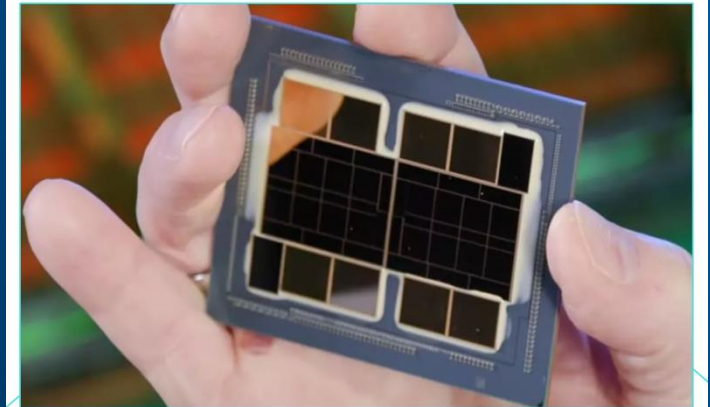# A Unified Network for HPC and AI

**Uri Elzur, Intel**

September 21st , 2023

Habana® Gaudi®2

Intel® Data Center GPU Max Series

# Notices and Disclaimers

Intel technologies may require enabled hardware, software or service activation. // No product or component can be absolutely secure. // Your cost and results may vary. // Performance varies by use, configuration and other factors. // See our complete legal Notices and Disclaimers.

Intel is committed to respecting human rights and avoiding complicity in human rights abuses. See Intel's Global Rights Principles. Intel's products and software are intended only to be used in applications that do not cause or contribute to a violation of internationally recognized human rights.
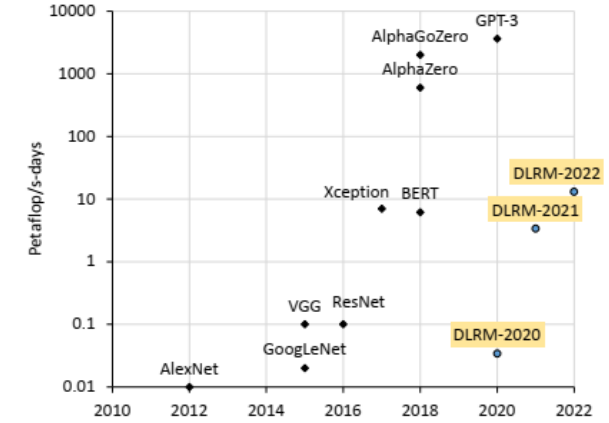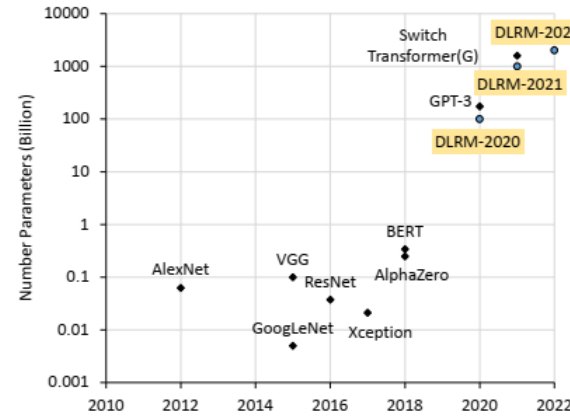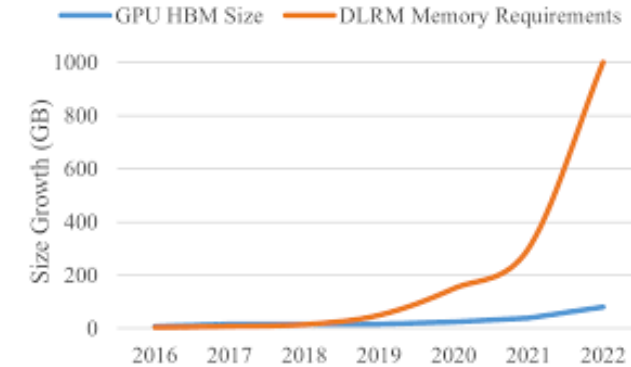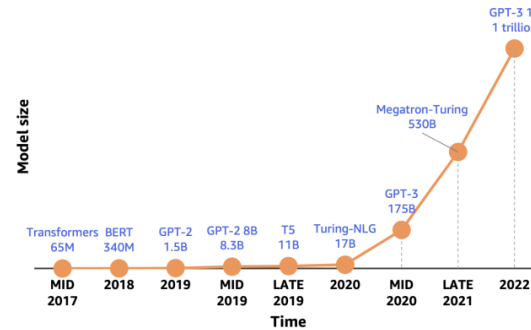
intel.

# AGENDA

- At a Crossroads

- Ultra Ethernet Consortium overview

- Summary

intel.

# At a Crossroads
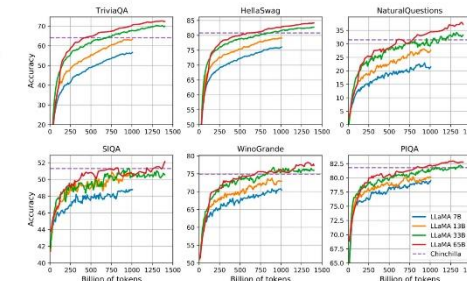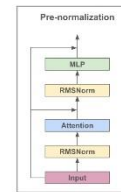
- **Workload large or small? Model vs Tokens**

- **Large Language Model and RecSys**

- Cloud approach vs dedicated Training cluster

- HPC embraces AI or AI using HPC? Convergence of sorts?

- GPU, TPU, Wafer Scale or other???

- Optics: NPO, CPO, Direct drive, Intra ASIC?

- Network: Lossy vs Lossless – religion or technology?



**15,000x increase in 5 years**

LLaMA: LLMs for Everyone!

# At a Crossroads

- WL large or small? Model vs Tokens

- LLM and RecSys

- **Cloud approach vs dedicated Training cluster**

- **HPC embraces AI or AI using HPC? Convergence of sorts?**

- GPU, TPU, Wafer Scale or other???

- Optics: NPO, CPO, Direct drive, Intra ASIC?

- Network: Lossy vs Lossless – religion or technology?



THE CONVERGENCE OF HPC * AI
Integrating the Third and Fourth Pillars of Scientific Discovery

| HPC | AI |
|---|---|
| 40+ years of algorithms based on first principles theory | New algorithms and models with potential to increase model size and accuracy |

Dramatically Improves Accuracy and /or Time-to-Solution at Large Scale

Commercially viable fusion energy

Improve or validate the Standard Model of Physics

Clinically viable precision medicine

Understanding cosmological dark energy and matter

Climate/weather forecasts with ultra-high fidelity

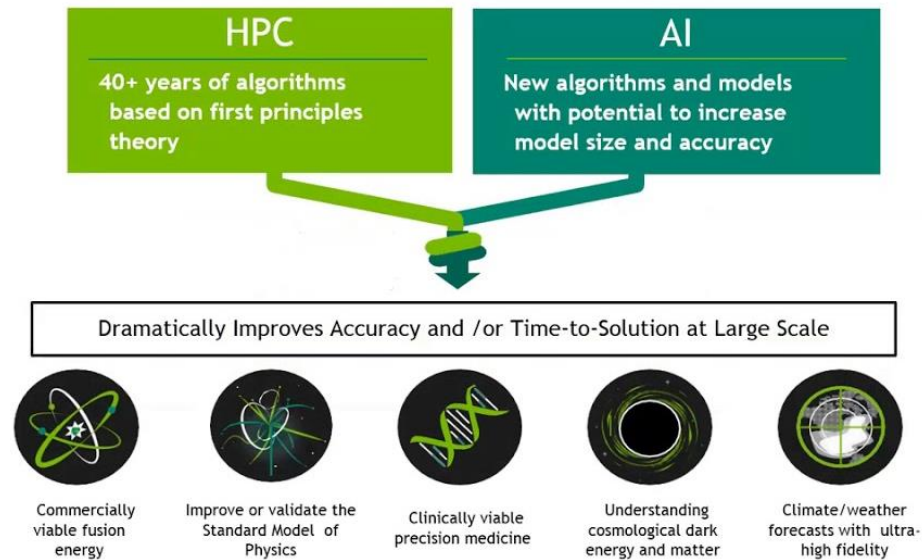# At a Crossroads

- WL large or small? Model vs Tokens

- LLM and RecSys

- Cloud approach vs dedicated Training cluster

- HPC embraces AI or AI using HPC? Convergence of sorts?

- **GPU, TPU, Wafer Scale or other???**

- Optics: NPO, CPO, Direct drive, Intra ASIC?

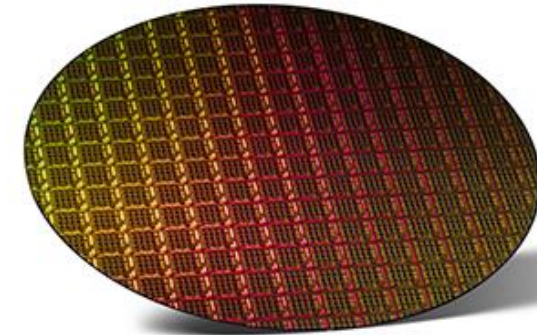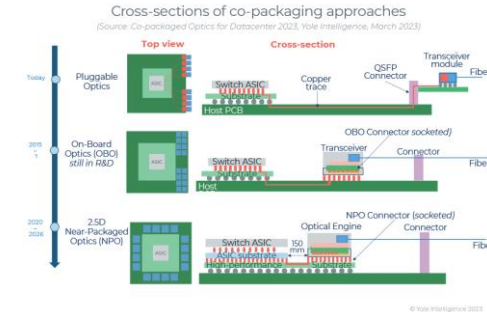- Network: Lossy vs Lossless – religion or technology?

intel.

# At a Crossroads

- WL large or small? Model vs Tokens

- LLM and RecSys

- Cloud approach vs dedicated Training cluster

- HPC embraces AI or AI using HPC? Convergence of sorts?

- GPU, TPU, Wafer Scale or other???

- **Optics: NPO, CPO, Direct drive, Intra ASIC?**

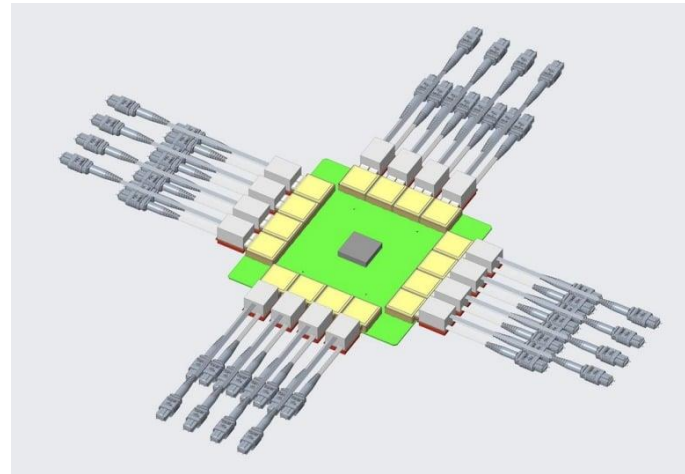- Network: Lossy vs Lossless – religion or technology?
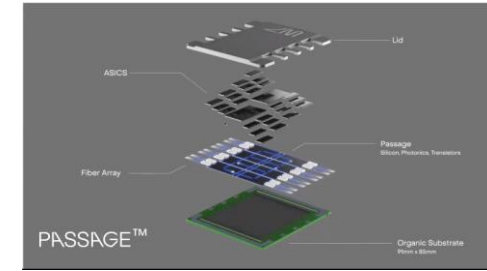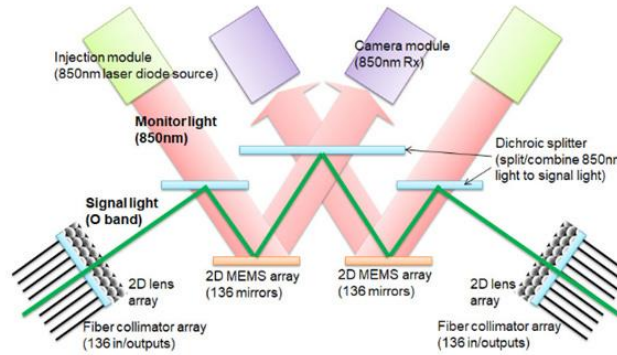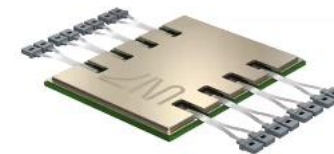
# At a Crossroads

- WL large or small? Model vs Tokens

- LLM and RecSys

- Cloud approach vs dedicated Training cluster

- HPC embraces AI or AI using HPC? Convergence of sorts?

- GPU, TPU, Wafer Scale or other???

- Optics: NPO, CPO, Direct drive, Intra ASIC?

- **Network: Lossy vs Lossless – religion or technology?**



| RDMA Application | Software |
|---|---|

RDMA API (Verbs)

| IBTA Transport Protocol |
|---|
| IBTA Network Layer | UDP |
| | IP |
| RoCE v1 | RoCE v2 |
| Ethernet Link Layer |

Typically Hardware



NIC — Bridge — Bridge — Congestion — Bridge — NIC

NIC RL — Bridge — CNM — Bridge — PFC

NIC — Bridge — NIC

From IEEE DCB tutorial

# At a Crossroads or maybe a Perfect Storm…?

# Networks of Interest: Basic Characteristics



PCIe/CXL
GPU Scale-up
GPU Scale-out
Ethernet

Network #1 – CSP or big lab - proprietary

Network #2 – Ultra Ethernet Consortium

Network #3 – Vendor specific? Network or Memory? ASIC/Node/Package/Optics Technology

## Primary DC network

- Used by all 3 deployment models
- Main network for some HPC At Scale
- Very large scale: up to 100K-1M Endpoints
- Distance: >150m ; RTT ~100 uS +; BW/GPU ~10GB/S
- Storage attached e.g., over RoCE RDMA
- Network semantics

## GPU/TPU Scale-Out Network

- DL/Inference Cluster -10k nodes and ↗
- Distance: <100m ; RTT <10 uS + ; BW ~50GB/S
- Main network for some HPC At Scale
- Network semantics

## GPU/TPU Scale-Up Network

- Within a node; small scale e.g., 256 XPU?
- Distance: ~1m ; RTT ~1 uS +; BW ~1000 GB/S
- Direct connect and/or switched
- Memory and Network semantics

# The 3 Key Deployment Models

## DL Pod

## AI/HPC Cloud

## HPC At Scale

### Characteristics

- Deployment Assumptions (network)
- Network traffic scheduling
- Network "primitives"
- Network Technology Requirements
- Traffic patterns
- Compute technologies

**DL Pod**
- Single Job
- Dedicated high BW network
- Framework, *CCL
- Moderate ⇔ low precision Algebra
- Optimized RDMA
- Collectives dominated
- Balanced System

**AI/HPC Cloud**
- Multi Tenant, multi Jobs
- Pockets of dedicated high BW network?
- Separate HPC and AI?
- eventually IaaS like?

**HPC At Scale**
- Single Job (multi jobs)
- Dedicated network
- Cross sectional BW
- Small message rate
- MPI, SHMEM etc.
- Verbs RDMA
- High precision Algebra

**One Network solution for all of the above, is feasible?**

# Common Requirements



Tail Latency

**Transport primitives for**

- Large Scale
- Multi pathing
- Relaxed ordering
- Modernized Congestion Control

- Optimized RDMA
- Performance – BW, latency, tail latency, Packets/S
- High network utilization
- Stability and Reliability

# The Network – direct workload performance influence!

## AI



- Framework coordinated – systolic
- High Bandwidth
- Large messages
- In Network Compute – 2x potential

https://youtu.be/miv5PExXTmc?t=782

## HPC



- MPI
- Small messages – Latency sensitive
- Existing application support - required

https://mvapich.cse.ohio-state.edu/static/media/talks/slide/kawthar-slingshot-osu-booth-sc22_2.pdf

intel

*Ultra Ethernet Consortium*

An Ethernet-based, open, interoperable, high performance, full-communications stack architecture to meet the growing network demands of AI & HPC at scale

**Uri Elzur**
**Technical Advisory Committee Chair, Ultra Ethernet Consortium**

# INTRODUCING: THE PROMISE OF ULTRA ETHERNET

https://*ultraethernet*.org/

**THE NEW ERA NEEDS A NEW NETWORK**

Ultra*Ethernet*

As **performant** as a supercomputing interconnect

As **ubiquitous** and **cost-effective** as Ethernet

As **scalable** as a cloud data center

Ultra*Ethernet*

# Steering Committee Members

Ultra Ethernet

# TARGET DEPLOYMENT MODELS / USE CASES

| | | AI | HPC |
|---|---|:---:|:---:|
| **Use Case** | Public cloud | ✓ | ✓ |
| | On-prem | ✓ | ✓ |

**Workloads At Scale**

*Profiles defined for AI and HPC use cases*

# APPROACH

The founding companies are seeding the consortium with highly valuable contributions in four working groups: **Physical Layer, Link Layer, Transport Layer and Software Layer**.

UEC will follow a **systematic approach with modular, compatible, interoperable layers** and tight integration of these layers to provide a holistic improvement for demanding workloads is paramount.

The consortium will work on **minimizing communication stack changes** while maintaining and **promoting Ethernet interoperability**.

Project under the Joint Development Foundation (JDF) of the Linux Foundation

# TECHNICAL GOALS

**Open** specifications, APIs, source code for optimal performance of AI and HPC workloads at scale.

**Consortium Focus**

- Electrical and optical signaling characteristics
- Link level and end-to-end transport protocols
- Link level and end-to-end congestion control
- Telemetry and signaling mechanisms
- Application programming interfaces and data structures
- Storage, management, and security constructs

Ultra Ethernet

# UEC TRANSPORT ADDRESSES GRAND CHALLENGES

- Future proof system scale with up to 1M endpoints

- Improved network utilization with multi-path routing

- Lower tail latency with flexible packet ordering

- Faster congestion control response times

- Modernized & optimized RDMA operations and APIs

- Security built-in from the beginning

- End-To-End telemetry provides improved network visibility



**Future Proof Scale**

**Network Utilization**

AI & HPC Networking

**Advanced Security**

**Congestion Control**

# FUTURE PROOF SYSTEM SCALE & NETWORK UTILIZATION

- Determinism and predictability become more difficult as systems grow

  - Network Stability, Fairness, re-convergence times, deadlock avoidance are part of the design

- "Packet spraying" - every flow to simultaneously uses all paths to the destination, vs flow using a single path
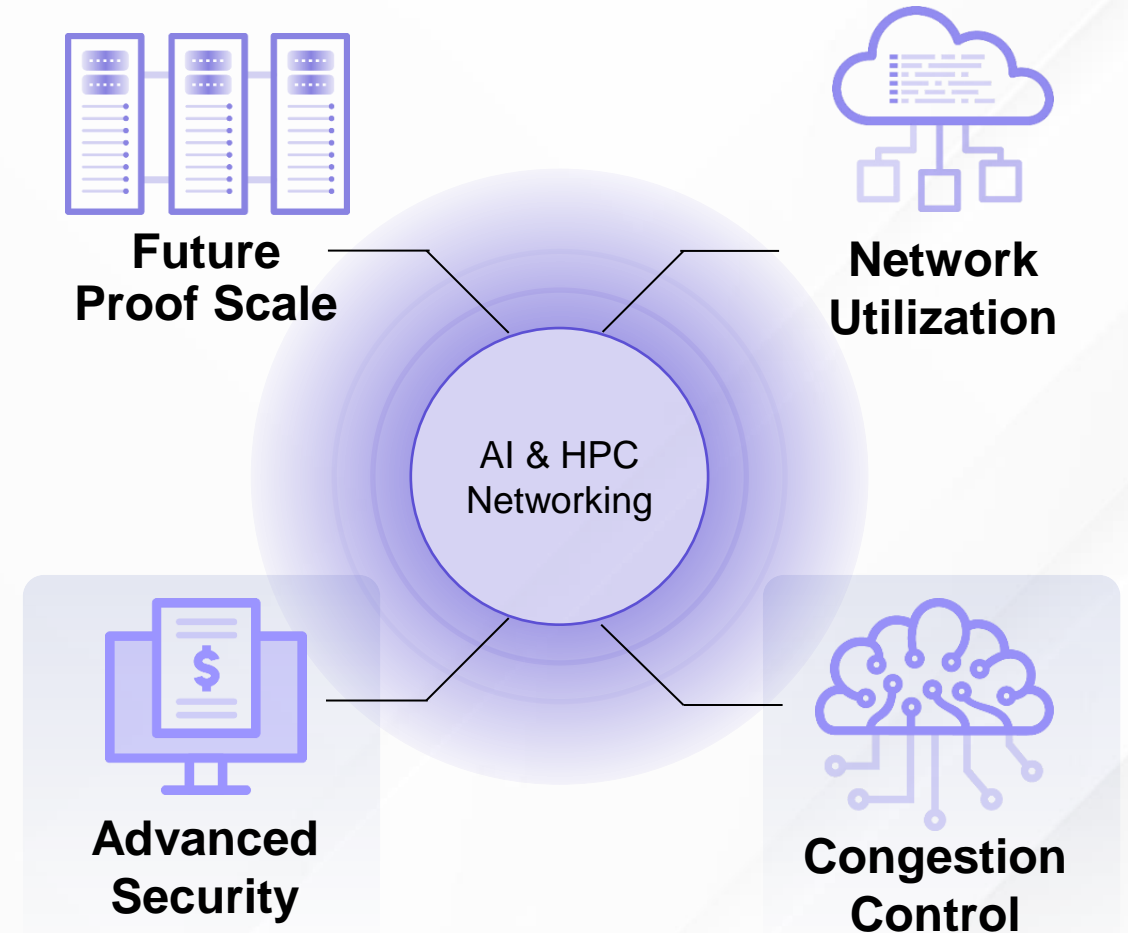
  - Achieves more balanced use of entire network

- From Rigid to Flexible Ordering

  - Rigid packet and message ordering uses "go-back-n" for loss recovery, but restricts network utilization and increases tail latencies

  - Flexible ordering enables packet-spraying in bandwidth-intensive large collective operations; without reorder buffers

  - Supports modernized RDMA operations and APIs, relaxing packet ordering while enabling maintenance of message ordering

  - Minimize state and complexities of Initiator and target

  - Critical to curtail tail latencies

**Future Proof Scale**

**Network Utilization**

AI & HPC Networking

**Advanced Security**

**Congestion Control**

*Ultra Ethernet*

# ADVANCED SECURITY, CONGESTION CONTROL & TELEMETRY

- Congestion

  - Optimized response time while maintaining high utilization

  - Support packet spraying

  - Address incast (e.g., as a result of All-to-All)

- Telemetry

  - Address wire and end-point congestion

  - Leverage shortened congestion signaling path, with more information to the endpoints to allow a more responsive congestion control

  - Information = location and cause of the congestion

- Advanced Security

  - Encryption support that doesn't balloon the session state in hosts and network interfaces

  - Similar conditions in AI and HPC

**Future Proof Scale**

**Network Utilization**

AI & HPC Networking

**Advanced Security**

**Congestion Control**

# Modern Transport and RDMA Services Needs for AI and HPC

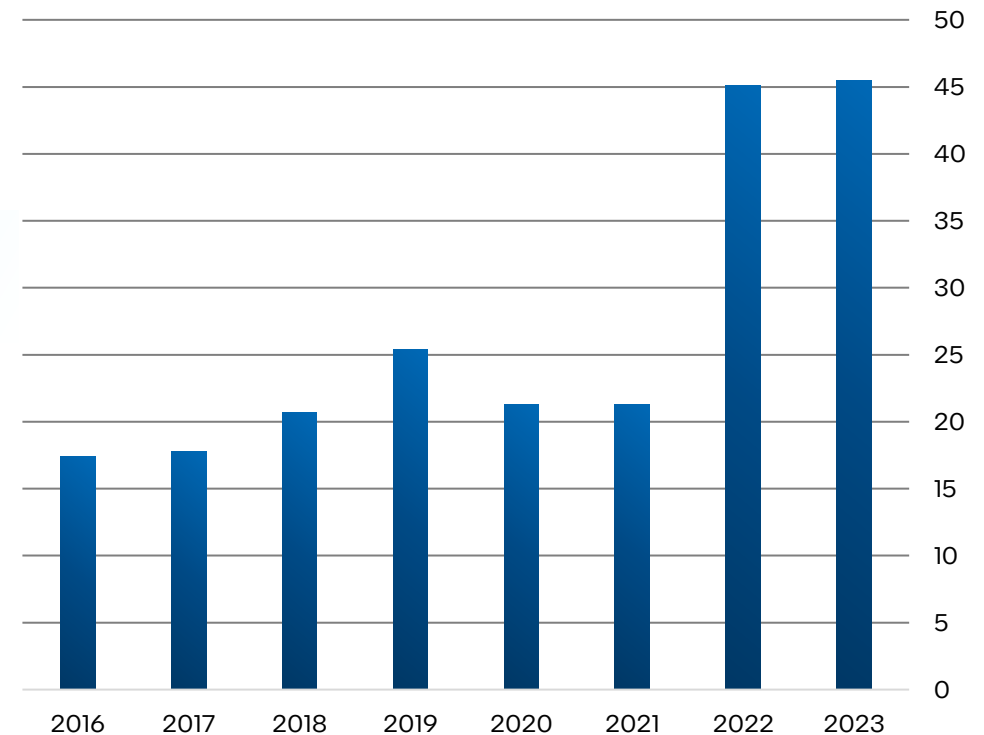| Requirement | UEC Transport | Legacy RDMA | UEC Advantage |
|---|---|---|---|
| Multi-Pathing | Packet spraying | Flow-level multi-pathing | Higher network utilization |
| Flexible Ordering | Out-of-order packet delivery with in-order message delivery | N/A | Matches application requirements, lower tail latency |
| AI and HPC Congestion Control | Workload-optimized, configuration free, lower latency, programmable | DCQCN: configuration required, brittle, signaling requires additional round trip | Incast reduction, faster response, future-proofing |
| E2E Telemetry | Sender or Receiver | ECN | Faster congestion resolution, better visibility |
| Simplified RDMA | Streamlined API, native workload interaction, minimal endpoint state | Based on IBTA Verbs | App-level performance, lower cost implementation |
| Security | Scalable, 1st class citizen | Not addressed, external to spec | High scale, modern security |
| Large Scale with Stability and Reliability | Targeting 1M endpoints | Typically, a few thousand simultaneous end points | Current and future-proof scale |

# Summary

- The Network as an island of stability amidst the storm

- Collaborate with us to move Ethernet to next level

  - Join UEC

  **www.ultraethernet.org**

- Industry benefits

  - Std high volume Ethernet based AI/HPC network products

  - AI/HPC convergence support/acceleration

**Ethernet Interconnect Family Performance Share**



Source: https://www.top500.org