![AsteraLabs logo]

# Scaling GPU Clusters & Low Latency Memory Fabrics With Active PCIe / CXL Cabling

Chris Blackburn

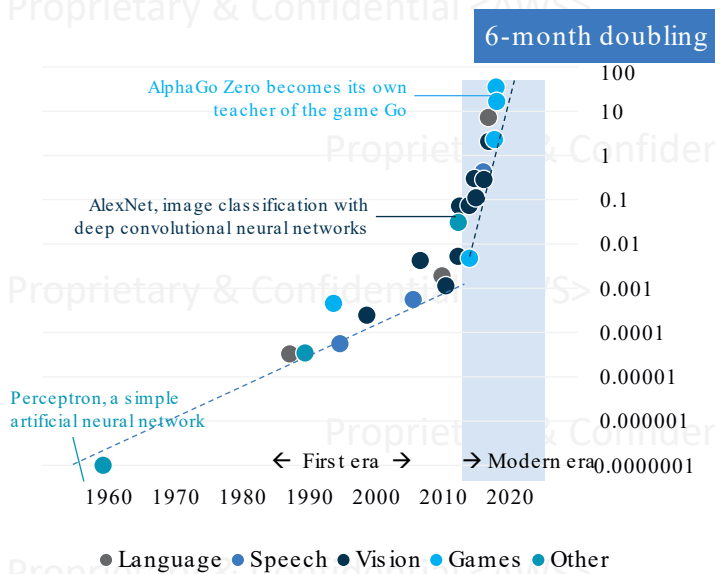System Architect & Director of Field Applications Engineering
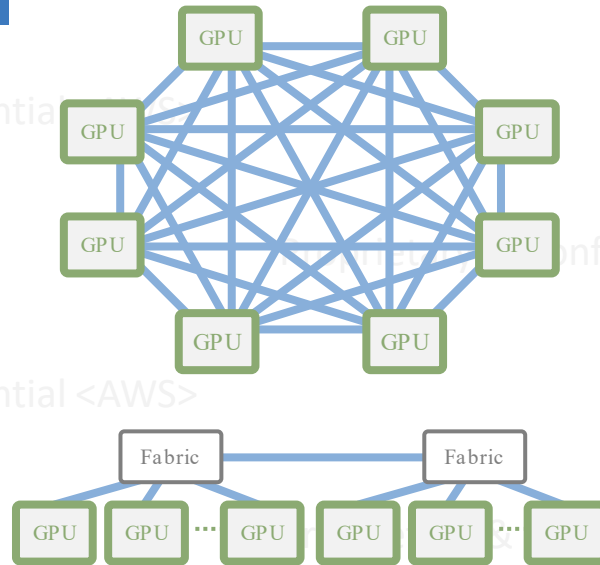
August 2024

**FMS** *the Future of Memory and Storage*

# AI Infrastructure Scale Challenges

## AI Models Continue to Expand

6-month doubling

AlphaGo Zero becomes its own teacher of the game Go

AlexNet, image classification with deep convolutional neural networks

Perceptron, a simple artificial neural network

← First era →   → Modern era

100
10
1
0.1
0.01
0.001
0.0001
0.00001
0.000001
0.0000001

1960 1970 1980 1990 2000 2010 2020

● Language ● Speech ● Vision ● Games ● Other

Model sizes have doubled after 6 months*

## Larger GPU Clusters Required

GPU GPU GPU GPU GPU GPU GPU GPU
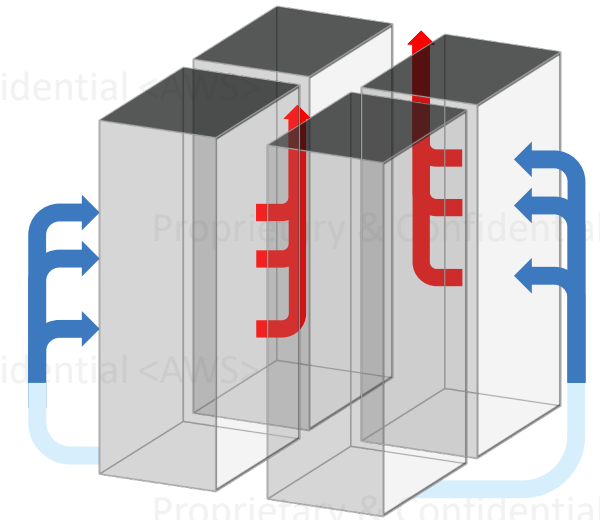
Fabric    Fabric

GPU GPU ... GPU GPU GPU ... GPU

Scale up fabrics connect hundreds of GPUs

## Limited Power per Rack

AI servers consume 8X more power than CPU servers**

## Thermal Constraints

GPUs transitioning from air to liquid cooling***

## AI infrastructure under heavy pressure to scale clusters across several racks

# Emerging Application: Multi Rack AI Fabric

**Current Generation**

Up to 3m

**Next Generation**

Up to 7m

700W GPUs

700W GPUs

1400W GPUs

1400W GPUs

1400W GPUs

1400W GPUs

**Challenges**
- Rack power density
- Rack thermal density
- Increasing number of high-bandwidth links
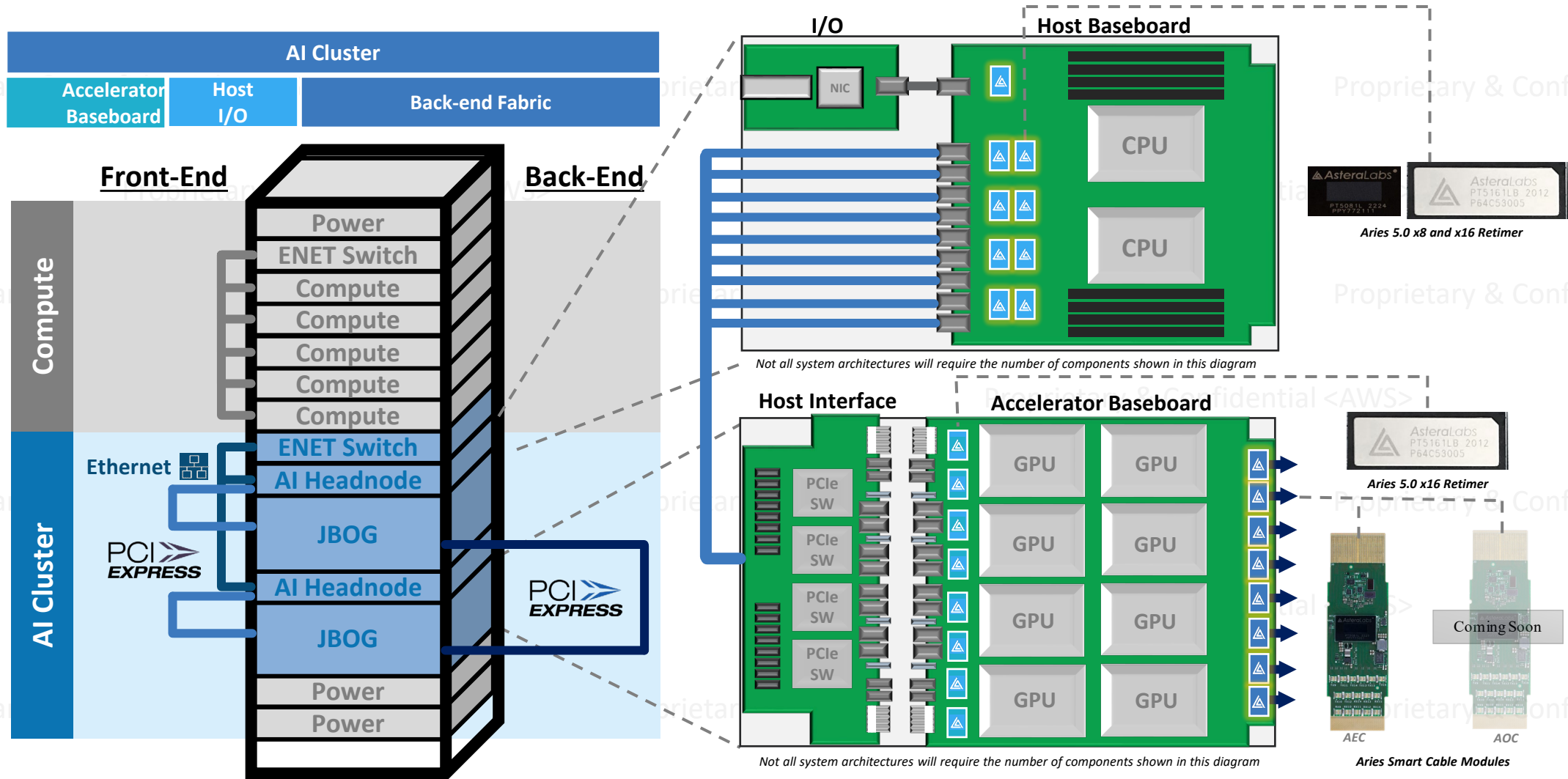
**Advantages**
- Accommodates higher-performant GPUs
- Allows for larger clusters while optimizing rack power/thermals

# Broadening PCIe Connectivity for the Era of AI



Not all system architectures will require the number of components shown in this diagram

Not all system architectures will require the number of components shown in this diagram

*Aries 5.0 x8 and x16 Retimer*

*Aries 5.0 x16 Retimer*

*Aries Smart Cable Modules*

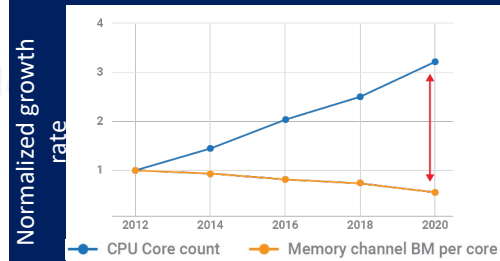# Memory Bottlenecks Due to AI / ML Workloads

**AI Model Complexity Doubling Every ~6 Months**



(Source: Open AI)

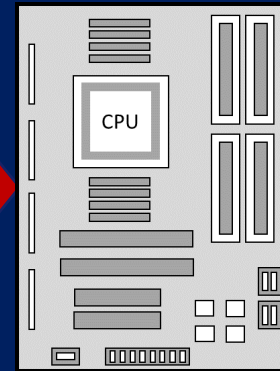**Workloads need higher CPU efficiency & memory expansion**

**Memory Bandwidth Per Core is Declining**



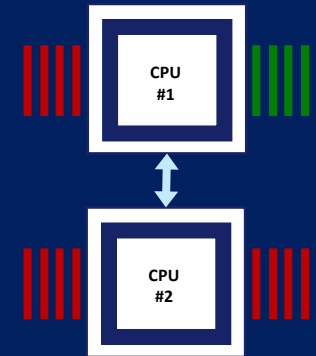(Source: Meta OCP Presentation, Nov '21)

**CPU efficiency is declining due to declining memory bandwidth per core**

**Server CPU Package & Thermal Constraints Limit Memory Channels**



**Memory bottlenecks caused by CPU Pin and thermal constraints**
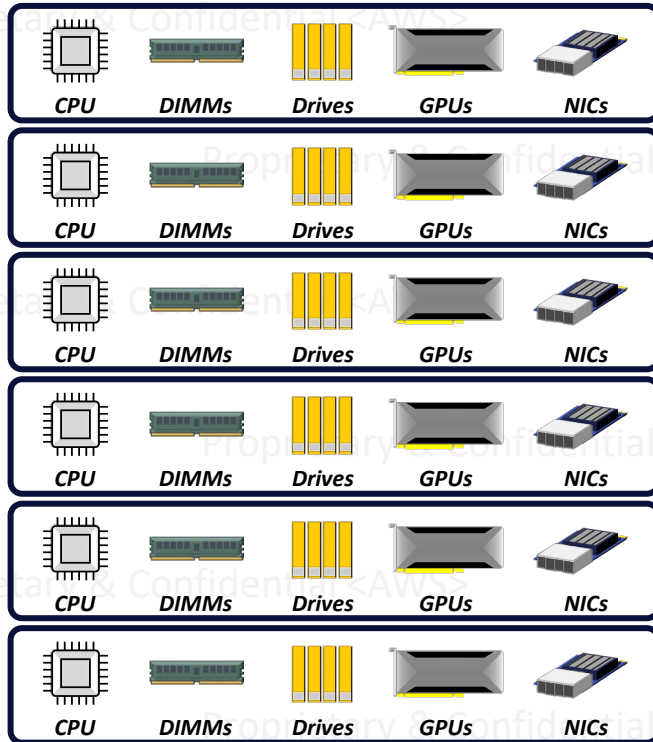
**Memory Capacity is Tied to Compute Nodes**



**Memory is stranded behind compute leading to over provisioning**

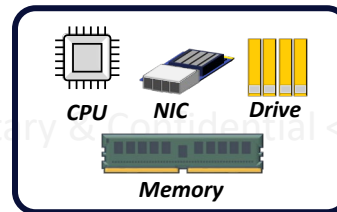# Emerging Application: Heterogeneous Infrastructure

## Converged Infrastructure



CPU — DIMMs — Drives — GPUs — NICs

| Static Configuration | Adaptive Environment |
|---|---|
| Stranded Resources | Efficient Performance |
| Rigid System Design | Composable Resources |
| Fixed Hosting Costs | Flexible Cost Model |

### Workload 1



CPU   NIC   Drive
Memory

### Workload 2



CPU   NIC
Memory   GPU

| OPEX Challenges | OPEX Advantages |
|---|---|
| High PUE | Low PUE |
| Thermal Management | Focused Cooling Zones |
| SW Performance Tuning | Bare-Metal Performance |

## Disaggregated/Composable Infrastructure



CPU   CPU   CPU   CPU   CPU   CPU

PCIe/CXL Switch

NIC   NIC   NIC   NIC   NIC   NIC

Drive   Drive   Drive   Drive   Drive   Drive   Drive

Memory   Memory   Memory

GPU   GPU   GPU   GPU

PCI EXPRESS   CXL Compute Express Link   **Box-to-box cabling**

# Broadening PCIe Connectivity for the Era of Compute

# External Cabling Reach Considerations

**PCIe/CXL Retimer**

## Within-the-rack
E.g., AI Server headnode to JBOG
*3m reach requirement*

**Active Riser Card**

**PCIe Passive DAC with Aries active riser cards**

**Active Riser Card**

Cable Lengths
**Up to 3m**

Full PCIe channel budget

Full PCIe channel budget
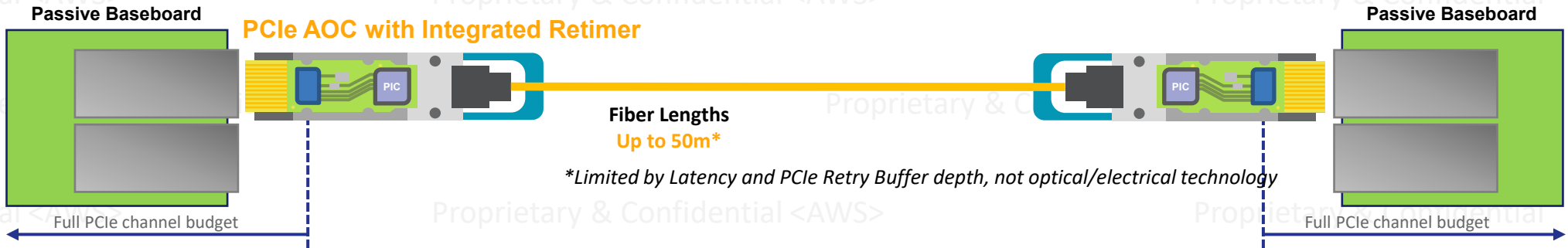
Full PCIe channel budget

## Across-racks
E.g., JBOG to JBOG
*5-7m reach requirement*

**Passive Baseboard**

**PCIe AEC with Integrated Retimer**

**Passive Baseboard**

Cable Lengths
**Up to 7m**

Full PCIe channel budget

Full PCIe channel budget

Full PCIe channel budget

## Across-rows
E.g., Switch to Switch (future)
*20-50m reach requirement*

**Passive Baseboard**

**PCIe AOC with Integrated Retimer**

**Passive Baseboard**

PIC

PIC

Fiber Lengths
**Up to 50m***

*\*Limited by Latency and PCIe Retry Buffer depth, not optical/electrical technology*

Full PCIe channel budget

Full PCIe channel budget

# PCIe/CXL AECs: Handling PCIe Side-Band Signals

**Three "required" side-band signals defined in PCI-SIG's Card Electomechanical (CEM) Specification:**

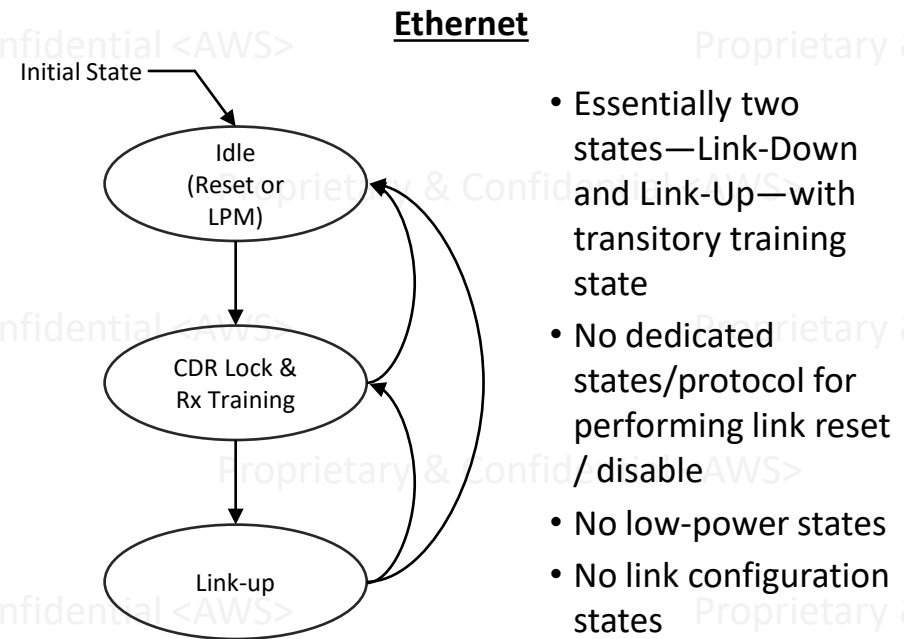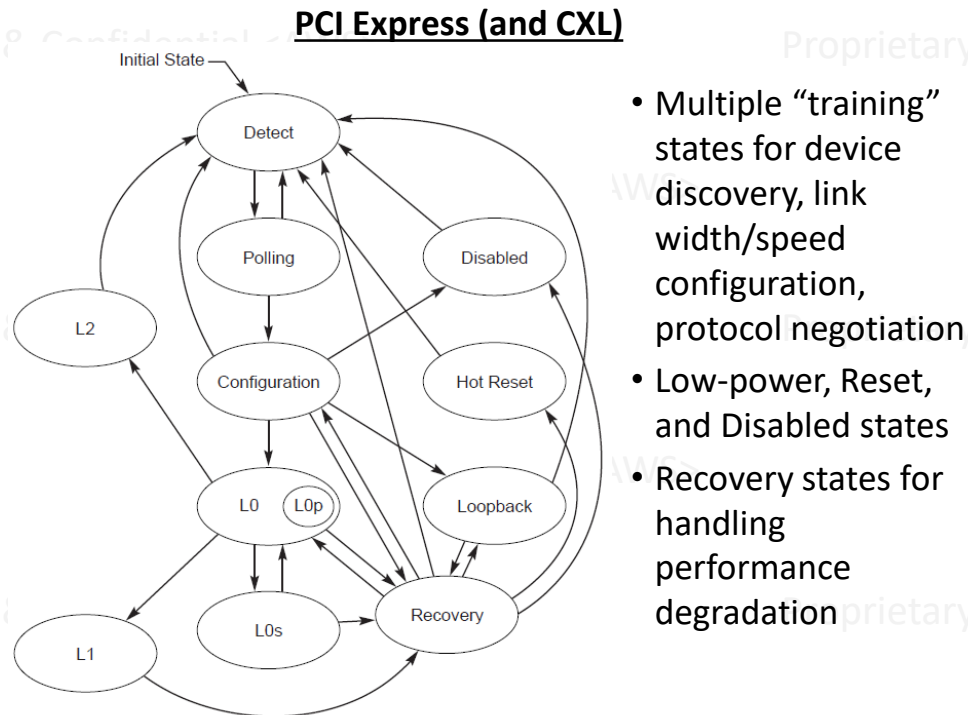| PCIe Side-Band Signal | Description | Option for handling within an AEC | Alternative |
|---|---|---|---|
| REFCLK | 100 MHz HCSL clock with or without spread-spectrum modulation | Dedicated differential pair to carry REFCLK from one side to the other.<br>**Pros**: Allows for common clock topologies<br>**Cons**: Extra cable cost, "asymmetric" cable design | No REFCLK transport in cable: SRNS/SRIS.<br>**Pros**: lower cost, "symmetric" cable; scalable to multi-link AECs<br>**Cons**: CC topology requires dedicated side-band cable between systems |
| PERST# | PCIe Protocol Reset | Dedicated single-ended line to carry PERST#.<br>**Pros**: Allows PERST# synchronization on a per-link basis<br>**Cons**: Extra cable cost, "asymmetric" cable design | No PERST# transport in cable. PCIe Reset events are handled through in-band Hot Reset, host-coordinated local reset, side-band management, and/or Hot Plug support.<br>**Pros**: Lower cost, "symmetric cable"; scalable to multi-link AECs<br>**Cons**: No dedicated per-link PERST# |
| PRSNT# | Cable (cable) present indicator | Pluggable cable MSAs (OSFP, OSFP-XD, etc.) include ModPrsL functionality already | N/A |

# AECs: PCIe VS. Ethernet

- Two main differences:

**Protocol complexity**: PCIe's backwards compatibility and link training requirements make AECs more complex for PCIe compared to Ethernet

**Interoperability**: The variety of device types and ecosystem players is significantly more for PCIe compared to Ethernet

**PCI Express (and CXL)**



**Ethernet**



- Multiple "training" states for device discovery, link width/speed configuration, protocol negotiation
- Low-power, Reset, and Disabled states
- Recovery states for handling performance degradation

- Essentially two states—Link-Down and Link-Up—with transitory training state
- No dedicated states/protocol for performing link reset / disable
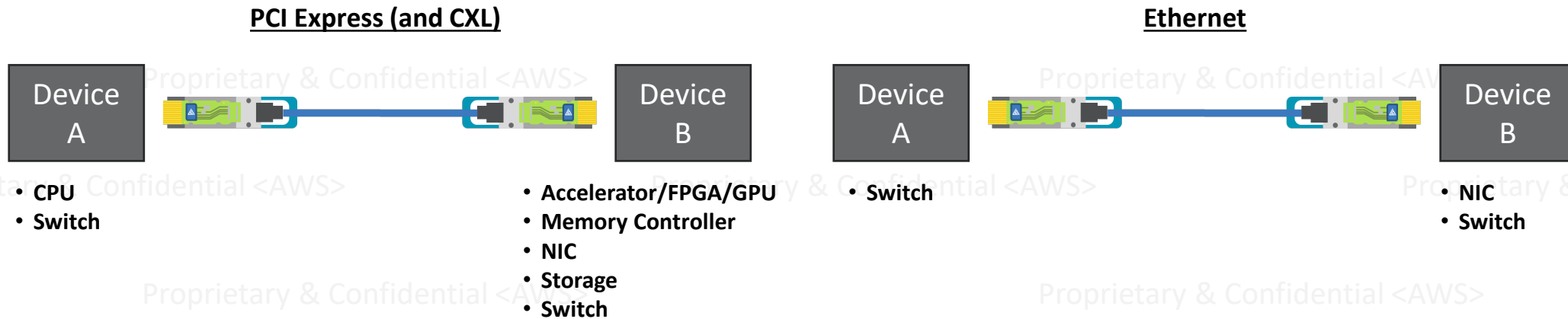- No low-power states
- No link configuration states

# AECs: PCIe VS. Ethernet

- Two main differences:

  **Protocol complexity**: PCIe's backwards compatibility and link training requirements make AECs more complex for PCIe compared to Ethernet

  **Interoperability**: The variety of device types and ecosystem players is significantly more for PCIe compared to Ethernet

**PCI Express (and CXL)**

**Ethernet**

| Device A | Device B |
|---|---|
| • **CPU**<br>• **Switch** | • **Accelerator/FPGA/GPU**<br>• **Memory Controller**<br>• **NIC**<br>• **Storage**<br>• **Switch** |

| Device A | Device B |
|---|---|
| • **Switch** | • **NIC**<br>• **Switch** |

# PCIe Cabling Form Factor Comparison

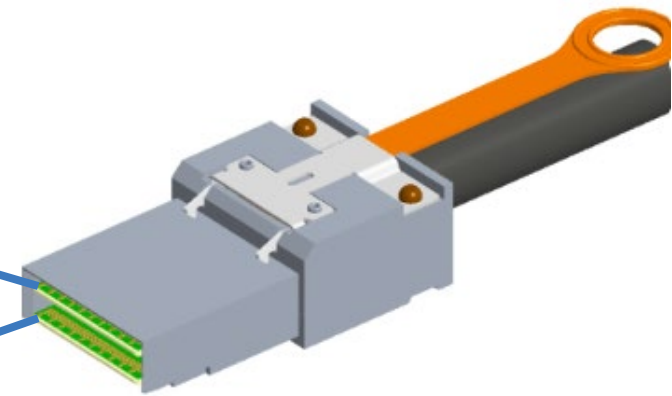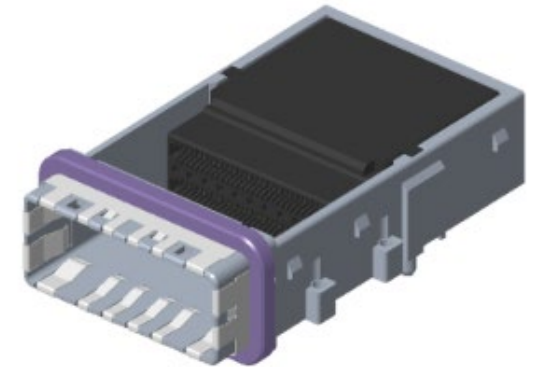| | OSFP-XD Under consideration for Optical | CDFP (x16) CopprLink | OSFP Under consideration for Optical | QSFP-DD | QSFP |
|---|---|---|---|---|---|
| High-speed lane count (full duplex) | 16 | 16 | 8 | 8 | 4 |
| X-Y PCB Size (normalized to x16) | 2292 mm² | 1460 mm² | 3989 mm² | 2472 mm² | 3933 mm² |
| Connector Contact Pitch | 0.60 mm | 0.75 mm | 0.60 mm | 0.80 mm | 0.80 mm |
| Cable Gauge Supported | 26-32 AWG | 28-32 AWG | 26-32 AWG | 27-32 AWG | 26-32 AWG |
| 32 GT/s Max DAC reach (at max gauge) | 4 m | 3.0 m | 4 m | 3.5 m | 4 m |
| 32 GT/s Max AEC reach (at max gauge) | 7 m | 5.5 m | 7 m | 6 m | 7 m |
| 64 GT/s Max DAC reach (at max gauge) | 3 m | 2.5 m | 3 m | 2.5 m | 3 m |
| 64 GT/s Max AEC reach (at max gauge) | 6 m | 5 m | 6 m | 5 m | 6 m |
| Active Copper cable | Yes | No | Yes | Yes | Yes |
| Active Optical cable | Yes | No | Yes | Yes | Yes |
| Power Capability per Lane | 8x2.5A@3.3V 66W/16 = 4.125W | 1x1.5A@12V + 1x1.5A@3.3V 23W/16= 1.44W | 4x2.5A@3.3V 33W/8= 4.125W | 6x1.5A@3.3V 30W/8= 3.75W | 3x1.0A@3.3V 10W/4= 2.5W |

**Assumptions**:
- Twinax losses: 28/27/26AWG=4.3/4.0/3.6 + 10% dB/m at 16 GHz.
- AEC: Retimer silicon to cable pads: 4 dB @ 16 GHz
- DAC: Retimer silicon (behind cage) to passive DAC cable pads: 9.5 dB @ 16 GHz
**Reference**: https://drive.google.com/file/d/12ZSTkIgkzESbf4fZzj7WQBi3y4oto-gB/view

# CopprLink and Active External Cables

- SFF-TA-1032 (CDFP) uses two physical paddle cards inside a cable assembly

- This presents a significant challenge: **How can you connect Tx and Rx signals from separate paddle cards into a Retimer component?**

- Rx and Tx both terminating in the Retimer is necessary for:
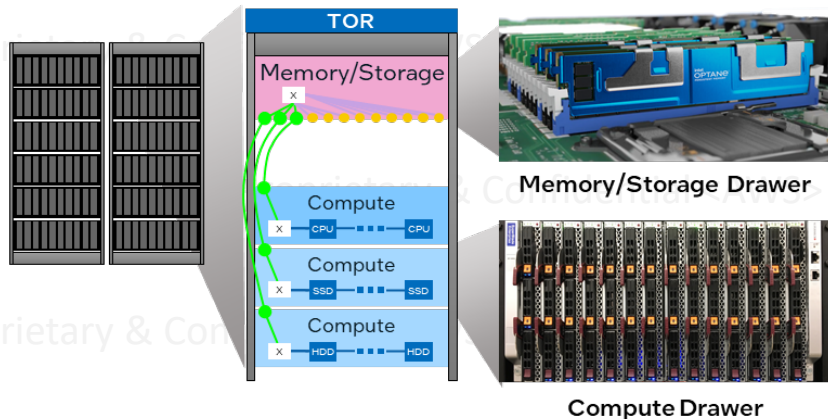  Equalization Phase 2/3 training
  In-band lane margining

Upper paddle card pinout



Lower paddle card pinout

# Wrap Up



Memory/Storage Drawer

Compute Drawer

PCIe Passive DAC with Aries active riser cards

PCIe AEC with Integrated Retimer

PCIe AOC with Integrated Retimer

- Evolving AI and disaggregated compute system topologies require more **external cabling**

- Reach requirements vary from **2m** (within the rack), to **7m** (rack to rack), and **beyond** (larger clusters)

- Retimer-based AEC and optical solutions enable reach extension while presenting an easy-to-design-to **PCIe compliance point** to the host/device

- Implementing PCIe AEC and optical involves higher design complexity in terms of **protocol and interoperability** as compared to Ethernet

- **OSFP-XD/OSFP** represents an attractive option for PCIe/CXL x16/x8 applications, allowing for passive DAC, AEC, and Optical solutions