# Modern Ethernet Networks for the AI Era

The AI era promises to innovate research and information sharing. The computational demands of AI applications are significant and push the limits of our digital infrastructure. To optimise AI workloads and reduce Job-Completion Time (JCT), we not only need more powerful CPUs, GPUs and TPUs but also smarter, faster networks.

**Words:** Christophe Compain, Arista Networks

To find out how Arista supports Higher Education and Research facilities visit **solutions.arista. com/arista-for-higher-education-and-research**

**Visit Arista Networks on booth no. 7 at TNC24**

As a steering member of the Ultra Ethernet Consortium (UEC), Arista is proud to be at the forefront of building the best networking infrastructure for resilient AI clusters, leveraging its industry-leading, standards-based Extensible Operating System (EOS®) software stack and Arista Etherlink™ platforms to facilitate the next technological evolution.

## Application Requirements

In training large AI models, parameters are distributed across thousands of GPUs, with each GPU communicating the results of its calculations to its neighbours via a dedicated network.

A typical AI training workload involves potentially billions of parameters, large sparse matrices, derivatives and scalar computations distributed across hundreds or thousands of specialised processors or 'peers'. Data from these peers undergoes reduction or merging with local data, initiating another cycle of processing.

**In this compute-exchange-reduce cycle, between 20% and 50% of the job time is spent communicating across the network.**

A high-specification, robust network infrastructure is essential to ensure efficient data transfer during the exchange-reduce cycle of AI applications by minimising network congestion. Additionally, quality networks facilitate seamless data import during the initiation of new AI sessions, streamlining operations and enhancing overall performance.

## The Scalability Challenge for AI Networking

As the distributed computing environment expands and the number of possible connections between nodes increases exponentially, efficient resource utilisation across a network needs meticulous optimization and carefully considered network architecture choices.

Efficiently distributing this computation is essential for delivering results within acceptable processing times. However, the distributed nature of AI application logic has three significant implications for the network.

Firstly, the primary goal is to synchronise all GPUs to process simultaneously and produce results collaboratively, necessitating the use of RDMA (Remote Direct Memory Access) transport to minimise latency and facilitate collaborative communication patterns.

Secondly, the nature of AI training involves moving large amounts of data with a small number of flows, requiring a network with substantial bandwidth and mechanisms to efficiently manage it.

Thirdly, it follows that in order to support RDMA, the network must have specific characteristics including:

- tight synchronisation to coordinate bursty traffic flows efficiently
- specialised handling including back-pressure mechanisms to prevent congestion in "many-to-one" flows or incast scenarios
- mechanisms for efficiently managing a diverse set of substantial data transfers including a small number of large-size flows

## Arista Ethernet-based AI Networking

Similar to high performance computing (HPC) and supercomputing deployments, building connectivity for AI clusters involves both front-end and back-end networks, each with specific requirements.

The front-end network facilitates general purpose connectivity externally, handling general-purpose tasks such as user access, control and administration, and connectivity to high-performance shared storage for staging parameters and cluster outputs.

The back-end network is an island that provides a high-capacity messaging bus for the cluster. Design goals for this network deviate from those of a typical data centre and are central to AI networking.

The main components of high-speed back-end AI networks are:

- **Performance**: With 400GbE and 800GbE network switches, Ethernet can provide low latency and scalability for AI workloads

- **Lossless behaviour**: Crucial for efficient data transport, ROCEv2 (RDMA Over Converged Ethernet version 2) addresses low latency and lossless requirements
- **Flow-distribution**: IP/Ethernet delivers load balancing and no collisions for low-entry AI flows across various architectural scales
- **Back-pressure**: PFC/ECN protocols efficiently handle large bandwidths and mitigate congestion impact due to 'incast' flow patterns
- **Telemetry**: Real-time traffic counting at microsecond intervals and monitoring interface congestion/queueing latency provide in-depth visibility of AI workloads
- **Security and Management**: Virtual LANs, access control lists, multi-tenancy with VxLAN and encryption maintain data centre security and compliance

Embraced by major AI users, open, standards-based IP/Ethernet infrastructure like Arista's is favoured for both front-end and back-end networks and, unlike proprietary networking technologies, can be readily redeployed into other parts of the enterprise if needed.

Arista's programmable and highly modular EOS software stack is unmatched in the industry, empowering customers to construct resilient AI clusters. With support for hitless upgrades, it ensures uninterrupted operation, avoiding downtime and thus maximising AI cluster utilisation.

Arista Etherlink™ supports dynamic load balancing, congestion control, and reliable packet delivery to all NICs supporting RoCE across a broad range of 800G systems and line cards based on Arista EOS.

As the UEC finalises its extensions to optimise Ethernet for AI workloads, Arista is poised to deliver UEC-compatible products. These offerings will be easily upgradable to the standards the UEC sets in 2025.

**ARISTA**