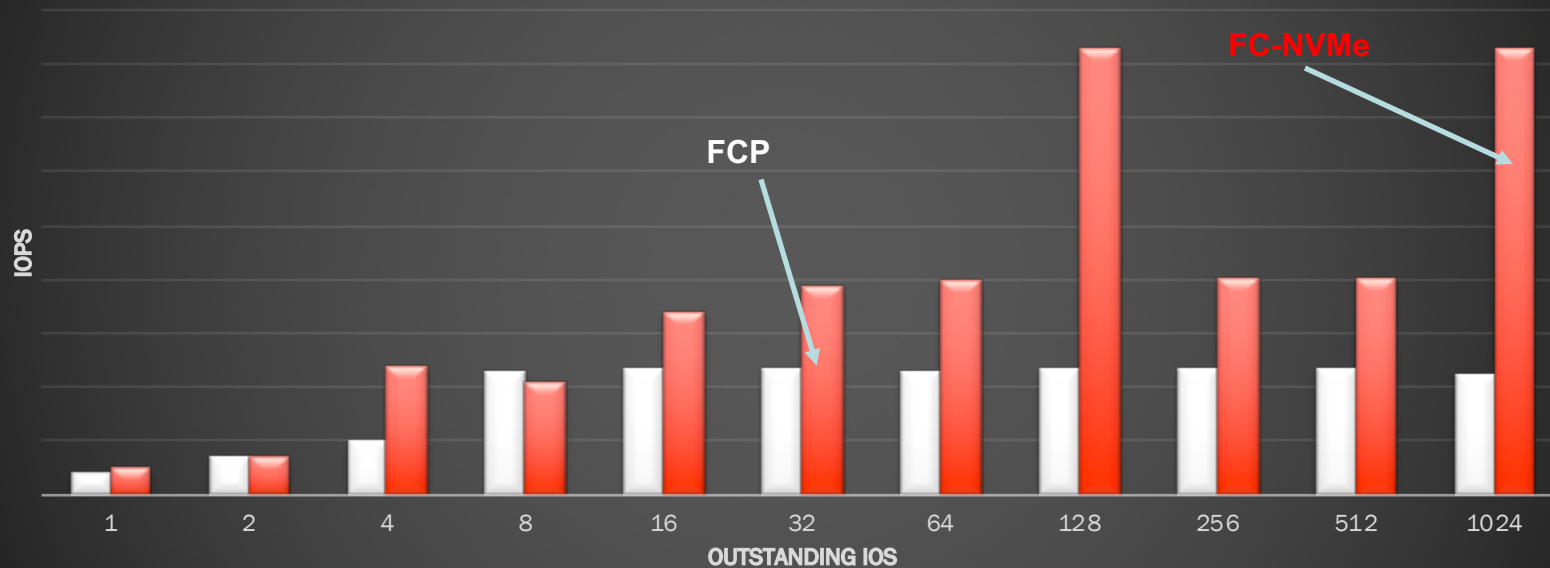
Leverage Existing Investments in Fibre Channel

FC-NVMe T11 Committee

Ecosystem Ready

# FCP vs. FC-NVMe

FCP vs. FC-NVMe: 4KB RD & 4 Jobs / DP to 1 LUN/NS per port



FC-NVMe Scales in performance

# Use Cases by Fabric

No one size fits all!

DAS, HPC, AI/ML



**NVMe/RDMA (Ethernet)**

Performance at the cost  
of complexity

Logos are indicative of workload characteristics only.

Enterprise Applications



**FC-NVMe (Fibre Channel)**

Leverage existing  
infrastructure. Reliability is  
key

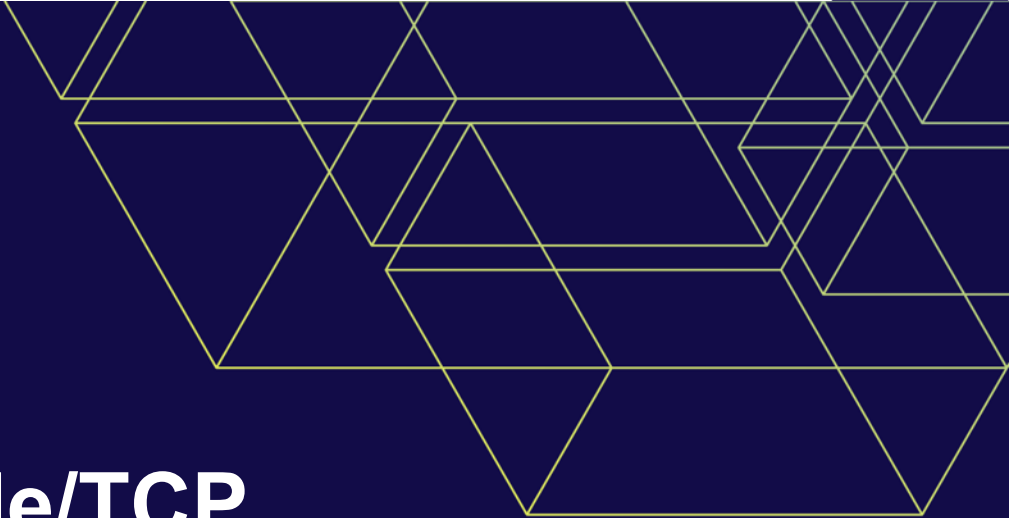
All Applications



**NVMe/TCP (Ethernet)**

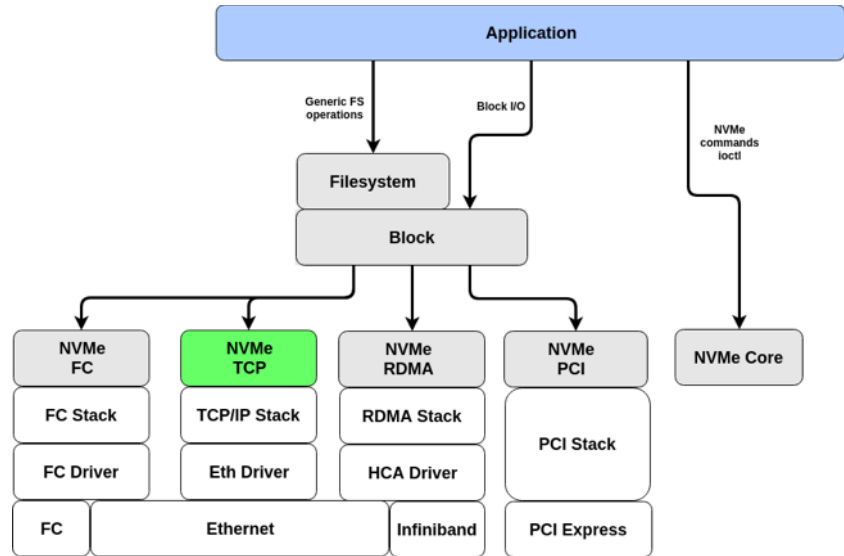
Simplicity is key. Balance of  
performance and cost

# NVMe/TCP



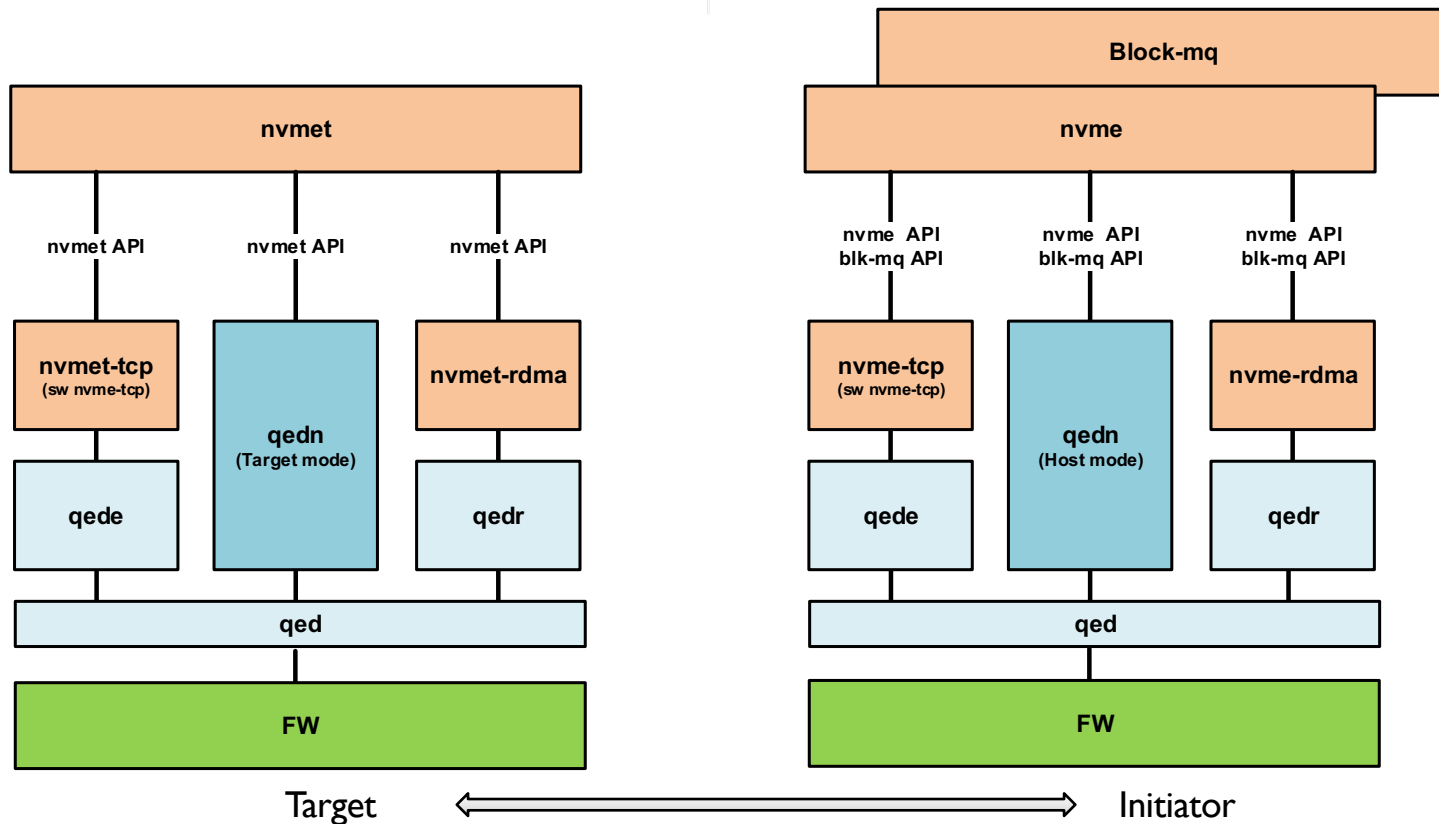
# NVMe-oF: NVMe/TCP

- What: Defines a TCP Transport Binding layer for NVMe-oF
- Promoted by Facebook, Google, Intel, Marvell etc.
- Not RDMA-based, Standardized on 15NOV18
- Why:
  - Enables adoption of NVMe-oF into existing datacenter IP network environments that are not RDMA-enabled



# NVMe-oF Driver Stack

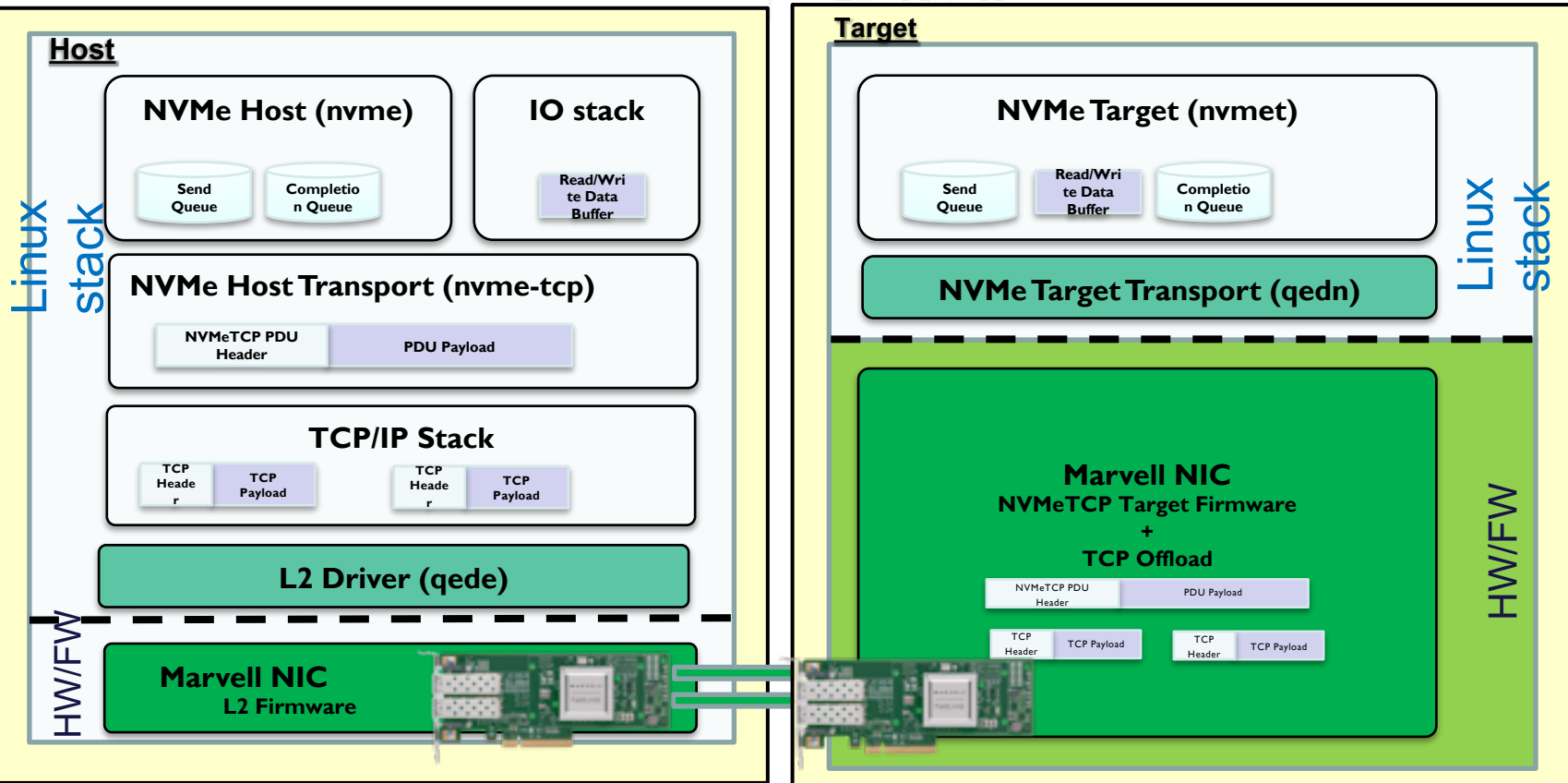
September 2019  
Santa Clara, CA



# Offloading NVMe/TCP

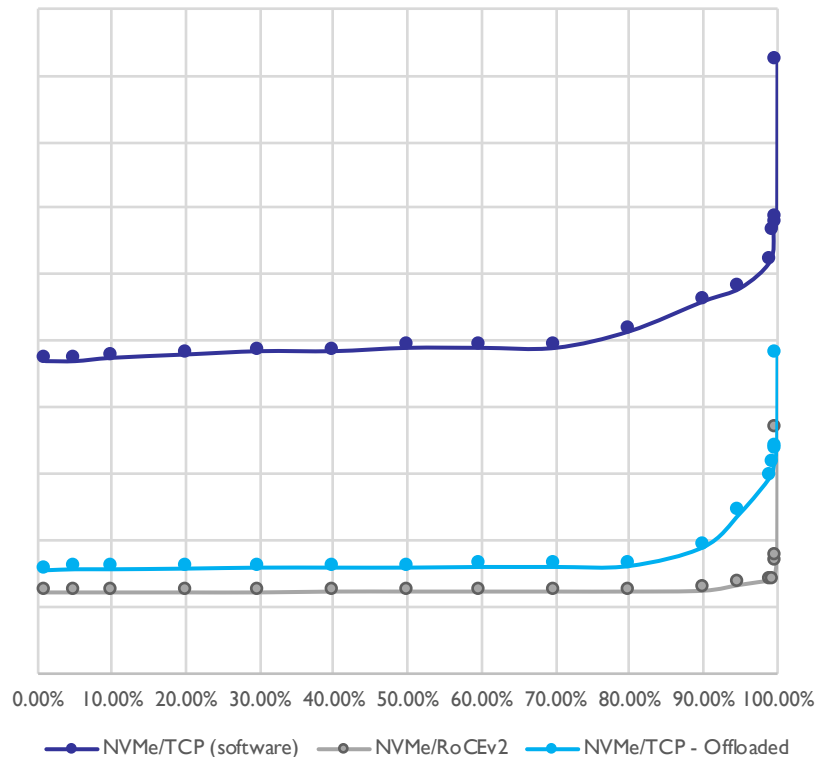
September 13-14, 2019  
Santa Clara, CA

SDC<sup>19</sup>

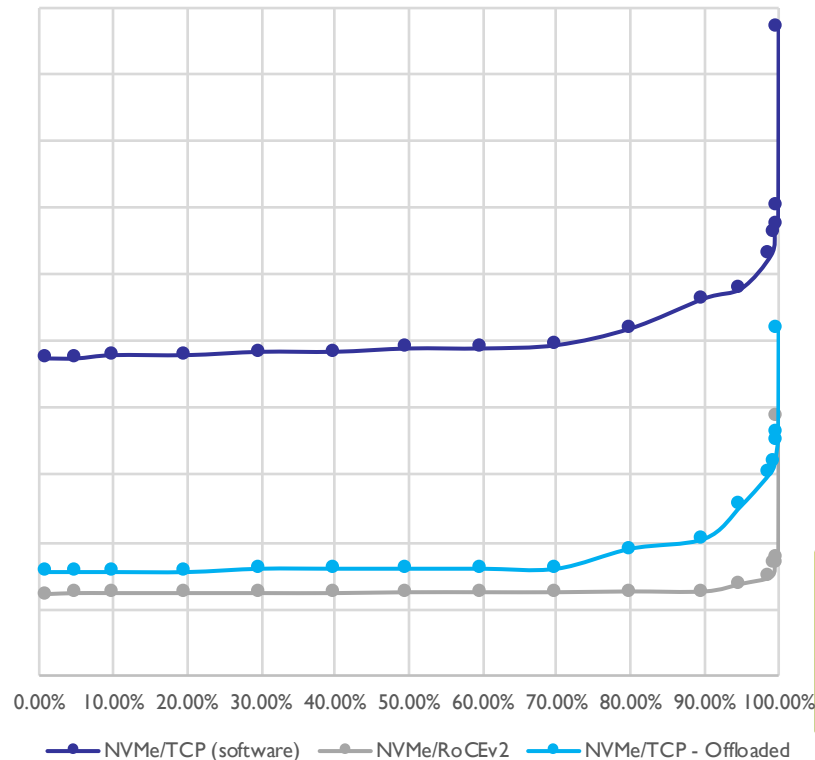


# Accelerating NVMe/TCP

4K Read IO - I pending latency [usec]



4K Write IO - I pending latency [usec]



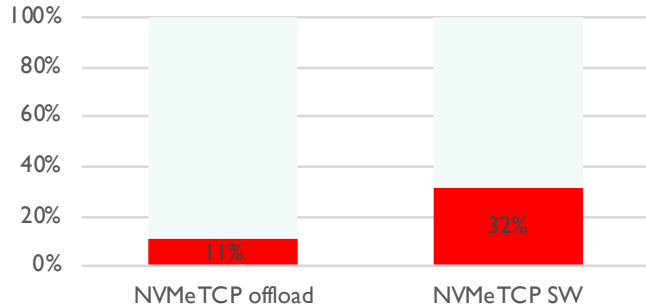


# Cost of I/O – NVMe/TCP

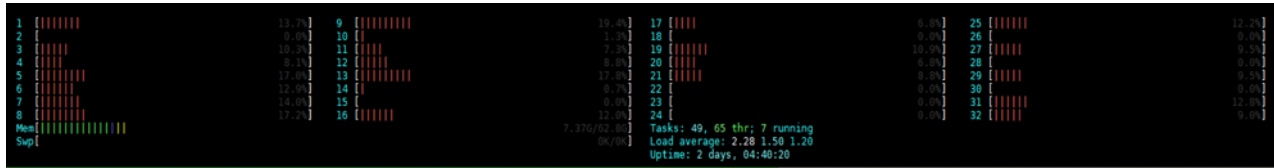
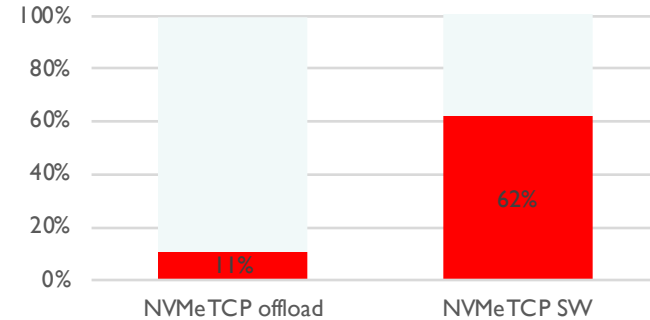
September 23-27, 2019  
Santa Clara, CA

SDC<sup>19</sup>

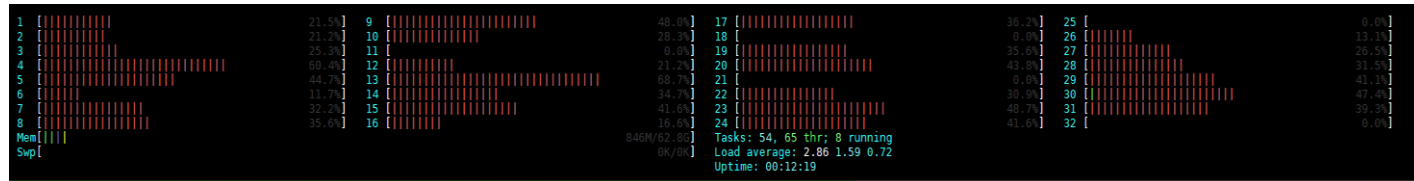
128K Read IOs - 100Gbps - CPU Utilization



128K Write IOs - 100Gbps - CPU Utilization



## Significant CPU Savings with NVMe/TCP Offload

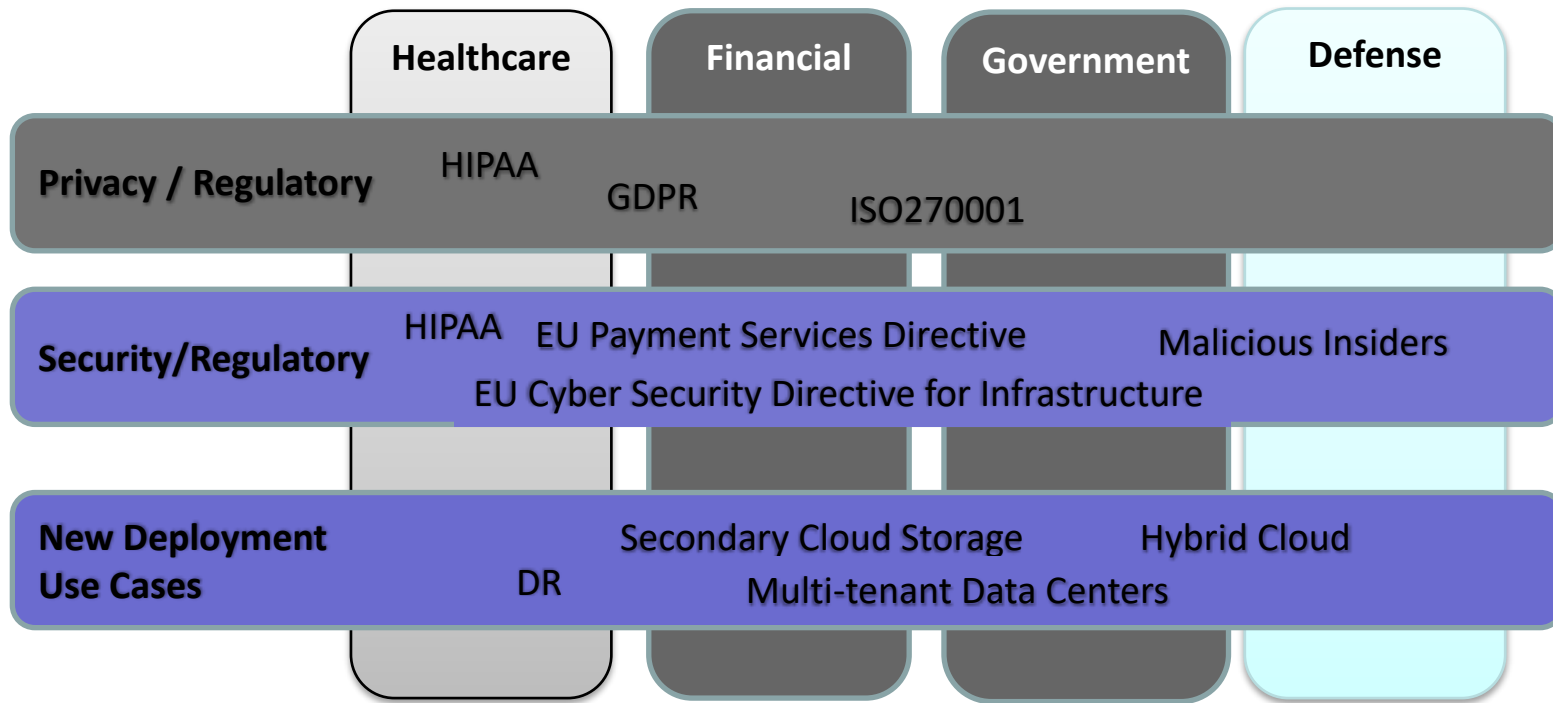





# **Security FC-NVMe**

# Drivers for FC-NVMe Security

## Security and Privacy Sensitive Verticals



## SDC 19



# Isn't FC Secure Already?

## Trusted Storage Interconnect for Decades

### Physical Security

- Data Centers are physically secured

### Segregation

- Fibre Channel SANs are segregated networks

### Partitioning

- FC Zoning ensures fabric partitioning

### Masking

- LUN masking restricts access to specific LUNs

### Management

- Out-of-Band Management (IP) is secure, OS Controls

# Yes, But...

September 23-26, 2019  
Santa Clara, CA

- New Data Center Architectures bring new threats
  - Distributed data centers - Remote replication and DR backups may be accessed by different users over Fabrics that span several sites
  - Multi Tenant data centers – Need to segregate and protect data traversing the same wire
- Increasing scale of FC SANs
  - Networks can be misconfigured
  - Fabric configuration databases are shared, have WKAs
- Existing mechanisms may not be enough
  - Switches are the sole entity that grant/deny access
    - Authorization based
  - “Segmentation” tools being used to implement “Security”
    - Soft zoning, LUN Masking

# Potential DC Storage Security Threats

## Sniffing Storage Traffic



## Storage Masquerading



## Data Corruption



## Session Hijacking



**Mitigated by Fibre Channel SAN Security**

# FC-SP-2: What and Why?

- **Why?** : Need to transition SANs from **Authorization and segmentation** based FC security to **authentication and encryption** based security!
- **What?** FC-SP-2 is a ANSI/INCITS standard (2012) that defines protocols to –
  - **Authenticate** Fibre Channel entities
  - **Setup** session **encryption keys**
  - Negotiate parameters to ensure per **frame integrity and confidentiality**
  - Define and **distribute security policies** over FC
- Designed to protect against several classes of threats



# Fabric Security Architecture

Components of FC-SP-2 Security Architecture

## Authentication Infrastructure

Secret, certificate, password and pre-shared key based architecture

## Authentication

Protocol to assure identify of communicating entities, negotiation of security requirement and protocol

## Security Associations

Protocol to establish Shared key between communicating entities, Based on IKEv2 (RFC4595)

## Crypto Integrity Confidentiality

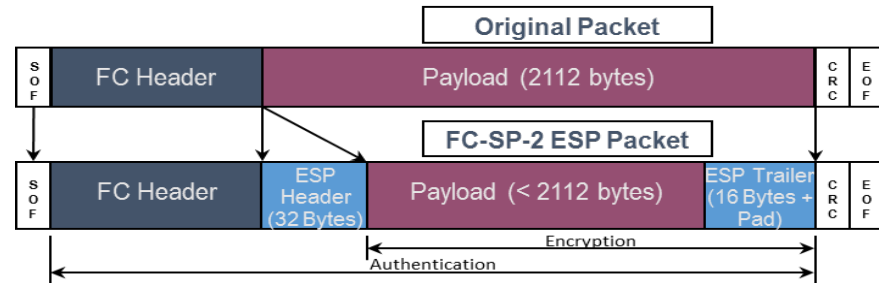
Frame by frame encryption, replay protection, origin authentication, ESP\_Header or CT\_Authentication

## Authorization

Fabric policies that control which entities can connect with each other, management access to the fabric

# FC-SP-2 *ESP\_header*

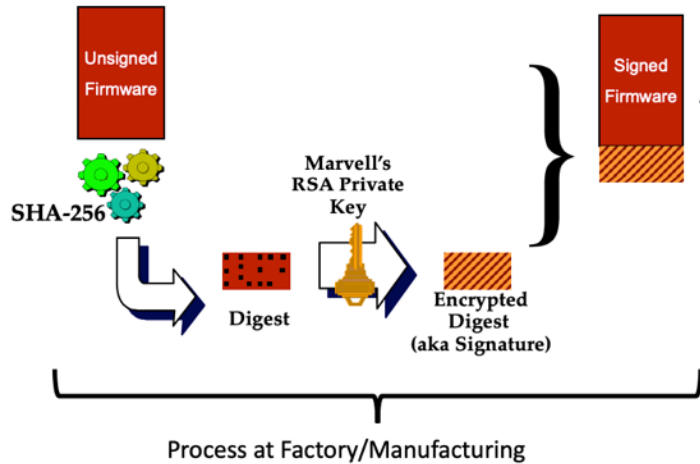
- *ESP\_header* (optional) is a layer 2 security protocol that provides
  - Origin authentication
  - Integrity
  - Anti-replay protection
  - Confidentiality
- Encapsulating Security Payload (ESP) is defined in RFC 4303
- FC-FS-3 defines optional headers for Fibre Channel, FC-SP defines how to use ESP in Fibre Channel
- Similar protections exist for CT\_Authentication



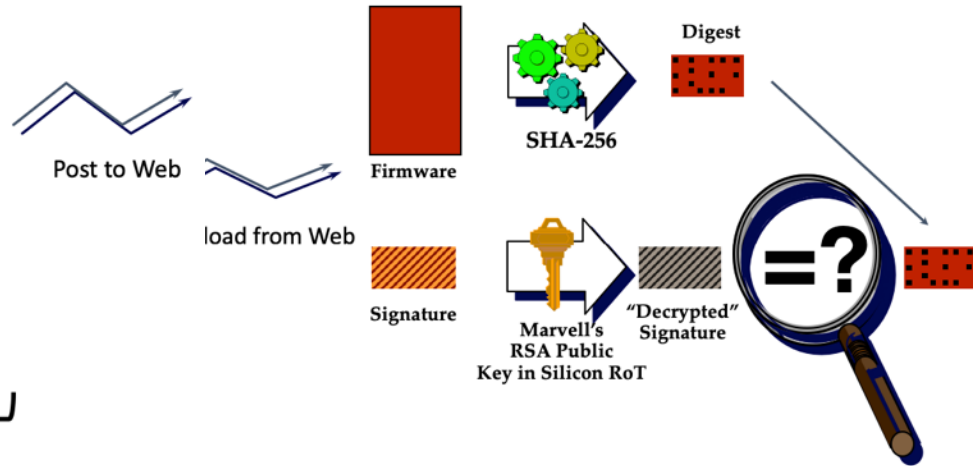
# Silicon Root of Trust

## Protecting the Integrity of Fibre Channel Firmware

### *Signing Fibre Channel firmware*



### Verifies firmware signature using Silicon Root of Trust



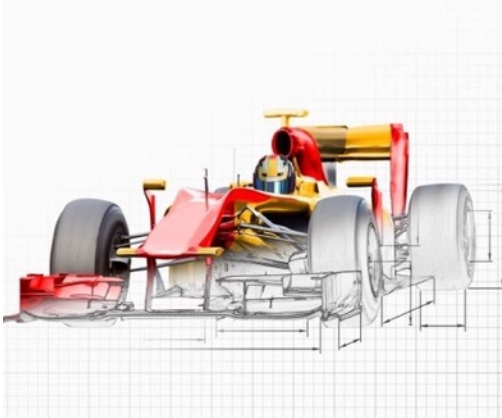


**Key Takeaway**

# Making the right “fabric” choice!

September 10-11  
Santa Clara, CA

SDC<sup>19</sup>



Not “just” about “fabrics”  
performance



Culture and Install Base



Use Cases and Security

**SDC**<sup>19</sup>

September 23-26, 2019  
Santa Clara, CA

**That's it!**

