# Mature at Scale Memory Fabrics for all Performance and Price Points
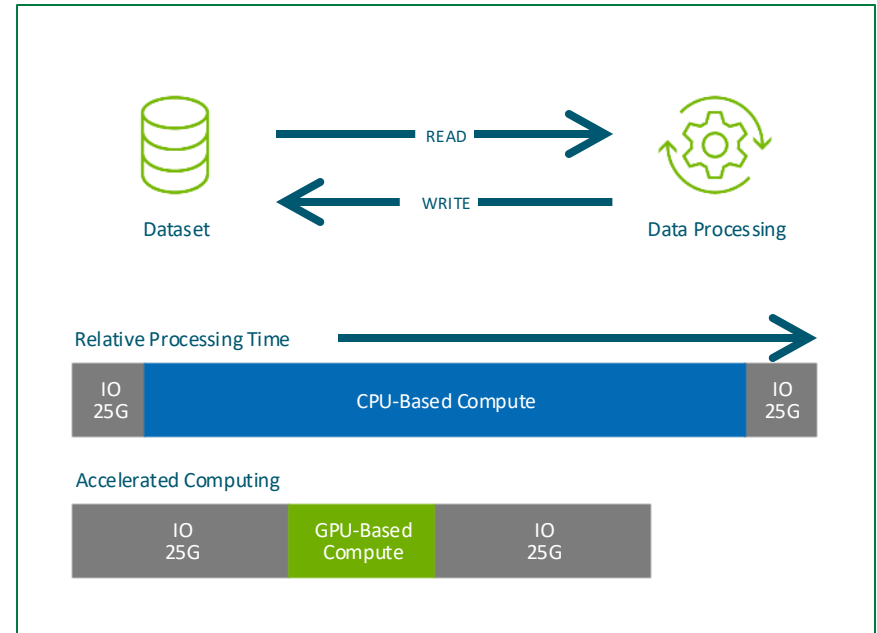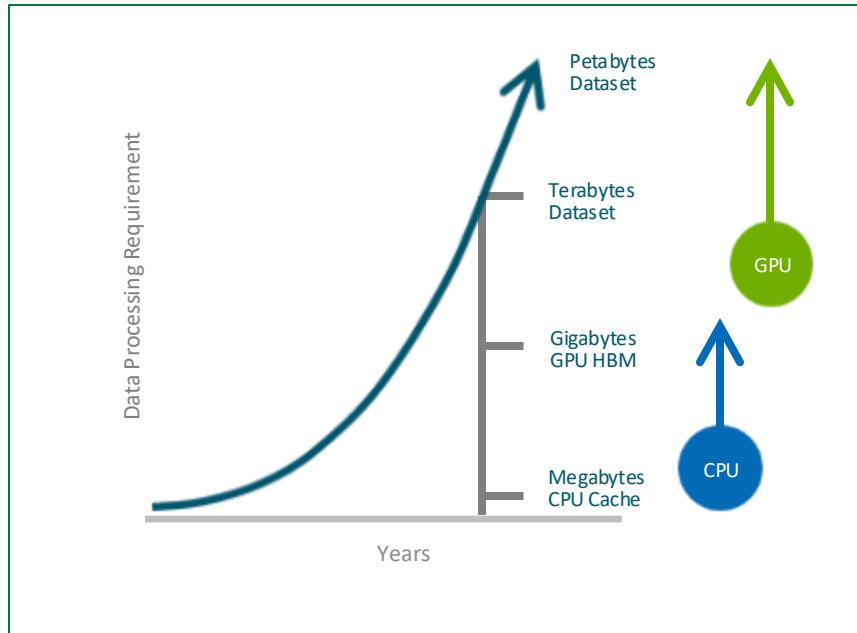
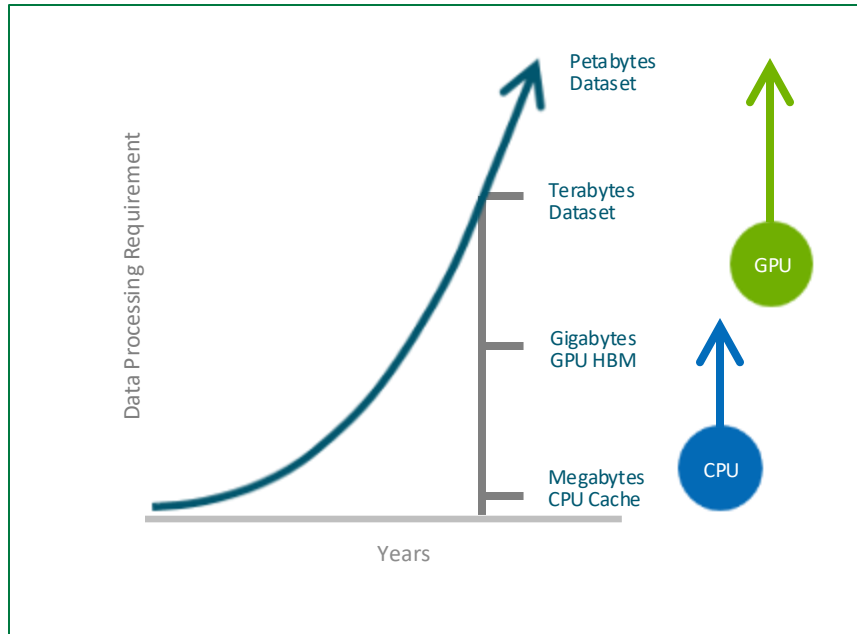Rob Davis, VP Storage Technology, NVIDIA

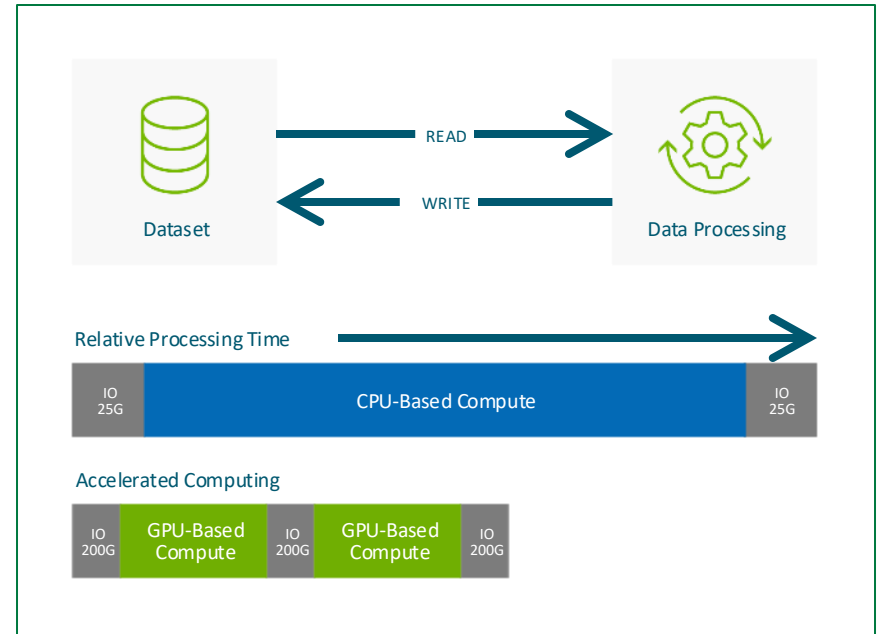# Networking Challenges for AI



GPU Application Data Sets

# Networking Challenges for AI



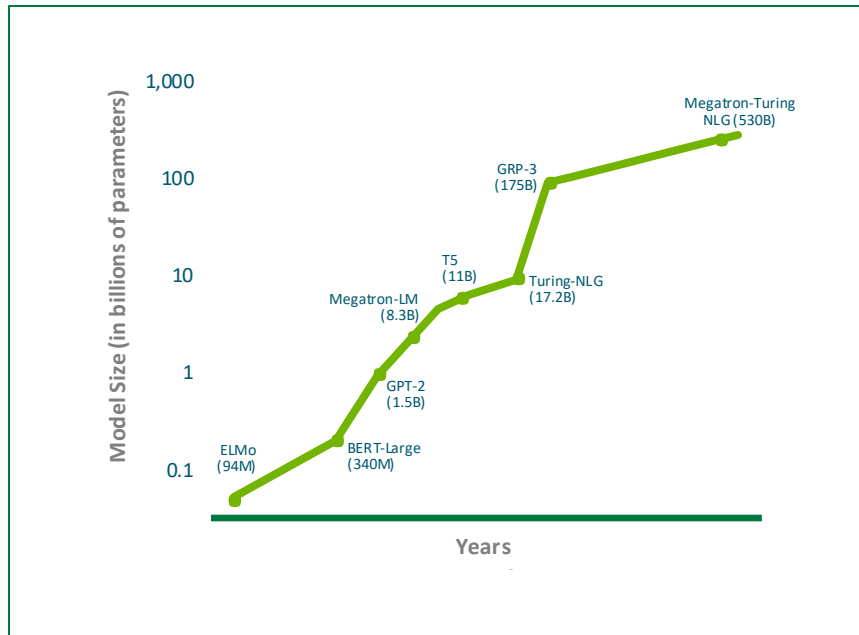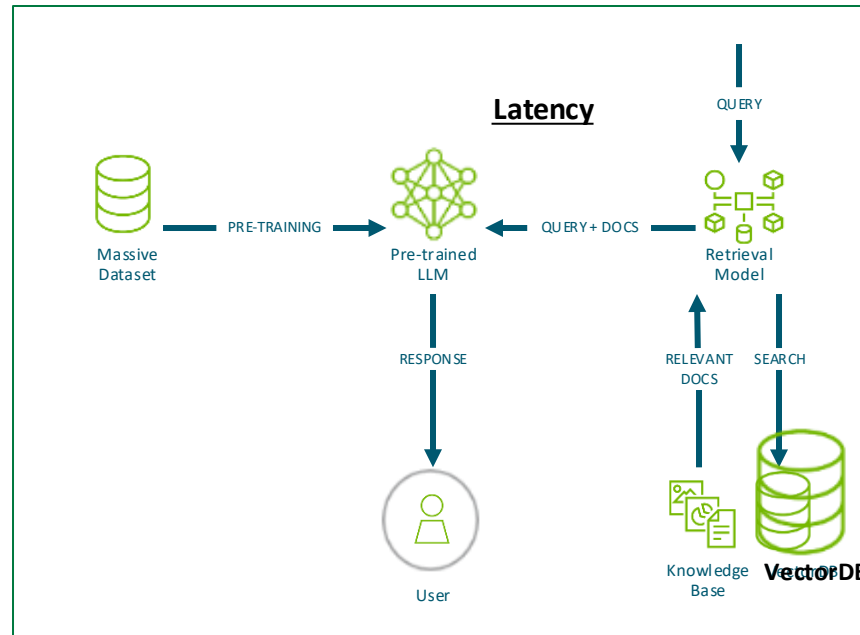GPU Application Data Sets



GPU Direct Storage (RDMA)

# Next Wave of Performance and Scalability
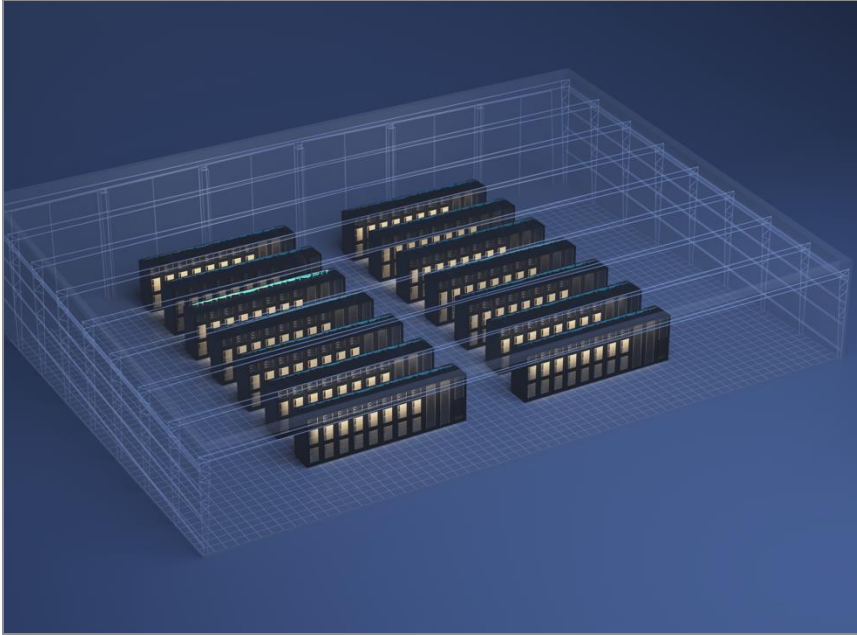
### Training



Exploding Model Size

### Inferencing



Retrieval Augmented Generation (RAG)

# Two Types of AI Data Centers



AI Factories

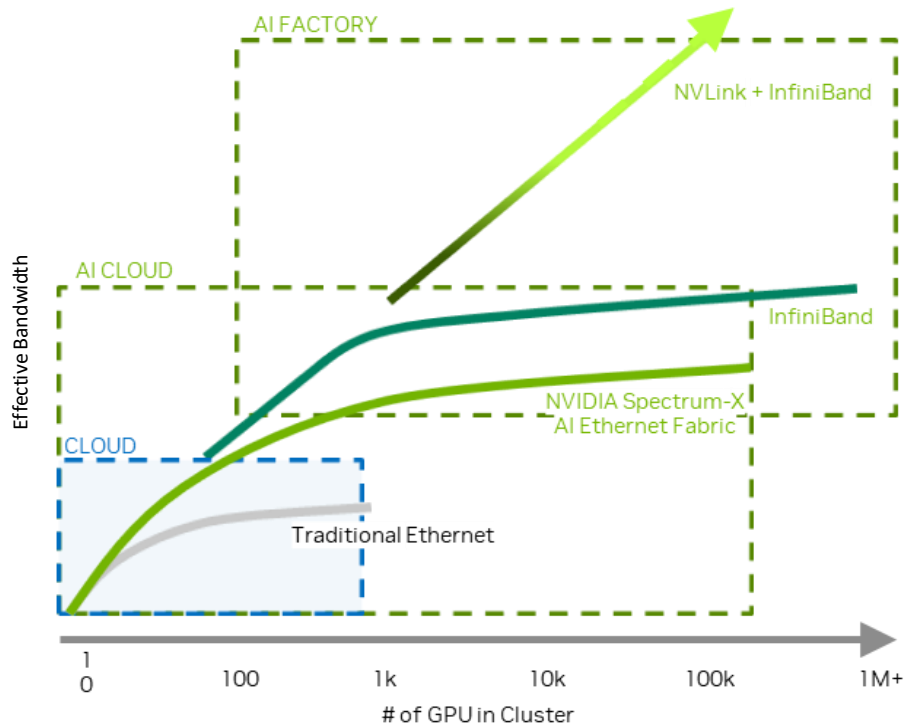Single or few users  |  Extremely large AI models  |  NVLink and InfiniBand AI fabric



AI Cloud

Multi-tenant  |  Variety of workloads  |  Ethernet network
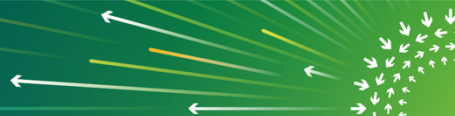
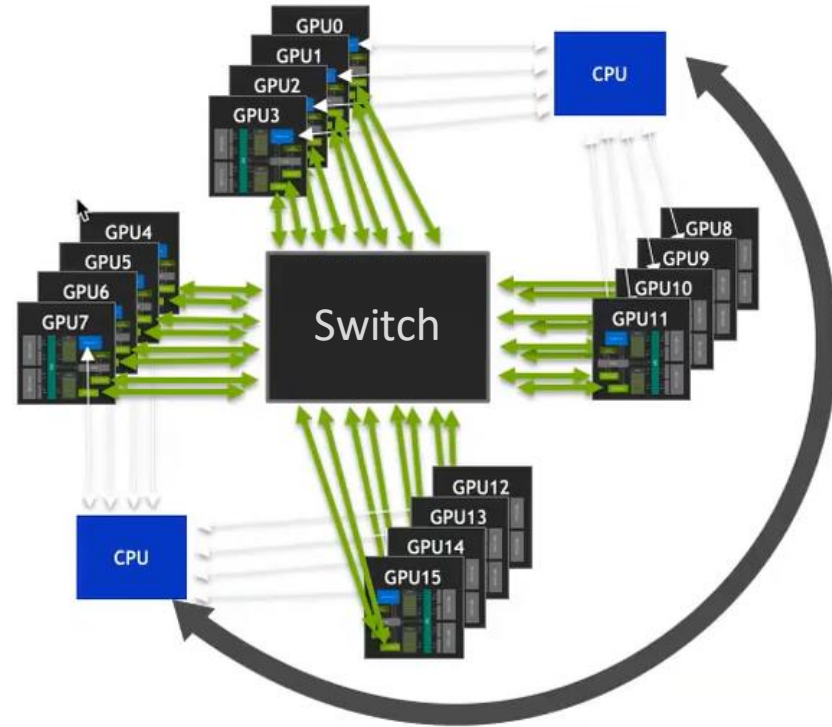# The Right Network for the Right AI Workload and Scale

# NVLINK Creates One Gigantic GPU Memory

- From the perspective of the GPUs, all HBMs can be accessed without intervention by other processes

- Load/Store, DMA

# NVLink is Tried and True and VERY Performant



NVLink At-Scale Performance

GB200 — 5th Generation NVLink
H200 — 4th Generation NVLink
A100 — 3rd Generation NVLink
V100 — 2nd Generation NVLink
P100 — 1st Generation NVLink

Y-axis: GB / sec (0, 300, 600, 900, 1,200, 1,500, 1,800, 2,100)
X-axis: Architecture Release (2014 → 2024)

# Fifth Generation NVLink Switch

**Single ASIC**

72 NVLink Ports

100GB/s port speed

7.2TB/s Total Bandwidth

3.6 TF In-Network Compute



I/O (Serdes)

…36 PORTS …

Management Logic
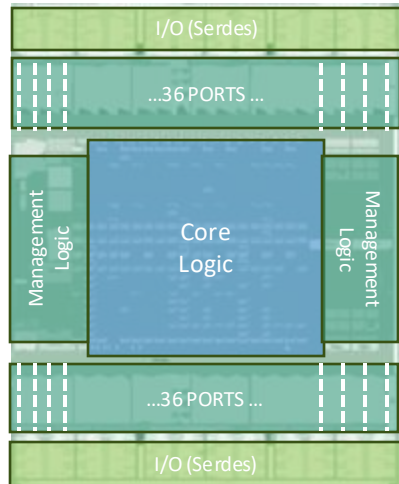
Core Logic

Management Logic

…36 PORTS …

I/O (Serdes)

1.8 TB/s Bidirectional Bandwidth

**Rack-Scale to Data Center Scale**

Up to 576 GPU NVLink Domain

Unified Fabric Management

# NVLink Deployed at Rack Scale



1,296 connections

4 GPUs per compute tray

10 compute trays

9 NVLink Switch trays

8 compute trays

18 compute trays

9 NVLink Switch trays
2 NVLink Switches with 72 ports each

# InfiniBand Creates Data Center Scale Memory with RDMA



**InfiniBand Roadmap**

| EDR | HDR | NDR | XDR | GDR | LDR |
|-----|-----|-----|-----|-----|-----|
| 100G* | 200G | 400G | 800G | 1600G | 3200G |
| (4X) | (4X) | (4X) | (4X) | (4X) | (4X) |

Link Bandwidth per direction, Gb/s

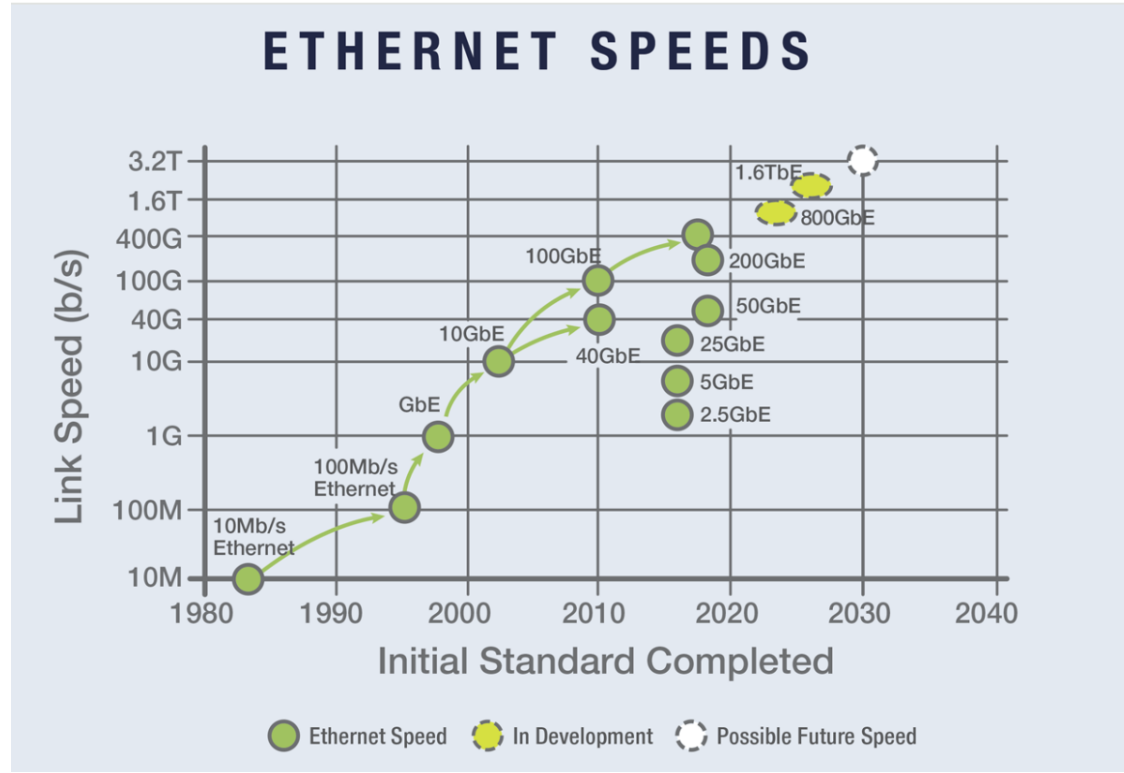Legend: 1X, 2X, 4X, 8X, 12X



144 port 800Gb/port InfiniBand Switch

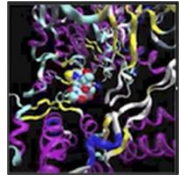# InfiniBand Scales to Hundreds of Thousands of Nodes

- 10,368 ports (2 levels)

- 746,496 ports (3 levels)

- Adaptive routing and congestion control

- Self-Healing

- Copper between switches (up to 1.5m)

# High Speed RoCE with Enhancements for AI



ETHERNET SPEEDS

Link Speed (b/s) vs. Initial Standard Completed

Legend: Ethernet Speed · In Development · Possible Future Speed

# Challenges to Running AI Workloads on Traditional Ethernet
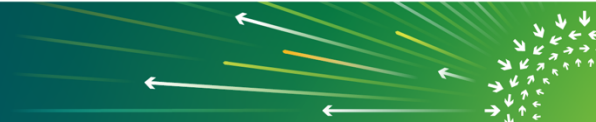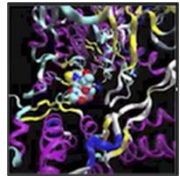


AI Workload
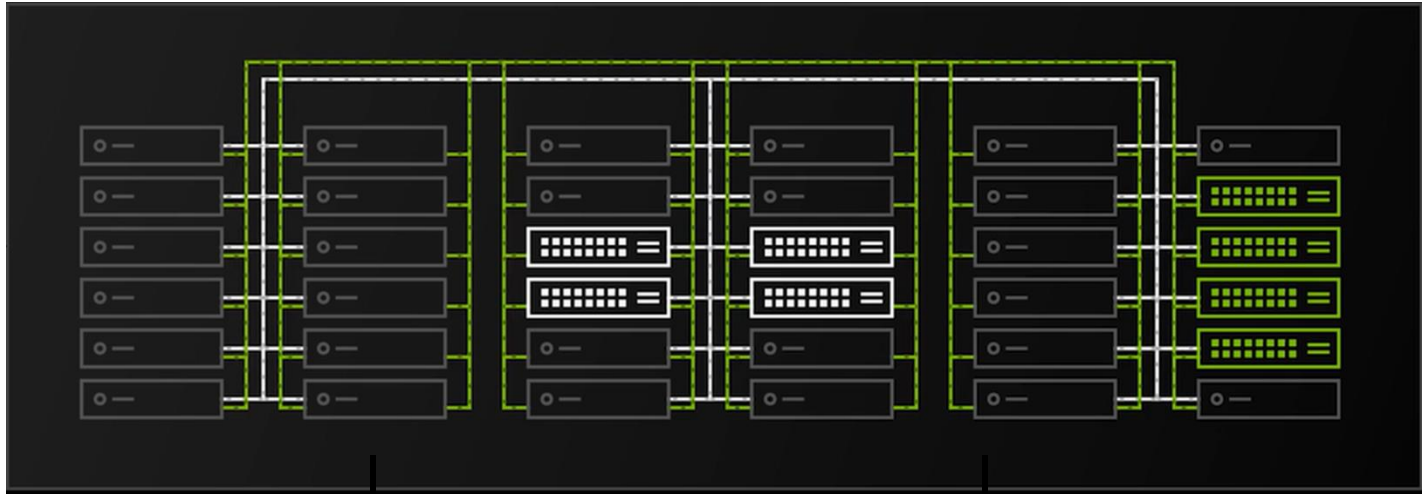
Significant Congestion

Increased Latency

Bandwidth Unfairness

# Ethernet Enhancements for AI – Spectrum-X



AI Workload

95%
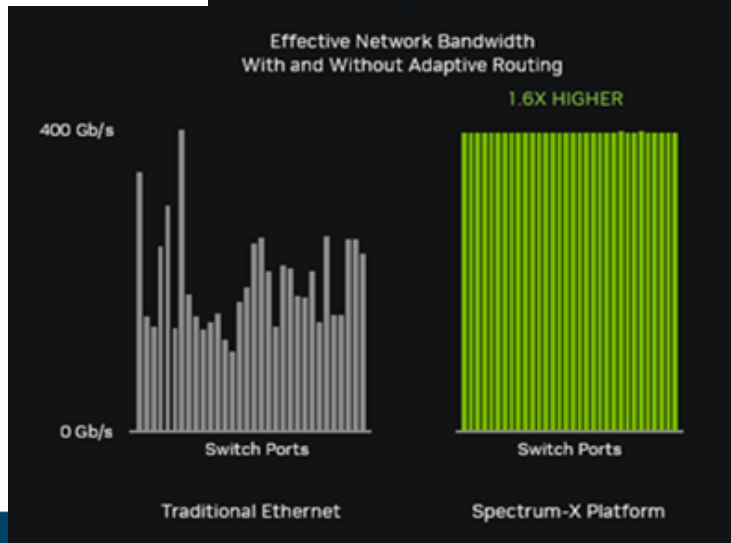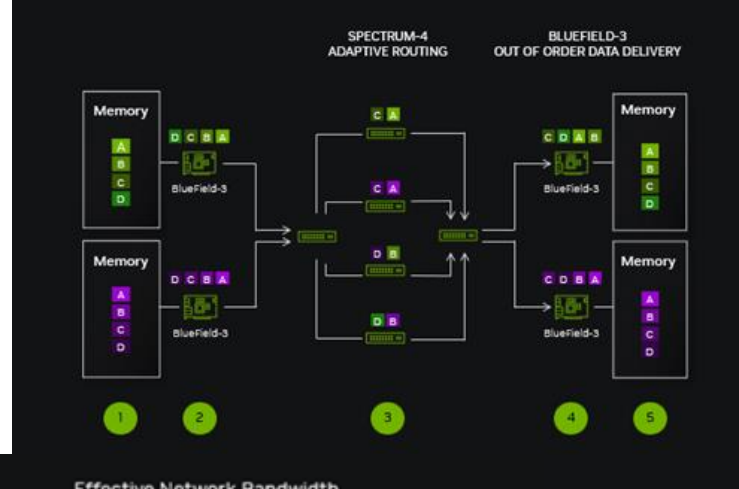Higher
Effective
Bandwidth

1.6X
Increased
AI Network
Performance
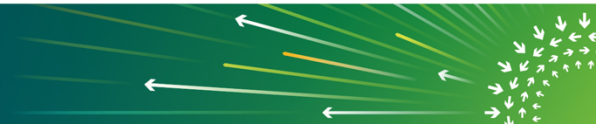
# Packet Level Adaptive Routing

- The NIC sends data into the switch network

- The switch spreads the data packets across all available routes

- The NIC ensures in-order data delivery

- Increase from typical 60% to 95% effective bandwidth

# Call to Action

- NVLink, InfiniBand and Spectrum-X Ethernet solutions are here today to improve GPU efficiency at different performance and price points

- NVLink, InfiniBand and Spectrum-X Ethernet products are available and welcome AI solution partners to test them and show the advantages

- Reach out to us for any questions

- Where to find additional information (URL links)

  - https://www.nvidia.com/en-us/data-center/nvlink/

  - https://www.infinibandta.org/infiniband-roadmap/

  - https://www.nvidia.com/en-us/networking/spectrumx/

Thank you!