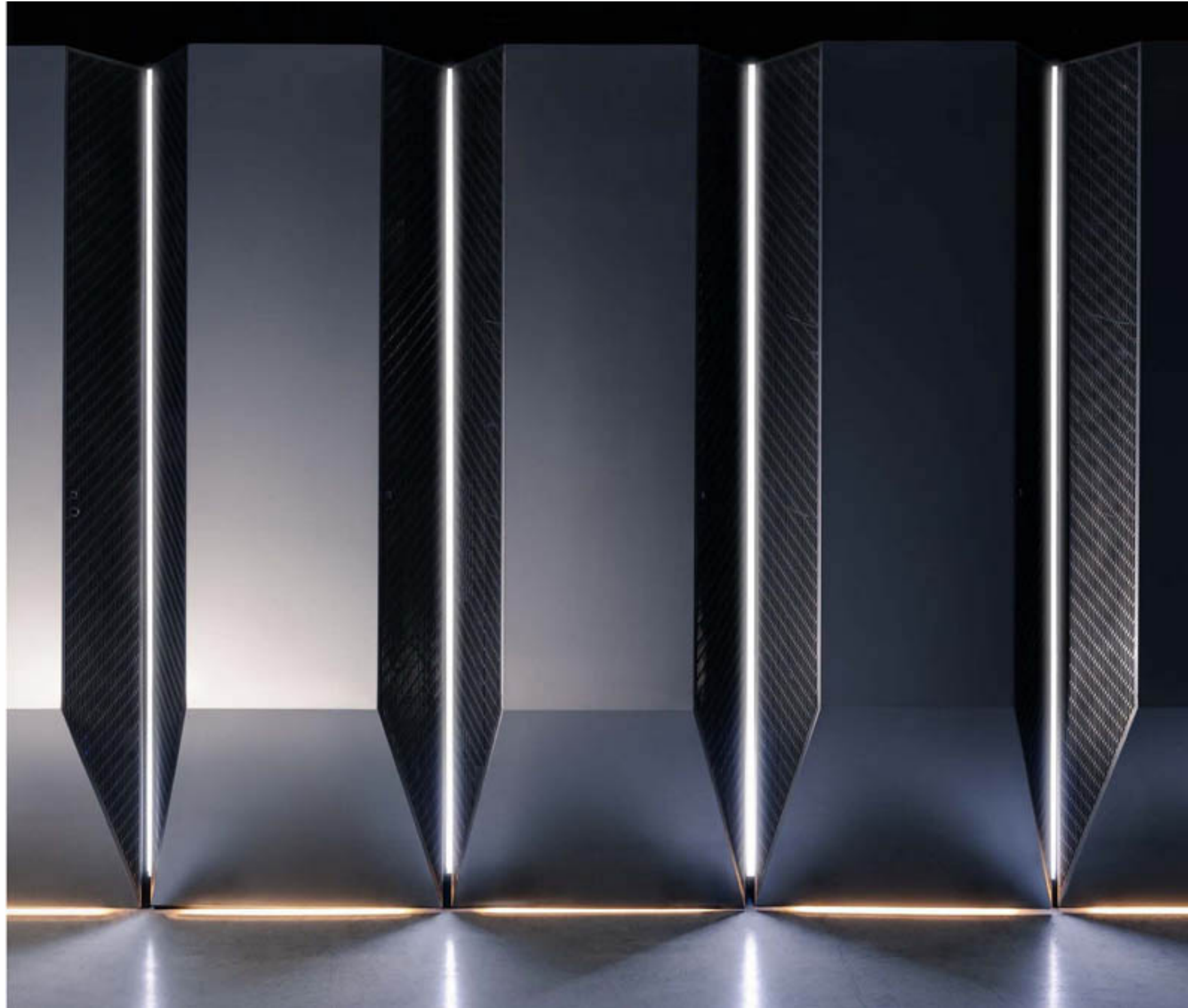TESLA

# Tesla Transport Protocol over Ethernet (TTPoE)

A new lossy, Exa-Scale fabric for the Dojo AI Supercomputer

Eric Quinnell, Ph.D.
Dojo Fabric Lead

# Problem Statement

TCP/IP is too slow for scaled AI interconnect
- Bound by CPU SW kernel

Lossless fabrics are complex and brittle
- Priority Flow Control (PFC) affects the global network

**Torsten Hoefler**
ETH Zürich and Microsoft

**Duncan Roweth, Keith Underwood, Bob Alverson**
Hewlett Packard Enterprise

**Mark Griswold, Vahid Tabatabaee, Mohan Kalkunte, Surendra Anubolu**
Broadcom

**Siyan Shen**
ETH Zürich

**Abdul Kabbani, Moray McLaren, Steve Scott**
Microsoft

**Ideal Fabric:**

- Lowest latency
- Highest bandwidth
- Simple Software

**For Tesla AI:**
- Layer 2 only
- Collective communications and ingest
- Low congestion, single application

T E S L A

# TTPoE

Tesla Transport Protocol over Ethernet (TTPoE)
is a peer-to-peer ethernet Transport Layer Protocol executed **entirely in hardware**.



Why a custom transport protocol?

1. **Vertical Integration** – *extend Dojo RDMA onto optical fabric*

2. **"Lossy" ethernet network** – *ease of scaling, cost, congestion mgmt.*

3. **Use 3rd party hardware** – *Ethernet II frames "Just Work"*

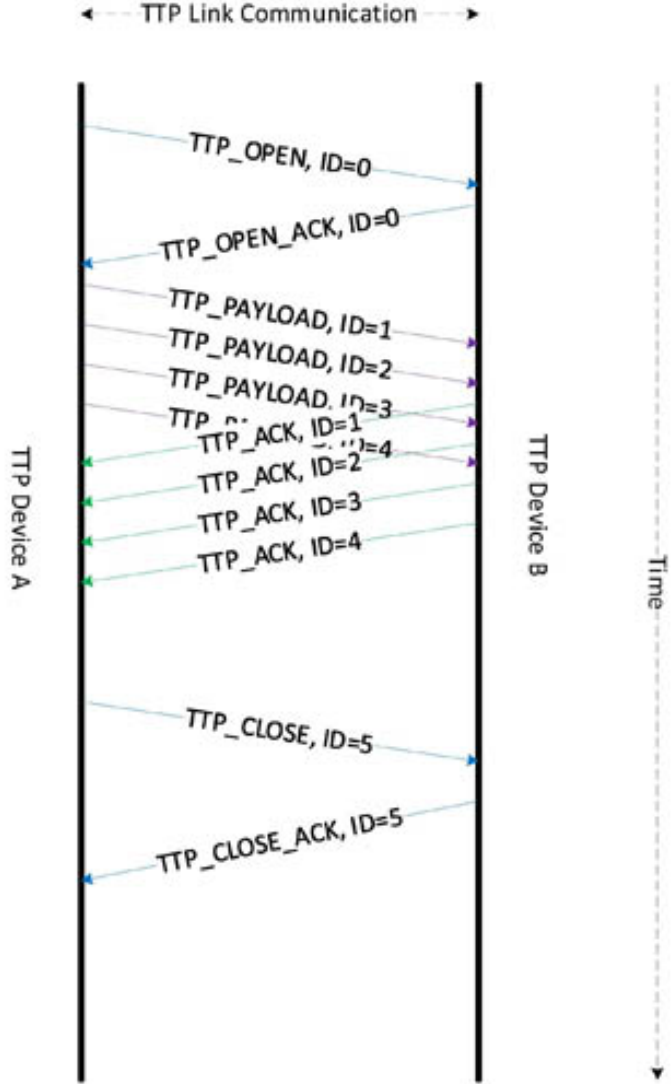*TCP got it right – just do it in hardware*

# Dojo OSI Layers

## Standard Stack

| OSI Layer | Example Protocols (TCP/IP) | TCP/IP Implementation |
|---|---|---|
| **Layer 7** Application | HTTP, Telnet, FTP | Software |
| **Layer 6** Presentation | JPEG, PNG, MPEG | |
| **Layer 5** Session | NFS, SQL | |
| **Layer 4** Transport | TCP, UDP | |
| **Layer 3** Network | IPv4/IPv6 | |
| **Layer 2** Data Link | Ethernet Frames, MAC addresses, VLAN | Hardware |
| **Layer 1** Physical | Data Encoding, Physical Specs | |

## Dojo Stack

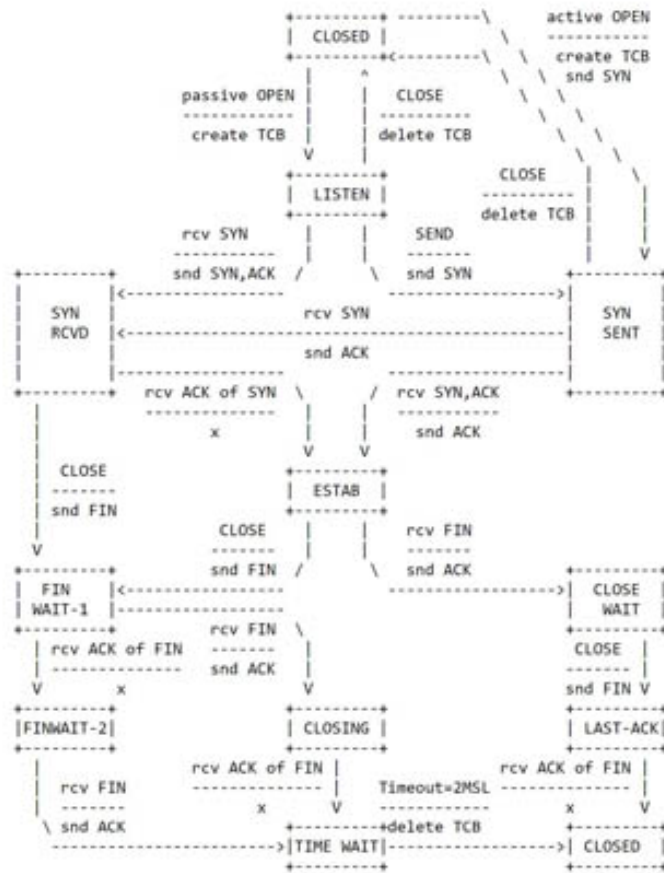| OSI Layer | Example Protocols | Dojo Implementation |
|---|---|---|
| **Layer 7** Application | Pytorch, Dojotorch | Software |
| **Layer 6** Presentation | FFMPEG, HEVC, YUV | |
| **Layer 5** Session | Dojo RDMA Descriptors | |
| **Layer 4** Transport | **TTP** | Hardware |
| **Layer 3 (Optional)** Network | IPv4/IPv6 (Optional) | |
| **Layer 2** Data Link | Ethernet Frames, MAC addresses, VLAN | |
| **Layer 1** Physical | Data Encoding, Physical Specs | |

# TTP transaction examples



**Clean TTP transfer Example**

**NACK TTP transfer Example.**
TTP_PAYLOAD, ID=3 is either lost or out of order

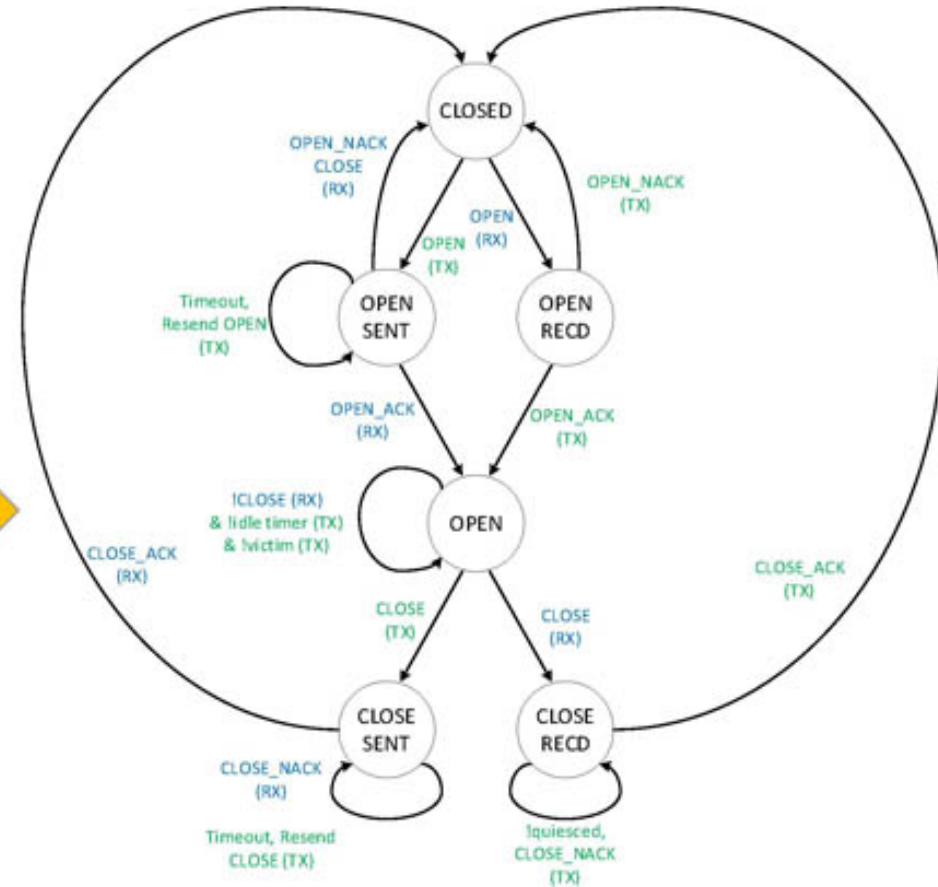T E S L A

# Transport Layer State Machines



**TCP STATE MACHINE**

TCP Connection State Diagram
Figure 6.
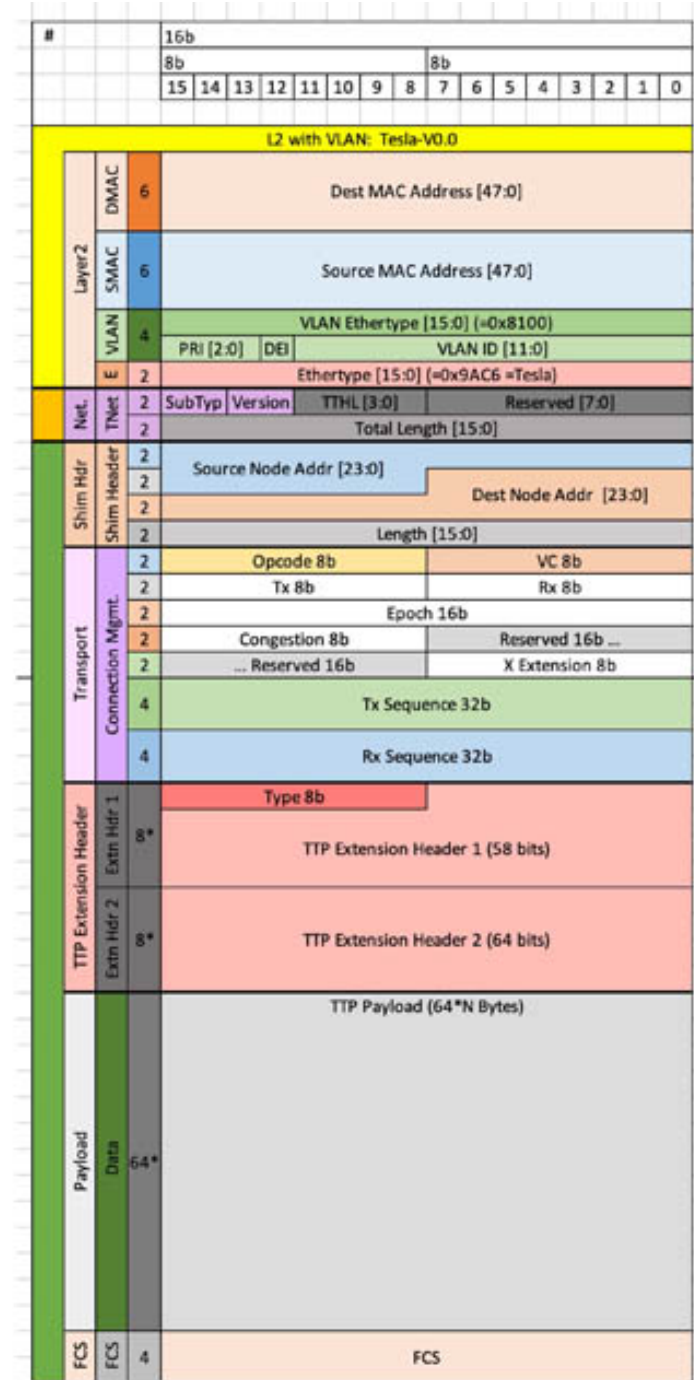
IETF RFC-793

**TTP STATE MACHINE**

HW CONSTRAINED

Modifications made for hardware-only execution

- 2 millisecond quiesce in a microsecond protocol is too long

- No reliance on virtual memory – physical memory only

- Automatic OPEN/CLOSE with no SW involvement

T E S L A

# TTP Header Frame

**TTP uses Ethernet-II simple formats with optional standard Layers**

- Dojo at scale uses only Layer 2, currently not using Layer 3

- MAC addresses are a hardware hash of the SOW Physical Address (PA)

- A TTP endpoint can concurrently handle 512 unique links, dynamically replaced via victimization and LRU

- Virtual channels (VCs) allow for non-blocking control, semaphore, completion, and data movement



T Ξ 5 L ā

# Lossy Protocol

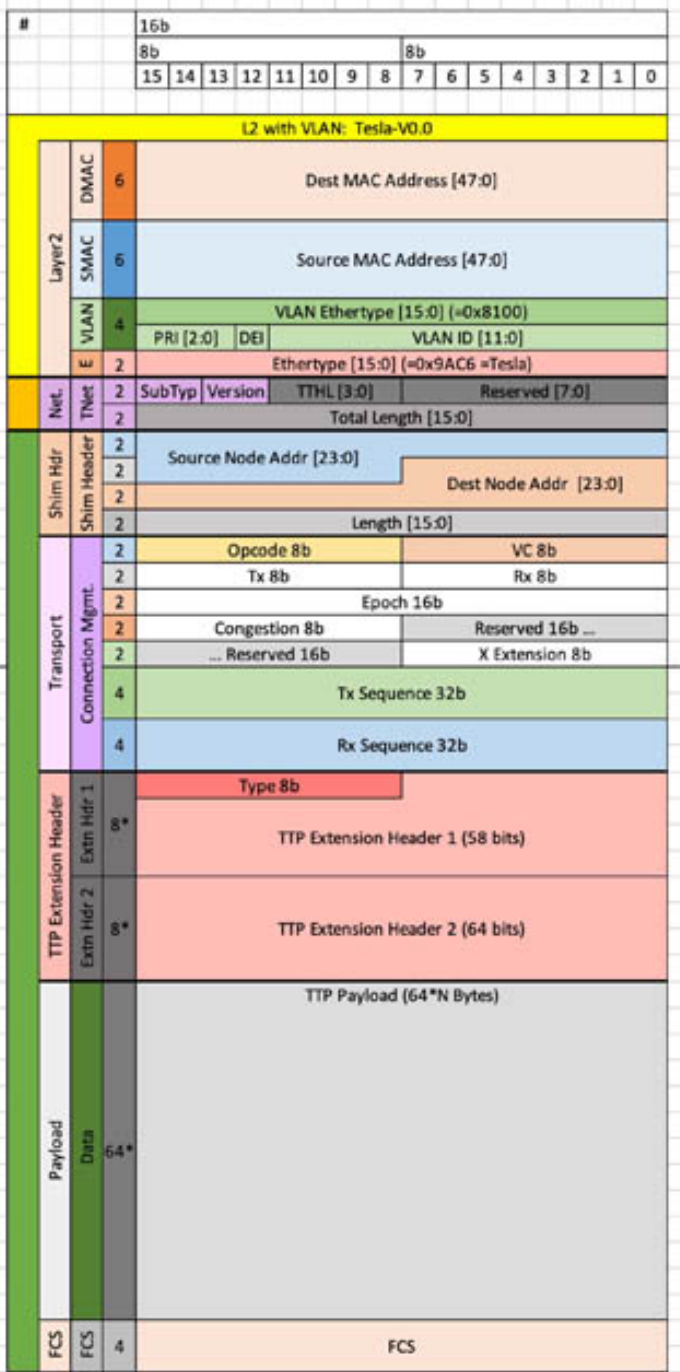## TTPoE is a "lossy" transport protocol

- "Lossy" transport meaning the underlying medium expects to lose packets and retry – full packet transmission is still guaranteed.
  - Similar to TCP and unlike UDP.

- TTP will default to packet drops and replays in corner cases of congestion, backpressure, or errors

- Speculative transmission is limited by SRAM size before a RTT ACK. This, in effect, forces a "TTP window size" beyond which bandwidth is lost

- Local SRAM lines are not retired/deallocated until the ACK comes back, allowing HW to replay the line.

- Replay amounts are also limited by SRAM, constraining the scale of replay storms



| # | | | 16b | | |
|---|---|---|---|---|---|
| | | | 8b | | 8b |
| | | | 15 14 13 12 11 10 9 8 | 7 6 5 4 3 2 1 0 | |

L2 with VLAN: Tesla-V0.0

| Layer 2 | DMAC | 6 | Dest MAC Address [47:0] | |
|---|---|---|---|---|
| | SMAC | 6 | Source MAC Address [47:0] | |
| | VLAN | 4 | VLAN Ethertype [15:0] (=0x8100) | |
| | | | PRI [2:0] DEI | VLAN ID [11:0] |
| | E | 2 | Ethertype [15:0] (=0x9AC6 =Tesla) | |

| Net. | TNet | 2 | SubTyp Version TTHL [3:0] | Reserved [7:0] |
|---|---|---|---|---|
| | | 2 | Total Length [15:0] | |

| Shim Hdr | Shim Header | 2 | Source Node Addr [23:0] | |
|---|---|---|---|---|
| | | 2 | | Dest Node Addr [23:0] |
| | | 2 | | |
| | | 2 | Length [15:0] | |

| Transport | Connection Mgmt. | 2 | Opcode 8b | VC 8b |
|---|---|---|---|---|
| | | 2 | Tx 8b | Rx 8b |
| | | 2 | Epoch 16b | |
| | | 2 | Congestion 8b | Reserved 16b … |
| | | 2 | … Reserved 16b | X Extension 8b |
| | | 4 | Tx Sequence 32b | |
| | | 4 | Rx Sequence 32b | |

| TTP Extension Header | Extn Hdr 1 | | Type 8b | |
|---|---|---|---|---|
| | | 8* | TTP Extension Header 1 (58 bits) | |
| | Extn Hdr 2 | 8* | TTP Extension Header 2 (64 bits) | |

TTP Payload (64*N Bytes)

| Payload | Data | 64* | |
|---|---|---|---|

| FCS | FCS | 4 | FCS |
|---|---|---|---|

T E S L A

# Congestion Management

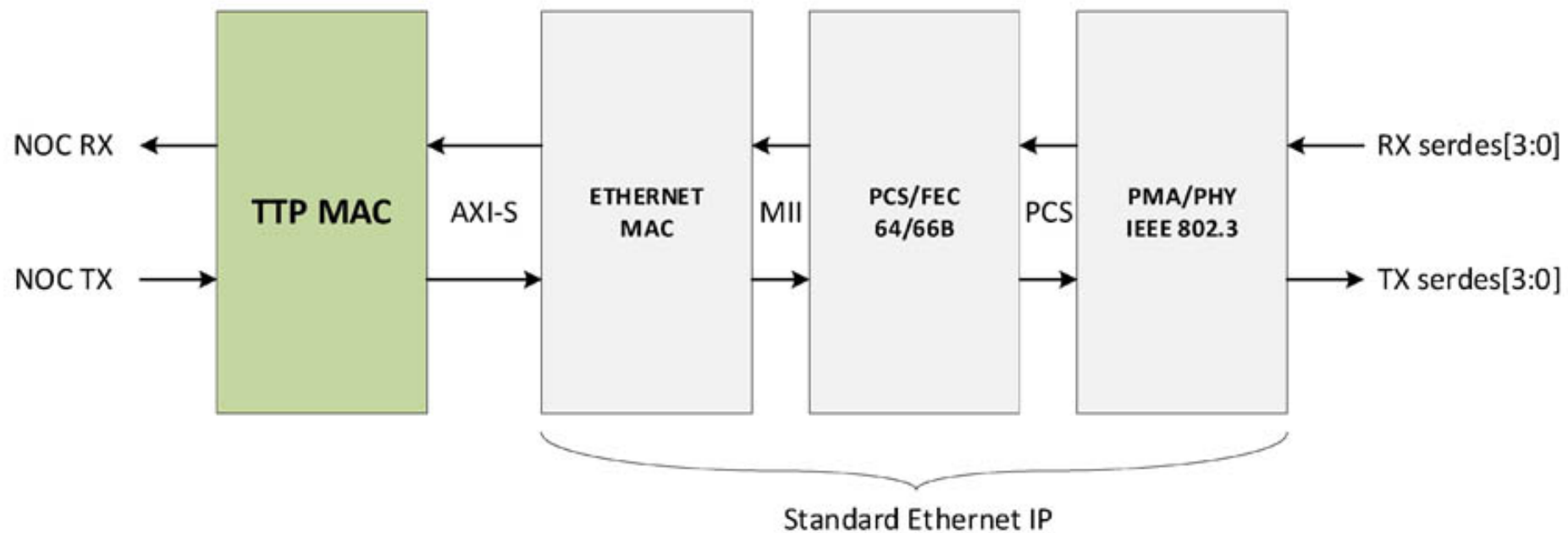## Congestion management is distributed

- Exponential backoff, rate control, and algorithms are handled by local link TX channels, not by central network or switch.

- Fault Tolerant flow "flushes" the TTP network and removes a bad link before continuing training

- No PFC, no Nagel Algorithm, no QoS, no tokens, no lossless artifacts

T Ξ S L Ā

# TTP MAC IP

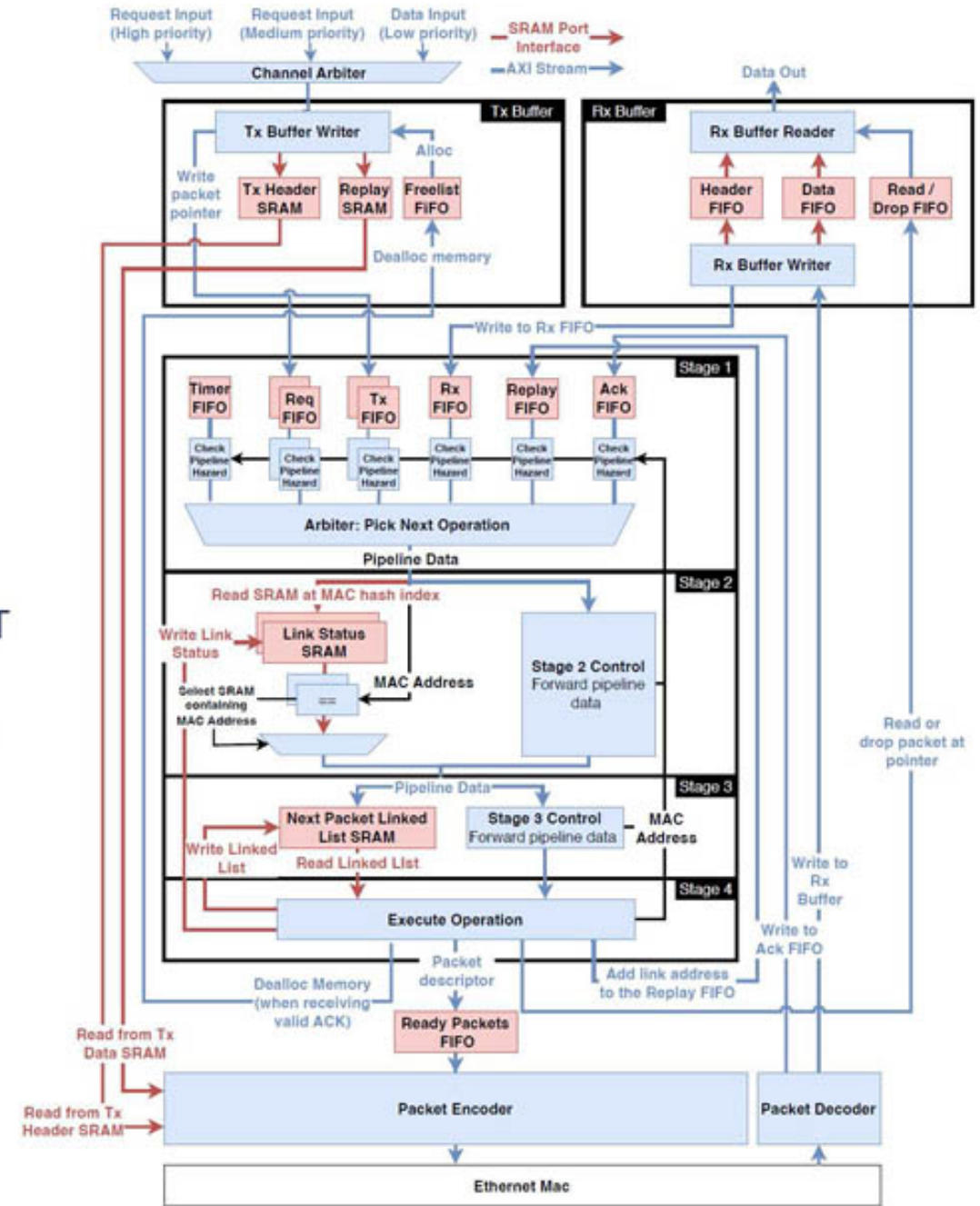The Transport Layer hardware is an IP block between a NOC and an Ethernet standard MAC
- Translates and coalesces 64B/cycle NOC packets into up to 1kB TTP Ethernet packets
- Speaks AXI-S or SOP/EOP formats
- Optionally activates standard MAC features – pause packets, counters, stats, LLDP
- IP block instantiated in FPGA and Silicon implementations

# TTP MAC Micro-Architecture

TTP's Micro-Architecture uses techniques from SMP Caches, Snoop Filters, CPUs
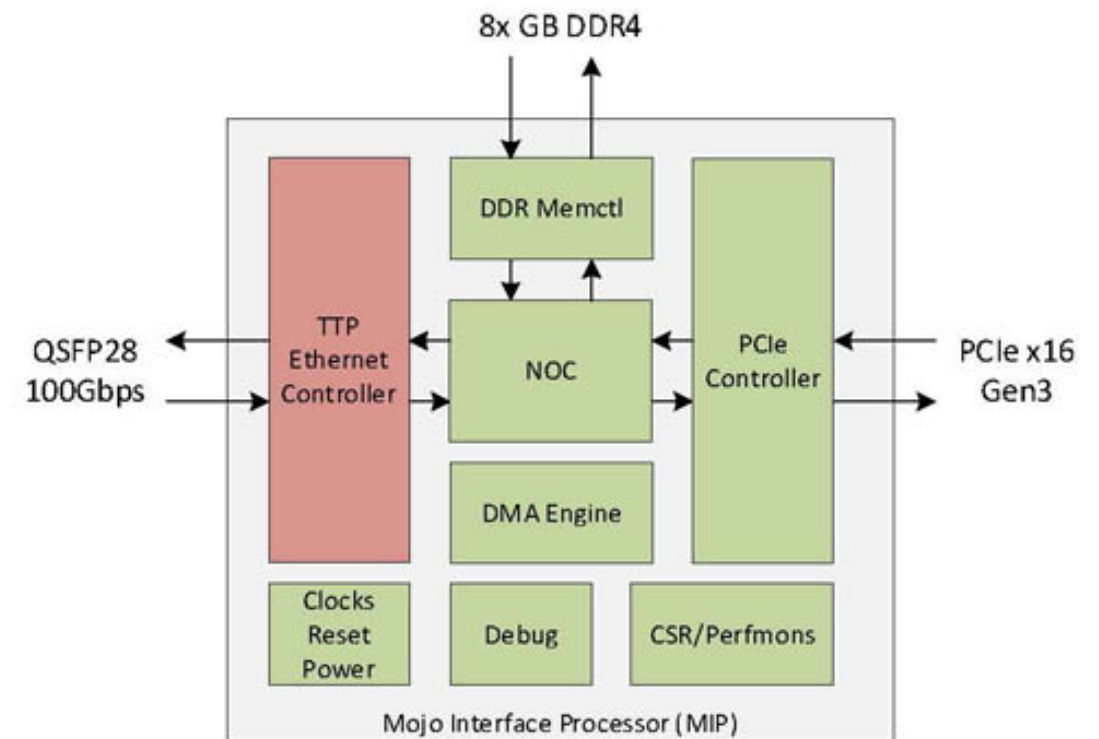
- 4-stage Read-Modify-Write (RMW) Pipeline
- TX Buffer size determines maximum outstanding packets before stall/backpressure
  - ACK packets "retire" a packet from the common buffer
  - 1MB TX Buffer allows for ~80 microseconds latency tolerance RTT
- Virtual Channels to prioritize and avoid livelock/deadlock
- Multi-channel "coherent" arbitration to update link and use the TX Physical Channel
- DMA descriptors issue to TTP MAC
  - Can be PUSH for implicit pass-thru local-to-remote
  - Can be explicit HBM2HBM fabric memcpy



TESLA

# "Mojo" 100Gbps Dumb-NIC



| Feature | Spec |
|---|---|
| Ethernet Speed | 100Gbps QSFP |
| PCI-e | Gen3 x16 |
| Memory | 8GB DDR4 |
| Power | <20W max |
| Reliability | 5-year tested |
| DMA engine | Dojo DMA |
| CPU+OS | None |
| Active Links | 512 unique, 2-way, LRU |



8x GB DDR4

QSFP28 100Gbps

PCIe x16 Gen3

DDR Memctl

TTP Ethernet Controller

NOC

PCIe Controller

DMA Engine

Clocks Reset Power

Debug

CSR/Perfmons

Mojo Interface Processor (MIP)

# Second integration box – Dojo Training Tile

**5x5 array of known good D1 chips**
- 4.5TB/s off-tile bandwidth per edge
  - Half of in-tile bandwidth

**Fully integrated module**
- Electrical + thermal + mechanical
- 15kW of power delivery

**Custom power delivery**
- Horizontal data communication plane
- Vertical power delivery and cooling
- 15kW per module

**Custom high-density connectors**
- Seamless connection to neighboring training tiles

# V1 Dojo Interface Processor

## 32GB High-Bandwidth Memory
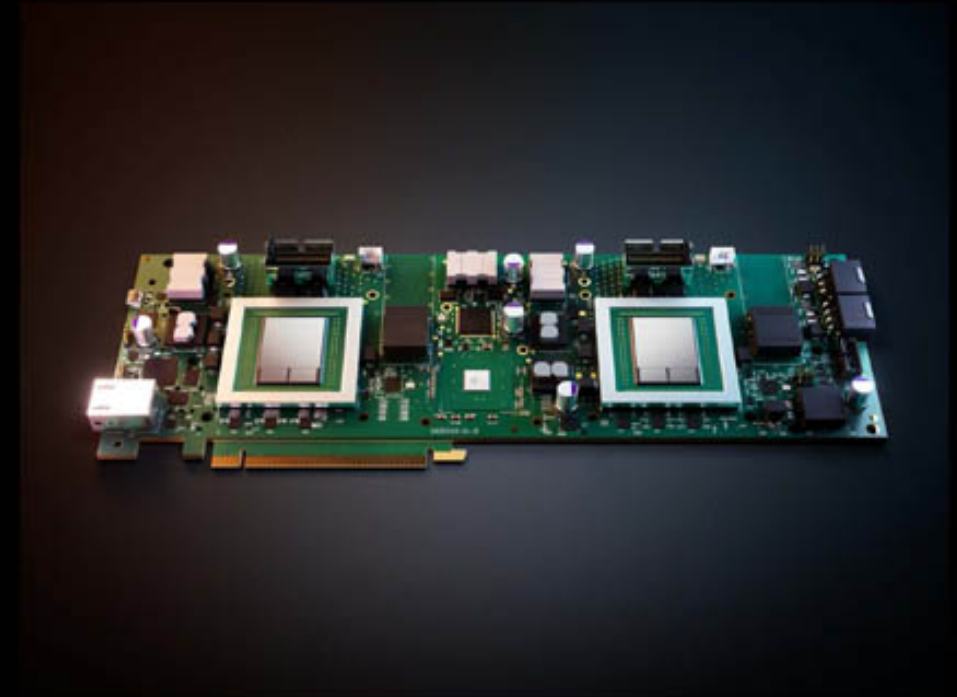
- 800 GB/s Total Memory Bandwidth

## 900 GB/s TTP Interface

- Tesla Transport Protocol (TTP) - Full custom protocol
- Provides full DRAM bandwidth to Training Tile

## 50 GB/s TTP over Ethernet (TTPoE)

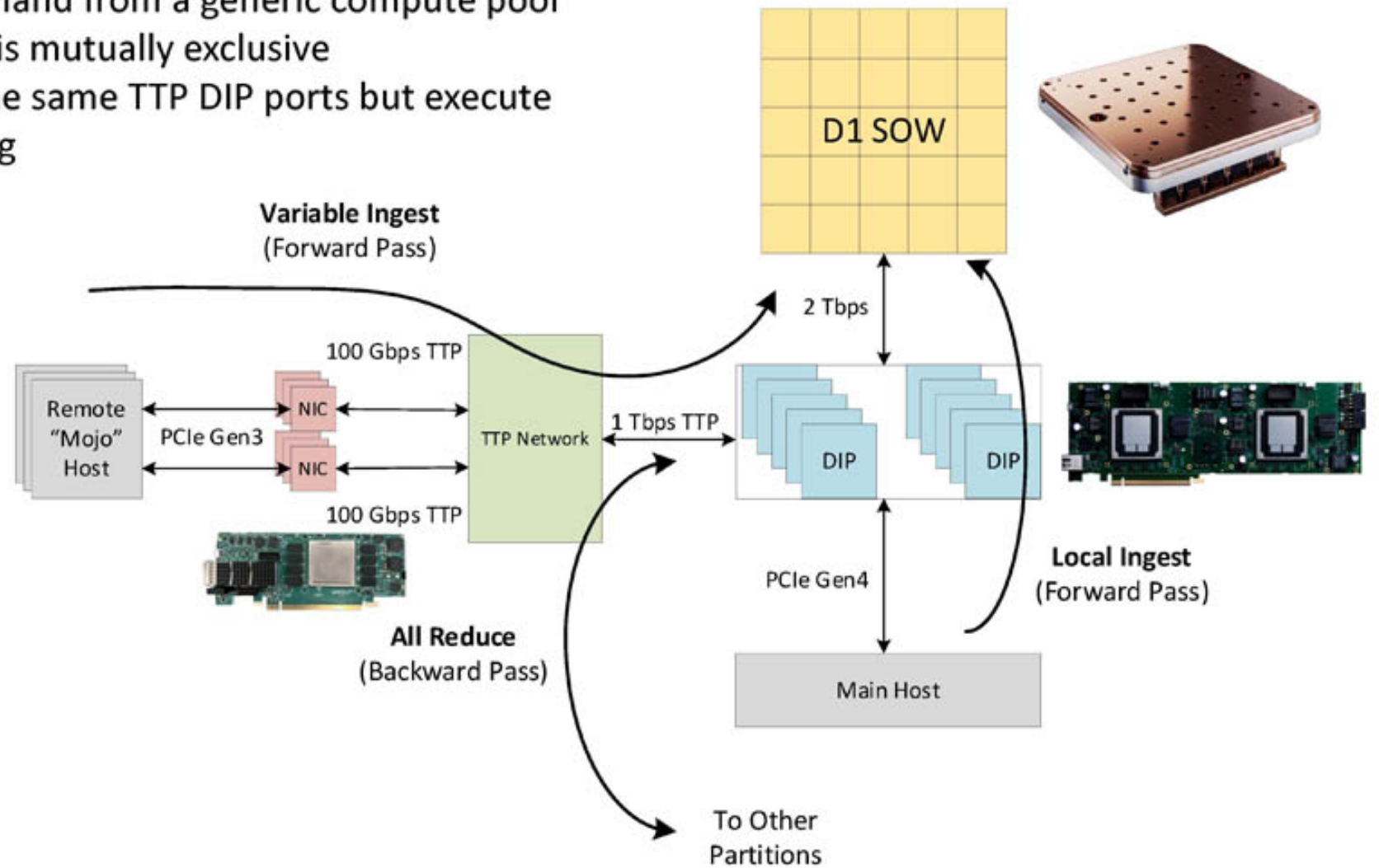- Enables extending communication over standard Ethernet
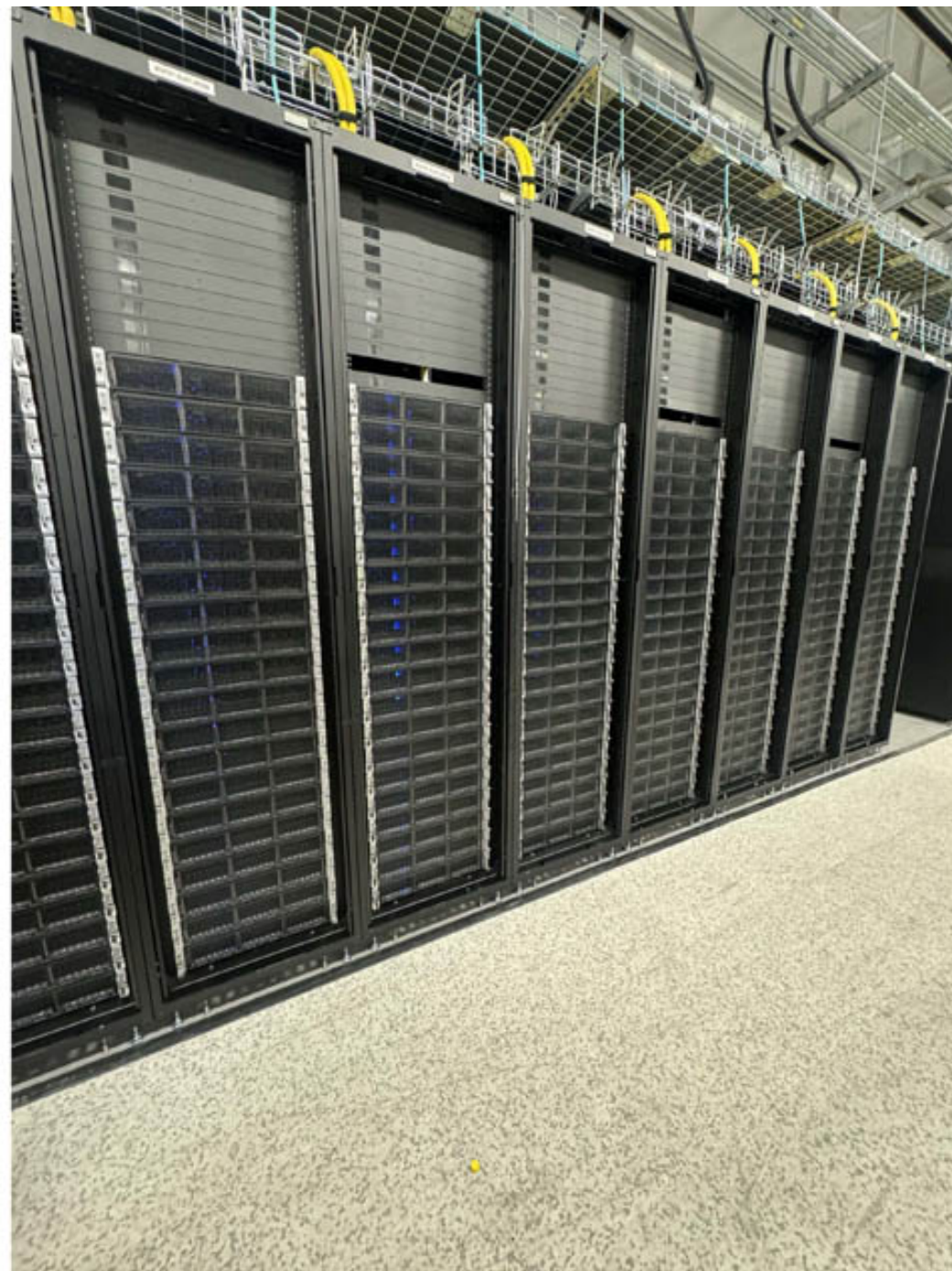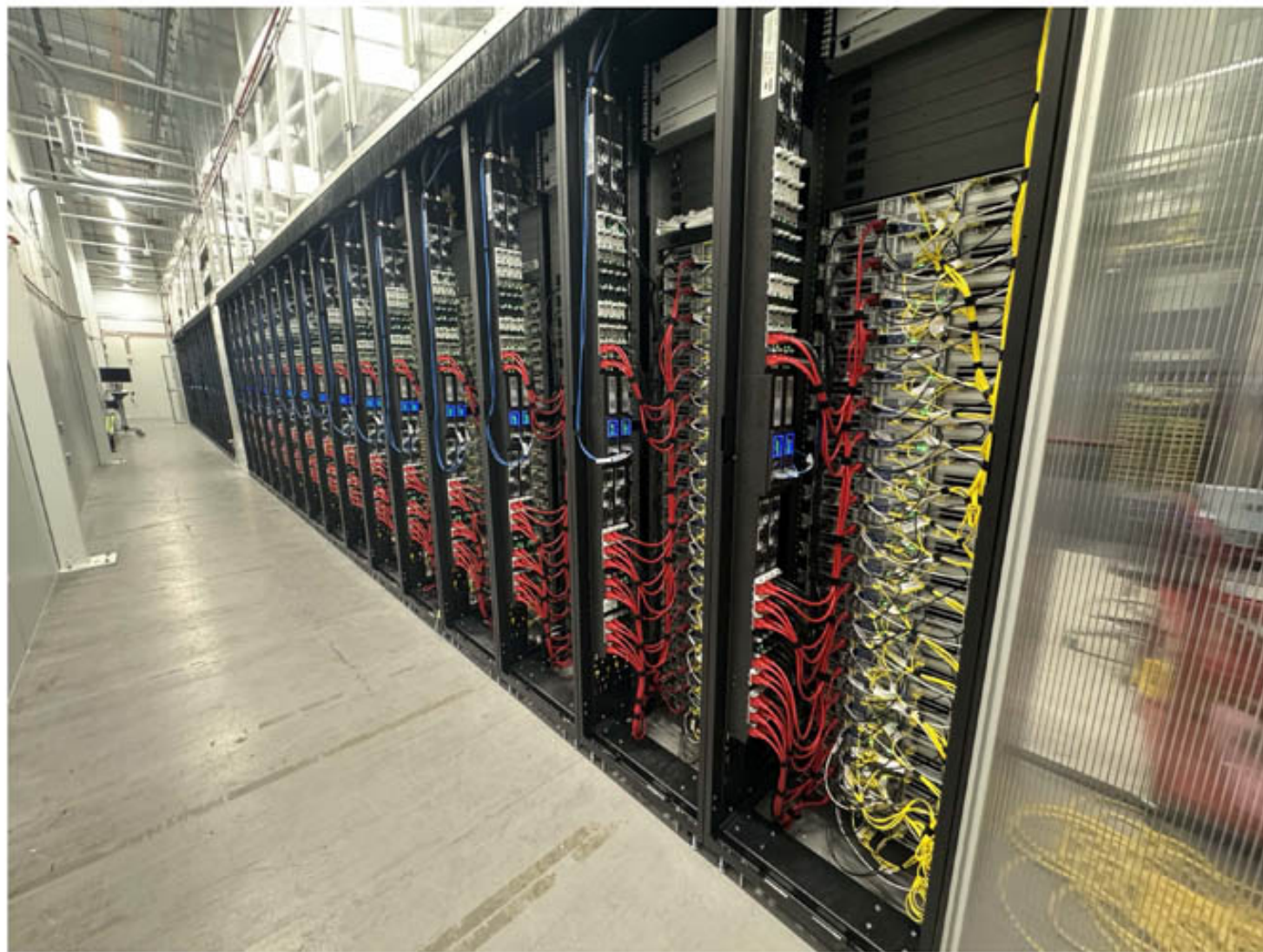- Native hardware support

## 32 GB/s Gen4 PCIe Interface

# "Mojo" Hosts – Variable Ingest via TTP Network

Vision networks can be heavily ingest limited
- Vision-based tensors and training clips in GBs
- "Mojo" Hosts are scheduled on demand from a generic compute pool
- Forward/Backward pass TTP traffic is mutually exclusive
  - i.e. ingest and all-reduce share the same TTP DIP ports but execute during different phases of training



**Variable Ingest**
(Forward Pass)

D1 SOW

2 Tbps

100 Gbps TTP

Remote "Mojo" Host

PCIe Gen3

NIC

NIC

TTP Network

1 Tbps TTP

DIP

DIP

100 Gbps TTP

PCIe Gen4

**Local Ingest**
(Forward Pass)

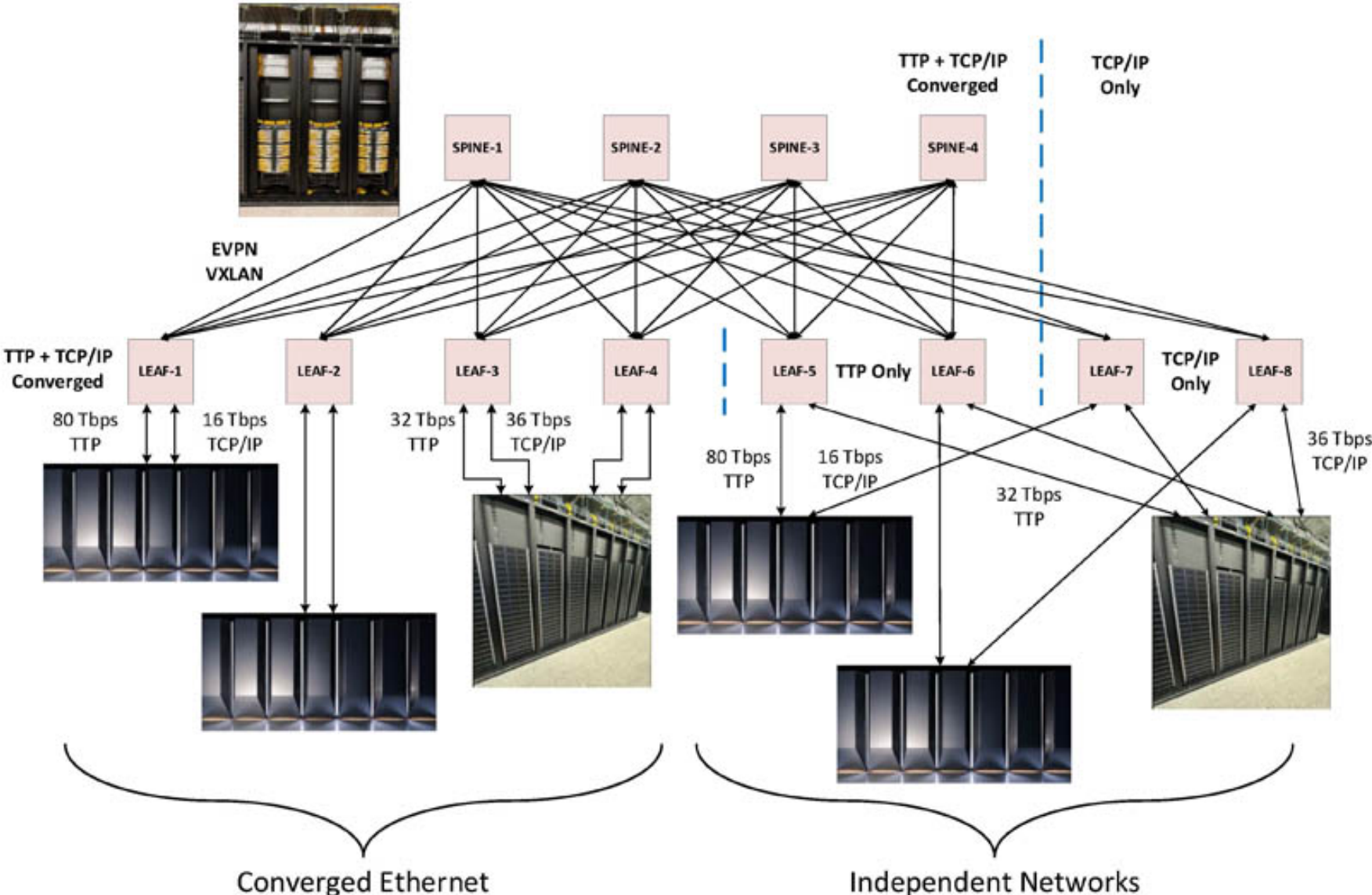**All Reduce**
(Backward Pass)

Main Host

To Other Partitions

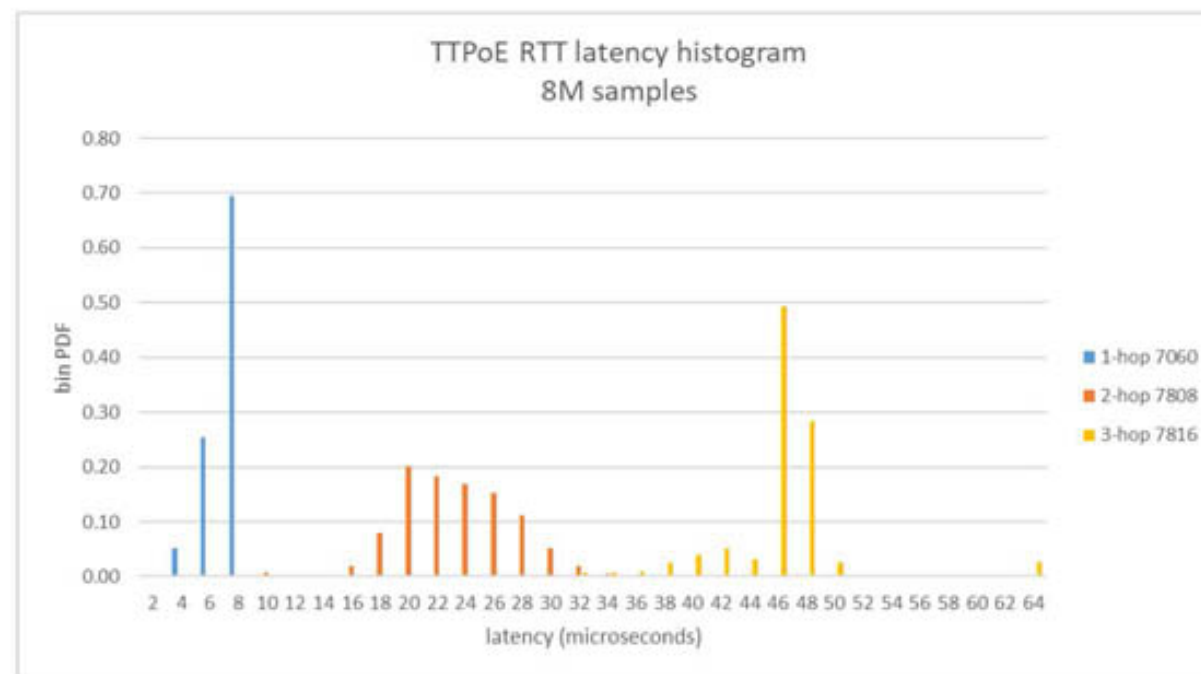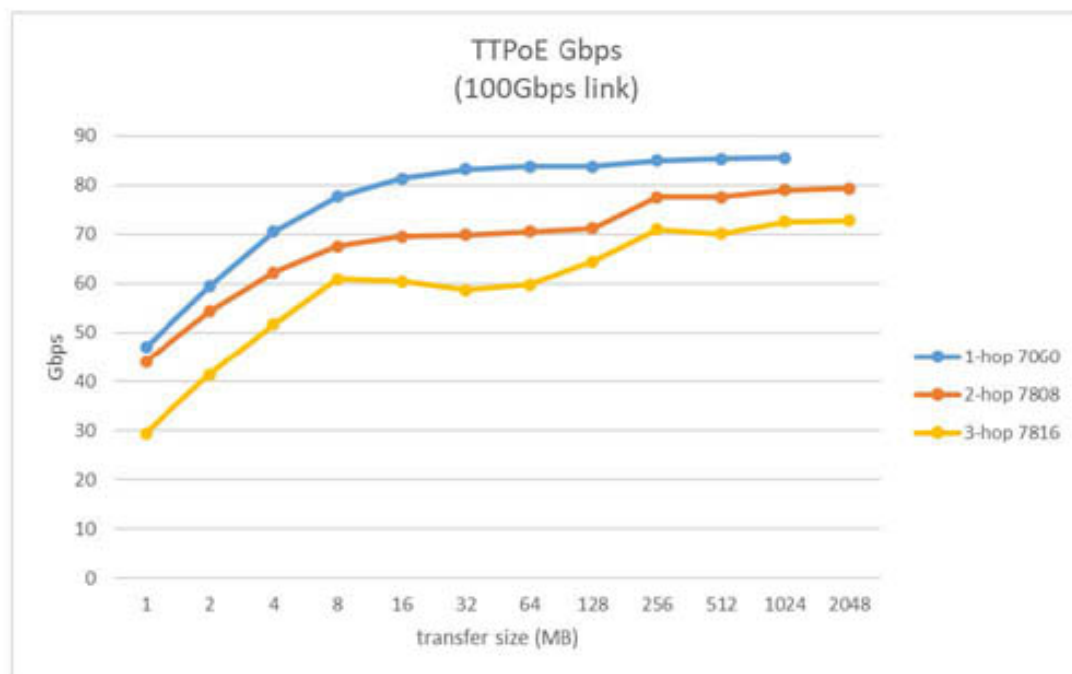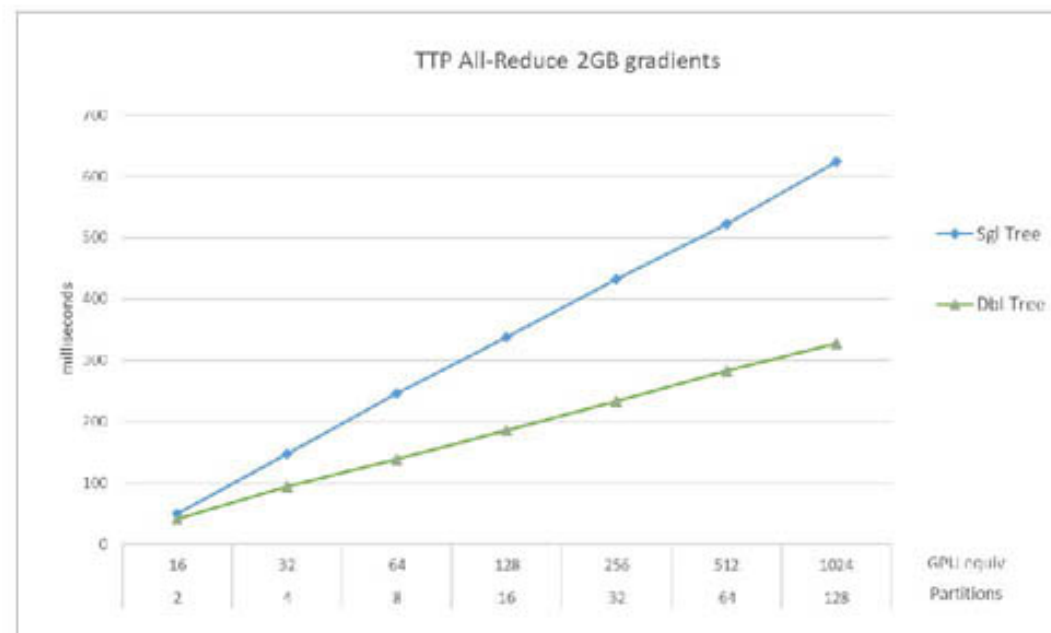TESLA

# MDCH – Mojo Dojo Compute Hall

# Dojo Engineering System

- 4xExaFLOP BF16/FP16 Cluster
- 40 PB Local Storage
- 40,960 Main Host Cores
- 61,440 Mojo Host Cores
- 320 Tbps TTP All-Reduce I/O (endpoint)
- 128 Tbps TTP Ingest I/O (endpoint)
- 208 Tbps TCP/IP (endpoint)
- Converged and non-Converged network experiments



EVPN
VXLAN

TTP + TCP/IP
Converged

TCP/IP
Only

SPINE-1   SPINE-2   SPINE-3   SPINE-4

TTP + TCP/IP
Converged

LEAF-1   LEAF-2   LEAF-3   LEAF-4   LEAF-5   TTP Only   LEAF-6   LEAF-7   TCP/IP   LEAF-8
                                                                          Only

80 Tbps   16 Tbps       32 Tbps   36 Tbps        80 Tbps   16 Tbps                    36 Tbps
TTP       TCP/IP        TTP       TCP/IP         TTP       TCP/IP                     TCP/IP

                                                                    32 Tbps
                                                                    TTP

Converged Ethernet                              Independent Networks

# Results

- Measured on Arista 7060, 7808, and 7816 switches
- RTT latency is random sampling of in-flight packets + ACK return
- Gbps is wall time real-data movement
- All-reduce measure is network only, non-pipelined
  - SOW has all-reduce not shown (pre-network)
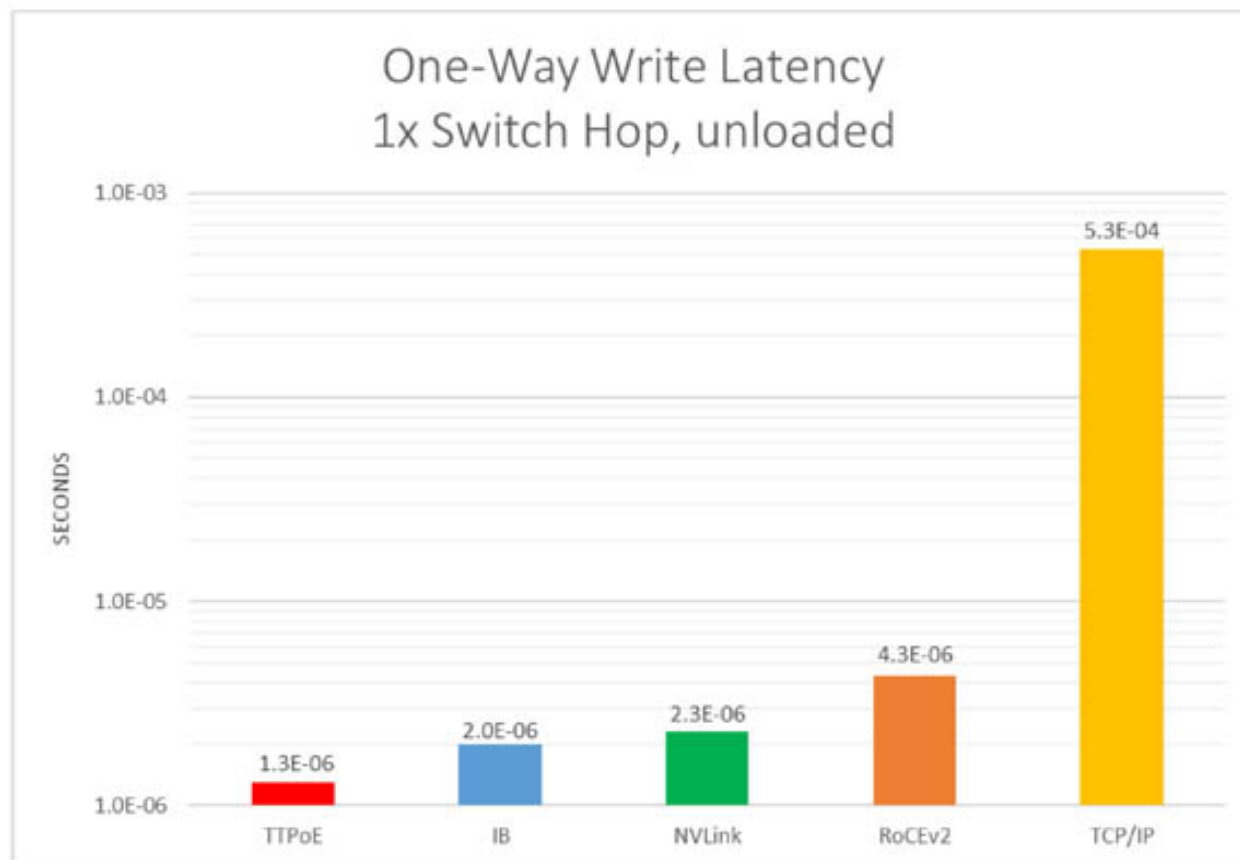- All-reduce throughput is determined by the slowest node in system



TTP All-Reduce 2GB gradients



TTPoE Gbps
(100Gbps link)



TTPoE RTT latency histogram
8M samples

T Ξ S L A

# Backup – Latencies

*Intended de-emphasis on synthetic latency measurements*

Differences of greater consequence:
- lossy vs lossless
- centralized vs distributed congestion
- proprietary vs open source
- sustained bandwidths at scale



**One-Way Write Latency**
**1x Switch Hop, unloaded**

TTPoE, TCP/IP – Spectrum3 SN4700
IB – Spectrum 9700 IB
Nvlink – DGX-H100 NvSwitch level1 (internal)
RoCEv2 – 7812 R3

Inconsistent methodology and hardware, not at scale

T E S L A

# TTPoE in Ultra Ethernet Consortium (UEC)

**Ultra Ethernet**
*Consortium*

https://ultraethernet.org/

### Steering Members

AMD · ARISTA · BROADCOM · CISCO · EVIDEN an atos business

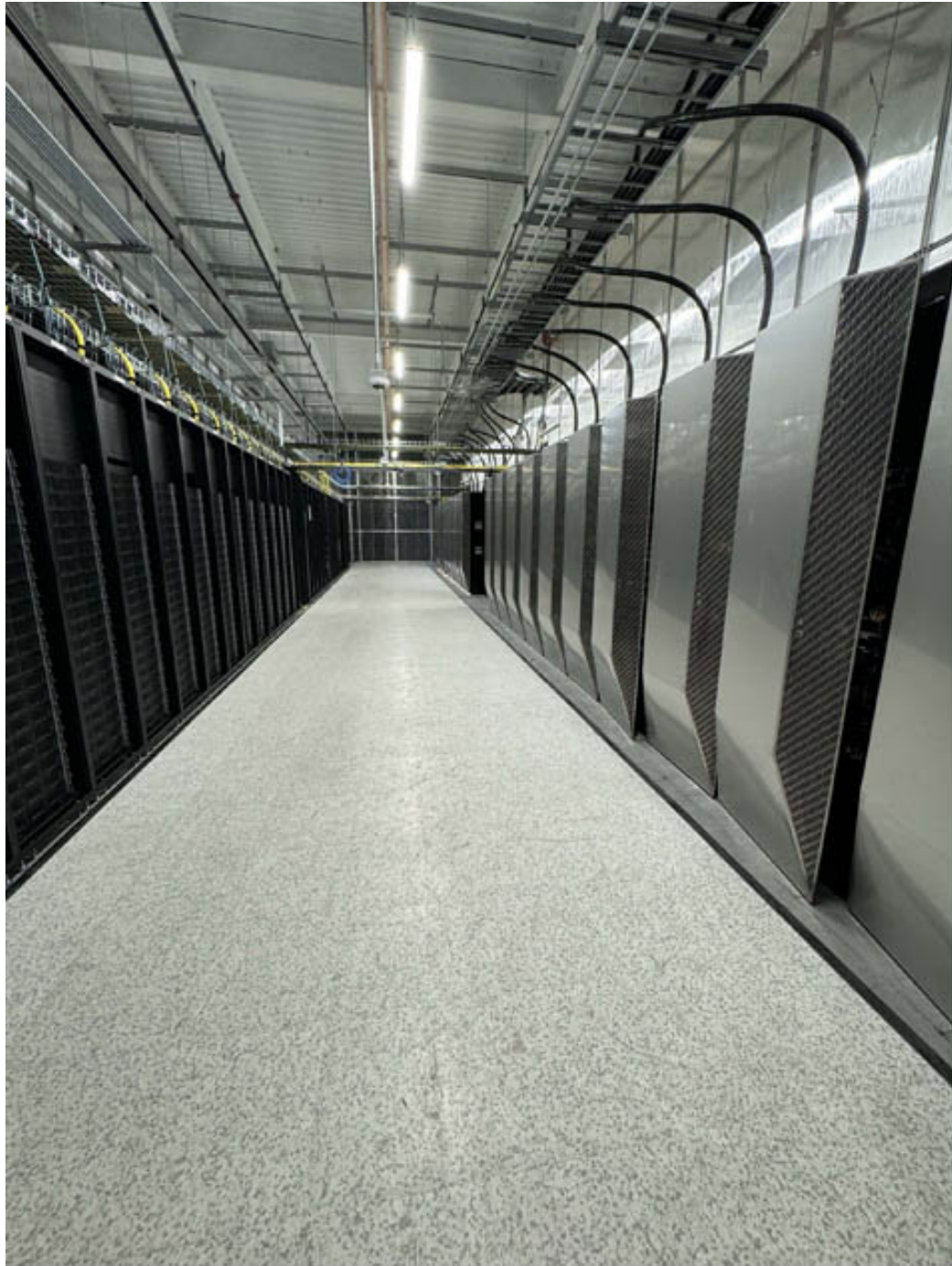Hewlett Packard Enterprise · intel · Meta · Microsoft · ORACLE

> While large lossless RoCE networks can and have been successfully deployed, they require careful tuning, operation, and monitoring to perform well without triggering these effects. This level of investment and expertise is not available to all network operators and leads to a high TCO. A transport protocol that does not depend on a lossless fabric is needed.

https://ultraethernet.org/wp-content/uploads/sites/20/2023/10/23.07.12-UEC-1.0-Overview-FINAL-WITH-LOGO.pdf

Tesla has achieved Exa-scale with a lossy fabric, executing real training runs deployed in FSD

## Tesla is joining the UEC and offering the TTPoE protocol publicly

TESLA

## Team Acknowledgements

**Prototyping is Easy. Scaling is Hard**

Thanks to the

TTPoE Original Inventors, Network Deployment Team, Silicon Design Team, System and Infrastructure Team, SW and Drivers Team, Linux Patch Team, SDN Team, DevOps Team, QA Team, DC Tech Team, Supply Team, and all TTP/Mojo Interns

T E S L A

Tesla Transport Protocol over Ethernet (TTPoE)