# Western Digital.

# A Journey into NVMe-oF™:
## Options, Trade-offs and Challenges

Ihab Hamadi

Fellow, Western Digital

August 8, 2019

8/9/2019

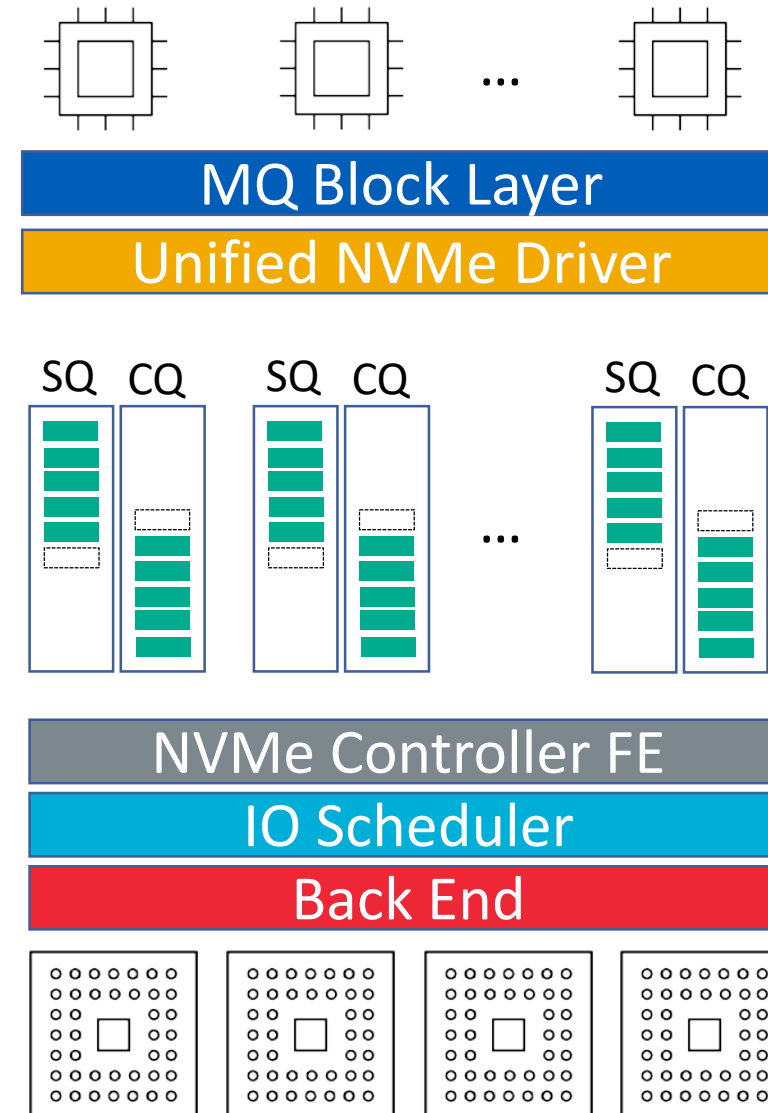# Agenda

**1** Background

**2** Landscape of NVMe™ Fabrics

**3** Lossless vs. Lossy

**4** Fabric Selection Criteria

**5** Case Study

Western Digital®
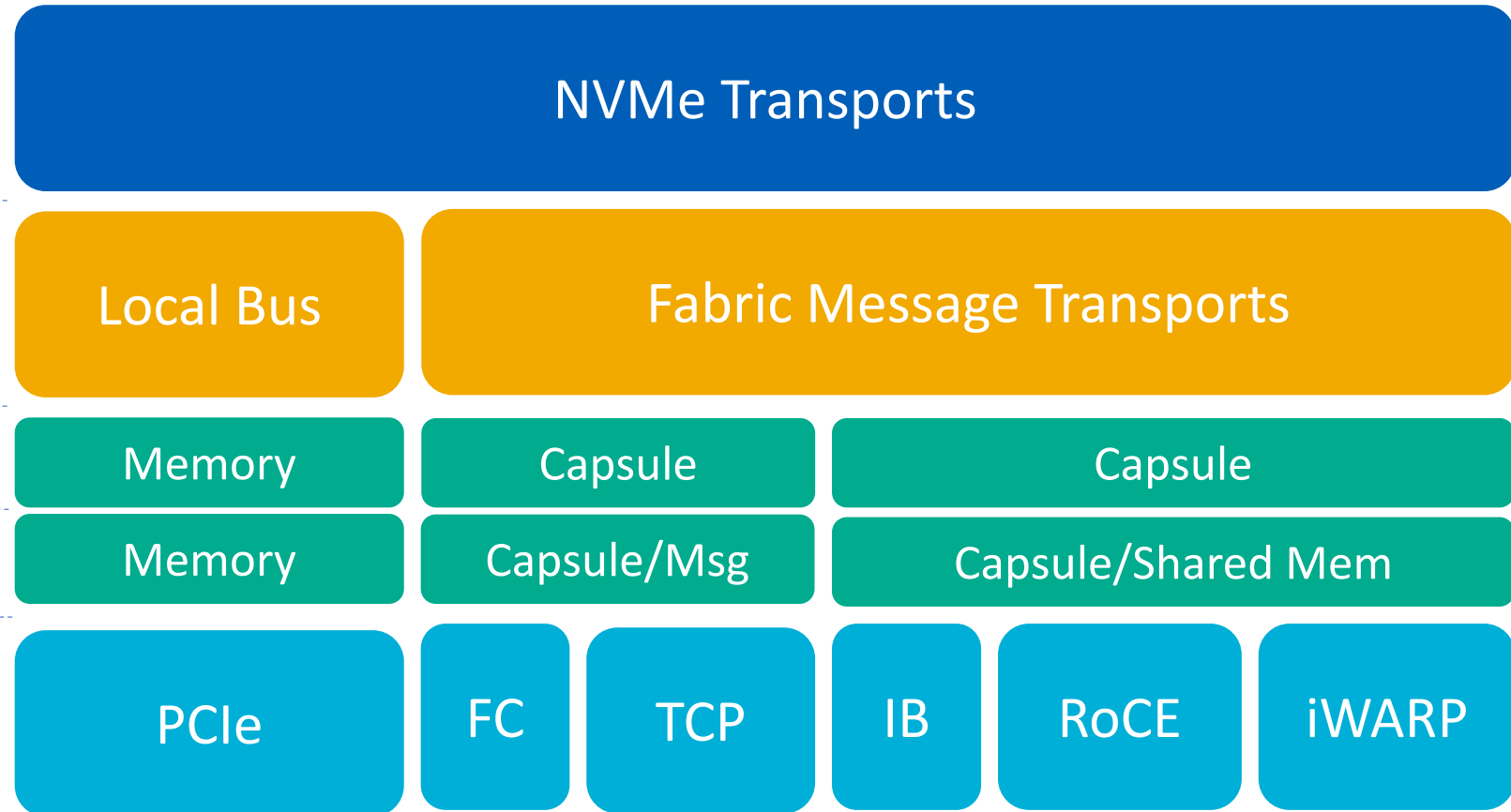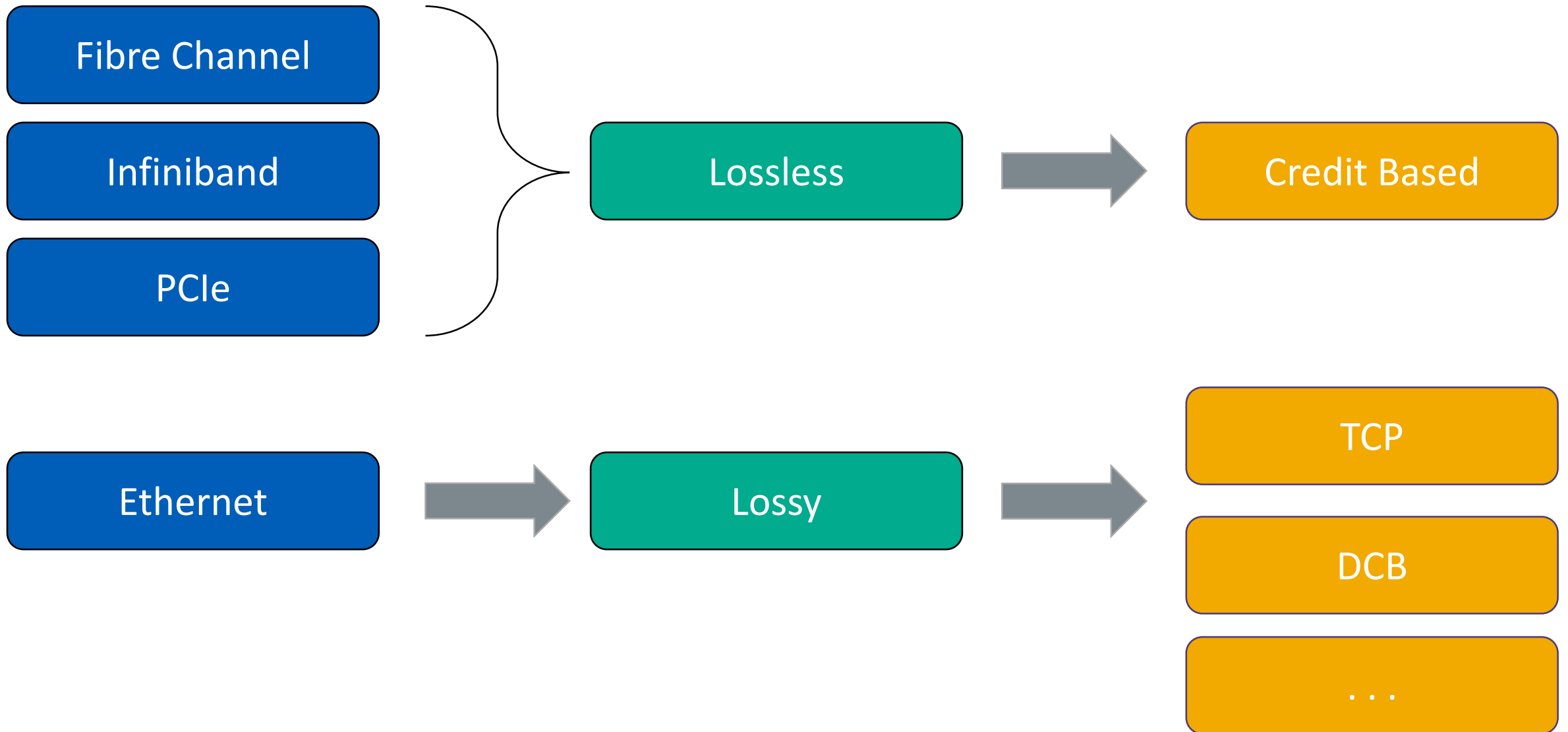
# Background: Why NVMe? Why NVMe-oF?

- **Parallelism fits multi-core CPUs**
  - Also reduces/spreads host CPU load

- **Removes some cost components**
  - Some common HW blocks
  - One driver

- **Storage System Benefits**
  - Lower latency (average & tail)
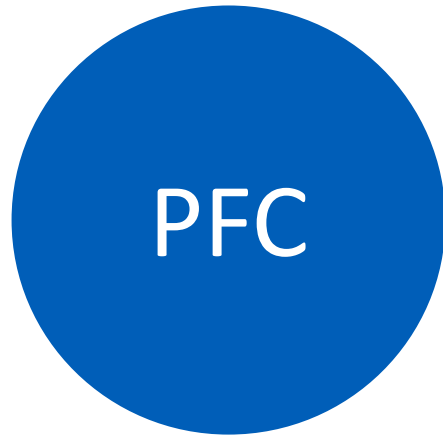  - Higher BW

- **NVMe-oF Motivation:**
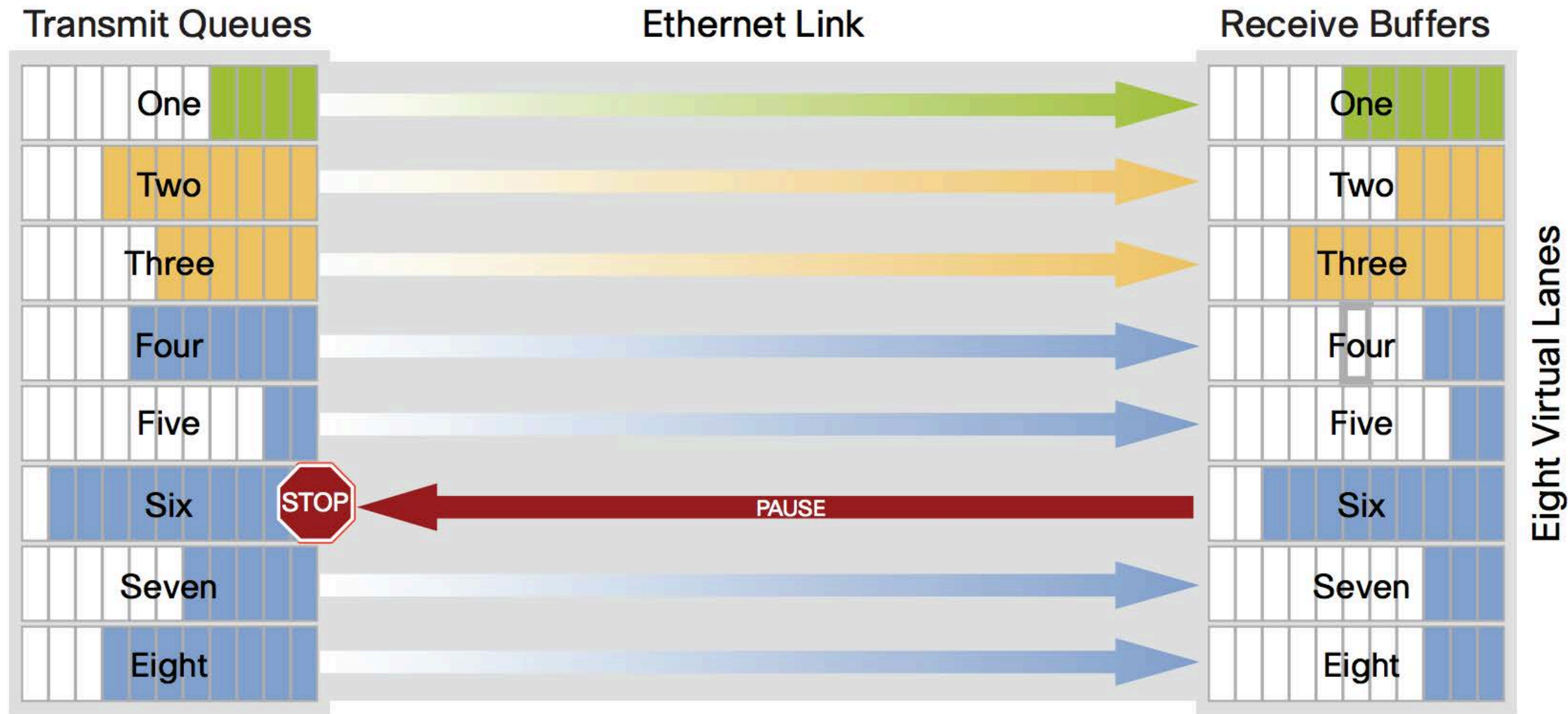  - Extend benefits end-to-end

# NVMe Transport Model

| | | | |
|---|---|---|---|
| | **NVMe Transports** | | |
| **Locality** | **Local Bus** | **Fabric Message Transports** | |
| **Model: Cmd/Rsp** | Memory | Capsule | Capsule |
| **Model: Data** | Memory | Capsule/Msg | Capsule/Shared Mem |
| **Fabric Type** | PCIe | FC    TCP | IB    RoCE    iWARP |

**Western Digital.**

# Fabric 101: Lossy vs. Lossless Fabrics

```
Fibre Channel ┐
              │
Infiniband    ├──→  Lossless  ──→  Credit Based
              │
PCIe          ┘

Ethernet  ──→  Lossy  ──→  TCP
                           DCB
                           . . .
```

**Western Digital.**

8/9/2019          5

# Data Center Bridging (DCB)
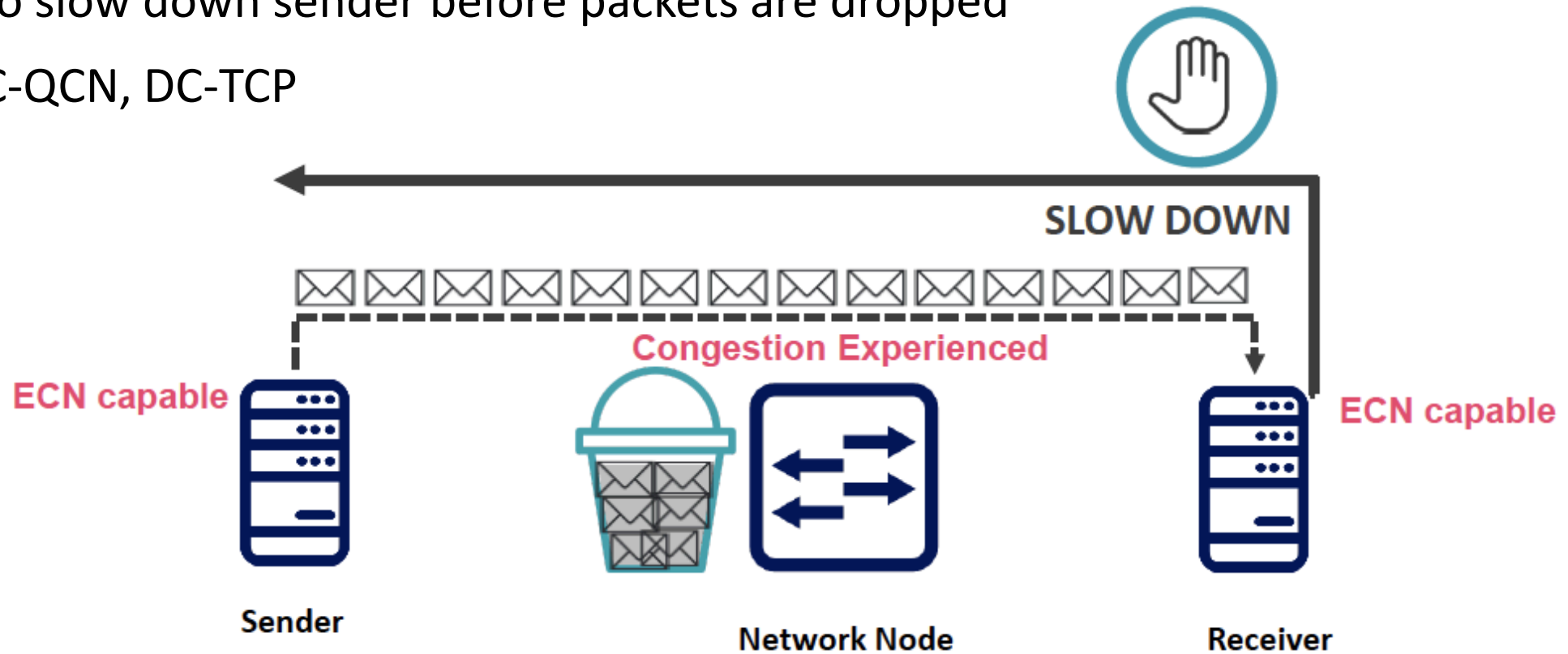
PFC

ECN

DCBx

ETS

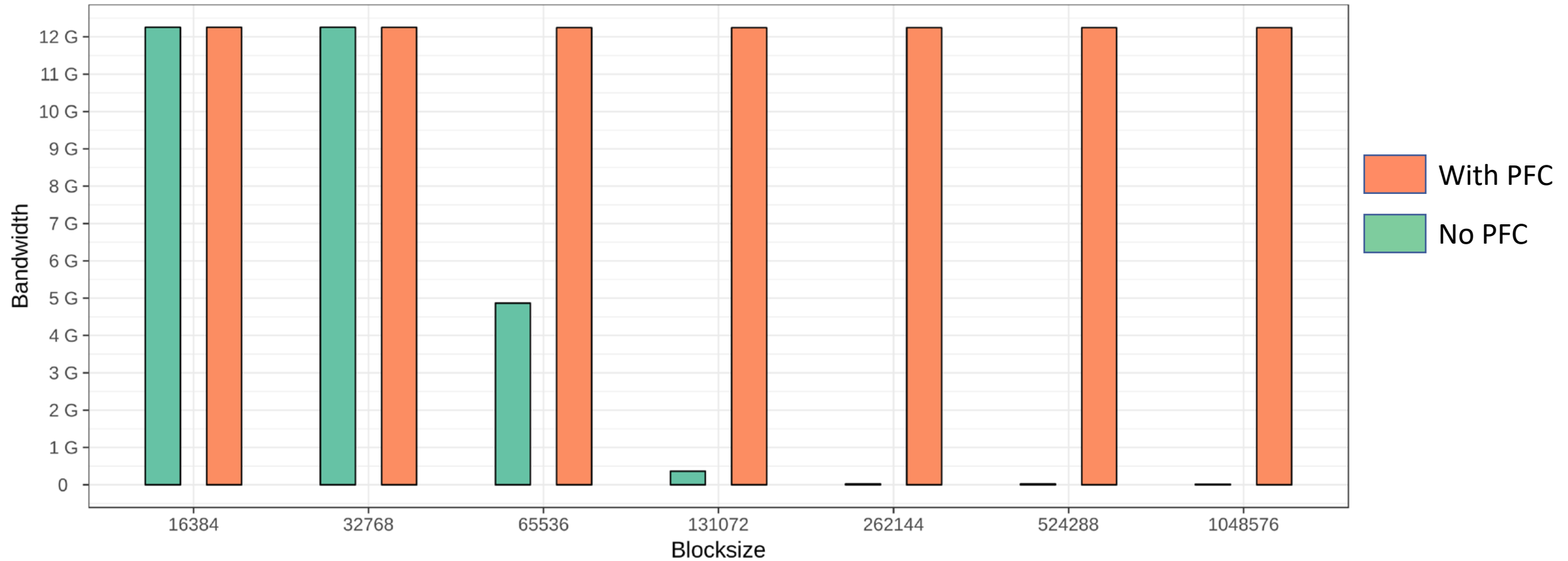**Western Digital.**

# DCB: PFC

# DCB: ECN

- ECN is end-to-end congestion management mechanism

- Three roles: Sender (RP), Switch (CP), Receiver (NP)

- Goal is to slow down sender before packets are dropped

- QCN, DC-QCN, DC-TCP



SLOW DOWN

ECN capable

Congestion Experienced

ECN capable

Sender

Network Node

Receiver

# Do I Really Need DCB (Lossless Net) with RoCE?

*BW vs. IO Size*



Source: Western Digital Performance Tests

**Western Digital.**

# Fabric Selection Criteria

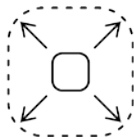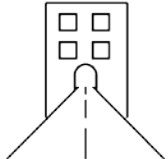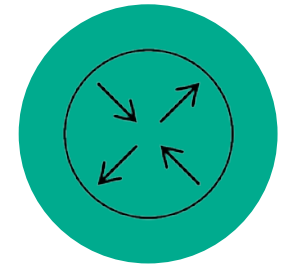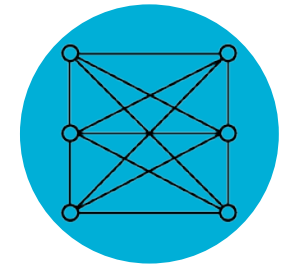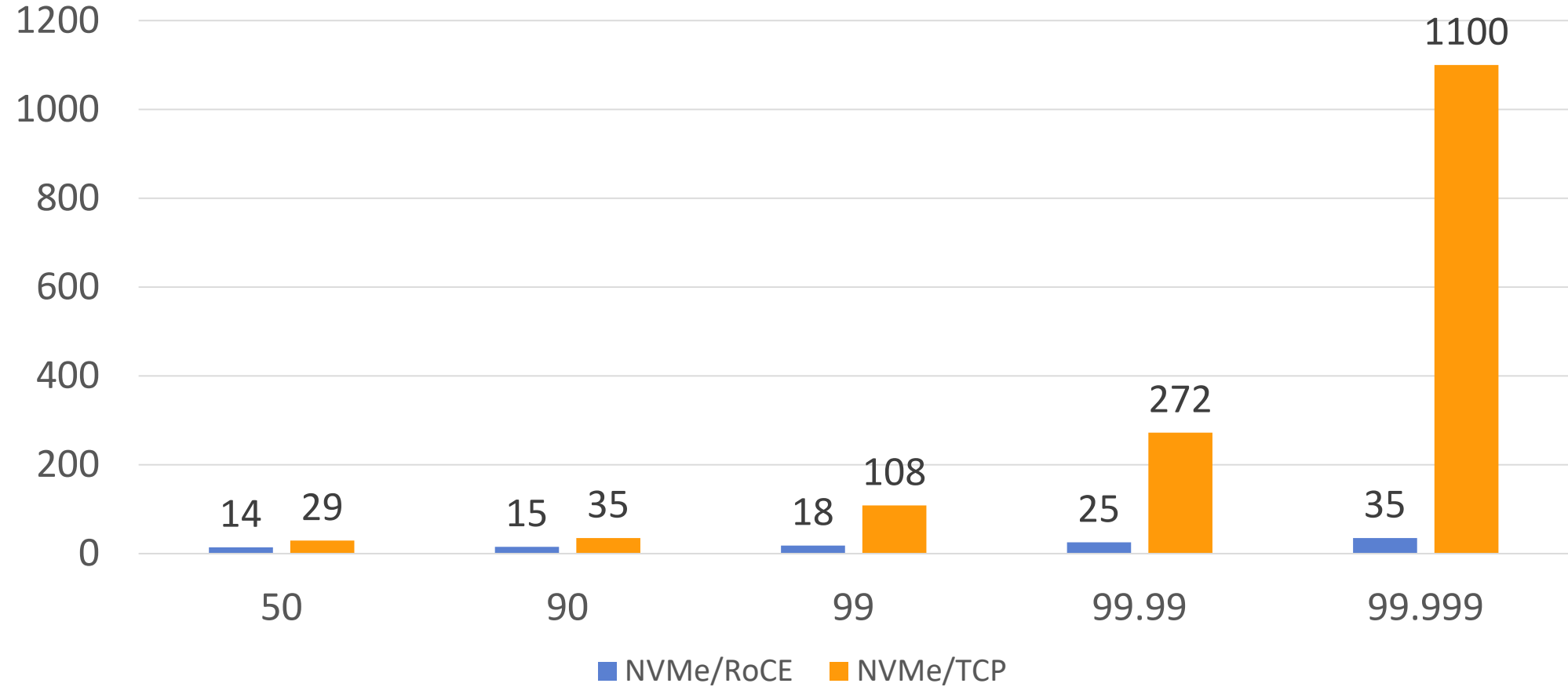| Environment | Metrics | Scale | Operations | Future |
|---|---|---|---|---|
| Target Loc<br><br>Accessibility<br><br>Distance<br><br>Existing Fabrics<br><br>Consumer Loc<br>Regulatory<br><br>Multi-tenancy | Perf: Latency<br>Perf: Predictability<br>Perf: Consistency<br>Perf: Bandwidth<br><br>Cost: $/Port<br>Cost: CPU/BW<br>Cost: CPU/IOPS | Single Rack<br>Multi-rack<br>Clos architecture<br>Oversubscription<br>Link aggregation<br>Redundancy | Onboarding<br>Configuration<br>Automation<br>Adv Telemtry<br>Intent Based<br>SW Defined <x> | Future Roadmap<br>Scale-up<br>Upgrade |

**Western Digital.**

# Case Study: Fabrics Comparison (partial sample)

| | NVMe/RoCE | NVMe/TCP |
|---|---|---|
| **Max Speed (current->next gen)** | 200G → 400G | 200G → 400G |
| **Link Aggregation** | Yes. HW based | Yes. HW based |
| **½ Round Trip Transport Latency** | 1.4us | 8-30us |
| **4k Write Latency (50th Percentile)** | 14us | 31us |
| **4k Write Latency – Tail/QoS (99.99th percentile)** | 25us | 272us |
| **Encapsulation** | UDP | TCP |
| **Routability** | Routable UDP based | Routable TCP based |
| **Scale** | Multi Rack | Multi Rack |
| **Convergence with other traffic** | Yes | Yes |
| **Switch ASIC (Merchant Silicon)** | Yes | Yes |
| **Disaggregated Switches** | Yes | Yes |
| **SDN** | Yes | Yes |

**Western Digital.**

# Latency Comparison



Latency (us) Percentiles

Source: Western Digital Performance Tests

# Latency vs. IOPS



Legend:
- RoCE P99
- RoCE P99.99
- TCP P99
- TCP P99.99

Y-axis: Latency (us) — 0, 100, 200, 300, 400, 500, 600, 700
X-axis: IOPS — 100k, 200k, 300k, 400k, 500k, 600k, 700k, 800k, 900k, 1M

Source: Western Digital Performance Tests

# Test Setup

- Linux kernel 5.0
- Mellanox ConnectX-5 NIC
- Mellanox 2700 32x100G switch
- Intel® Xeon® Gold 6150 CPU @ 2.70GHz
- 100G RAM disk

# Summary

- NVMe/RoCE and NVMe/TCP are complimentary technologies

- RoCE has lower and more consistent latency

- RoCE needs DCB

- RoCE uses less CPU cycles

- TCP does not need DCB

- TCP appears less optimized for performance and efficiency

- No "One Size Fits All"

# Western Digital.®

8/9/2019