

WORKSHOP GARR 2025

**NET MAKERS**



Finanziato  
dall'Unione europea  
NextGenerationEU



Ministero  
dell'Università  
e della Ricerca



Italiadomani  
PIANO NAZIONALE  
DI RIPRESA E RESILIENZA



RESTART

# Inside the AI Fabric: Architectures, Trends, and the Role of SRv6

Stefano Salsano

Università di Roma Tor Vergata

# Why Networking Is a Bottleneck for Large-Scale AI

## Massive Bandwidth Demand

- All-to-all parameter exchanges in data- and model-parallel training can easily saturate 200/400 Gbps links

## Ultra-Low Latency Requirements

- Tight synchronization across hundreds of GPUs means even microsecond-scale delays degrade convergence speed

## Congestion & Tail-Latency Spikes

- East-West traffic patterns create transient hotspots and microbursts that traditional queues and ECN can't fully tame

# Why Networking Is a Bottleneck for Large-Scale AI

## Scalability & Topology Constraints

- Multi-pod/DC fabrics introduce multiple hops, oversubscription ratios, and complex traffic engineering needs

## Limited Visibility & Troubleshooting

- Lack of in-network measurement makes pinpointing path-level stalls or misrouted flows extremely challenging

## Dynamic & Bursty Traffic

- Shifting traffic patterns during gradient updates demand on-the-fly steering without heavy control-plane overhead

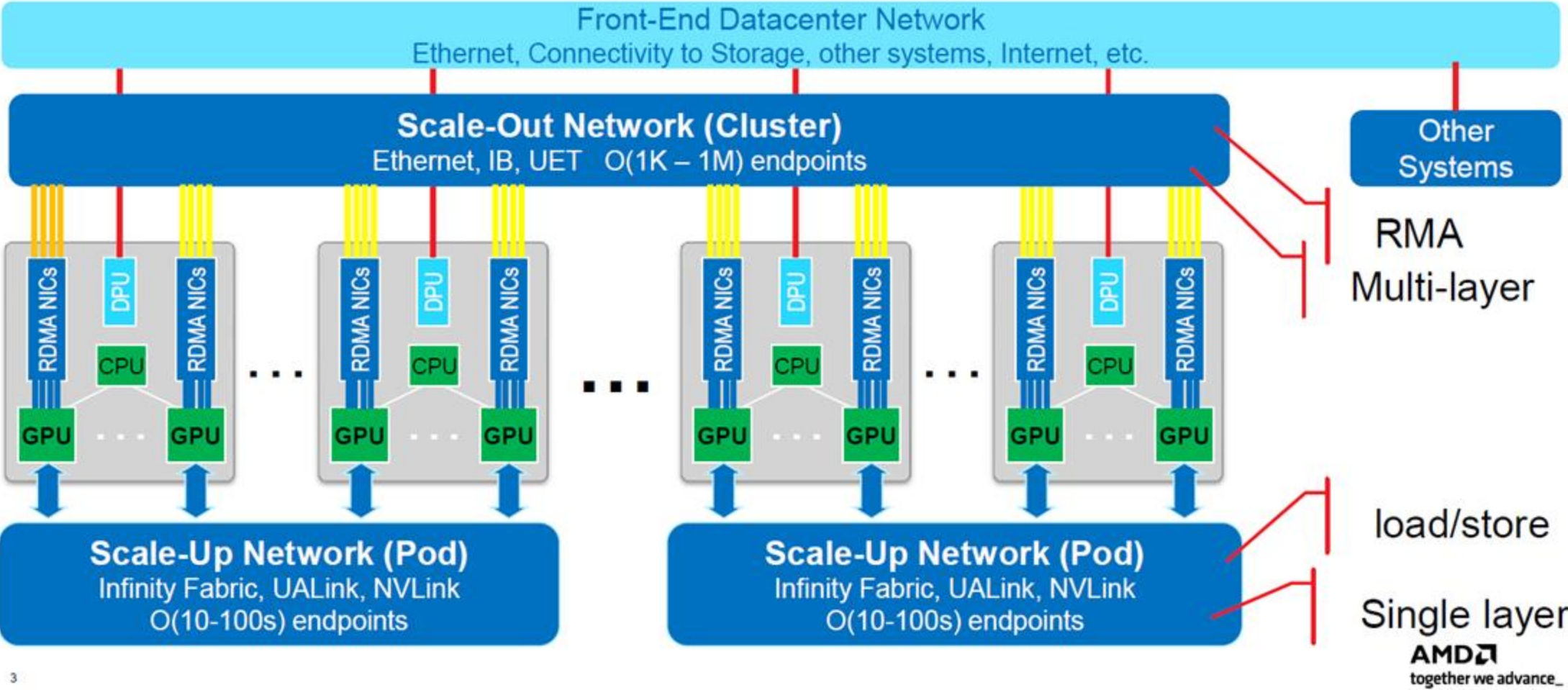
# What Characterizes Networking for AI?

- Extreme Scale
  - Low cost
  - Low power consumption
- Extreme performance: high bandwidth and low latency
- Recurring and synchronized traffic patterns
- Collective communications: all or nothing

AMD  
together we advance\_

© Mario Baldi [2]

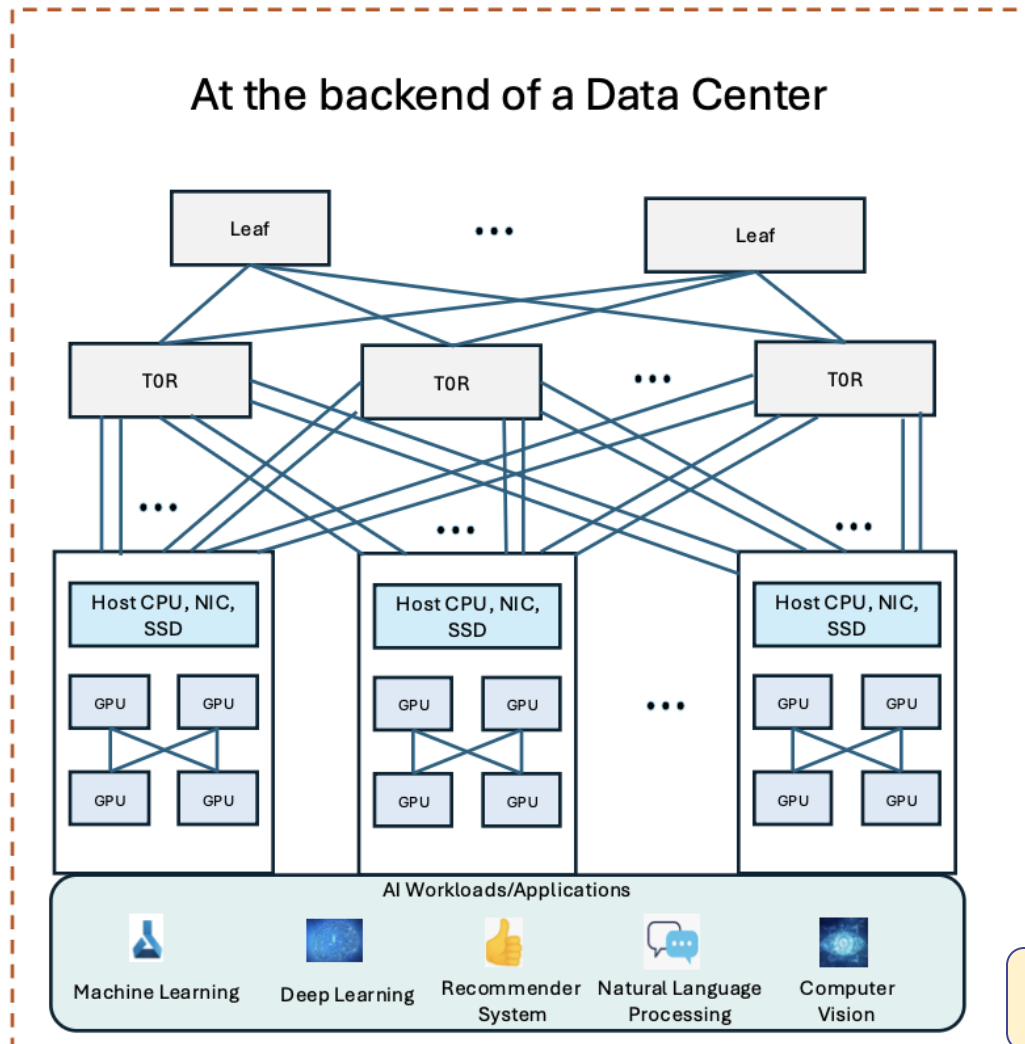
# Network architecture for AI training workloads



3

© Mario Baldi [2]

# The AI fabric



## Artificial Intelligence in the Cloud

Raising the Bar for Hyperscale Datacenter Networks

- Long lasting and large flows of training data
- Large bursts of data sent synchronously
- Long training time demands reliable networks
- Retries of failed jobs increased costs
  - Efficient traffic management, monitoring and visibility
- AI applications need fast processing and responses
  - Lossless traffic with low latency
- Traffic using RoCEv2 has low entropy for ECMP
  - Traffic engineering technology for AI backend network

© Rita Hui – Microsoft @ MPLS & SRv6 World Congress 2025

# The Ultra Ethernet Consortium (UEC)

## Founders (2023)

- AMD, Arista, Broadcom, Cisco, HPE, Intel, Meta, Microsoft, and others

## Goals

- Redefine Ethernet for AI and HPC fabrics — achieving the ultra-efficient, deterministic, and scalable transport required by large-scale training
- Deliver an *open, Ethernet-based alternative* to proprietary interconnects (e.g., InfiniBand)

## Approach

- Co-design of the full stack — hardware, transport, and software orchestration layers
- Focus on:
  - High-performance, lossless data transport
  - Congestion control and telemetry
  - Fabric management and scheduling for AI workloads

# UEC for Scale-Up (Intra-Node & Pod-Level)

**Optimize communication within a single node or pod (hundreds of GPUs)**

**Goals: Minimize latency and synchronization delay in collective operations**

## Technical directions

- New Ultra Ethernet Transport (UET) optimized for AI collectives.
- In-network acceleration for reduction/scatter operations.
- NIC and switch enhancements for deterministic latency and fine-grained flow control.
- Tight integration with training frameworks (PyTorch, TensorFlow).

# UEC for Scale-Out (Multi-Pod / Multi-Cluster)

**Extend the Ethernet fabric beyond a single rack**

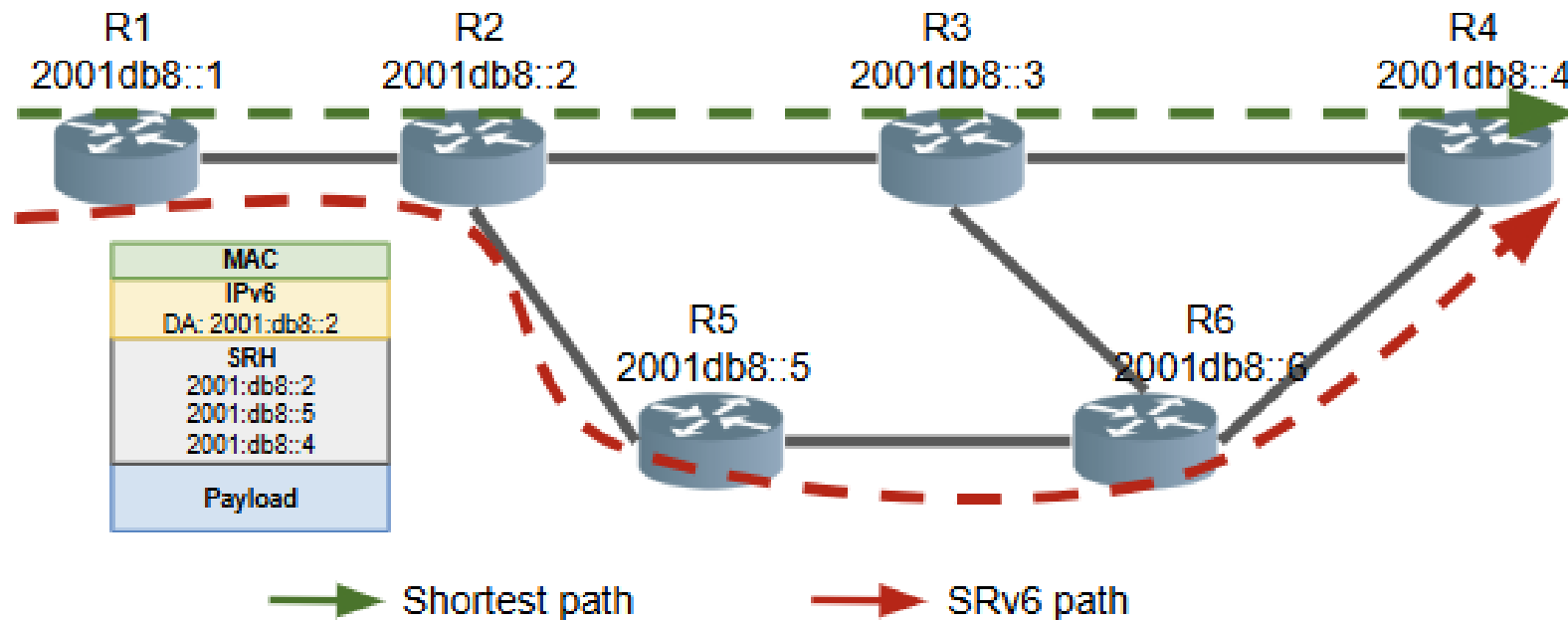
**Goals: Achieve efficient, congestion-resilient scaling of AI jobs across pods or sites**

## Technical directions

- **Hierarchical congestion control and telemetry feedback** to coordinate fabric-wide flow adjustments across multiple tiers.
- **Traffic isolation between concurrent AI jobs** ensuring predictable performance and fairness under mixed workloads.
- **Multipath RDMA transport** enabling parallel data flows across diverse routes for higher throughput and load balancing, while maintaining lossless semantics.
- **Path-aware routing for burst mitigation**, exploiting dynamic path diversity to spread traffic and minimize microbursts.

# Segment Routing for IPv6 (SRv6)

Segment Routing over IPv6 is a loose source routing technology. A list of segments is added to the packet headers. Each segment is an IPv6 address (128 bit / 16 byte)

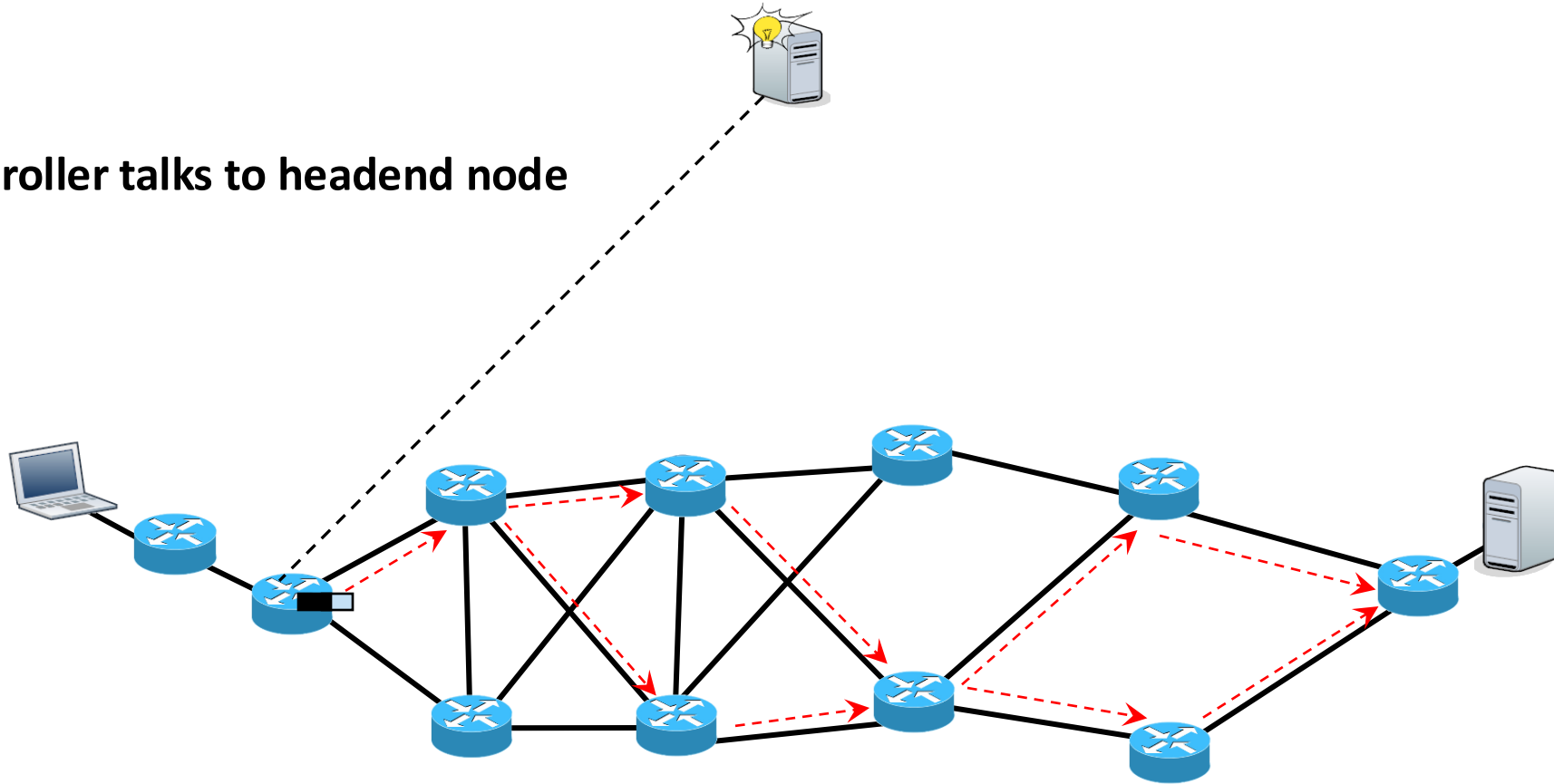


# SRv6 as a Transport (Replacement for MPLS, VXLAN...)

- Encapsulates existing IP, Ethernet, or RDMA traffic using an IPv6 outer header.
- Uses the Segment Routing Header (SRH) to encode the full path or function chain — no per-hop state.
- Provides the same traffic engineering and service chaining capabilities as MPLS, fully over IPv6.
- Eliminates MPLS control-plane complexity (LDP, RSVP-TE) while integrating with standard IPv6 and BGP.
- Enables a unified, programmable transport for backbone and AI fabric connectivity.

# SRv6 and SDN

SDN controller talks to headend node



# SRv6 and UEC: Complementary Approaches for Scale-Out AI Fabrics

- **SRv6 provides transport programmability** – path steering, traffic isolation, and in-network telemetry over standard IPv6.
- **Multipath and path-aware routing** in SRv6 can enhance UEC's congestion control and load balancing mechanisms.
- **SRv6 network programming (SIDs)** enables dynamic placement of AI workloads or control functions across domains.

# References

**[1] Industrial dissemination workshop on AI and Programmable Networks**

<https://tinyurl.com/idw-aipn25>

**[2] Mario Baldi, Networking for Distributed AI – What’s so special about it**  
available at: <https://tinyurl.com/idw-aipn25>

# **DOMANDE?**

**<http://wooclap.com>  
codice WSGARR25**

**Stefano Salsano**

**Università di Roma Tor Vergata**



# Sezione o slide conclusiva

Informazioni, approfondimenti, contatti

WORKSHOP GARR 2025

**NET MAKERS**

# Abstract

The talk will focus on the “AI fabric.” The participants will learn the architectures of current AI fabrics (i.e., data center infrastructure) for connecting tens or hundreds of GPUs with 400G or 200G links.

The current developments proposed by the vendors also within the Ultra Ethernet consortium will be reviewed.

After this overview, the application of SRv6 to an AI fabric will be discussed.