# From Fast to Faster

**An Overview of Network Architecture for AI Workloads**

Romeo Lazar
Sales Manager Eastern Europe
Corning Optical Communications
Mai 2024

CORNING

**The Ethernet and InfiniBand (Technology) Roadmap**

CORNING

# What Are We Going to Talk About?

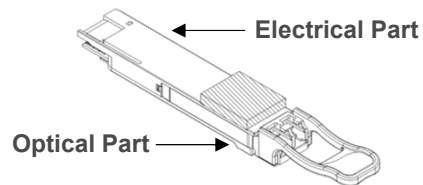## How much data can you transfer in 1 second?

- 🔴 10G
- 🟢 40G
- ⚪ 100G
- 🔵 200G
- 🟡 400G
- 🟢 800G
- 🟣 1.6T

## In simple words, how do we do it?

**Switches**

**Transceiver**
- Electrical Part
- Optical Part
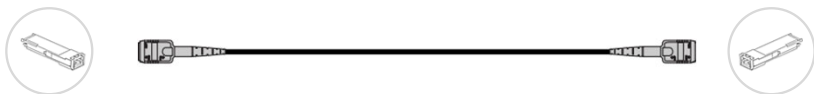
**Optical Fiber Components**
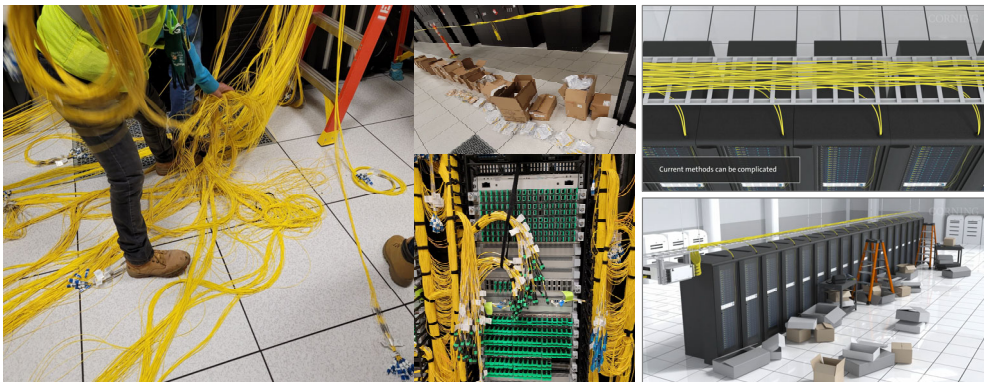
**Transmit / Receive**

CORNING

# How Can The Cabling be Done?

## Point-to-Point Cabling (Unstructured Cabling)
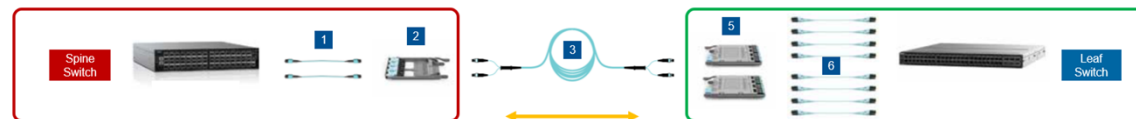
**The cabling starts with a few connections . . .**



**And this is how it ends up ...**



Current methods can be complicated

- ❯ No defined cable paths
- ❯ Problem determination difficult
- ❯ Changes made at active equipment
- ❯ System growth can be impacted

## Structured Cabling



| ■ MDA | ■ Switch | ■ EDA | ■ Horizontal Cabling | ■ Housing |

**1** **MTP Patch Cords**
MTP patch cords with MTP PRO to allow field management of pinning and polarity. MTP patch cords support parallel optics like QSFP, QSFP-DD and OSFP

**2** **MTP Adapter Panel**
Reverse polarity adapter for field polarity management

**3** **Trunk**
MTP trunk with 100 lb pulling grip to simplify installation

**5** **Module**
MTP-LC cassette to support port breakout functionality

**6** **LC Uniboot Patch Cords**
Reverse polarity uniboot patch cords minimize patch cord density and optimize routing

- ❯ Maximizes space and reduces installation time and cost
- ❯ Moves, adds, and changes (MACs) can be made easily
- ❯ A structured cabling system will provide the extra space needed for future growth
- ❯ A well-planned infrastructure can last 15-20 years and remain operational through multiple generations of system equipment and data-rate increases

CORNING

# How Can The Cabling be Done?
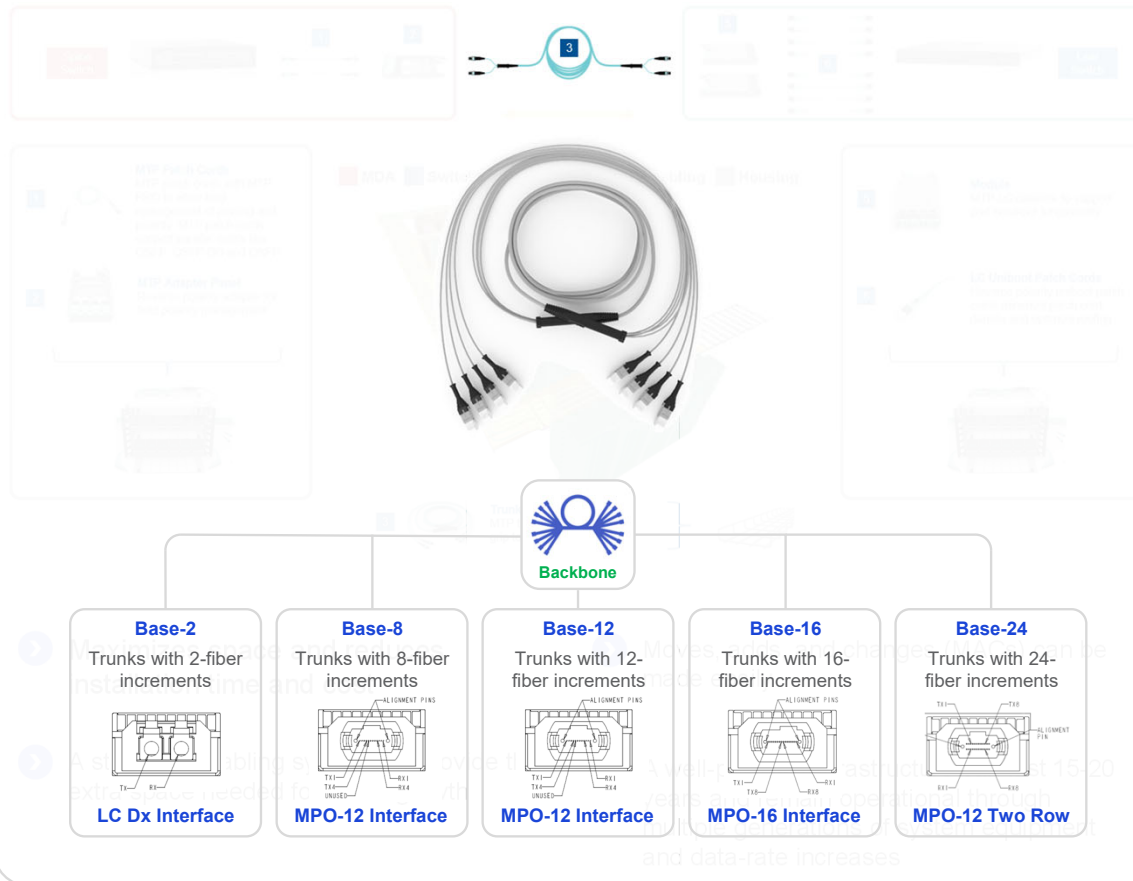
## Point-to-Point Cabling (Unstructured Cabling)

The cabling starts with a few connections . . .

**Switch Connectivity**

| Duplex LC Interface | MPO-8/12 Interface | MPO-16 APC Interface | Quad SN Interface |
|---|---|---|---|
| Dual Duplex LC Interface | Dual MPO-12 Interface | MPO-12 Two-Row Interface | 8x MDC and SN Interface |
| Dual Mini-LC Interface | Dual CS Interface | Quad MDC Interface | |

● Transceiver footprint available in the market
● Transceiver footprint not yet available

No defined cable paths

Changes made at active equipment

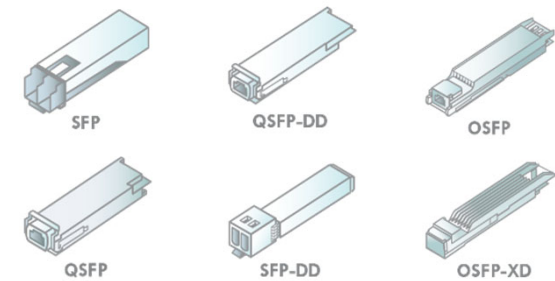Problem determination difficult

System growth can be impacted

## Structured Cabling

**Backbone**

| Base-2 | Base-8 | Base-12 | Base-16 | Base-24 |
|---|---|---|---|---|
| Trunks with 2-fiber increments | Trunks with 8-fiber increments | Trunks with 12-fiber increments | Trunks with 16-fiber increments | Trunks with 24-fiber increments |
| LC Dx Interface | MPO-12 Interface | MPO-12 Interface | MPO-16 Interface | MPO-12 Two Row |

CORNING

# Transceiver Roadmap and Backbone of Choice

| Transceiver Speed | 10G | 25G | 40G | | 50G | 100G | | | 200G | | | 400G | | | | 800G | | | | 1.6T |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Pluggable Module** | SFP | SFP | SFP / QSFP | | SFP / QSFP | SFP / SFP-DD / QSFP / QSFP-DD / OSFP | | | QSFP / QSFP-DD / SFP-DD | | | QSFP / QSFP-DD / OSFP | | | | QSFP / QSFP-DD / OSFP | | | | QSFP / QSFP-DD / OSFP / OSFP-XD |
| **SMF** | LR | LR | LR4 FR4 | PLR4 PLRL4 | LR FR | LR FR DR LR4 CWDM4 | N/A | PSM4 | LR4 FR4 FR DR | N/A | DR4 | LR8 FR8 FR4 LR4-6 LR4-10 | 2FR4 | DR4 DR2 DR4-2 | N/A | LR8 FR8 | 2LR4 2FR4 FR4 | DR4 DR4-2 | 2DR4 2PLR4 8FR DR8 DR8-2 | DR8 DR8-2 |
| **MMF** | SR | SR | BiDi SWDM4 | SR4 eSR4 | SR | BiDi SWDM4 VR SR | SR2 | SR4 eSR4 | N/A | VR2 SR2 | SR4 | N/A | N/A | SR4.2 VR4 SR4 | SR8 | N/A | N/A | VR4.2 SR4.2 | SR8 VR8 2VR4 2SR4 | VR8.2 SR8.2 |
| **Fibers per transceiver** | 2 | 2 | 2 | 8 | 2 | 2 | 4 (2x2) | 8 | 2 | 4 (2x2) | 8 | 2 | 4 (2x2) | 8 | 16 (16x1) | 2 | 4 (2x2) | 8 | 16 (8x2 or 16x1) | 16 (8x2 or 16x1) |
| **Base-2** | ● | ● | ● | ○ | ● | – | – | – | – | – | – | – | – | – | – | – | – | – | – | – |
| **Base-8** | ● | ● | ● | ● | ● | ● | ● | ● | ● | ● | ● | ● | ● | ● | ● | ● | ● | ● | ● | ● |
| **Base-12** | ● | ● | ● | ○ | ● | ● | ● | ○ | ● | ● | ○ | ● | ● | ○ | ○ | ● | ● | ○ | ○ | ○ |
| **Base-16** | ⊙ | ⊙ | ⊙ | ⊙ | ⊙ | ⊙ | ⊙ | ⊙ | ⊙ | ⊙ | ⊙ | ⊙ | ⊙ | ⊙ | ● | ⊙ | ⊙ | ⊙ | ● | ● |
| **Base-24** | ● | ● | ● | ○ | ● | ● | ● | ● | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ |

**Legend:**

- ● Allow full scalability, 100% fiber utilization and migration
- ⊙ Allow scalability and migration. Limited backward compatibility with existing Base-8 and Base-12 backbones / installations
- ○ Scalability and migration complexity in some degree (base conversion components, partial fiber utilization)
- – Not recommended due to scalability limitations and high complexity

Backbone

| Base-2 | Base-8 | Base-12 | Base-16 | Base-24 |
|---|---|---|---|---|
| Trunks with 2-fiber increments | Trunks with 8-fiber increments | Trunks with 12-fiber increments | Trunks with 16-fiber increments | Trunks with 24-fiber increments |
| LC Dx Interface | MPO-12 Interface | MPO-12 Interface | MPO-16 Interface | MPO-12 Two Row |

**The connector in the backbone is relevant for: Flexibility, Migration to new technologies, Scalability, TCO**

SFP  QSFP-DD  OSFP
QSFP  SFP-DD  OSFP-XD

Picture source: Ethernet Alliance

CORNING

# Corning's Way of Working

**Voice of Customer** + **Voice of Technology** + **R&D** = **Future-ready**

| | | | |
|---|---|---|---|
| Hyper | CISCO | EDGE™ | TCO |
| MTDC | ARISTA | EDGE8® | ● 10G |
| Enterprise | DELL EMC | Clean Advantage™ | ● 40G |
| | BROADCOM® | EDGE™ Rapid Connect | ● 100G |
| | | EDGE™ MDC | ● 200G |
| | | | ● 400G |
| | | | ● 800G |
| | | | ● 1.6T |

CORNING

# OSFP Optical Interfaces

## 2-Fiber Transceivers

**Duplex LC Optical Interface**



## 4-Fiber Transceivers

**Dual Duplex LC Optical Interface**



**Dual Mini-LC Optical Interface**

LC Dup    Mini-LC



**Dual CS Optical Interface**



## 8-Fiber Transceivers

**MPO-8/12 Optical Interface**



**Quad SN Optical Interface**



**Quad MDC Optical Interface**



## 16-Fiber Transceivers

**Dual MPO-12 Optical Interface**



**MPO-16 Optical Interface**



**MPO-12 Two Row Optical Interface**



**8x MDC and SN Optical Interface**

MDC    SN



---

● Footprint available in the market    ● Footprint available and high adoption expected    ● Transceiver footprint **not** yet available

CORNING

# Base - 8 Backbone

**Backbone**

## MPO-8/12 Interface

ALIGNMENT PINS
TX1 RX1
TX4 RX4
UNUSED

**Trunks with 8-fiber increments**

**Switch Connectivity**

| Duplex LC Interface | MPO-8/12 Interface | MPO-16 APC Interface | Quad SN Interface |
|---|---|---|---|
| Dual Duplex LC Interface | Dual MPO-12 Interface | MPO-12 Two-Row Interface | 8x MDC and SN Interface |
| Dual Mini-LC Interface | Dual CS Interface | Quad MDC Interface | |

LC Dup   Mini-LC

MDC   SN

● **Transceiver footprint available in the market**

● **Transceiver footprint not yet available**

**Where Used**

- **Backwards compatible** with existing Base-8 and Base-12 architectures.

- Used in **small to large data centers**, enabling **migration to new transceiver technologies** with minimal to no change in existing structured cabling

- Widely believed to be the **most flexible option** to accommodate future industry trends, supporting deployments of new varieties of connectors at the transceiver, with full fiber utilization

**Migration**

### Base-8 supports the following data rates

● 10G  ● 40G  ● 100G  ● 200G  ● 400G  ● 800G  ● 1.6T*

**Cabling Infrastructure**

### Examples of components used with different optical interfaces for different data rates

**40G | 100G | 400G | 800G | 1.6T**

| Single or Dual MTP-8 Interface | MTP-8 Patch Cord | MTP Panel | MTP-8 Trunk | MTP-8 to LC Module | LCDx Patch Cord | LCDx Interface |

**400G | 800G | 1.6T**

| MTP-16 APC Interface | MTP-16 to MTP-8 Harness | MTP Panel | MTP-8 Trunk | MTP Panel | MTP-8 Patch Cord | MTP-8 Interface |

**10G | 40G | 100G | 200G | 400G | 800G**

| VSFFC (MDC or SN) Interface | LCDx to VSFFC Patch Cord | MTP-8 to LC Module | MTP-8 Trunk | MTP-8 to LC Module | LCDx to VSFFC Patch Cord | VSFFC (MDC or SN) Interface |

*1.6T Transceivers using LC Duplex are also expected to be launched to the market

CORNING

# Understanding the Numbers of Connectivity

# Base - 8 Backbone

**Backbone**

## MPO-12 Interface

Trunks with 8-fiber increments

**Switch Connectivity**

| Duplex LC Interface | MPO-12 Interface | MPO-16 APC Interface | Quad SN Interface |
| Dual Duplex LC Interface | Dual MPO-12 Interface | MPO-12 Two-Row Interface | 8x MDC and SN Interface |
| Dual Mini-LC Interface | Dual CS Interface | Quad MDC Interface | |

LC Dup   Mini-LC

MDC   SN

● Transceiver footprint available in the market
● Transceiver footprint not yet available

**Where Used**

- **Backwards compatible** with existing Base-8 and Base-12 architectures.

- Used in **small to large data centers**, enabling **migration to new transceiver technologies** with minimal to no change in existing structured cabling

- Widely believed to be the **most flexible option** to accommodate future industry trends, supporting deployments of new varieties of connectors at the transceiver, with full fiber utilization

**Migration**

Base-8 supports the following data rates

● 10G   ● 40G   ● 100G   ● 200G   ● 400G   ● 800G   ● 1.6T

**Cabling Infrastructure**

Rotatable Strain-Relief Plate | EDGE8 TAP Modules | Tool-Less, Snap-On Integration Clip | Base-8 Pinned Trunk

Adjustable Mounting Bracket | Patch Cord Routing Guides

Top and Bottom Parking Locations | EDGE8® Harness

LC Duplex Uniboot Patch Cords | Grey Colour and Imprinted "8" | Port Breakout Module | 4-Port MTP® to MTP Panel with shuttered Adapters | 8-Fiber MTP Patch Cord

**EDGE8 Solutions**

- The best option supporting **migration** from 10G to 1.6T

- Supports Base-2, Base-8 and Base-16 connectivity with **duplex and parallel architectures**

- Support port **breakout solutions** to save space, power and cooling

- Supports **network monitoring** without adding separate space consuming hardware

- Supports keyed connectivity for **Secure Solutions**

- Supports **latency sensitive** applications

- **High Density** supporting **144F** per RU using **LC Dx** or **576F** per RU using **MTP-8**

- **Optical frames** available in single and dual versions: **5,760 duplex** or **23,040 parallel fibers**

*1.6T Transceivers using LC Duplex are also expected to be launched to the market

CORNING

# An Overview of Network Architecture for AI Workloads

# Interconnecting MDA to EDA with EDGE8®

**Spine Switch** · **1** · **2**

**3**

**4** · **5** · **Leaf Switch**

## MTP Patch Cords
**1** MTP patch cords with MTP PRO to allow field management of pinning and polarity. MTP patch cords support parallel optics like QSFP, QSFP-DD and OSFP

## MTP Adapter Panel
**2** Reverse polarity adapter for field polarity management

🔴 **MDA**  ⚪ **Switch**  🟢 **EDA**  🟡 **Horizontal Cabling**  ⚫ **Housing**

## Module
**5** MTP-LC cassette to support port breakout functionality

## LC Uniboot Patch Cords
**6** Reverse polarity uniboot patch cords minimize patch cord density and optimize routing

**3** **Trunk**
MTP trunk with 100 lb pulling grip to simplify installation

CORNING

ChatGPT

Examples
"Explain quantum computing in simple terms" →
"Got any creative ideas for a 10 year old's birthday?" →
"How do I make an HTTP request in Javascript?" →

Capabilities
Remembers what user said earlier in the conversation
Allows user to provide follow-up corrections
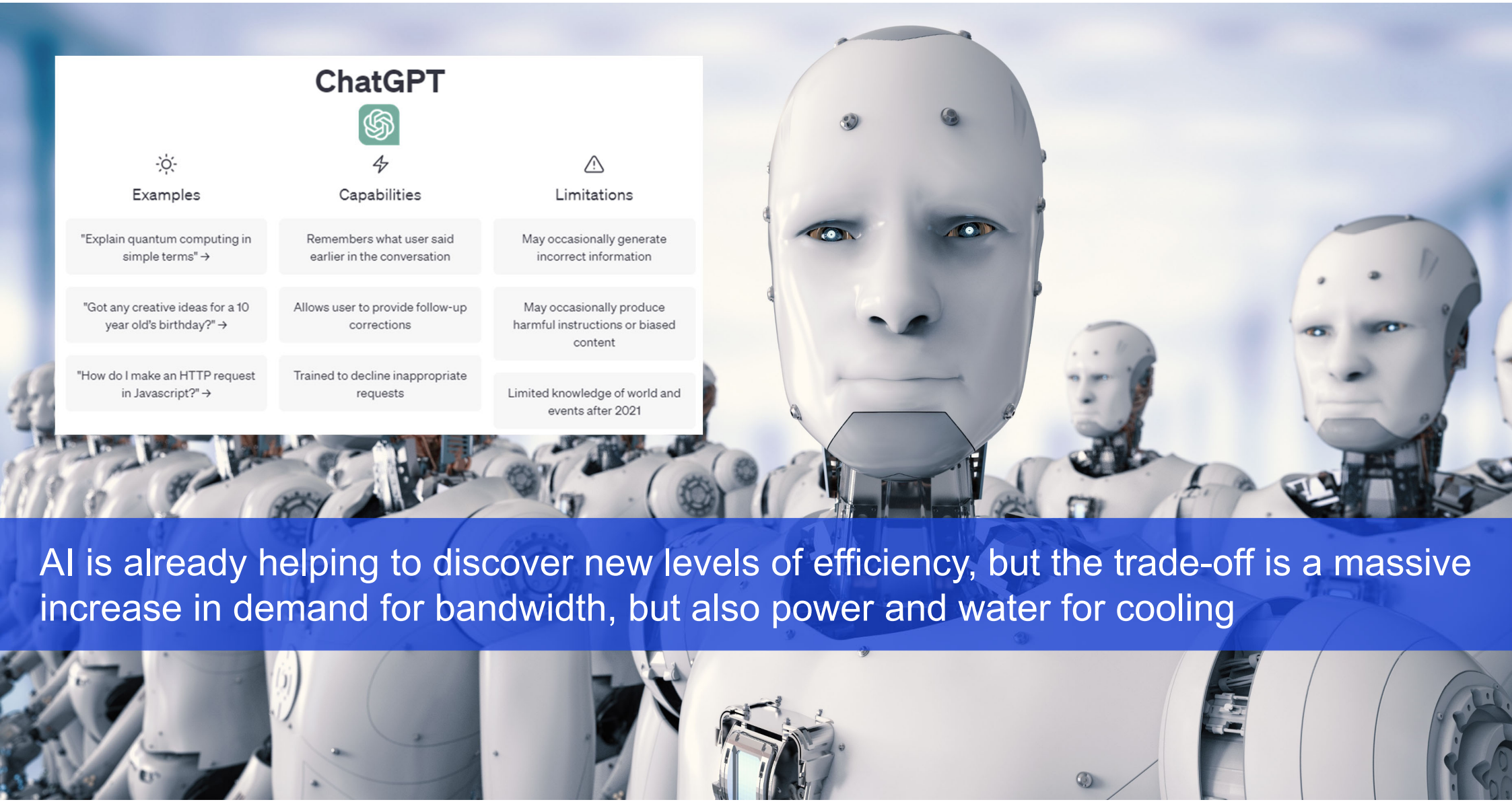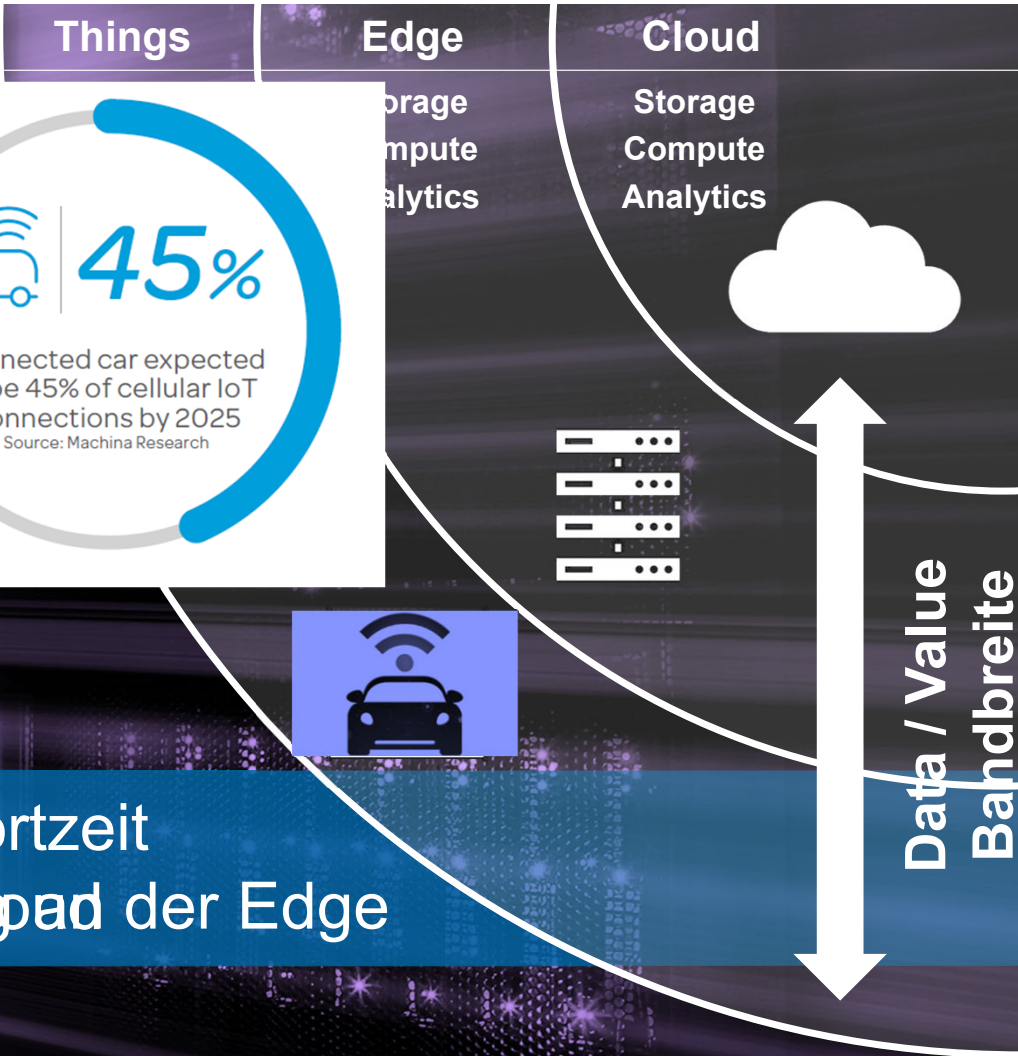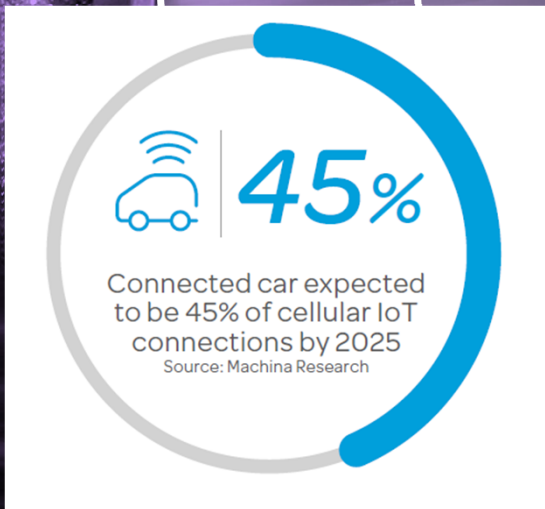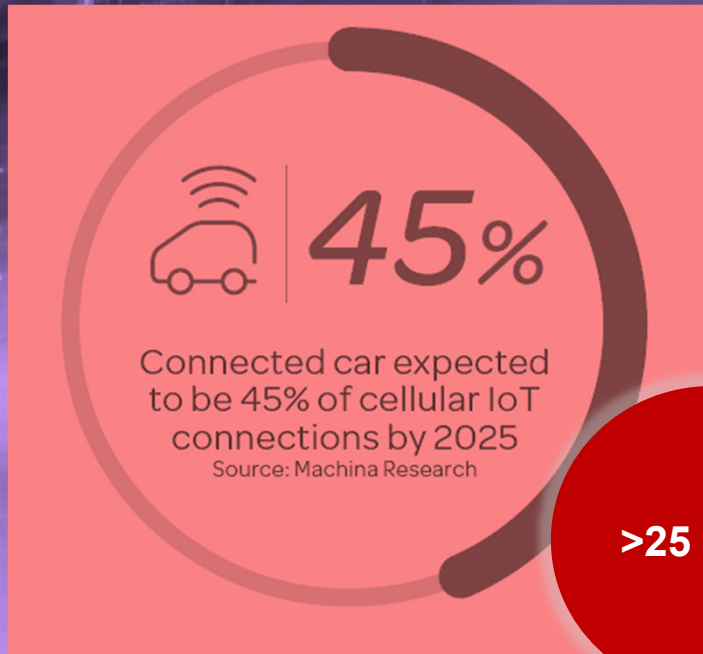Trained to decline inappropriate requests

Limitations
May occasionally generate incorrect information
May occasionally produce harmful instructions or biased content
Limited knowledge of world and events after 2021

AI is already helping to discover new levels of efficiency, but the trade-off is a massive increase in demand for bandwidth, but also power and water for cooling

CORNING

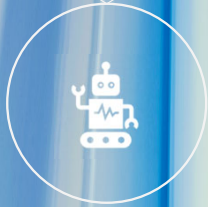# High Performance Computing, Artificial Intelligence & Machine Learning



*Image Source: ChatGPT*

ChatGPT was trained using **10,000 of Nvidia's GPUs** clustered together in a supercomputer on Azure.

Moreover, there plans for significantly increased GPU usage, with speculation that their **upcoming AI model** may require as many as **10 million GPUs**.

NVIDIA dominates the market for chips used in AI systems, with about 90% of the GPU market for ML.

CORNING

Hyper/Cloud

MTDC

Enterprise

Current AI Training workloads require large GPU clusters (32k) driving need for power and cooling efficiency and high bandwidth in the MTDC and Cloud

CORNING

# Two Different Approaches to AI/ML
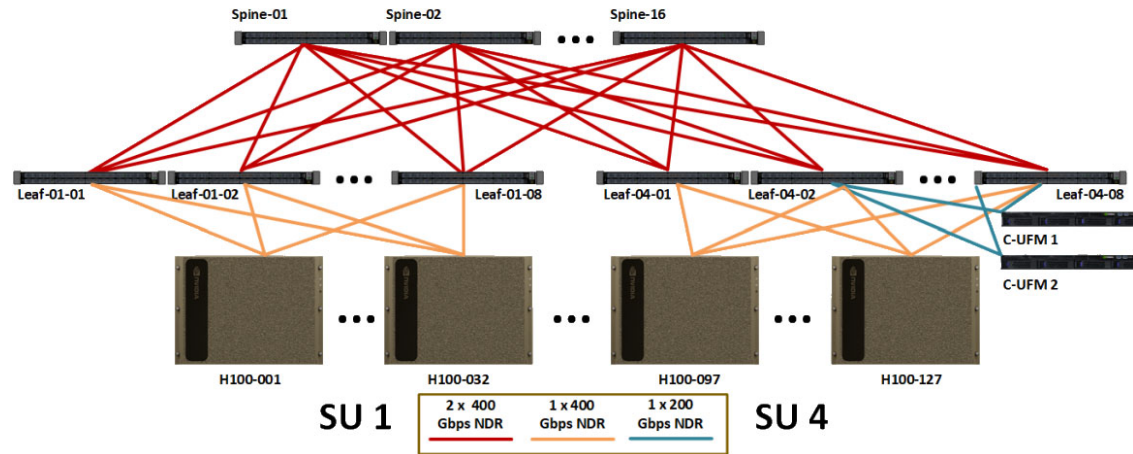
**E**

**Ethernet (Ultra Ethernet Consortium starting)**

**InfiniBand**

**IB**

**Dedicated Network for AI/ML**

**NVIDIA design as example:**



*Source: NVIDIA DGX SuperPOD. Reference Architecture Featuring NVIDIA DGX H100 Systems*

CORNING

# An Overview of Network Architecture for AI Workloads

**NVIDIA DGX-H100 SuperPOD**

**Compute Network Fabric**

**Storage Network Fabric**

**In-Band Network Fabric**
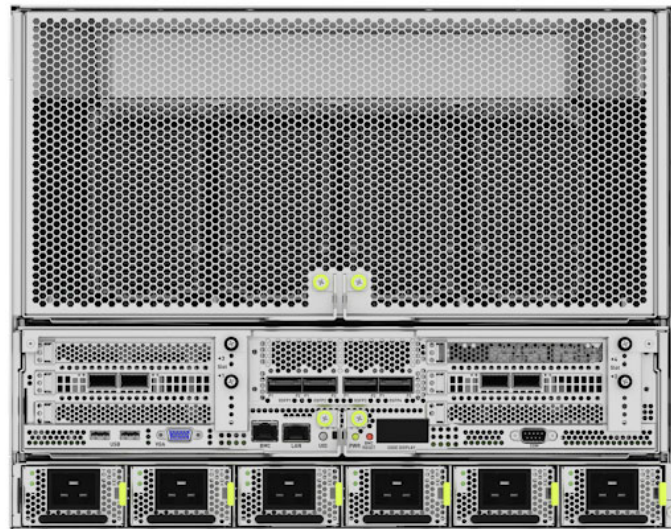
**Out-of-Band Network Fabric**

# NVIDIA DGX-H100 Compute Node (DGX-Data Center GPU Accelerated)

- For most machine learning workloads, both GPU and CPU work together to maximize performance.

- The CPU performs data cleaning on raw datasets before training models.

- Once this data is pre-processed, the CPU sends it to the GPU for parallel training/inference.

- After which, the GPU accelerates parallelizable math operations during training.

- Both are necessary for high-performance machine learning.
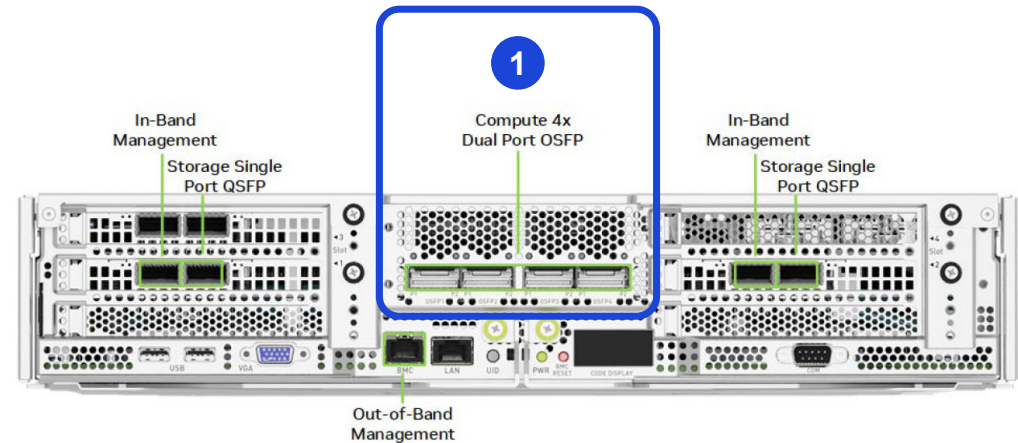
**8x GPU**

**2x CPU**



**DGX-H100 Node/Server**

| Feature | Description |
|---|---|
| Form Factor | 8U Rack mount |
| System Weight (max) | 287.6 lbs 130.45 kg |
| Input (200–240-volt AC) (max) | **10.2 kW** |

**~$482,000 USD at release**

# Compute Network Fabric



**DGX-H100 Node (Sever) Networking**



In-Band Management
Storage Single Port QSFP

**1** Compute 4x Dual Port OSFP

In-Band Management
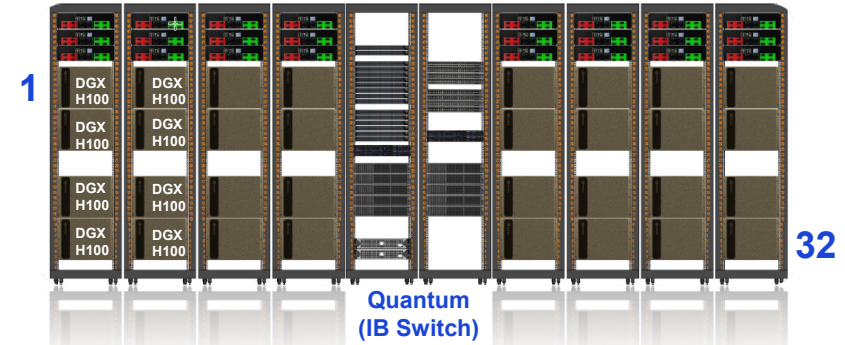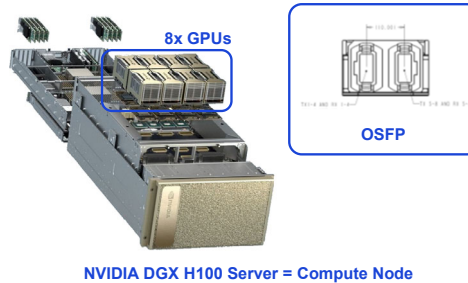Storage Single Port QSFP

Out-of-Band Management

1. **Compute Fabric:** There are 4xDual Port OSFP on each Node (8x400G connections)

2. There are 8 Leaf switches for each Scalable Unit (POD)

3. The fabric is rail-optimized, meaning that all the same Host Channel Adaptors (HCA) from each node are connected to the same leaf switch.

4. The fabric is built using Quantum 9700 Infiniband switches using 800Gbps/ Twin port OSFP transceivers



RX1 RX2 RX3 RX4 TX4 TX3 TX2 TX1
TX5 TX6 TX7 TX8 RX8 RX7 RX6 RX5

**1** 800Gb/s Twin-port OSFP, 2x400Gb/s

CORNING   *Images Source: H100 User Guide https://docs.nvidia.com/*

# NVIDIA's Reference Architecture as Example

The system is built upon building blocks of **scalable units (SU)**, each containing **32 DGX H100** systems, which provides for rapid deployment of systems of multiple sizes.

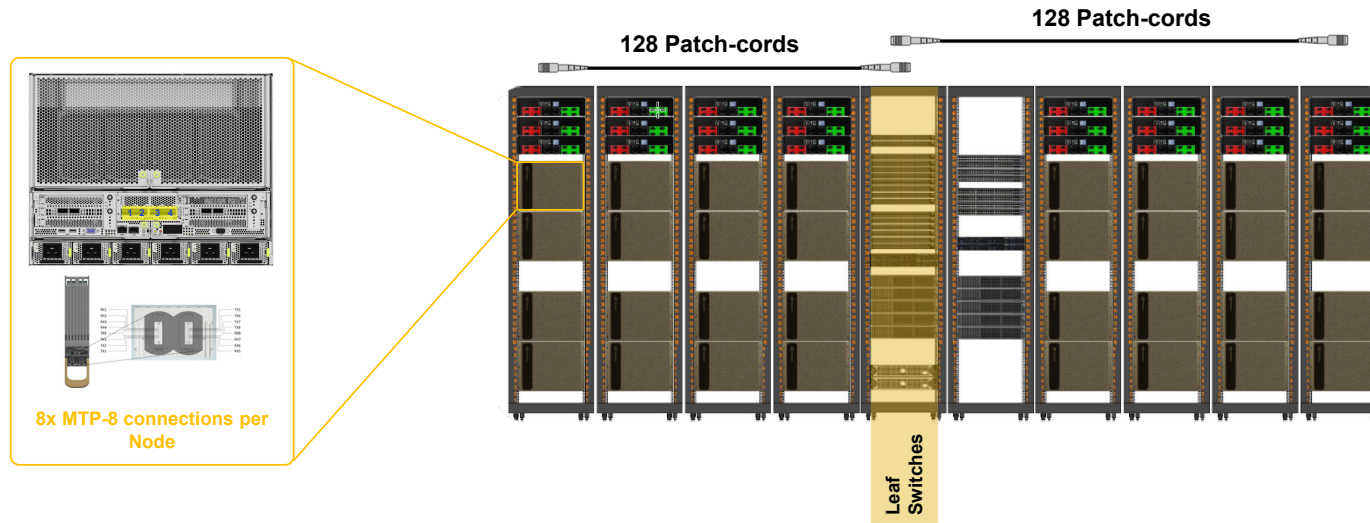Each SU has **256 GPUs**, this 32 DGX H100 in 8 racks



**8x GPUs**

**OSFP**

NVIDIA DGX H100 Server = Compute Node

**1**

**32**

**Quantum (IB Switch)**

## Example of Dedicated Network for AI/ML utilizing NVIDIA InfiniBand

| | SU Count | Node Count | GPU Count | Switch Counts | | | Cable Counts | | |
|---|---|---|---|---|---|---|---|---|---|
| | | | | InfiniBand Leaf | InfiniBand Spine | InfiniBand Core | Node-Leaf | Leaf-Spine | Spine-Core |
| **Small Cluster** | 4 | 128 | 1024 | 32 | 16 | -- | 1024 | 1024 | 1024 |
| | 8 | 256 | 2048 | 64 | 32 | -- | 2048 | 2048 | 2048 |
| **Medium Cluster** | 16 | 512 | 4096 | 128 | 128 | 64 | 4096 | 4096 | 4096 |
| | 32 | 1024 | 8192 | 256 | 256 | 128 | 8192 | 8192 | 8192 |
| **Large Cluster** | 64 | 2048 | 16384 | 512 | 512 | 256 | 16384 | 16384 | 16384 |

CORNING

# Cabling a Scalable Unit (POD)

**Each POD requires 256 MTP-8 connections (8 per H100 Node) to the Leaf Switches**

128 Patch-cords

128 Patch-cords

8x MTP-8 connections per Node

Leaf Switches

**Innovative Solutions in Progress: Partnering with Customers for Future Success**

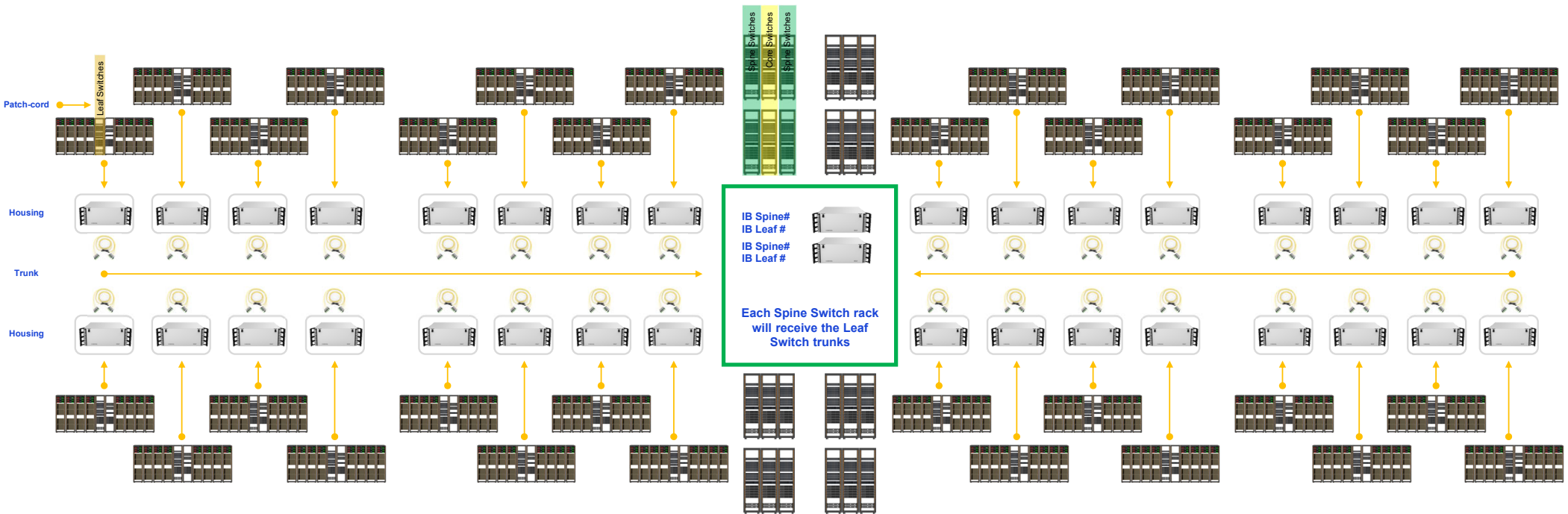**Jumpers (Patch-cords)**
**SMF, MMF**

**Bundle Jumpers**
**SMF, MMF**

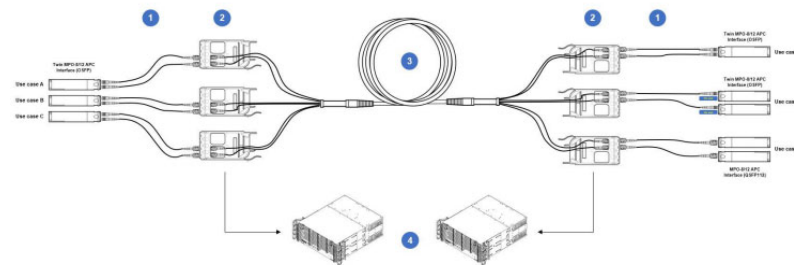**Trunks (w/pulling grip)**
**SMF, MMF**

**EDGE Distribution System**
**SMF**

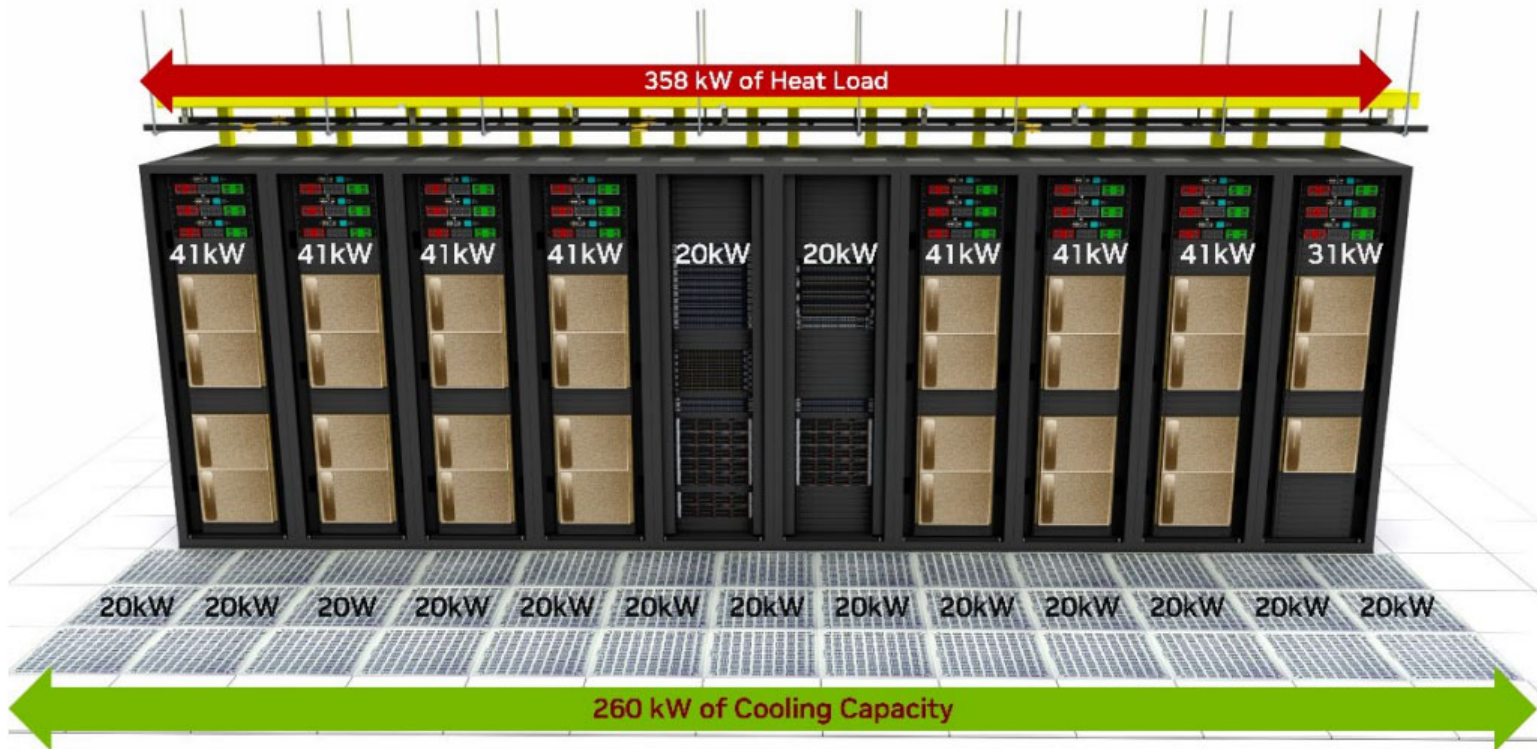CORNING

# Cabling a "Medium Cluster"

**Backbone cabling (512 trunks – 16 trunks per Leaf) will substitute 8,192 individual patch-cords, managing complexity across the data center**



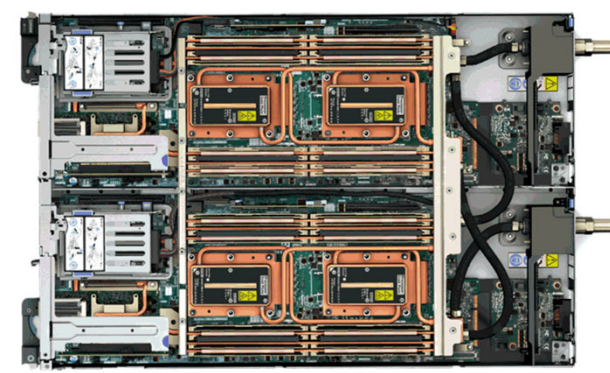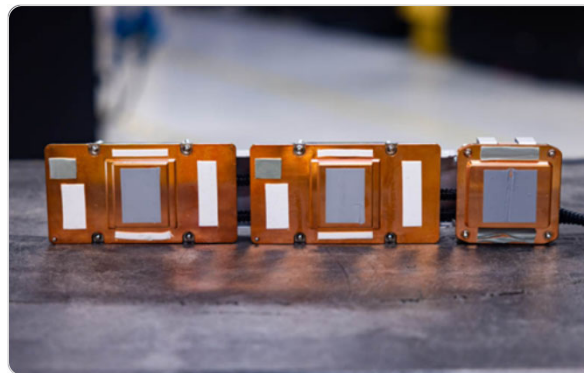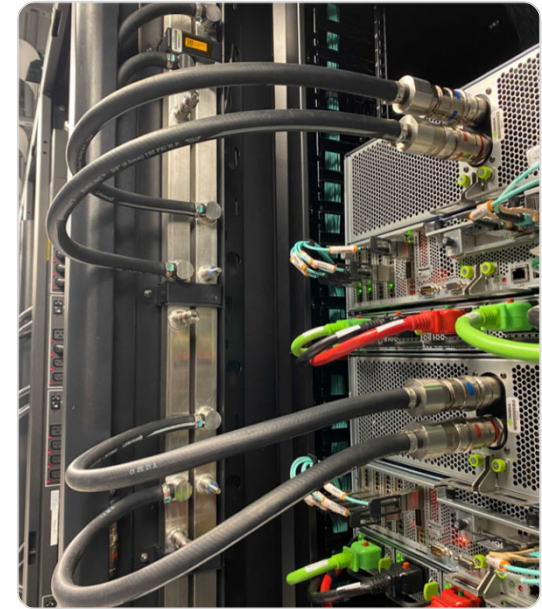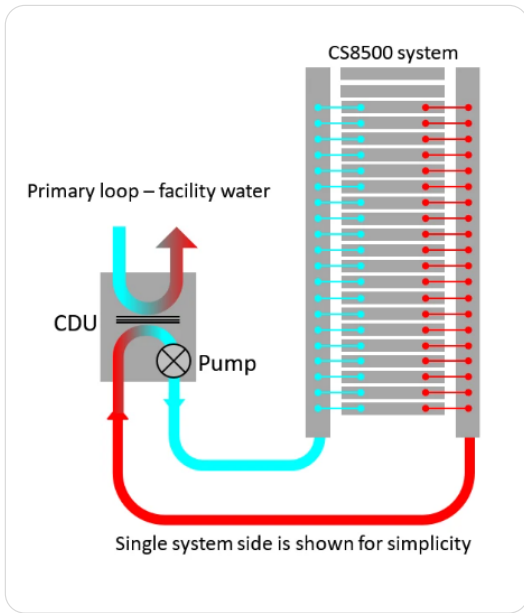| SU Count | Node Count | GPU Count | Switch Counts | | | Cable Counts | | |
|---|---|---|---|---|---|---|---|---|
| | | | InfiniBand Leaf | InfiniBand Spine | InfiniBand Core | Node-Leaf | Leaf-Spine | Spine-Core |
| 4 | 128 | 1024 | 32 | 16 | -- | 1024 | 1024 | 1024 |
| 8 | 256 | 2048 | 64 | 32 | -- | 2048 | 2048 | 2048 |
| 16 | 512 | 4096 | 128 | 128 | 64 | 4096 | 4096 | 4096 |
| 32 | 1024 | 8192 | 256 | 256 | 128 | 8192 | 8192 | 8192 |
| 64 | 2048 | 16384 | 512 | 512 | 256 | 16384 | 16384 | 16384 |

CORNING

# Power Requirement



- ✓ Due to power consumption of server rack (41 kW each) a full row of servers cannot be located in an existing DC design

- ✓ Base design is 256 GPUs/SU (Scalable Unit)

- ✓ There are 8 Production racks in a POD, meaning 8 racks x 32 GPUs

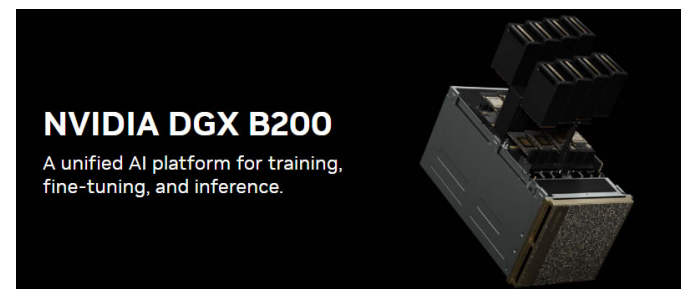# Cabling and Direct Liquid Cooling

CORNING

# Latest NVIDIA's DGX B200 (Blackwell GPU) Architecture

To achieve the most scalability, DGX SuperPOD is powered by several key NVIDIA technologies, including:
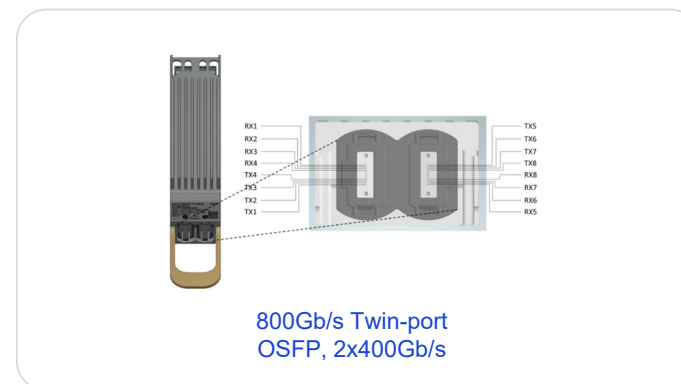
> NVIDIA DGX B200 system—to provide the most powerful computational building block for AI and HPC.

> NVIDIA NDR (400 Gbps) InfiniBand—bringing the highest performance, lowest latency, and most scalable network interconnect.

> NVIDIA NVLink® technology—networking technologies that connect GPUs at the NVLink layer to provide unprecedented performance for most demanding communication patterns.
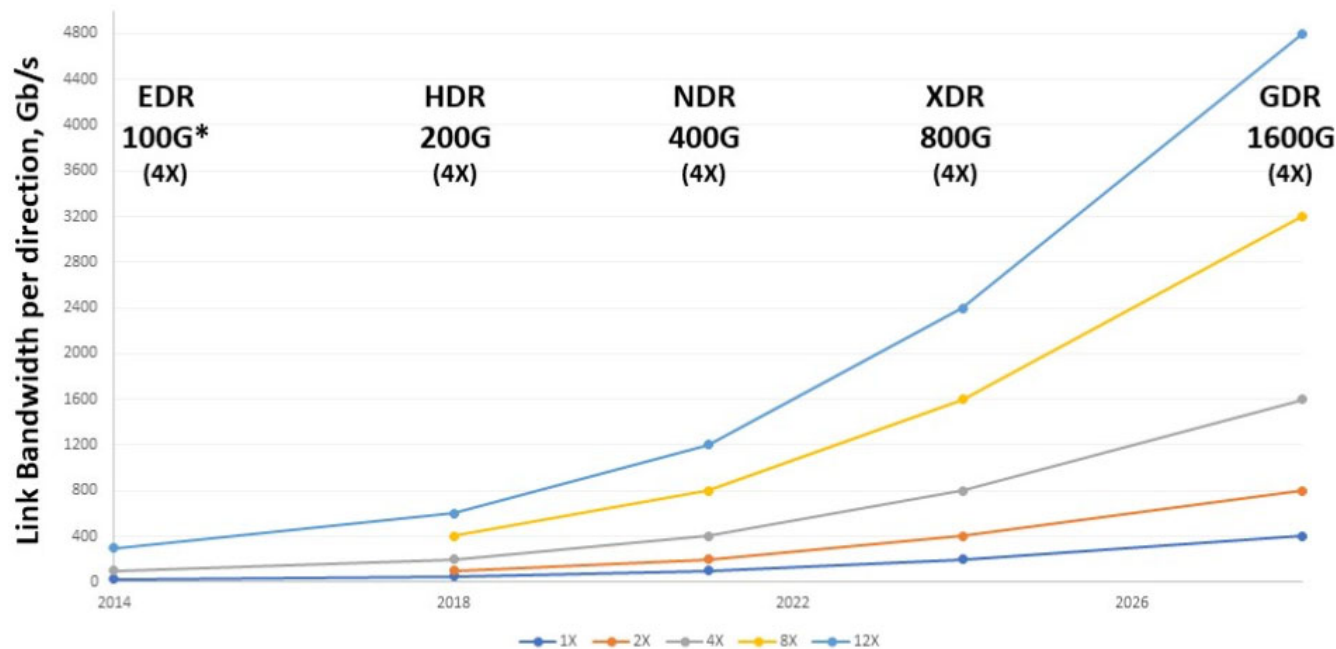


| Component | Technology | Description |
|-----------|-----------|-------------|
| Compute nodes | NVIDIA DGX B200 system with eight B200 GPUs | The world's premier purpose-built AI systems featuring NVIDIA B200 Tensor Core GPUs, fifth-generation NVIDIA NVLink, and fourth-generation NVIDIA NVSwitch™ technologies. |
| Compute fabric | NVIDIA Quantum QM9700 NDR 400 Gbps InfiniBand | Rail-optimized, non-blocking, full fat-tree network with eight NDR400 connections per system |

### NVIDIA DGX B200
A unified AI platform for training, fine-tuning, and inference.

| Feature | Description |
|---------|-------------|
| Form Factor | 10U Rack mount |
| Input (200–240-volt AC) (max) | **14.3 kW** |



800Gb/s Twin-port
OSFP, 2x400Gb/s

# InfiniBand Roadmap



| Full name | 1X (lane) | 4X (lanes) |
|---|---|---|
| Enhanced Data Rate (EDR) | 25G | 100G* |
| High Data Rate (HDR) | 50G | 200G |
| Next Data Rate (NDR) | 100G | 400G |
| Extreme Data Rate (XDR) | 200G | 800G |
| Gigantic Data Rate (GDR) | 400G | 1600G |

*Table 1. Summarized InfiniBand Roadmap*

*100G per Lambda*
*4x lanes = 8 Fibers*

*Link speeds specified in Gb/s at 4X (4 lanes)

*Source: InfiniBand Trade Association*

Current InfiniBand switches utilize **800G OSFP ports**, employing **dual 400G Next Data Rate (NDR) ports**. This configuration uses **8 fibers per port**, resulting in **64x400G ports per switch**. It's highly likely that the forthcoming generation of switches, whatever name they carry, will adopt **Extreme Data Rate (XDR) speeds**. This translates to **64x800G** ports per switch, also utilizing **8 fibers per port** – mostly **single mode fiber**. This 4-lane (8-fiber) pattern seems to be a recurring motif in the InfiniBand roadmap, summarized in Table-1, utilizing even faster speeds in the future.

CORNING

# AI/ML and Structured Cabling

## AI Cluster Investment



Legend: ■ AI Server  ■ Networking/Switches  ■ Structured Cabling

- More than **95% of the network cost** for the AI Server cluster is related to the **active gear**

- Power and Cooling investments will also be material

- **Fiber connectivity** is **facilitating** the networking **speed** and processing scalability.

- Fiber requirements for High Performance compute and AI networks is 5 times more than the traditional data center production network

- **Point-to-Point** bundled jumper assemblies can accommodate **smaller clusters, larger** scale-outs need **structured cabling** solutions

- **Roadmaps for Ethernet and InfiniBand transceivers** will scale with Base-8 fiber backbones

- Elements of the structured cabling system (Passive TAPs, Port-Breakout) will enable data center operators to **gain more value from the fiber infrastructure**

CORNING

# Planning for Migration



- **The path to higher speeds** will always depend on your unique needs.

- You may be happy with 40G now but planning to **upgrade to 100G** four years from now. Or maybe you are working with 400G and have your **eyes set on 800G** in five years: Migration will always vary based on your timeline and the available technologies in the market.

- But in most cases, **Base-8 will provide the ideal level of flexiblity to meet your needs throughout your transition**.

- Corning's **EDGE8 structured cabling solutions** will support your transition needs, doesn't matter if we talk about **Ethernet or InfiniBand**

CORNING

**Connect with us:**

@   lazarr@corning.com

in   linkedin.com/in/RomeoLazar

in   Corning Optical Communications

▶   Corning Optical Communications

𝕏   @CorningOpComm

CORNING

# CORNING