

Building AI Network Fabrics

Agenda

- 1. Why is networking for AI different than traditional networking?
- 2. What is the Arista strategy for AI networking ?
- 3. What do these networks look like ?
- 4. What are the Arista platforms for AI networking ?





Why is networking for AI different than traditional networking?

Confidential. Copyright © Arista 2025. All rights reserved.

AI workloads main characteristics

- It requires specialized Hardware and Software
 - XPU (GPU, ICU, TPU ...)
 - Library and development frameworks
- Principles
 - Al applications are "optimization" applications (with potentially billions of parameters to look after ...). These are executed on XPUs.
 - These XPU's execute algorithms/calculations in parallel and then share results amongst each other ("sharding" vs "replicating").
 - There are several ways of doing this, the commonly used term is "Collective Operations"
 - All XPU's have to wait until all XPU's have shared their messages and passed on results, the workload is stalled until this happens, commonly used term is "Barrier"
- Job Completion Time (JCT) is the key metric



What are the consequences ?

AI traffic characteristics

- This is not TCP but RDMA (Remote Direct Memory Access)
- RDMA NICs push at line rate of 100G/200G/400G
 - Traffic has a small number of large flows
 >> Very poor entropy
- The effects from the collective communications
 - Bursty traffic
 - Latency accumulation for "ring-based" collectives
 - Congestion in "tree-based" collectives
 - Fails "together"
- Current DC's are on demand, generalist applications and can recover easily from failure, in the AI world failures affect performance
 - Typical Data Center applications can easily recover from failure, in AI any type of failure will affect performance. Then uses checkpoints is critical in AI applications.









Arista Strategy for AI networking

Confidential. Copyright © Arista 2025. All rights reserved.

6

Arista Next Generation AI Etherlink Portfolio

| 3 | New Al Optimized 800GbE System Maximizing choice | s |
|-------------------------|---|---|
| \sim | Accelerator, Workload and NIC Agnostic | |
| N | Best in Class Performance for AI Workloads | |
| 5nm | Latest Generation 5nm Geometry 2x Density and 25% Less Power | |
| <u>~~~</u> ∗ ∲) -50% | >50% Power Reduction with LPO and Extended DAC Cables | |



Scaling AI from 10s to > 100,000 Accelerators



Key architecture building blocks

| Small to Medium cluster | Medium to Big cluster | Massive Cluster |
|-------------------------|-----------------------|-----------------|
| | | |
| | | |
| | | |
| | | |
| | | |

Accelerator & NIC Agnostic, Open Standards, Smart Al Features



Key architecture building blocks



Accelerator & NIC Agnostic, Open Standards, Smart Al Features





Arista Al Portfolio

Small to Medium cluster



Fixed & Modular Systems 51.2T ➡ 460.8T 64-576 x 800G Nodes 128-1152 x 400G Nodes Lowest Cost, Power & Complexity

Medium to Big cluster



2- or 3-Tier Leaf-Spine or Plane-based Multi-Petabit Bisectional B/W DLB, PFC, ECN, Tens of Thousands of Nodes

Massive Cluster



Single Hop Interconnect Fully Scheduled Lossless 100% Efficient, Cell Spraying Multi-Petabit Bisectional B/W Tens of Thousands of Nodes

Accelerator & NIC Agnostic, Open Standards, Smart Al Features





What do these AI networks look like?

Confidential. Copyright © Arista 2025. All rights reserved.

11

Where to connect XPU systems ?



ARISTA¹²

NVIDIA NCIS - H100 motherboard





What are the Arista platforms for Al networking ?

Confidential. Copyright © Arista 2025. All rights reserved.

Connecting High Speed NIC / XPU Port Types

NIC/XPU Port Form Factor

Example HW (not exhaustive, examples only)



Copyright @ Arista 2025. All Rights Reserved

Platforms for AI Networking - 2024







Next Generation – 7060X6 Series for 800G



64 x 800G AI Optimized Leaf



32 x 800G AI Optimized Leaf

High Performance and Radix:

- Up to 64 x 800G wire-speed ports
- 51.2T and 25.6T Bandwidth

Advanced AI Features:

- Advanced load balancing DLB, RDMA Aware I B
- SSU, Network Telemetry, DCQCN
- EVPN VXLAN for multi-tenancy
- LPO Support

2-tier scaling to thousands of ports



Next Generation 7800R4 Series



High Performance Modular System:

- Up to 576 x 800G / 1152 x 400G
- Non-blocking 460 Tbps (Full Duplex)
 - Foundation for Scale Out AI:
 - AI optimized pipeline
 - Lossless fully scheduled VOQ
 - 100% fair cell-based fabric
 - Non-blocking architecture
 - Integrated over-provisioning
 - Deep buffering for sustained traffic
- Ideal for single tier clusters or 2-tier scaling to tens of thousands of ports



AI Center

AI Center

Key requirements

400GbE access ports No oversubscription Optimized flow distribution Lossless Advanced Telemetry

Key Variables

Total # of AI NIC ports AI NIC Transceivers cap. Rack physical layout and fiber plant Cost

Single-tiered AI Fabric XPU system XPU system XPU system XPU system Small and Moderate AI applications (10s and 1k of xPUs) **Multi-tiered AI Fabric** -----



ARISTA



Single-tiered fixed



Up to 128 xPUs at 400Gbps



7060X6-64E 64 port 800G OSFP

- Fixed Configuration Switch
 - Up to 64x 800G
- No flow collisions
 - Single-asic line-rate forwarding
- ECN and/or PFC to handle incasts
 - Low buffers Requires tuning



Single-tiered modular





- 7800R4 Modular chassis offering high port density 4, 8, 12 or 16 slots / 36x 800GbE Linecards
- Non-blocking distributed forwarding Leaf (Linecards) & Spine (Fabrics) in a single chassis
- No flow collisions between line card and fabric Scheduled Cell-based Fabric Built in overprovisioning between line card and fabric 100% Fair and Efficient Load Balancing within the chassis



- High Availability Fabric, fan, power supply, sup redundancy
- ECN and/or PFC to handle incasts Deep buffers - Requires minimal tuning







- High Potential for Flow collision Mitigated with Optimal load-balancing (DLB, Source LB)
- Uplinks over-provisioning on AI-leaf (1:1,2) No oversubscription in all circumstances Address per-flow ECMP imbalance and link/spine failure
- ECN and/or PFC to handle incasts Low buffers on AI-leaf - Requires tuning



Scaling out the R4 AI Spine

7800R4 AI Spine



7700R4 Distributed Etherlink Switch



Lossless Fully Scheduled VOQ 100% Fair And Efficient Traffic Spraying Integrated Redundancy and Resilience Optimized Pipeline for AI Workloads Lossless Fully Scheduled VOQ 100% Fair And Efficient Traffic Spraying Integrated Redundancy and Resilience Optimized Pipeline for AI Workloads

Common Architecture – Re-packaged for Al Scale Out



Copyright @ Arista 2025. All Rights Reserved

Distributed Etherlink Switch – 4k+ 400GbE XPU ports



- Single-tier distributed switching system
 - Single logical switch
 - No tuning
- No flow collisions
 - 100% efficient VOQ
 - Cell Architecture
- Rich EOS and CloudVision
 - Single point of management
 - Independent device upgrade and replacement





Ethernet AI success story

Confidential. Copyright © Arista 2025. All rights reserved.

25

Ethernet AI success at Meta

.. we have successfully used both **RoCE** and InfiniBand clusters for large, GenAI workloads (including our ongoing training of Llama 3 on our **RoCE** cluster) **without** any network bottlenecks..



- 2 Clusters with 24,576 NVIDIA Tensor Core H100 XPUs each:
 - one based on NVIDIA IB
 - one based on Ethernet (Arista 7800 + Wedge and Minipack aka Arista 7388X5)
- Similar Results, Similar performance, Similar End to End Support
- Only one technology is open and not only supported by one vendor. **Ethernet**!

https://engineering.fb.com/2024/03/12/data-center-engineering/building-metas-genai-infrastructure/





Conclusion

Confidential. Copyright © Arista 2025. All rights reserved.

27

Why Arista is leading in AI networking?



XPU agnostic, best performance, vetted at Scale Already deployed in XPU clusters of 8K, 16K, 32K and growing



Copyright @ Arista 2025. All Rights Reserved



Arista Al White Paper

Arista Al Networking

Arista in RDMA Networks



Arista 7800R3 Datasheet

Arista 7388X5 Datasheet

Arista 7060DX5 Datasheet

Arista 7060X6 Datasheet

Tech Library RoCEv2 Deployment Guide

Arista 800GbE FAQ

Arista 400GbE FAQ Arista 200GbE FAQ





ARISTA

Thank You

www.arista.com



Confidential. Copyright © Arista 2025. All rights reserved.