

T. HOEFLER

Ultra Ethernet for Next-Generation AI and HPC Workloads

Invited Talk at Hot Interconnects, August 2025

Ultra Ethernet's Design Principles and Architectural Innovations

TORSTEN HOEFLER, ETH Zurich, Switzerland & Microsoft, USA
KAREN SCHRAMM, Broadcom, USA
ERIC SPADA, Broadcom, USA
KEITH UNDERWOOD, Hewlett Packard Enterprise, USA
CEDELL ALEXANDER, Broadcom, USA
BOB ALVERSON, Hewlett Packard Enterprise, USA
PAUL BOTTORFF, Hewlett Packard Enterprise, USA
ADRIAN CAULFIELD, OpenAI, USA
MARK HANDLEY, OpenAI, USA
CATHY HUANG, Intel, USA
COSTIN RAICIU, Broadcom, USA
ABDUL KABBANI, Microsoft, USA
EUGENE OPSASNICK, Broadcom, USA
RONG PAN, AMD, USA
ADEE RAN, Cisco, USA
RIP SOHAN, AMD, USA

<https://arxiv.org/abs/2508.08906>

The recently released Ultra Ethernet (UE) 1.0 specification defines a transformative High-Performance Ethernet standard for future Artificial Intelligence (AI) and High-Performance Computing (HPC) systems. This paper, written by the specification's authors, provides a high-level overview of UE's design, offering crucial motivations and scientific context to understand its innovations. While UE introduces advancements across the entire Ethernet stack, its standout contribution is the novel Ultra Ethernet Transport (UET), a potentially fully hardware-accelerated protocol engineered for reliable, fast, and efficient communication in extreme-scale systems. Unlike InfiniBand, the last major standardization effort in high-performance networking over two decades ago, UE leverages the expansive Ethernet ecosystem and the 1,000x gains in computational efficiency per moved bit to deliver a new era of high-performance networking.

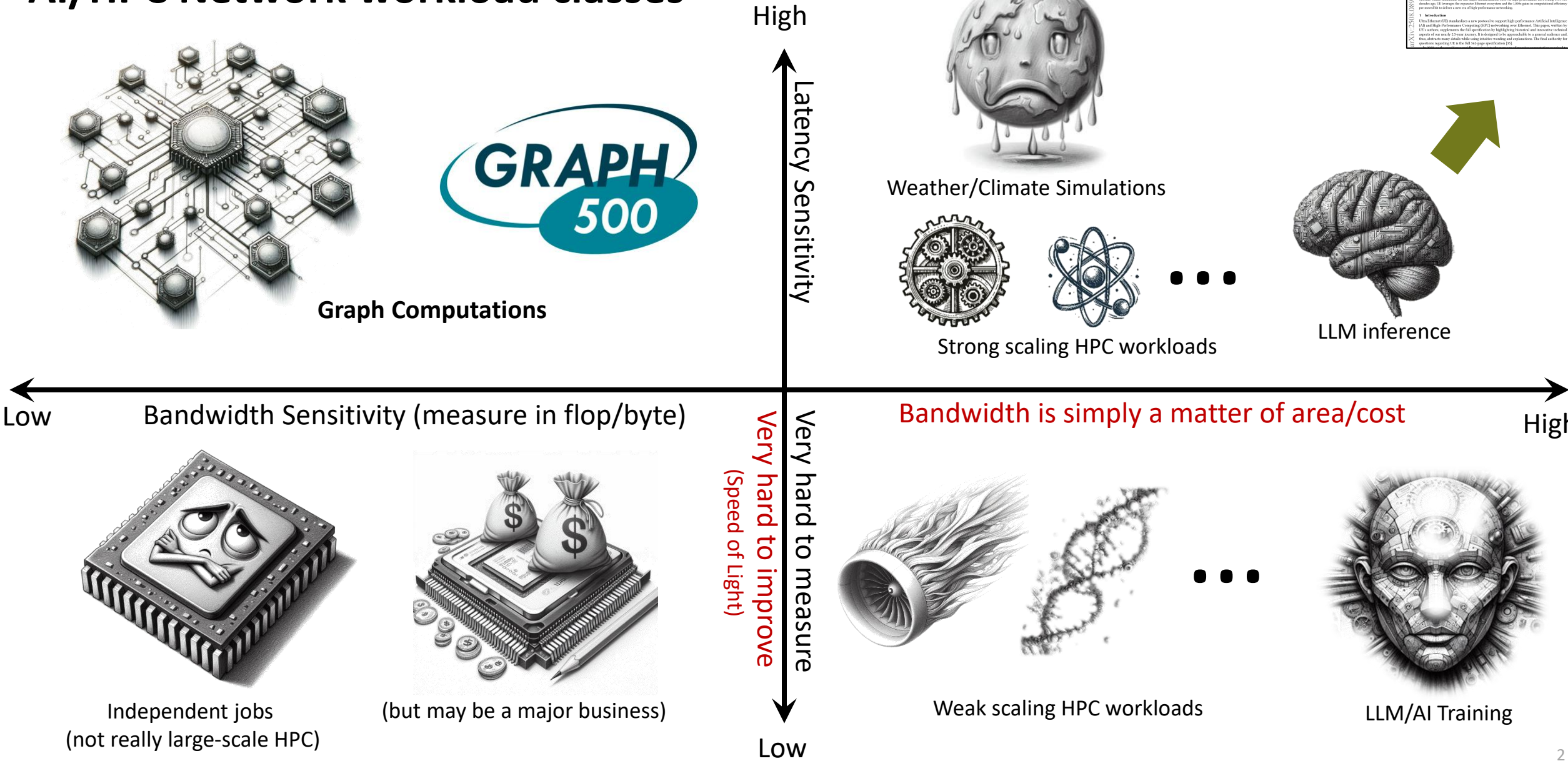
1 Introduction

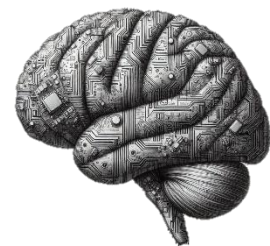
Ultra Ethernet (UE) standardizes a new protocol to support high-performance Artificial Intelligence (AI) and High-Performance Computing (HPC) networking over Ethernet. This paper, written by UE's authors, supplements the full specification by highlighting historical and innovative technical aspects of our nearly 2.5-year journey. It is designed to be approachable to a general audience and, thus, abstracts many details while using intuitive wording and explanations. The final authority for questions regarding UE is the full 562-page specification [35].

arXiv:2508.08906v1 [cs.NI] 12 Aug 2025



AI/HPC Network workload classes



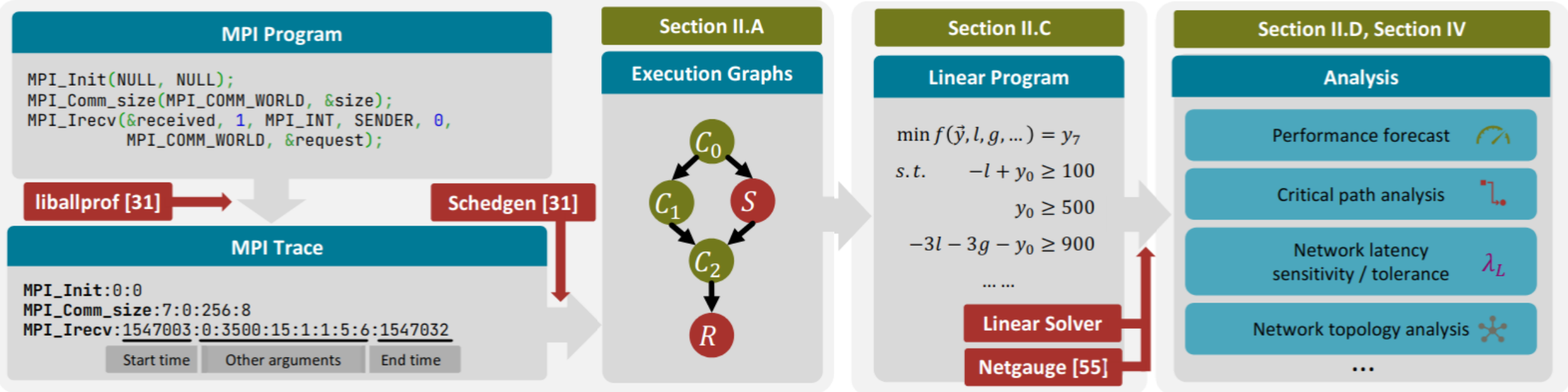


Measuring the **Latency Sensitivity** of Real Applications

How can we determine the latency sensitivity of an application?

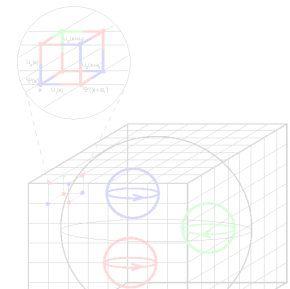
1. Large messages are less latency sensitive than small messages
2. Pipelined/overlapped small messages are less latency sensitive than small messages on the critical path

→ We need to analyze the whole execution DAG of a parallel (MPI, NCCL) application – message depth!



Latency sensitivity of applications varies widely!

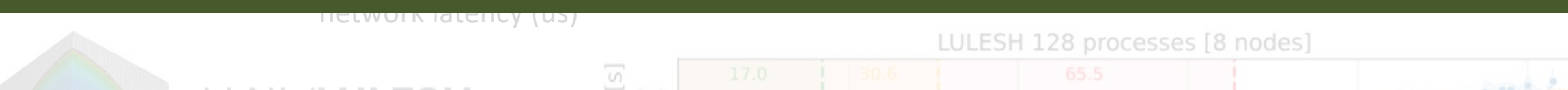
App. Slowdown
<1%, <2%, <5%



Section II.D, Section IV
Analysis
Performance forecast
Critical path analysis
Network latency sensitivity / tolerance
Network topology analysis
...

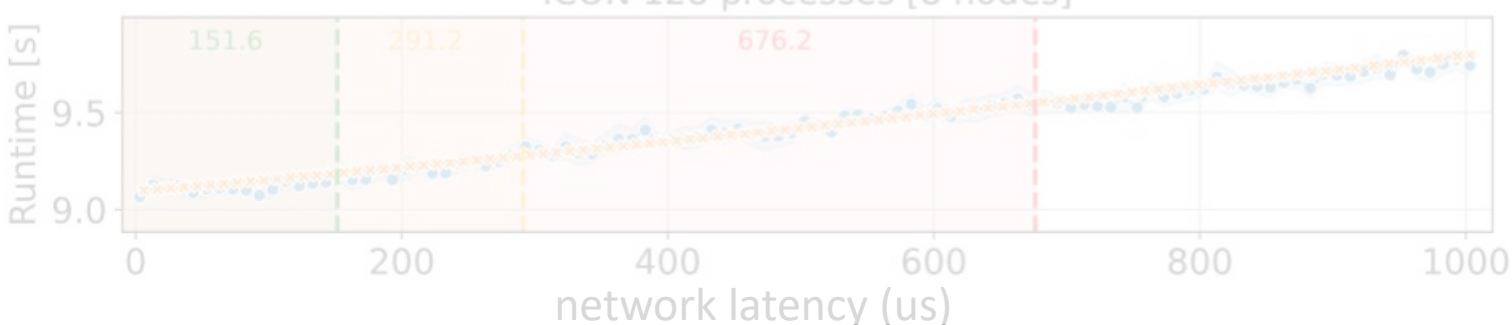
Latency Sensitivity Can be Accurately Measured Using Linear Programming

- <1% sensitivity **varies by 32x** between MILC and ICON -



Watch for ATLAHS traces and simulator tool-chain at SC25!

<https://arxiv.org/abs/2505.08936>



A new chapter was opened when ICON was released as open source code in January 2024. We are happy to be able to celebrate this and are grateful for the enormous amount of work that many colleagues have put in.

ICON Open Source Release





The Convergence of

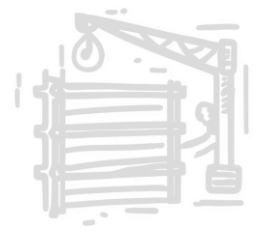
The network is the **cardiovascular system** of the datacenter
Convergence must happen around **Ethernet** !

Computing Networks

Torsten Hoefler, ETH Zurich
Ariel Hendel, Scala Computing
Duncan Roweth, Hewlett Packard Enterprise

We discuss the differences and commonalities between network technologies used in supercomputers and data centers and outline a path to convergence at multiple layers. We predict that emerging smart networking solutions will accelerate that convergence.

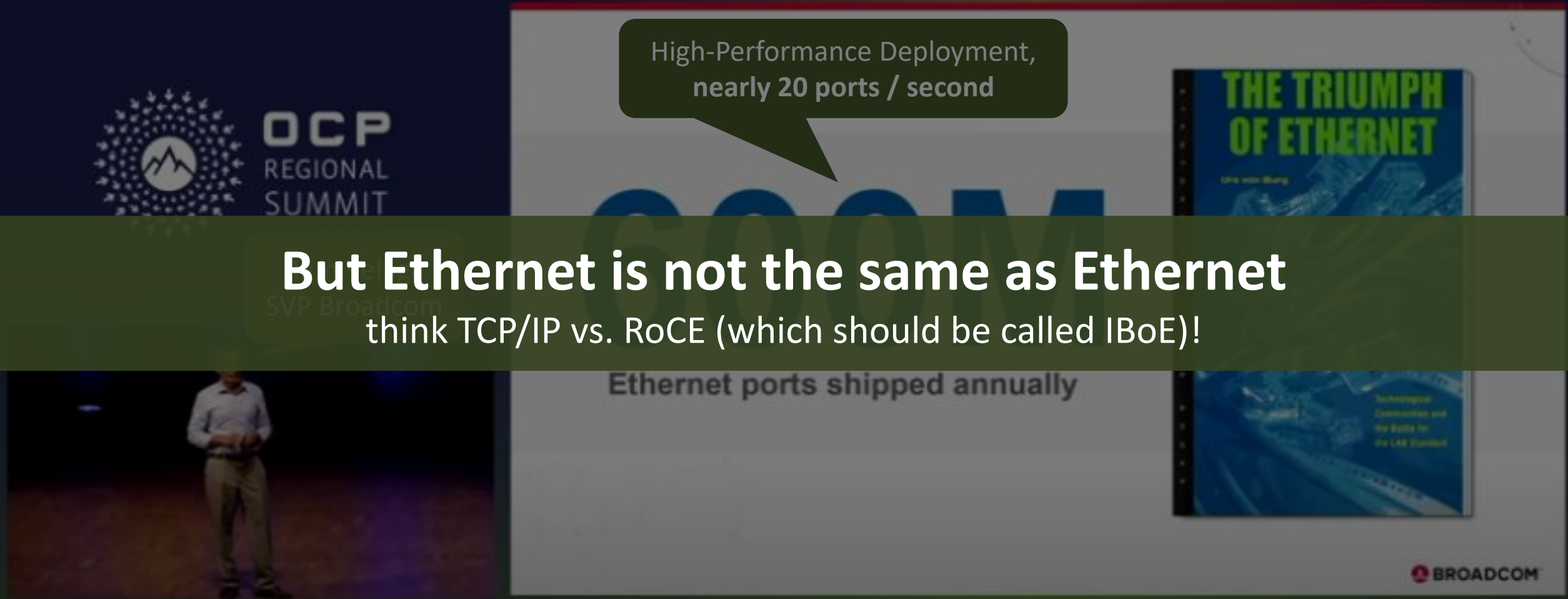
- **Design and Deployment**
 - One-off vs. incremental
 - Proprietary networks vs. Ethernet
 - ✓ AI supercomputers in the cloud
- **Operations philosophy**
 - Run-to-completion jobs vs. high-reliability services
 - Checkpoint/restart vs. replicated instances



- ✓ Most will be AI-driven to serve LLMs
- **Protocol stacks and layers**
 - Proprietary vs. task-adapted flow control
 - Simple protocols vs. multi-traffic protocols
 - Lossless vs. lossy
- **Utilization and applications**
 - High peak low noise vs. low peak high noise
 - High bandwidth low latency vs. normal bandwidth high latency
 - ✓ AI demands highest bandwidths and reasonable latency



The Ethernet Ecosystem – Is the **right one**!



High-Performance Deployment,
nearly 20 ports / second

But Ethernet is not the same as Ethernet
think TCP/IP vs. RoCE (which should be called IBoE)!

OCP REGIONAL SUMMIT

Gigamon

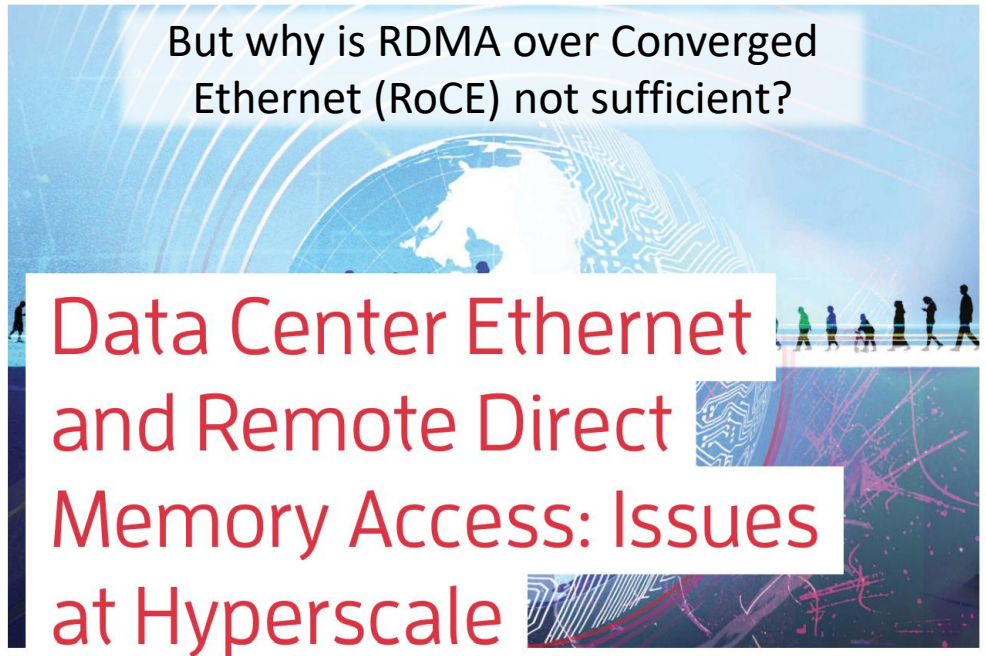
THE TRIUMPH OF ETHERNET
John W. Barry

Ethernet ports shipped annually

BROADCOM

Converging our HPC Networking Mess into a Unified Ethernet Standard

COVER FEATURE **TECHNOLOGY PREDICTIONS**



But why is RDMA over Converged Ethernet (RoCE) not sufficient?

Data Center Ethernet and Remote Direct Memory Access: Issues at Hyperscale


Torsten Hoefler¹, ETH Zürich
Duncan Roweth, Keith Underwood, and Robert Alverson, Hewlett Packard Enterprise
Mark Griswold, Vahid Tabatabaee, Mohan Kalkunte, and Surendra Anubolu, Broadcom
Siyuan Shen, ETH Zürich
Moray McLaren, Google
Abdul Kabbani and Steve Scott, Microsoft

Remote direct memory access (RDMA) over converged Ethernet (RoCE) was an attempt to adopt modern RDMA features into existing Ethernet installations. We revisit RoCE's design points and conclude that several of its shortcomings must be addressed to fulfill the demands of hyperscale data centers.



Founding Members





Ultra Ethernet™
Specification v1.0
June 11, 2025

<https://ultraethernet.org/uec-1-0-spec>

Ultra Ethernet's Design Principles and Architectural Innovations

TORSTEN HOEFLER, ETH Zurich, Switzerland & Microsoft, USA
KAREN SCHRAMM, Broadcom, USA
ERIC SPADA, Broadcom, USA
KEITH UNDERWOOD, Hewlett Packard Enterprise, USA
CEDELL ALEXANDER, Broadcom, USA
BOB ALVERSON, Hewlett Packard Enterprise, USA
PAUL BOTTORFF, Hewlett Packard Enterprise, USA
ADRIAN CAULFIELD, OpenAI, USA
MARK HANDLEY, OpenAI, USA
CATHY HUANG, Intel, USA
COSTIN RAICIU, Broadcom, USA
ABDUL KABBANI, Microsoft, USA
EUGENE OPSASNICK, Broadcom, USA
RONG PAN, AMD, USA
ADEE RAN, Cisco, USA
RIP SOHAN, AMD, USA

The recently released Ultra Ethernet (UE) 1.0 specification defines a transformative High-Performance Ethernet standard for future Artificial Intelligence (AI) and High-Performance Computing (HPC) systems. This paper, written by the specification's authors, provides a high-level overview of UE's design, offering crucial motivations and scientific context to understand its innovations. While UE introduces advancements across the entire Ethernet stack, its standout contribution is the novel Ultra Ethernet Transport (UET), a potentially fully hardware-accelerated transport protocol that leverages the expansive Ethernet ecosystem and the 1,000x gains in computational efficiency per moved bit to deliver a new era of high-performance networking.

08906v1 [cs.NI] 12 Aug 2025

<https://arxiv.org/abs/2508.08906>

Ecosystem is quickly growing

Today 10 steering companies, 17 general member companies, 51 contributor members



Chair's view of the Transport WG Meeting in March'24 (60+ members on site, 1,300+ total)

Ultra Ethernet Members – Join our Journey!

Ultra Ethernet Consortium
<https://arxiv.org/abs/2508.08906>













*not all members listed

75+ member companies
1,300+ individual participants

Modernizing RDMA for HPC and AI

Classic RDMA

Lossless (PFC or CBFC) operation

In-order transport and delivery

Inefficient go-back-n

Proprietary congestion control (e.g., DCQCN)

Single-path routing

No load balancing and “link polarization”

Large state per queue pair

kb NIC memory per peer

Security added at higher layers

IPSec, N^2 contexts, known attacks

UltraEthernet
Consortium

Lossy (& lossless) operation

Out-of-order data and message delivery

(Un)Reliable (Un)Ordered - ROD, RUD/RUDI, and UUD

Open, configurable, and flexible CC

Per-packet multipathing and load balancing

Including (close-to) zero state REPS

Connection-less API

Ephemeral zero-RTT reliability state

Built-in security

Cluster-wide keying, zero state replay protection



sRDMA – Efficient NIC-based Authentication and Encryption for Remote Direct Memory Access

Konstantin Taranov, Benjamin Rothenberger, Adrian Perrig, and
Torsten Hoefler, *ETH Zurich*

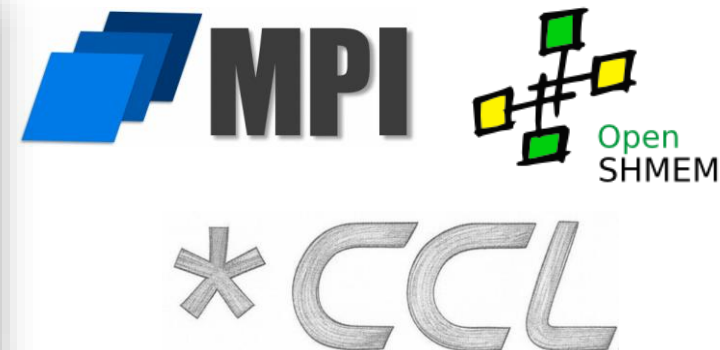
<https://www.usenix.org/conference/atc20/presentation/taranov>

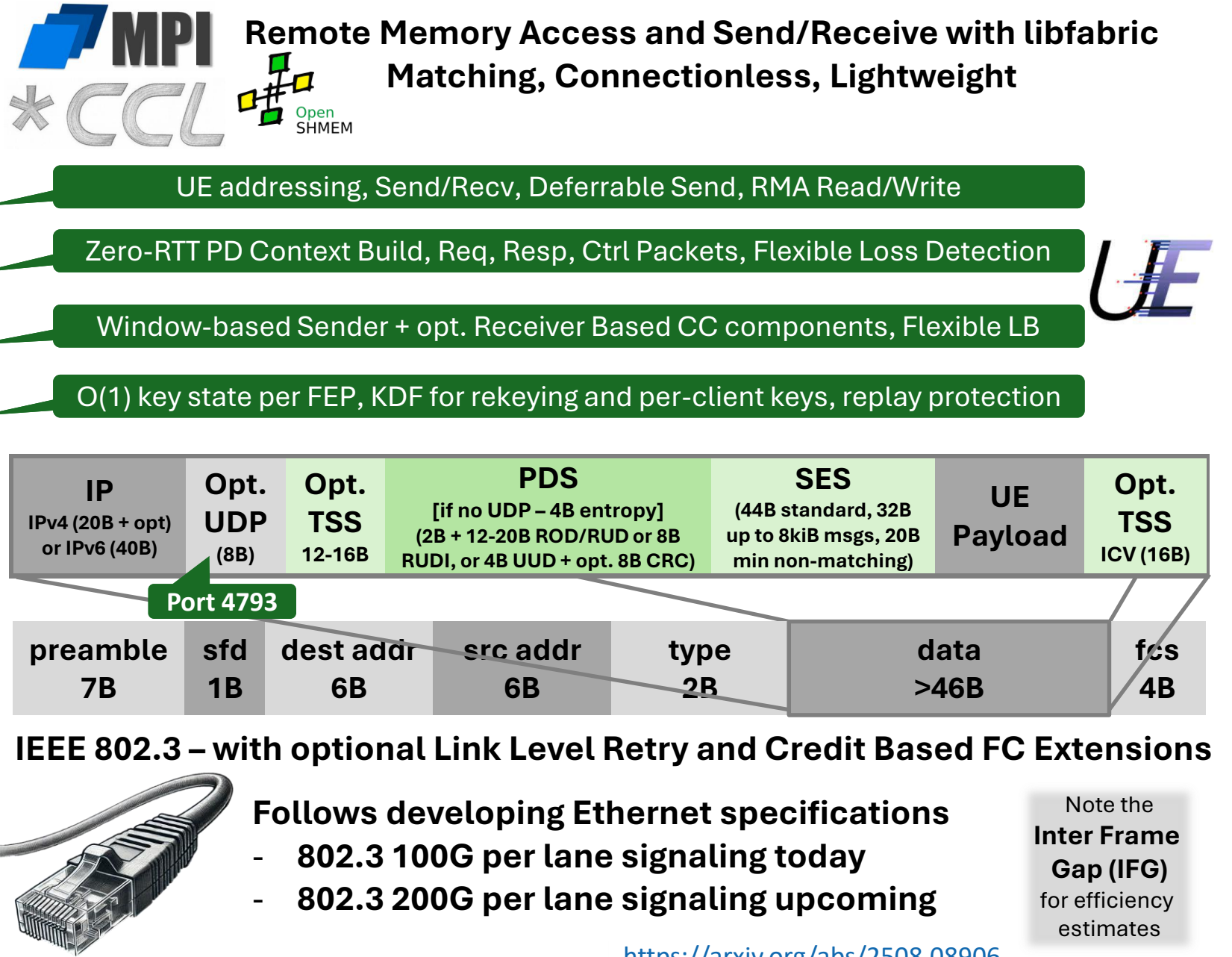
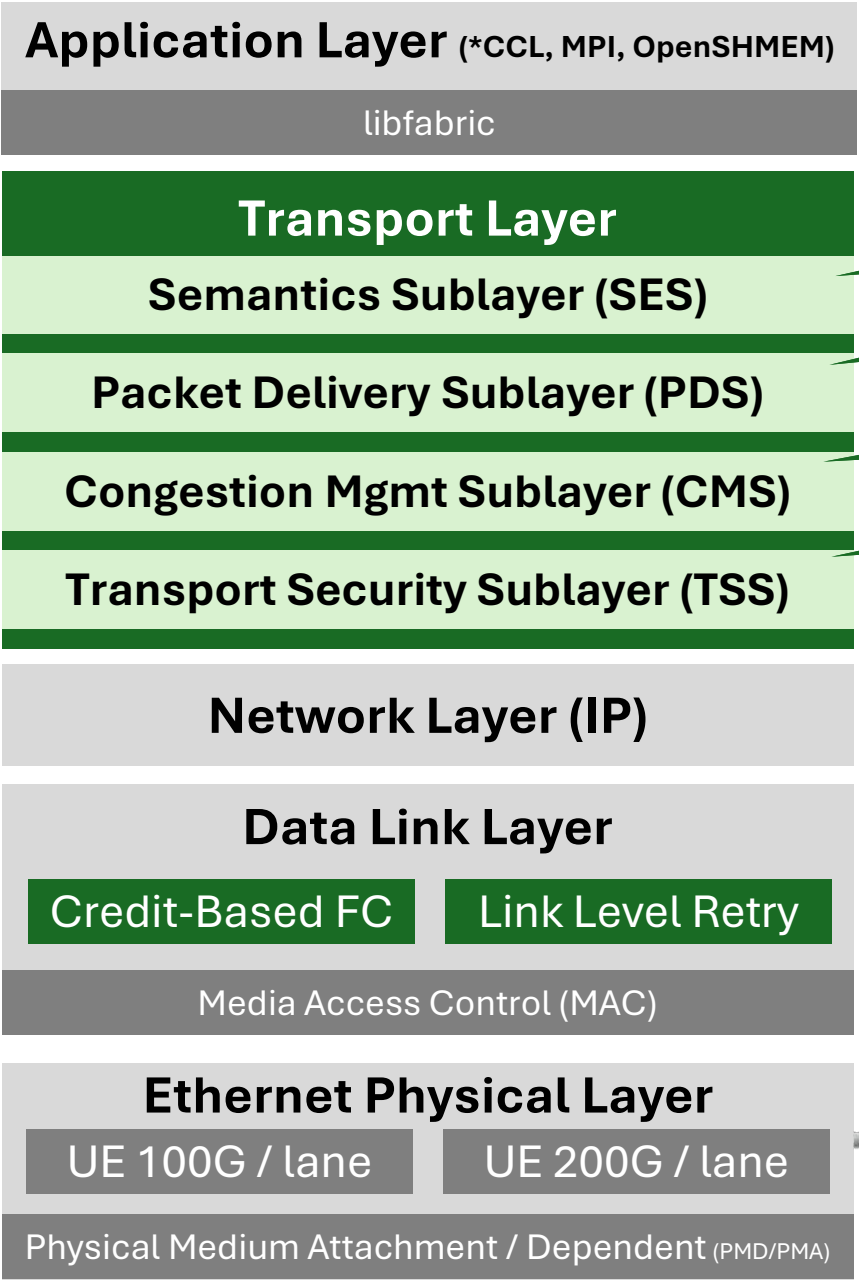


ReDMark: Bypassing RDMA Security Mechanisms

Benjamin Rothenberger, Konstantin Taranov, Adrian Perrig, and
Torsten Hoefler, *ETH Zurich*

<https://www.usenix.org/conference/usenixsecurity21/presentation/rothenberger>





Application Layer (*CCL, MPI, OpenSHMEM)

libfabric



Remote Memory Access and Send/Receive with libfabric
Matching, Connectionless, Lightweight

Transport Layer

Semantics Sublayer (SES)

UE addressing, Send/Recv, Deferrable Send, RMA Read/Write

Packet Delivery Sublayer (PDS)

Zero-RTT PD Context Build, Req, Resp, Ctrl Packets, Flexible Loss Detection

Congestion Mgmt Sublayer (CMS)

Window-based Sender + opt. Receiver Based CC components, Flexible LB

Transport Security Sublayer (TSS)

O(1) key state per FEP, KDF for rekeying and per-client keys, replay protection



Network Layer (IP)

IP	Opt. UDP	Opt. TSS	PDS	SES	UE Payload	Opt. TSS
IPv4 (20B + opt) or IPv6 (40B)	(8B)	12-16B	[if no UDP – 4B entropy] (2B + 12-20B ROD/RUD or 8B RUDI, or 4B UUD + opt. 8B CRC)	(44B standard, 32B up to 8kiB msgs, 20B min non-matching)		ICV (16B)

Port 4793

Data Link Layer

Credit-Based FC

Link Level Retry

preamble	sfd	dest addr	src addr	type	data	fcs
7B	1B	6B	6B	2B	>46B	4B

Media Access Control (MAC)

IEEE 802.3 – with optional Link Level Retry and Credit Based FC Extensions

Ethernet Physical Layer

UE 100G / lane

UE 200G / lane

Physical Medium Attachment / Dependent (PMD/PMA)



Follows developing Ethernet specifications

- 802.3 100G per lane signaling today
- 802.3 200G per lane signaling upcoming

Note the
Inter Frame
Gap (IFG)
for efficiency
estimates

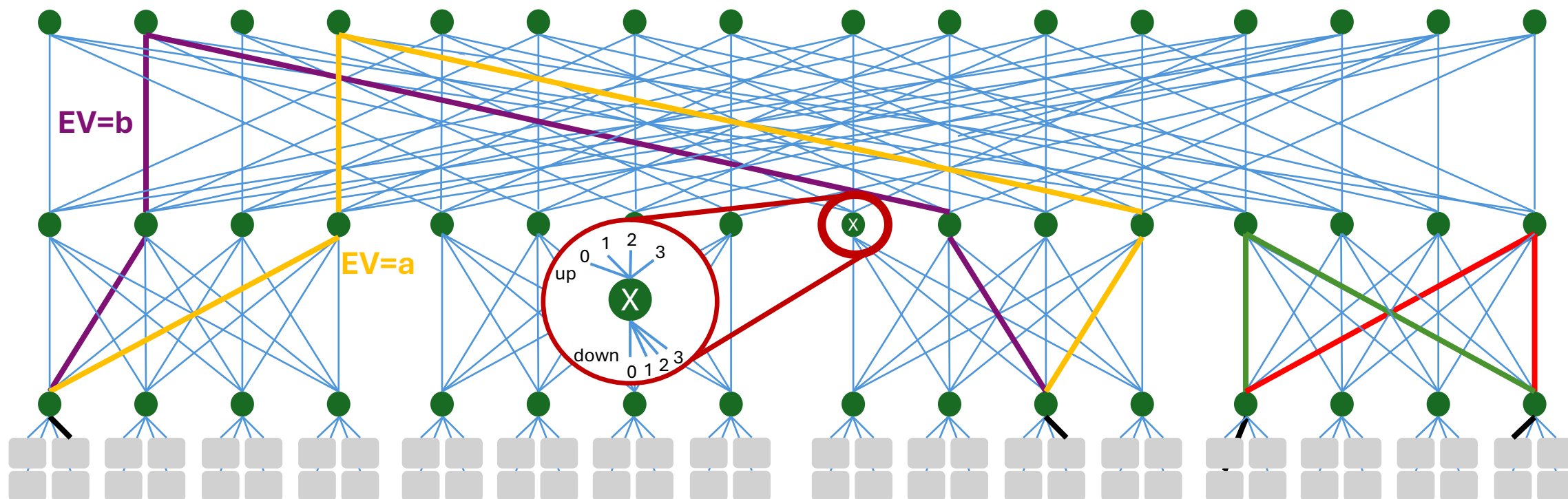
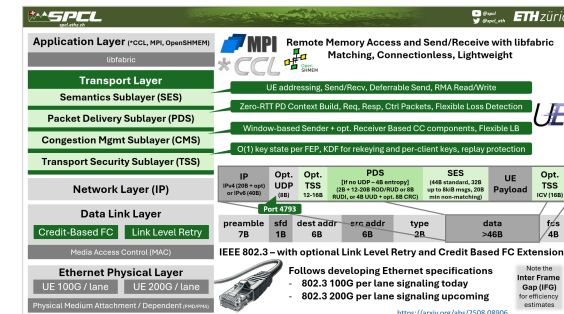
<https://arxiv.org/abs/2508.08906>

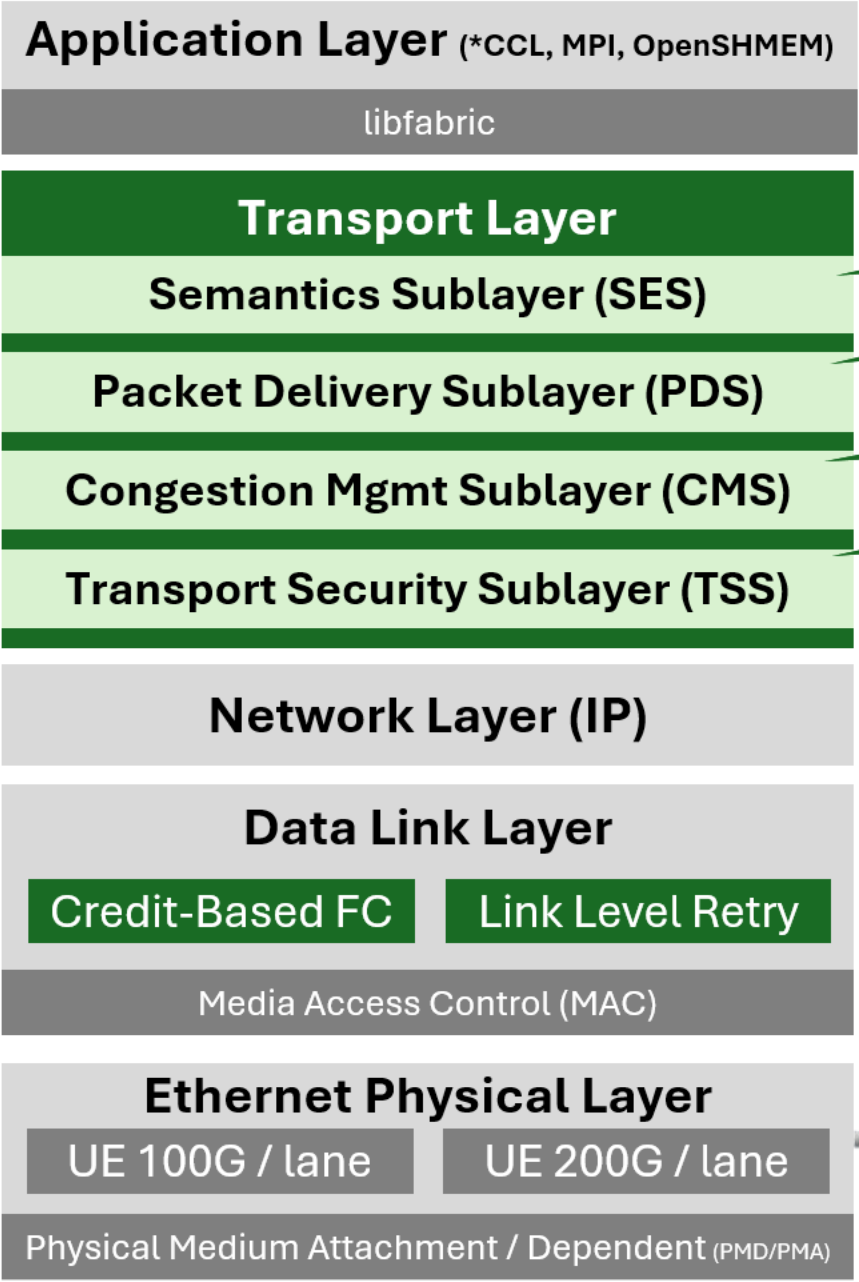
Ultra Ethernet Transport - Load Balancing Philosophy



Based on standard Ethernet Equal Cost Multi Pathing (ECMP)

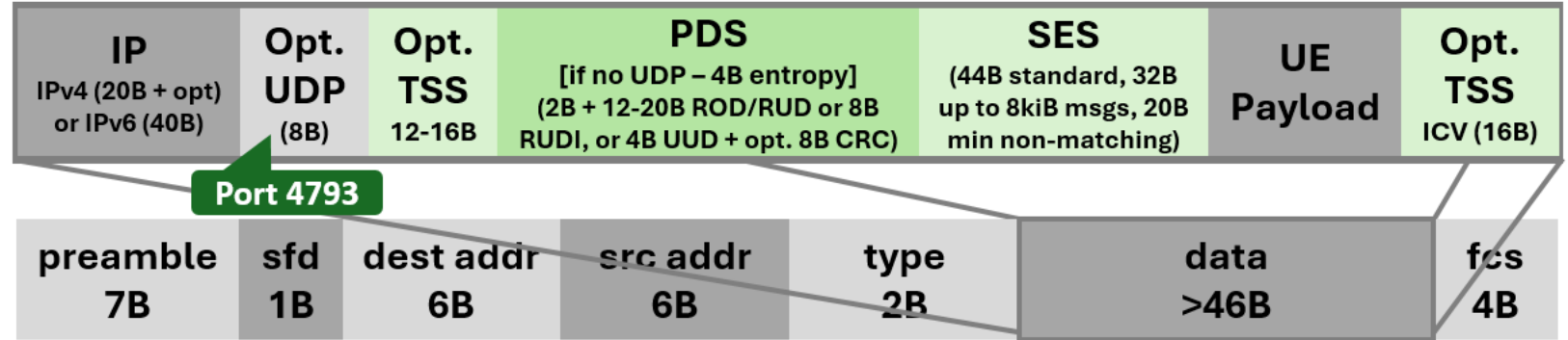
- Uses an “entropy value” (EV) to select from a set of output ports (encoded as UDP source port)
- Each EV selects a path (not necessarily unique)
- Same EV means same path (without failures)





Remote Memory Access and Send/Receive with libfabric
Matching, Connectionless, Lightweight

- UE addressing, Send/Recv, Deferrable Send, RMA Read/Write
- Zero-RTT PD Context Build, Req, Resp, Ctrl Packets, Flexible Loss Detection
- Window-based Sender + opt. Receiver Based CC components, Flexible LB
- O(1) key state per FEP, KDF for rekeying and per-client keys, replay protection



IEEE 802.3 – with optional Link Level Retry and Credit Based FC Extensions



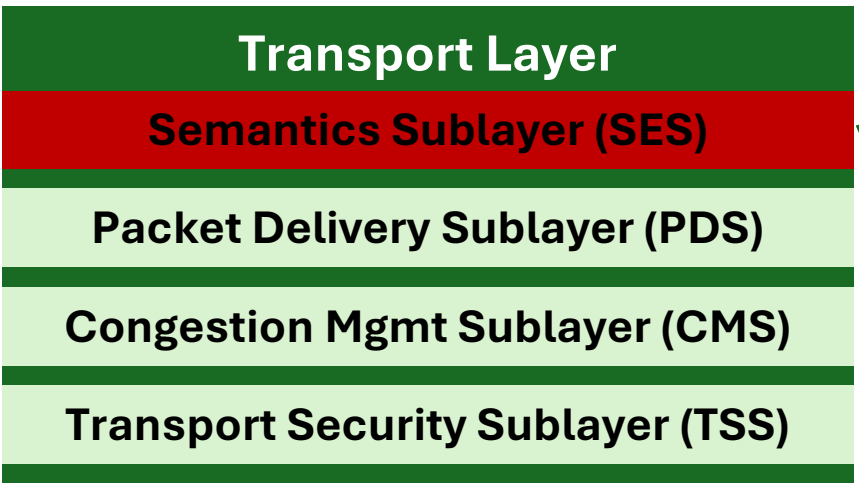
Follows developing Ethernet specifications

- 802.3 100G per lane signaling today
- 802.3 200G per lane signaling upcoming

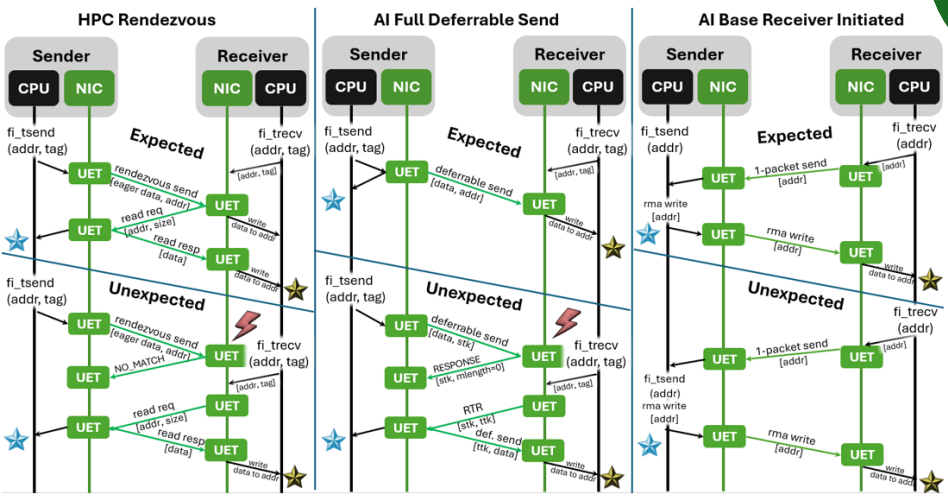
Note the **Inter Frame Gap (IFG)** for efficiency estimates

<https://arxiv.org/abs/2508.08906>

Transport layer - sublayers



- Compatible with existing applications (libfabric) – **no change!**
- RDMA services: Send/Recv + RMA (Write, Read, Atomics)
 - Focus on MPI and *CCL semantics
- Scalable addressing to millions of endpoints
- Optimized extensions:
 - Deferrable Send for optimized HW (aimed at AI)
 - Rendezvous using Send/Read (aimed at HPC)
 - Exact match tags for HW offload of ordering between endpoints using shared receive queues

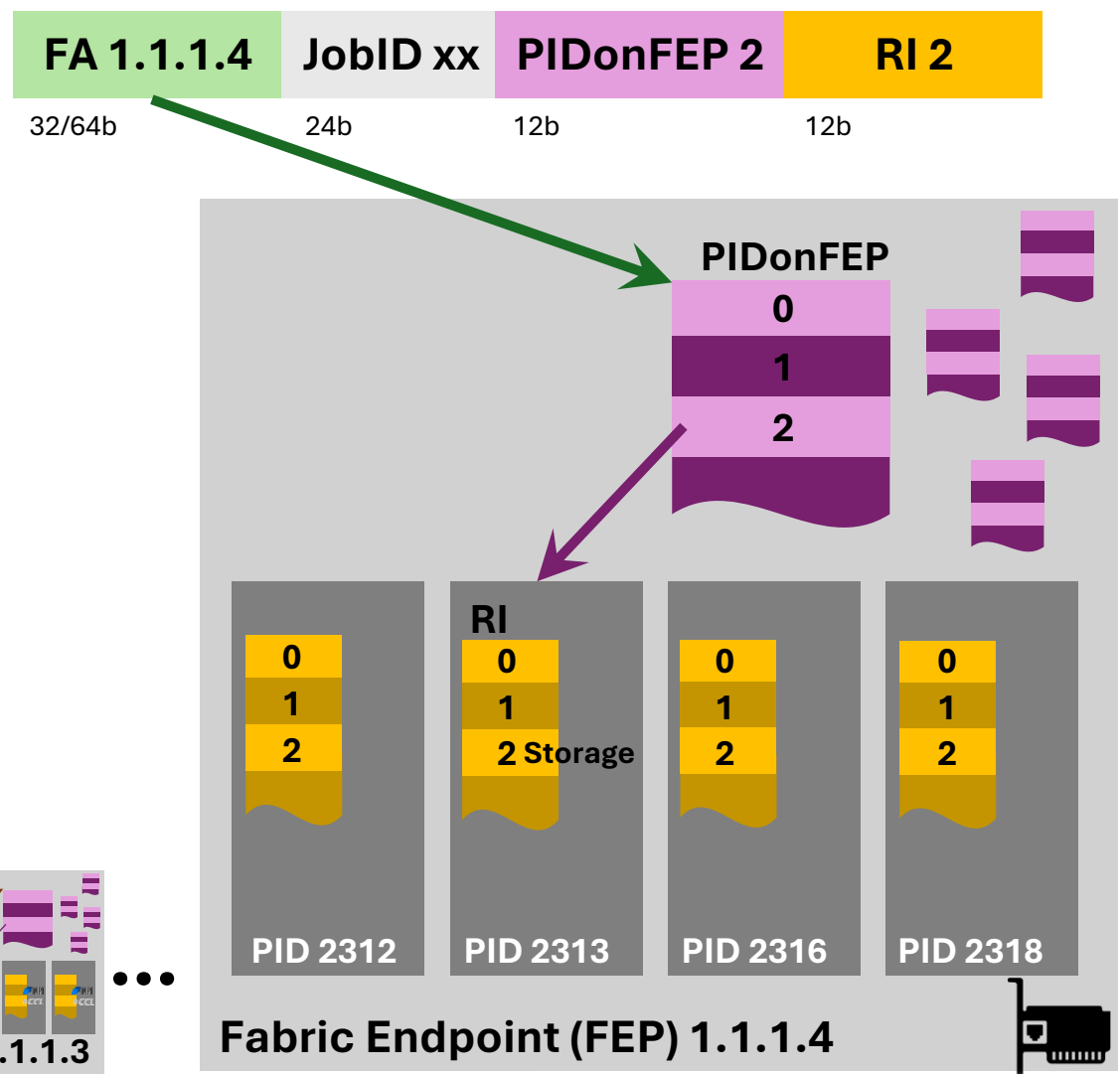
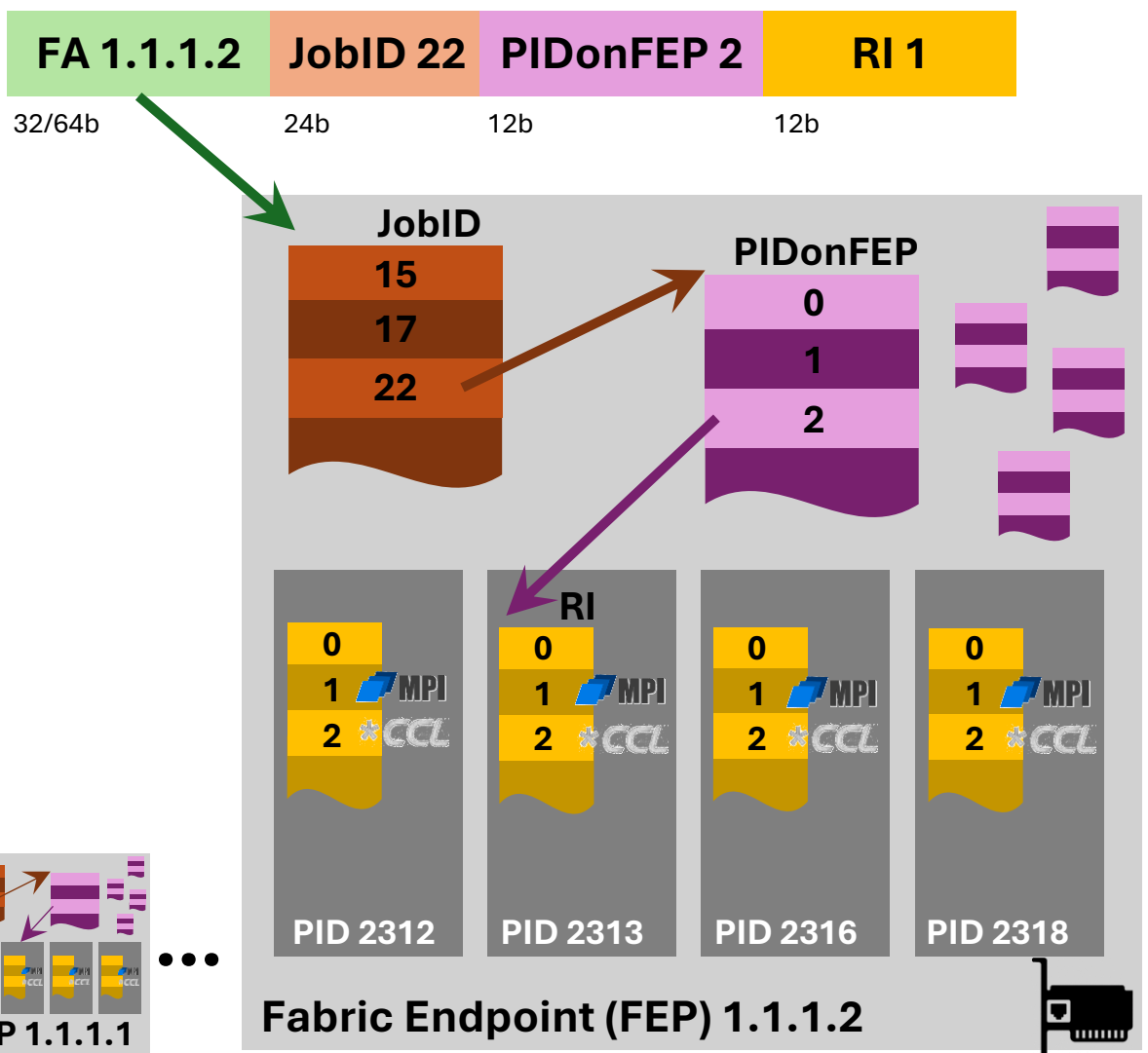


Use-case optimized communication profiles (AI Base, AI Full, HPC)

Ultra Ethernet Transport Scalable Addressing

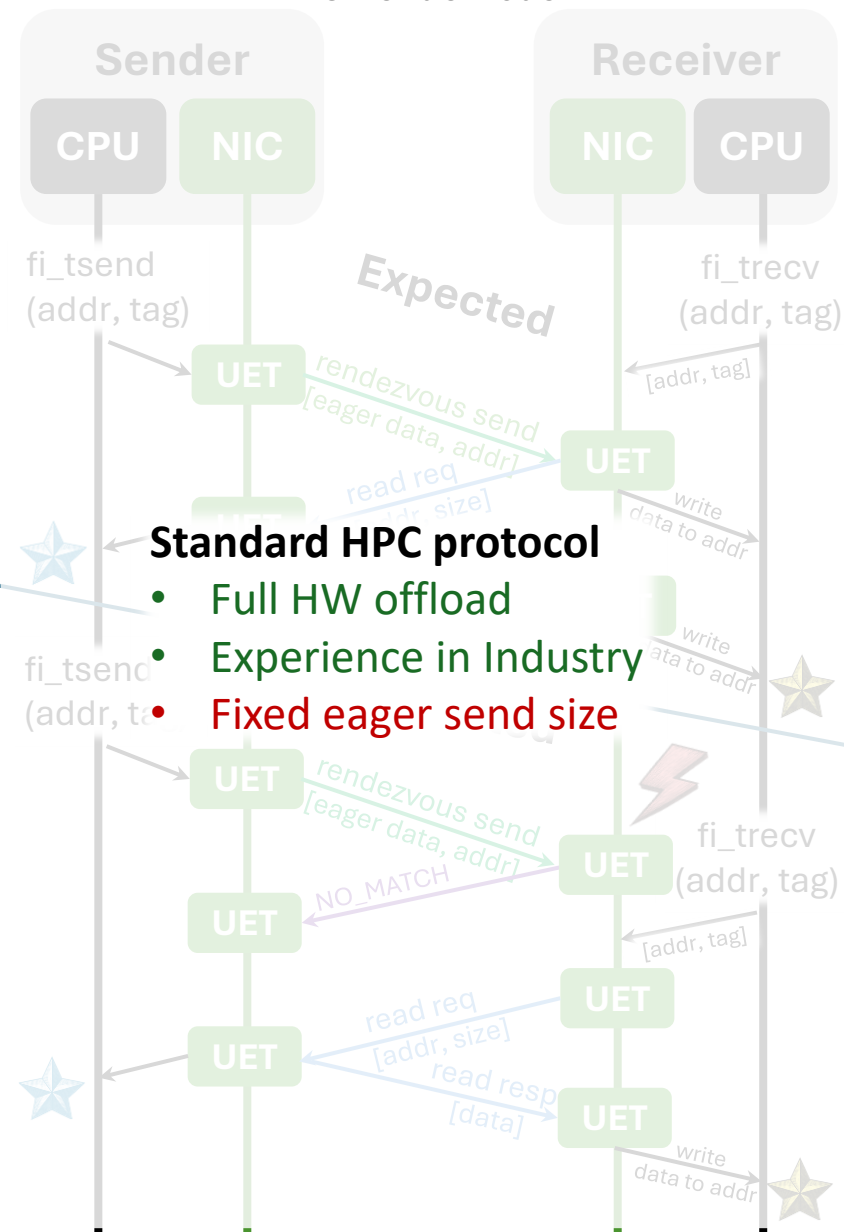
Relative (parallel job) Addressing

Absolute (client/server) Addressing



Ultra Ethernet Transport Large Message Delivery Options and Profiles

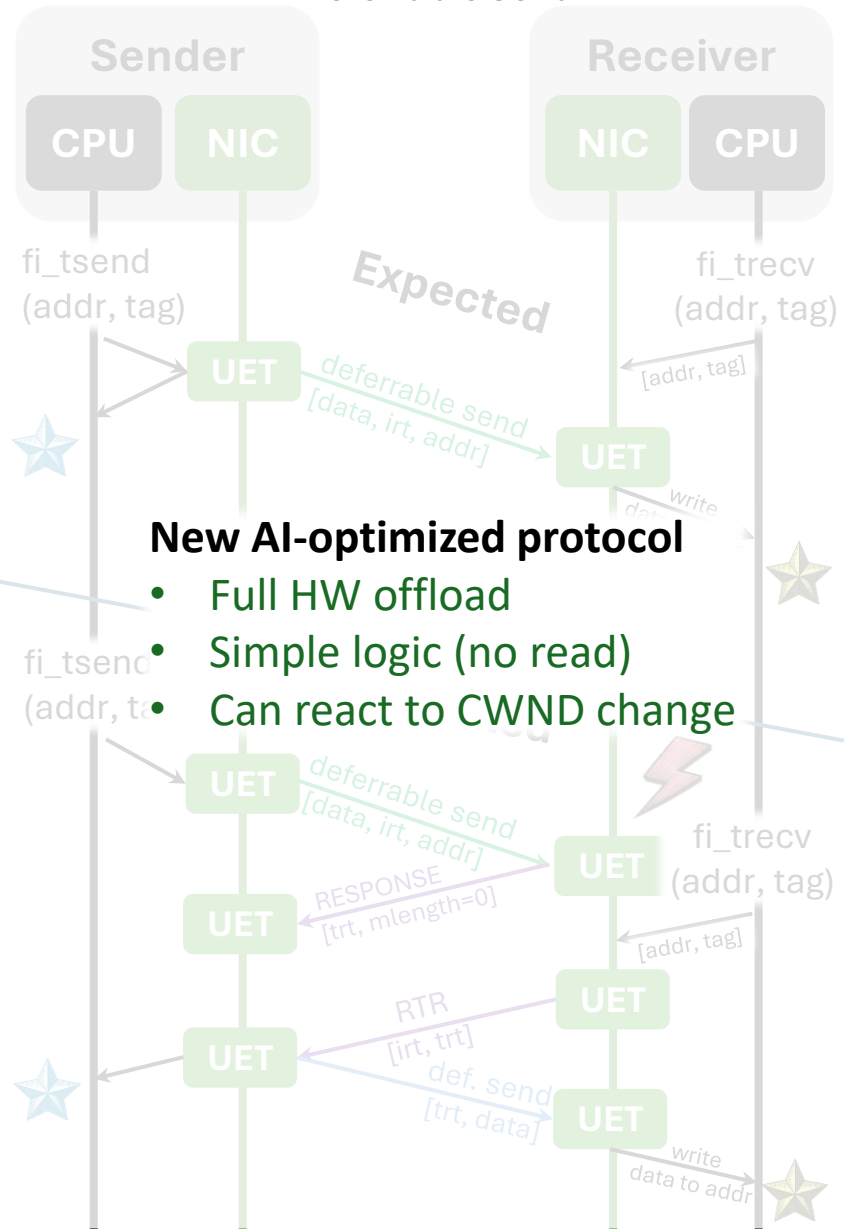
HPC Rendezvous



Standard HPC protocol

- Full HW offload
- Experience in Industry
- **Fixed eager send size**

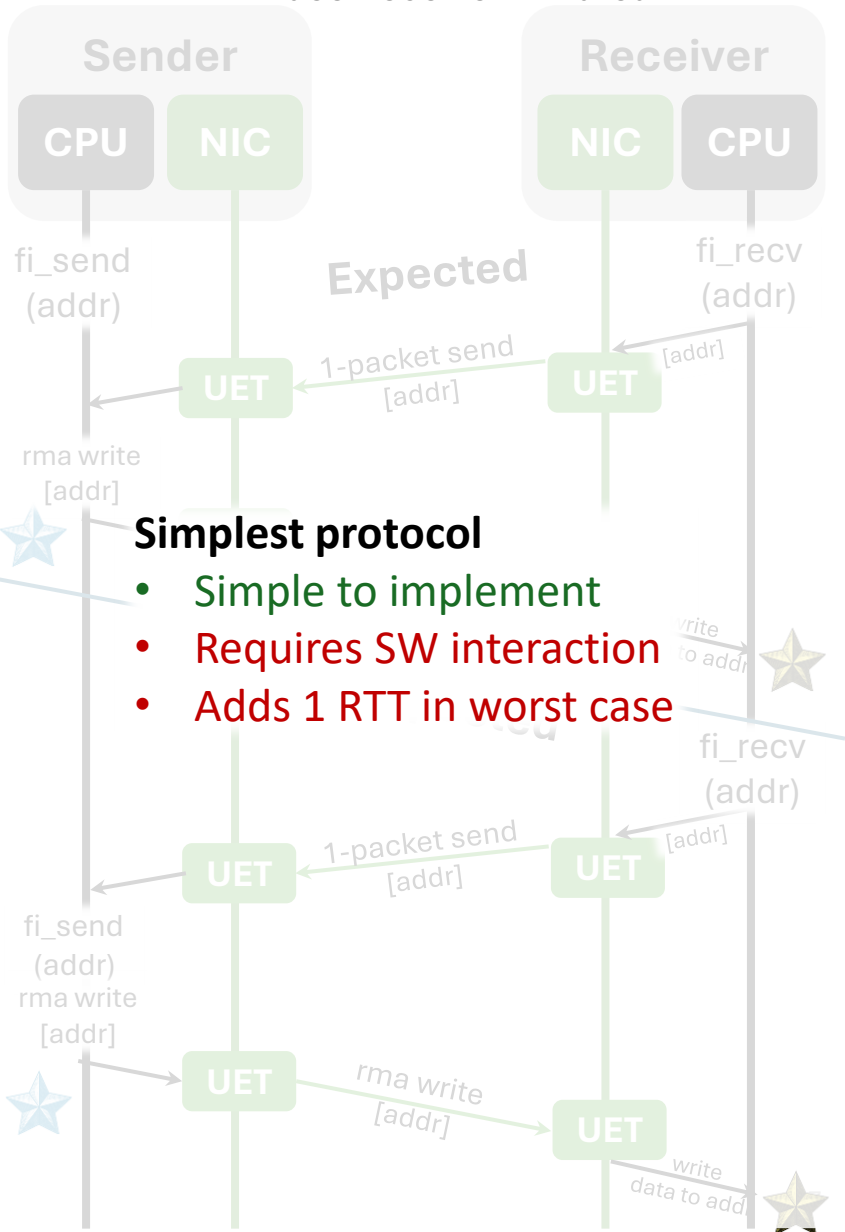
AI Deferrable Send



New AI-optimized protocol

- Full HW offload
- Simple logic (no read)
- Can react to CWND change

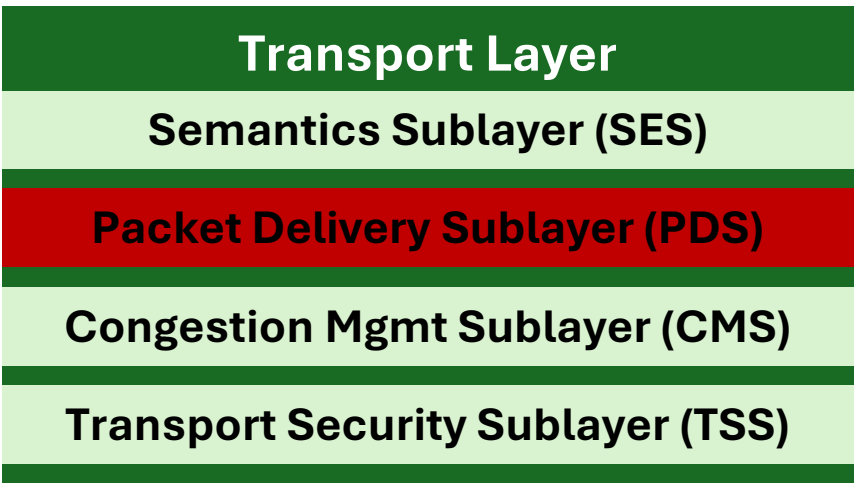
AI Base Receiver Initiated



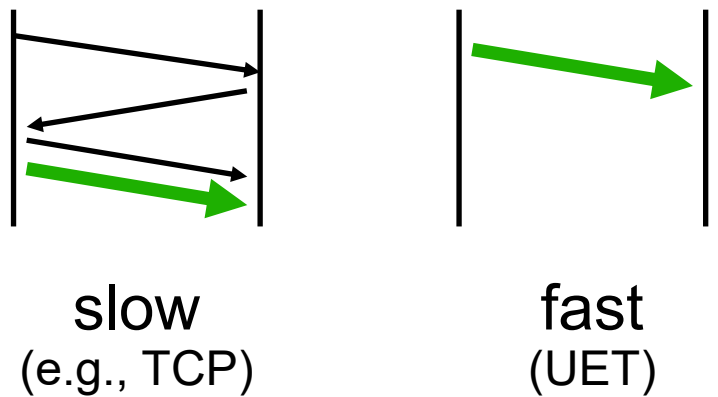
Simplest protocol

- Simple to implement
- **Requires SW interaction**
- **Adds 1 RTT in worst case**

Transport layer - sublayers

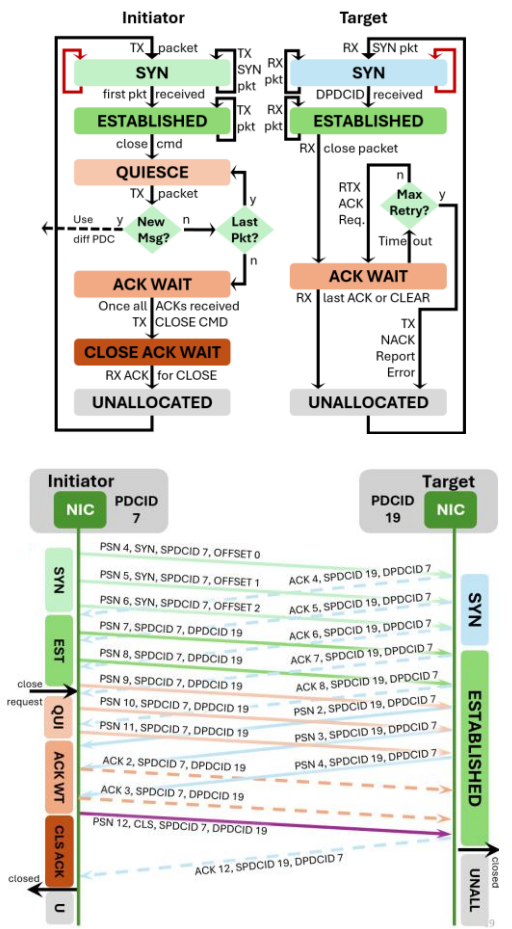


Zero-RTT Startup



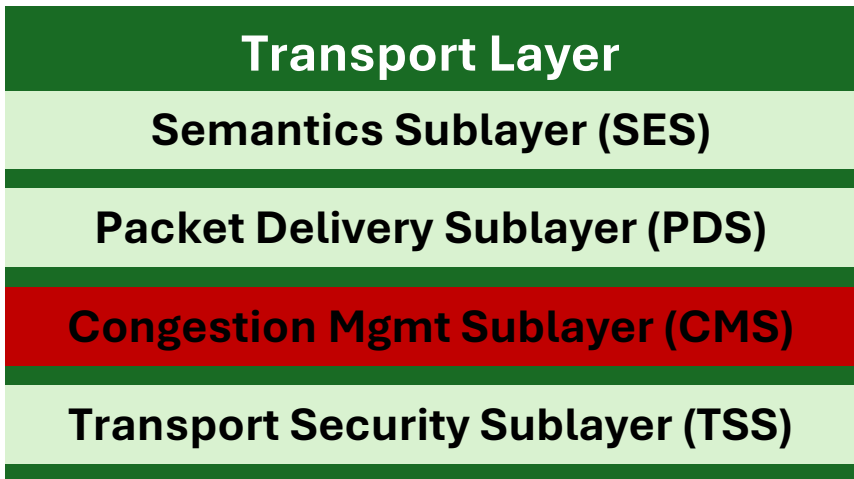
- Dynamic, ephemeral connections
 - Zero start up time, 1-RTT close
- 4 delivery services
 - ROD – Reliable, ordered
 - RUD – Reliable, unordered
 - RUDI – Reliable, unordered, idempotent (Write/Read)
 - UUD – Unreliable, unordered
- Shared receive queues
- Out-of-order packet arrival
- Selective acknowledgement and retransmission for RUD
 - ROD uses Go-Back-N

Ultra Ethernet Consortium



Fastest startup, drop state when convenient, rebuild it quickly!

Transport layer - sublayers



- Multipath with congestion avoidance
 - Leveraging ECMP
- Trimming with NACK signal
- Network Signaled CC (NSCC)
 - Window based at sender using RTT and ECN
- Receiver Controlled CC (RCCC)
 - Credit based at receiver

Network Signal Based CC (Sender-controlled)

- Available in all UE products
- Can be disabled
- Flexible for most deployments

Receiver Controlled CC

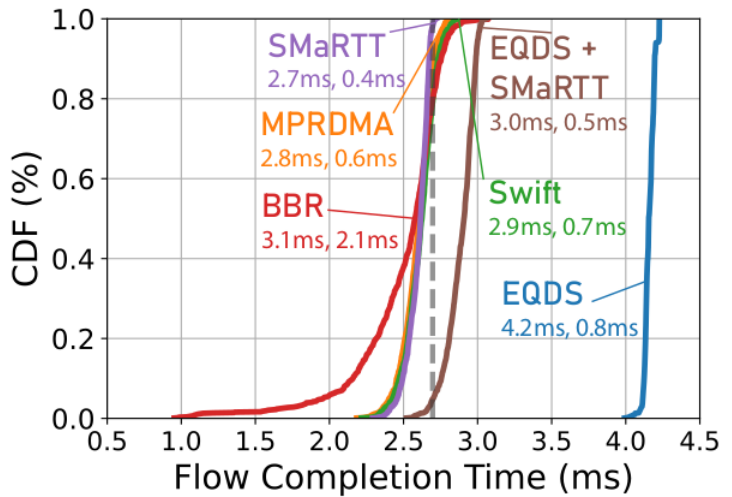
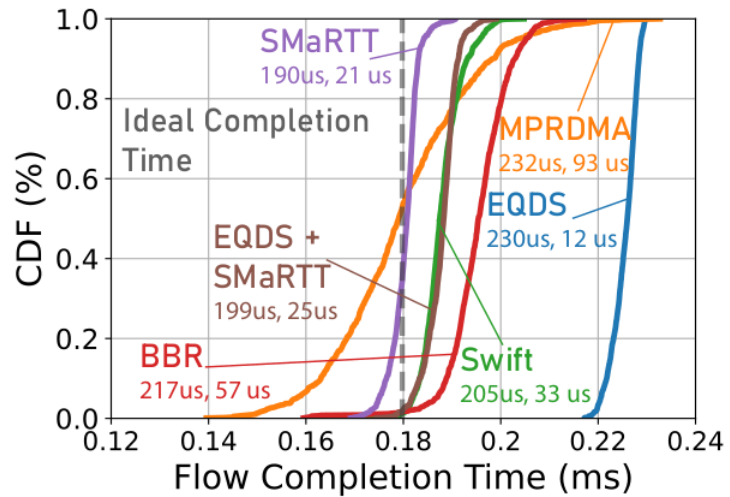
- Available in some UE products
- Receiver hands out credits
- Ideal for incast patterns

Work together for HPC+AI multi-pathing

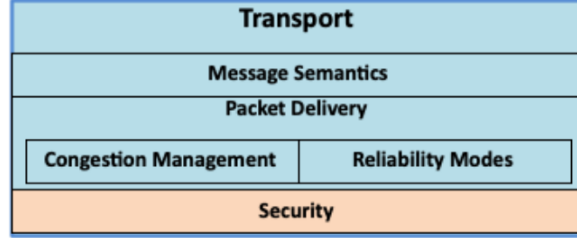
SMaRTT-REPS enables Modern Packet Spraying



- “State of the art” (2024), easily configured congestion control mechanisms

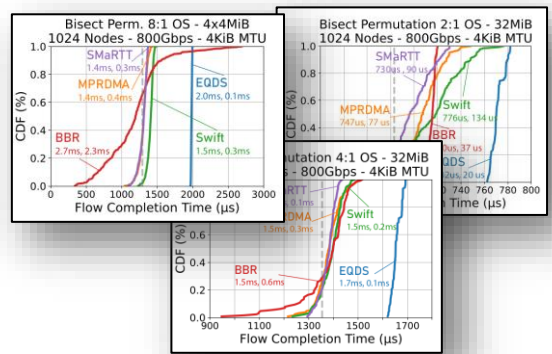
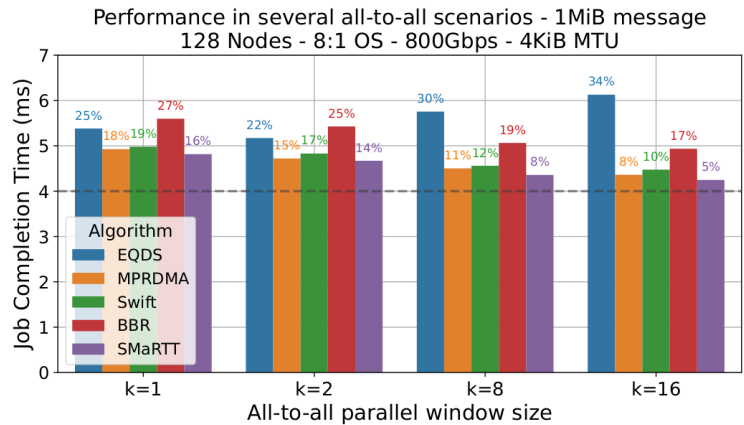


```
Algorithm 1 SMaRTT Pseudocode
1: acked, bytes_ignored = 0
2:
3: procedure congestion_loop_loop(p)
4:   acked += p.size
5:   bytes_ignored += p.size
6:
7:   if p.is_ack then
8:     if bytes_ignored < bytes_to_ignore then
9:       return
10:    else
11:      return
12:   end if
13:
14:   can_decrease = wait_to_decrease(p)
15:   adp = quick_adapt(p)
16:   fast = fast_increase(p)
17:   if adp or fast then
18:     return
19:   end if
20:
21:   if p.con and p.rt <= rtt and can_decrease then
22:     fair_decrease(p)
23:   else if p.con and p.rt > rtt and can_decrease then
24:     multiplicative_decrease(p)
25:   else if p.con and p.rt > rtt then
26:     fair_increase(p)
27:   else if p.con and p.rt <= rtt then
28:     multiplicative_increase(p)
29:   end if
30:
31:   if p.is_ack and p.timeout_triggered then
32:     cond = p.size
33:     trigger_ga = true
34:     retransmit_packet(p)
35:   end if
36:   if bytes_ignored >= bytes_to_ignore then
37:     quick_adapt(p)
38:   end if
39:   cond = max(cond, bdp, mru)
40: end procedure
```



37 lines simple pseudo-code

2 MiB Flows
32 MiB Flows
Permutation traffic on 8:1 oversubscribed fat tree



SMaRTT-REPS: Sender-based Marked Rapidly-adapting Trimmed & Timed Transport with Recycled Entropies

- | | | |
|---|--|---|
| Tommaso Bonato
ETH Zürich
Microsoft | Abdul Kabbani
Microsoft | Daniele De Sensi
Sapienza University of Rome |
| Rong Pan
AMD | Yanfang Le
AMD | Costin Raiciu
Broadcom Inc. |
| Mark Handley
Broadcom Inc. | Timo Schneider
ETH Zürich | Nils Blach
ETH Zürich |
| Ahmad Ghalayini
Microsoft | Daniel Alves
Microsoft | Michael Papamichael
Microsoft |
| Adrian Caulfield
Microsoft | Torsten Hoefler
ETH Zürich
Microsoft | |

Transport layer features

Transport Layer

Semantics Sublayer (SES)

Packet Delivery Sublayer (PDS)

Congestion Mgmt Sublayer (CMS)

Transport Security Sublayer (TSS)

- End-to-end AES encryption
- Key derivation for additional security
- Replay protection
- Scalable security domains
- Optional within UET

- **Builds on state of the art of IPSec and PSP – fixes all known attacks on RDMA**
 - AES-GCM, KDFs, IVs, Key Rotation, Anti-Replay
 - Protect data, connection establishment, replay in all scenarios
- **High scalability**
 - Group (re)keying
 - Secure Domains
 - Strong isolation (also wrt. in-network computation)

Key Points and Conclusions

More of SPCL's research:

youtube.com/@spcl 210+ Talks

twitter.com/spcl_eth 1.4K+ Followers

github.com/spcl 2K+ Stars

... or spcl.ethz.ch



Converging our HPC Networking Mess into a Unified Ethernet Standard

COVER FEATURE TECHNOLOGY PREDICTIONS

But why is RDMA over Converged Ethernet (RoCE) not sufficient?

Ultra Ethernet Consortium

Founding Members: AMD, ARISTA, BROADCOM, CISCO, EVIDEN, intel, Meta, Microsoft

Ultra Ethernet™ Specification v1.0
June 11, 2025

<https://ultraethernet.org/uec-1-0-spec>

IEEE Computer, June 2023

Ecosystem is quickly growing

Today 10 steering companies, 17 general member companies, 51 contributor members

Chair's view of the Transport WG Meeting in March '24 (60+ members on site, 1,300 members online)

Application Layer (CCL, MPI, OpenSHMEM)

Transport Layer

- libfabric
- UE addressing, Send/Recv, Deferrable Send, RMA Read/Write

Semantics Sublayer (SES)

- Zero-RTT PD Context Build, Req, Resp, Ctrl Packets, Flexible Loss Detection

Packet Delivery Sublayer (PDS)

- Window-based Sender + opt. Receiver Based CC components, Flexible LB

Congestion Mgmt Sublayer (CMS)

- O(1) key state per FEP, KDF for rekeying and per-client keys, replay protection

Transport Security Sublayer (TSS)

Network Layer (IP)

- IP (IPv4 (208 + opt) or IPv6 (408))
- Opt. UDP (68)
- Opt. TSS (12-168)
- PDS (if no UDP - 4B entropy) (28 + 12-258 RUD/RUD or 88 RUD, or 48 UUD + opt. 88 CRC)
- SES (44B standard, 32B up to 80B mgs, 208 min non-matching)
- UE Payload
- Opt. TSS (168)

Data Link Layer

- Credit-Based FC
- Link Level Retry

Ethernet Physical Layer

- UE 100G / lane
- UE 200G / lane

Physical Medium Attachment / Dependent (PMA/PCS)

Follows developing Ethernet specifications

- 802.3 100G per lane signaling today
- 802.3 200G per lane signaling upcoming

<https://arxiv.org/abs/2508.08906>

SPCL

SMaRTT-REPS enables Modern Packet Spraying

"State of the art" (2024), easily configured congestion control mechanisms

Permutation traffic on 8:1 oversubscribed fat tree

Performance in several all-to-all scenarios - 1MB message

128 Nodes - 8:1 OS - 800Gbps - 4KIB MTU

<https://arxiv.org/abs/2508.08906>

Ultra Ethernet's Design Principles and Architectural Innovations

- TORSTEN HOEFLER, ETH Zurich, Switzerland & Microsoft, USA
KAREN SCHRAMM, Broadcom, USA
ERIC SPADA, Broadcom, USA
KEITH UNDERWOOD, Hewlett Packard Enterprise, USA
CEDELL ALEXANDER, Broadcom, USA
BOB ALVERSON, Hewlett Packard Enterprise, USA
PAUL BOTTORFF, Hewlett Packard Enterprise, USA
ADRIAN CAULFIELD, OpenAI, USA
MARK HANDLEY, OpenAI, USA
CATHY HUANG, Intel, USA
COSTIN RAICIU, Broadcom, USA
ABDUL KABBANI, Microsoft, USA
EUGENE OPSASNICK, Broadcom, USA
RONG PAN, AMD, USA
ADEE RAN, Cisco, USA
RIP SOHAN, AMD, USA
- <https://arxiv.org/abs/2508.08906>

The recently released Ultra Ethernet (UE) 1.0 specification defines a transformative High-Performance Ethernet standard for future Artificial Intelligence (AI) and High-Performance Computing (HPC) systems. This paper, written by the specification's authors, provides a high-level overview of UE's design, offering crucial motivations and scientific context to understand its innovations. While UE introduces advancements across the entire Ethernet stack, its standout contribution is the novel Ultra Ethernet Transport (UET), a potentially fully hardware-accelerated protocol engineered for reliable, fast, and efficient communication in extreme-scale systems. Unlike InfiniBand, the last major standardization effort in high-performance networking over two decades ago, UE leverages the expansive Ethernet ecosystem and the 1,000x gains in computational efficiency per moved bit to deliver a new era of high-performance networking.

1 Introduction

Ultra Ethernet (UE) standardizes a new protocol to support high-performance Artificial Intelligence (AI) and High-Performance Computing (HPC) networking over Ethernet. This paper, written by UE's authors, supplements the full specification by highlighting historical and innovative technical aspects of our nearly 2.5-year journey. It is designed to be approachable to a general audience and, thus, abstracts many details while using intuitive wording and explanations. The final authority for questions regarding UE is the full 562-page specification [35].

arXiv:2508.08906v1 [cs.NI] 12 Aug 2025

More detailed requirements for HPC and AI networks



- Low latency / RTT
- Small message efficiency / message rate
- Tag matching (MPI, complex)
- Large # of connections (>10k for some apps)



- Extreme bandwidth requirements at endpoint
- No tags, in-order delivery though
- Connecting to few (<1k) endpoints
- Regular (oblivious) patterns (pre-plannable)

Bulk Synchronous Application – Last Message / Flow that finishes determines performance!

