**facebook**
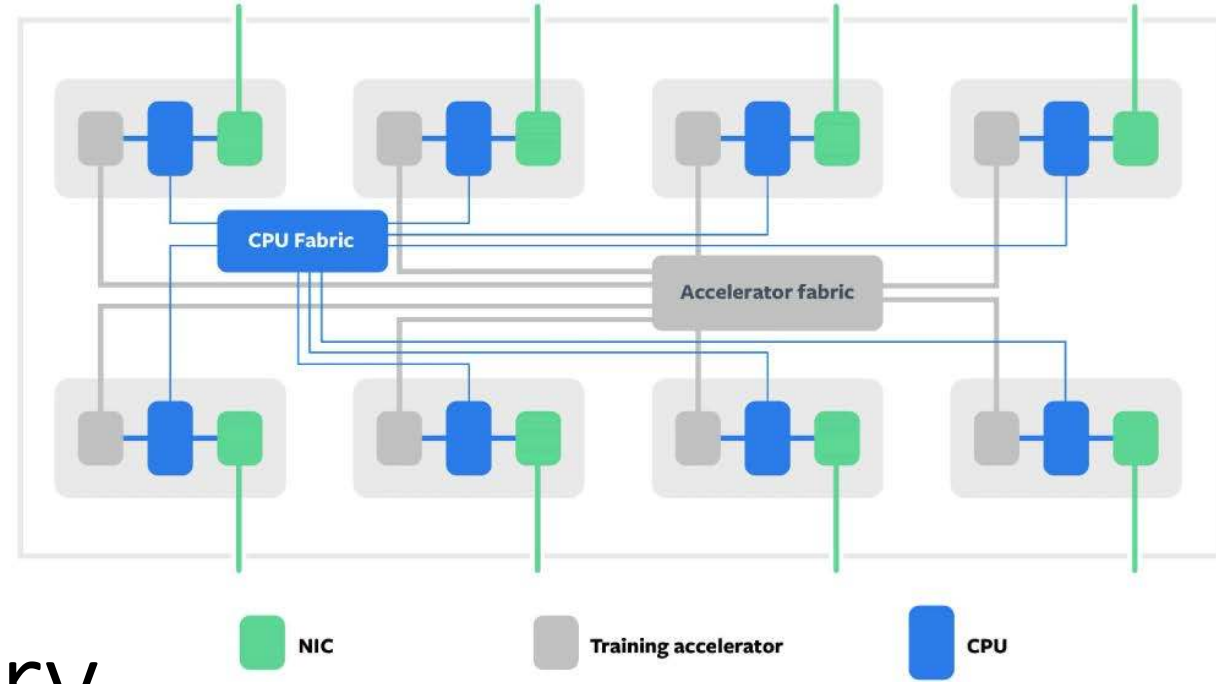
# Zion: Facebook Next-Generation Large Memory Training Platform

**Misha Smelyanskiy**

Hot Chips 31, August 19, 2019

# The Growth of ML at Facebook

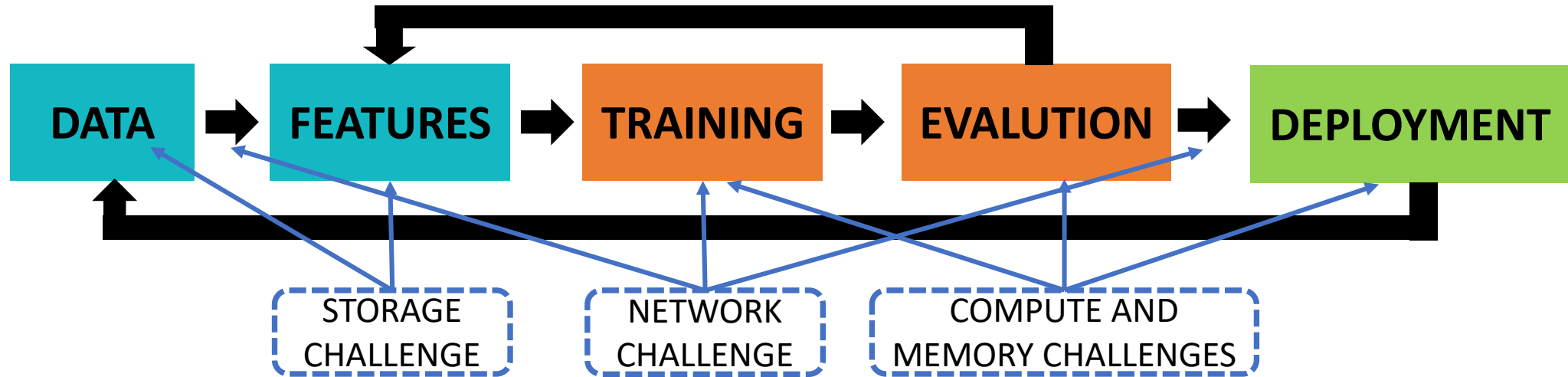DATA → FEATURES → TRAINING → EVALUTION → DEPLOYMENT

- **ML pipeline data growth**
  - Usage in **2018**: **30%**
  - Usage **today**: **50%**
  - ML data growth in **one year**: **3X**

- **12-month ML Training growth**
  - # of ranking engineers: **2X**
  - Workflows trained: **3X**
  - Compute consumed: **3X**

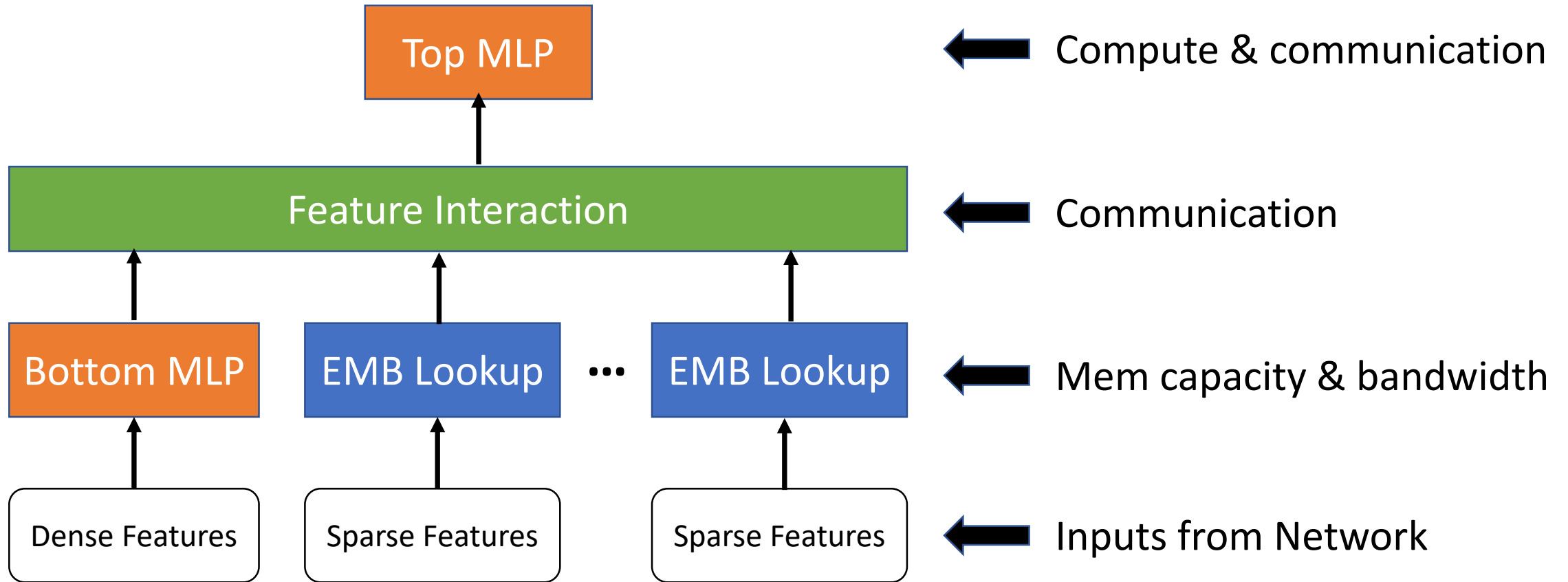facebook

# Training Infrastructure Challenges



- Strains memory, compute, storage, and network
- ML engineers expect developer efficiency and flexibility
- **Motivated SW/HW co-design of training platform**

facebook

# Major AI Services @ Facebook

- Ranking and recommendation
  - news feed, and search

- Computer vision
  - image classification, object detection, and video understanding

- Language
  - translation, content understanding

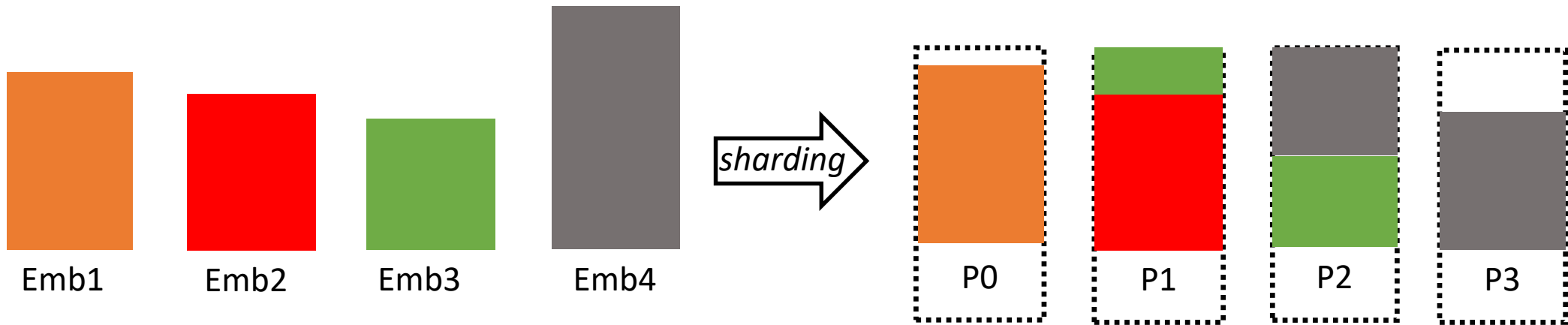- Recommendation models are among most important models

**facebook**

# Deep Learning Recommendation Models



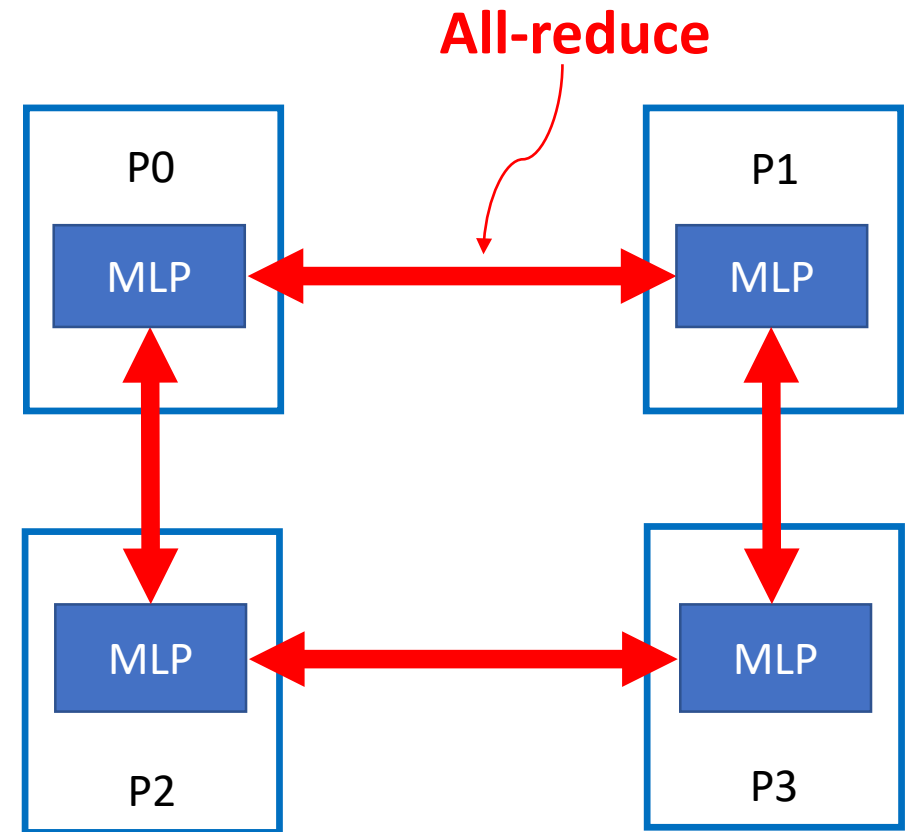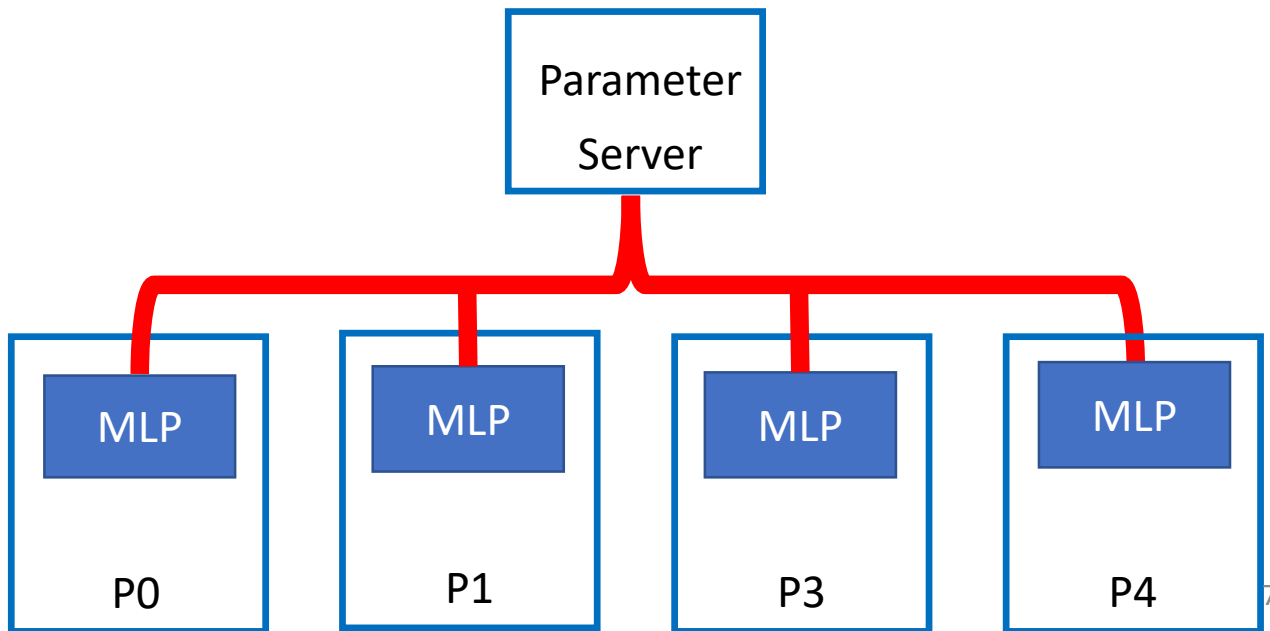- [Open-sourced](#) as a deep learning recommendation model benchmark

# Training Embedding Tables

- Very large embedding tables – O(10+) GBYTES

- Low arithmetic intensity, irregular memory accesses

- Model Parallelism
  - Map embedding tables to different compute devices
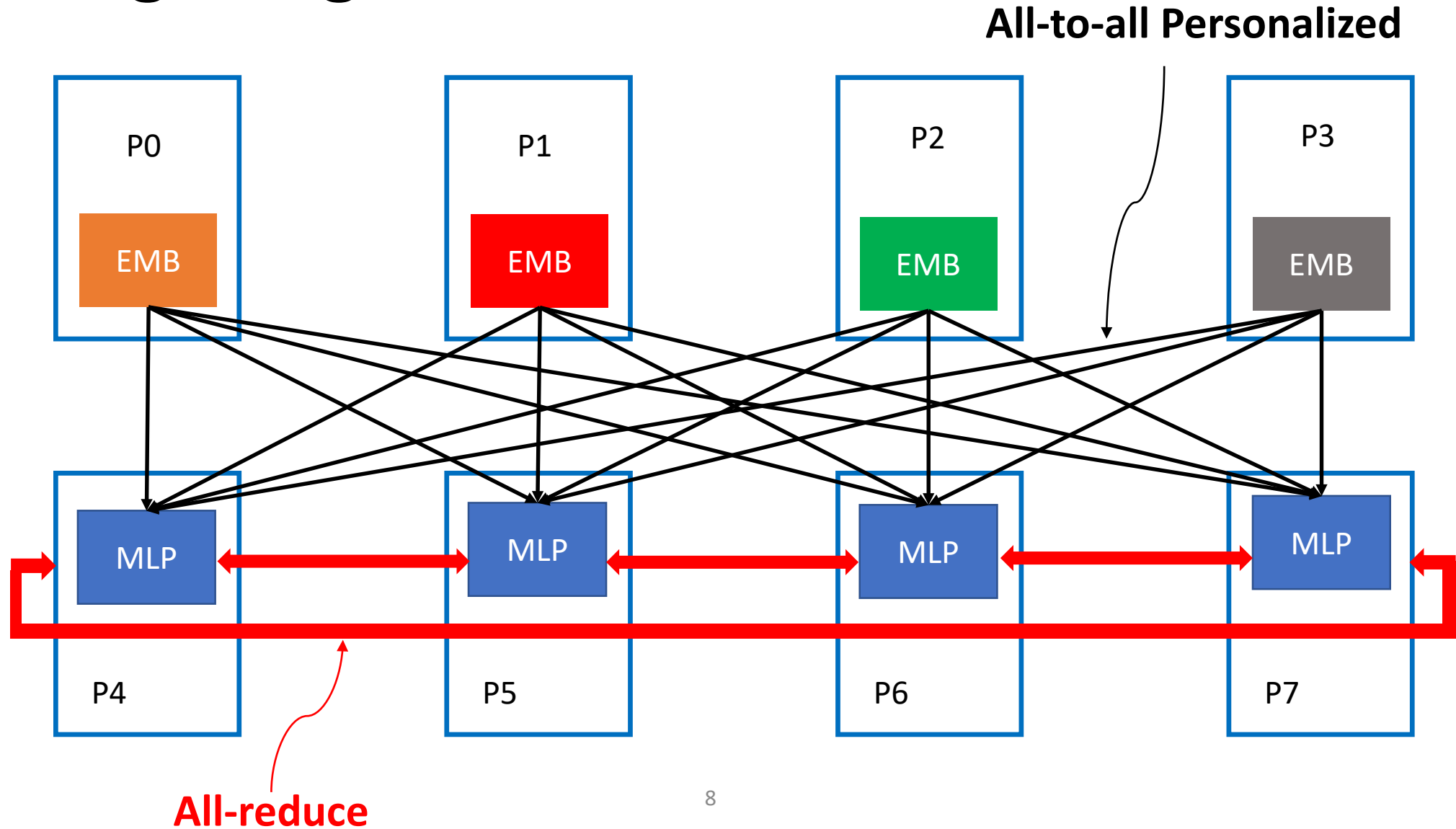  - Shard to balance out utilization given memory constraints

Emb1  Emb2  Emb3  Emb4  *sharding*  P0  P1  P2  P3

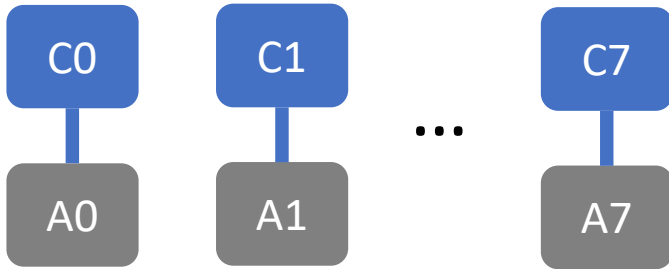facebook

# Training MLP

- Parallelism: model or data
- Updates: <u>asynchronous</u> or <u>synchronous</u> (via all-reduce)
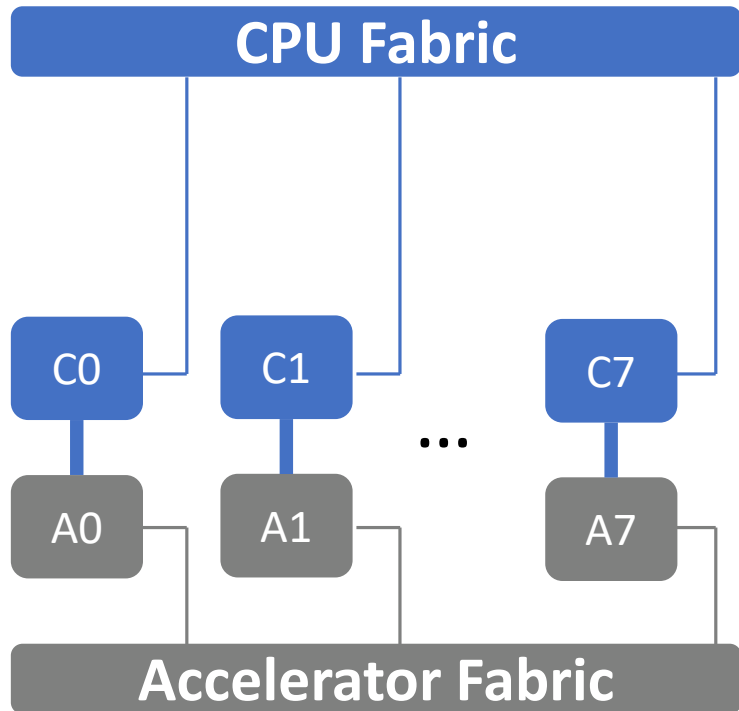- Dense regular compute, tall-skinny GEMMs

# Putting it together



**All-to-all Personalized**

P0    EMB

P1    EMB

P2    EMB

P3    EMB

MLP    MLP    MLP    MLP

P4    P5    P6    P7

**All-reduce**

facebook

8

# Zion: MLP and Embedding Support

C0 — A0  C1 — A1  ...  C7 — A7

| | CPU | Accelerator |
|---|---|---|
| # of devices | 8 | 8 |
| Total BF16 Compute (TFLOPS) | O(1) | O(10) |
| Power per device | ~100w | ~200w |

| | CPU | Accelerator |
|---|---|---|
| Mem Type | DDR4 | HBM2 |
| Total Capacity (GBYTES) | O(1000) | O(100) |
| Total BW  (TB/s) | O(1) | O(10) |

- Unified BFLOAT16 format with CPU and accelerators
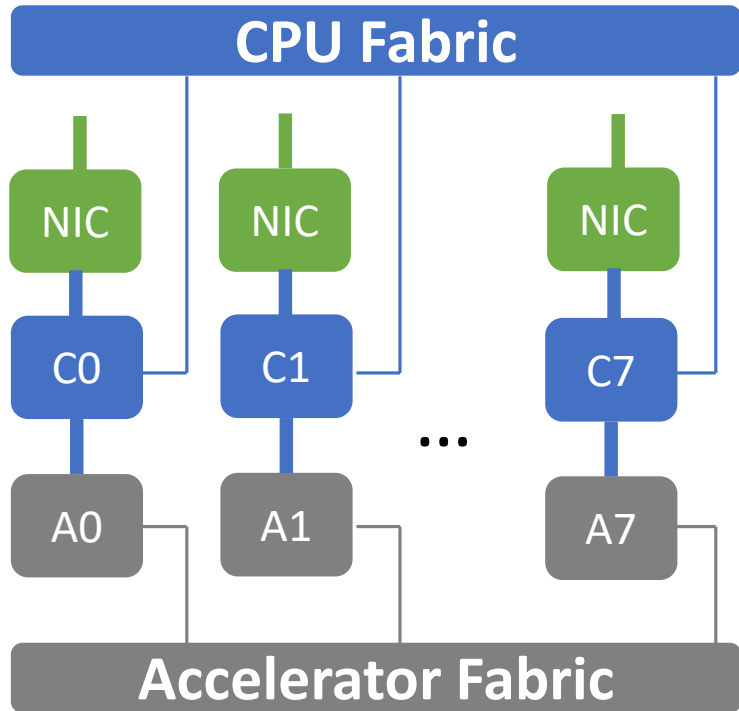- High capacity, low BW DDR;  low capacity, high BW HBM

facebook

# Zion: Communication Support



| | CPU | Accelerator |
|---|---|---|
| *Fabric Type* | cache-coherent UPI | vendor |
| *Fabric Topology* | Twisted Hypercube | varies |
| *Total BW (TB/s)* | O(1) | O(1) |

- Supports all-reduce and all-to-all
- Twisted hypercube has lower diameter than hypercube
- Use non-temporal stores on CPU to reduce coherent traffic

facebook

# Zion: Scaling Out



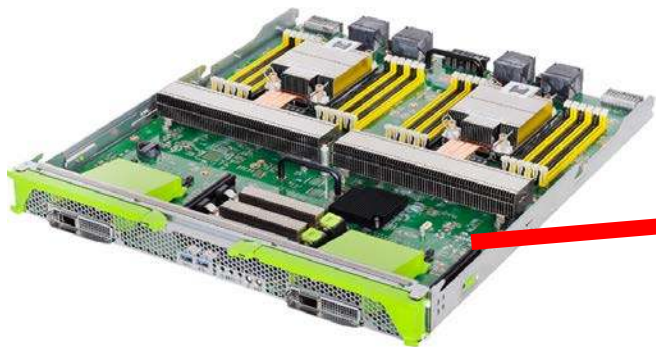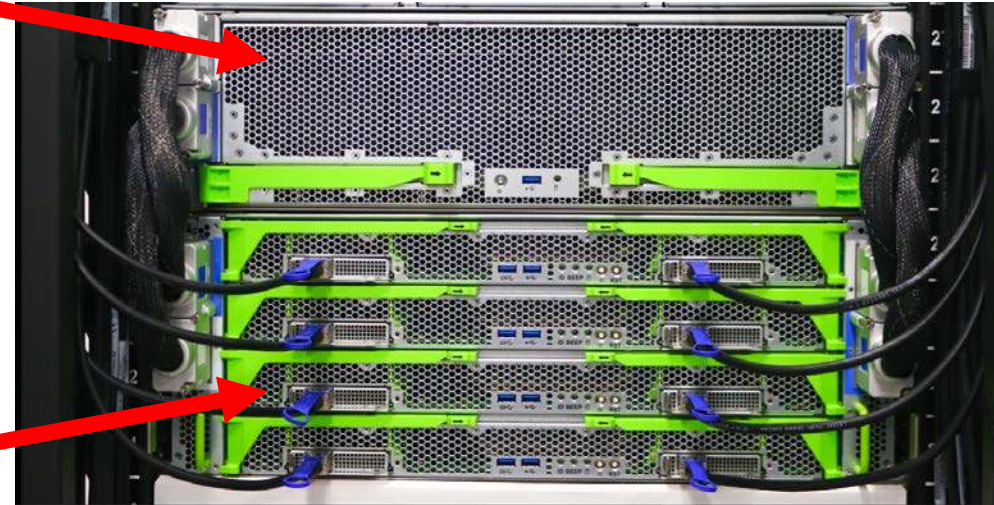| | CPU | Accelerator |
|---|---|---|
| NIC (Gbps) | 8 x 100 | n/a |
| PCIe (Gen3 or 4) | X16 | n/a |

- Via host NIC, P2P, RDMA, PCI-SWITCH

facebook

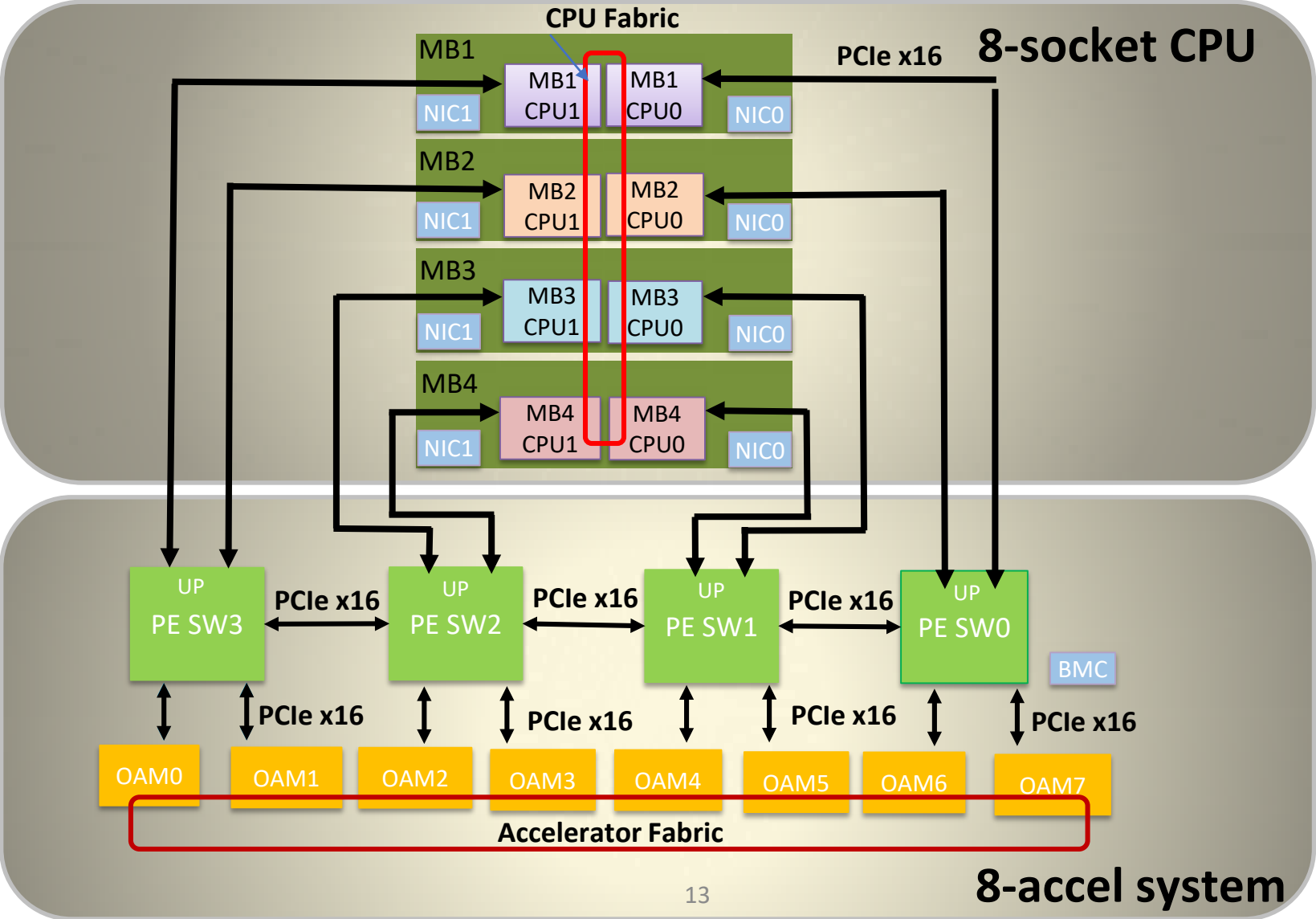# Modular Physical Design

OCP Accelerator Module

8-Accelerator System

**Zion System**

Dual-socket MB Module

**8-Socket System**

facebook
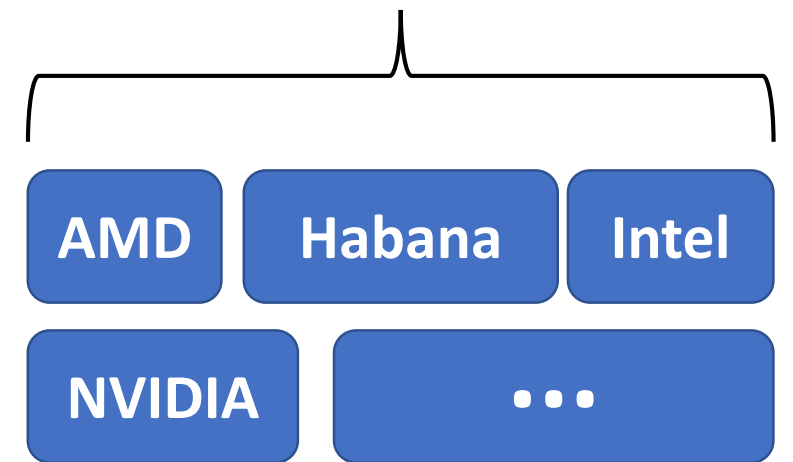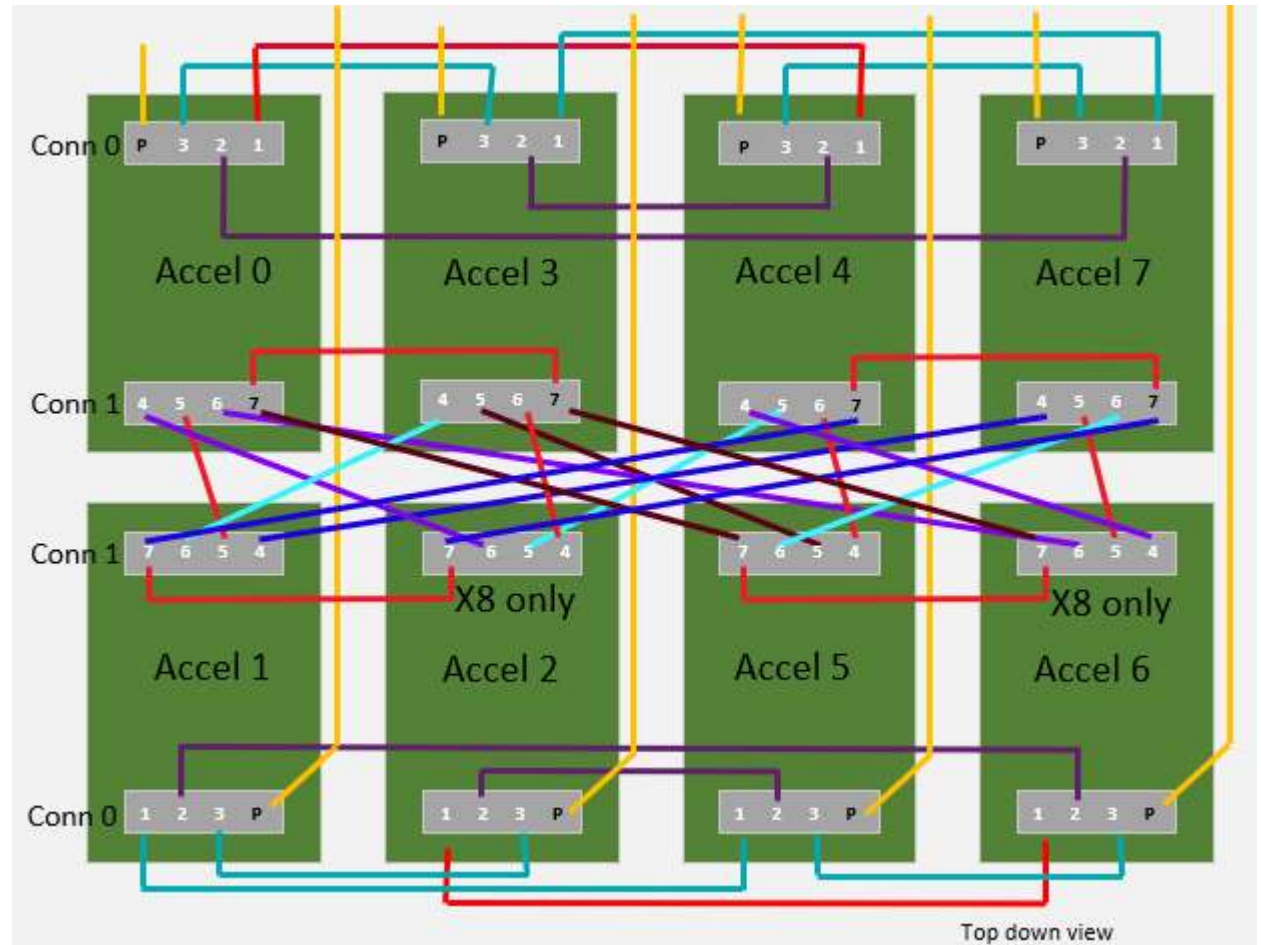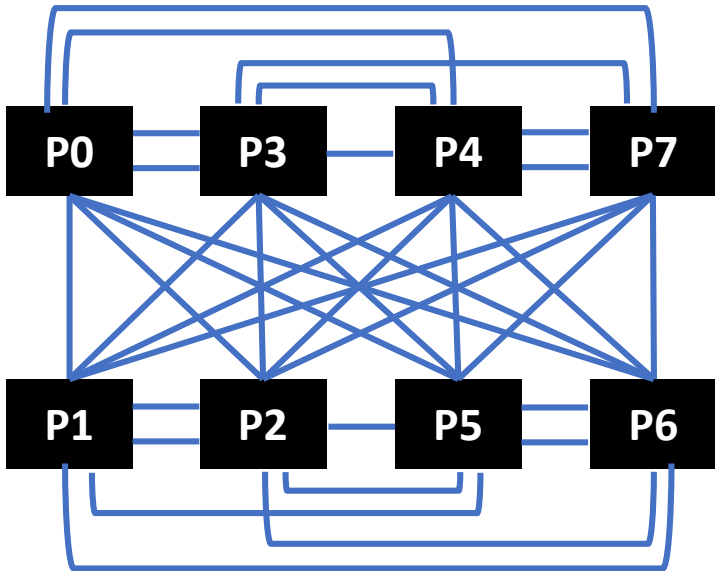
# Platform Architecture

# OCP Accelerator Module (OAM)

- Challenge: which accelerator do we use?
    - Very large number of accelerators
    - Limited resource to enable multiple systems

- Solution: OCP Accelerator Module(OAM)
    - Facebook led efforts
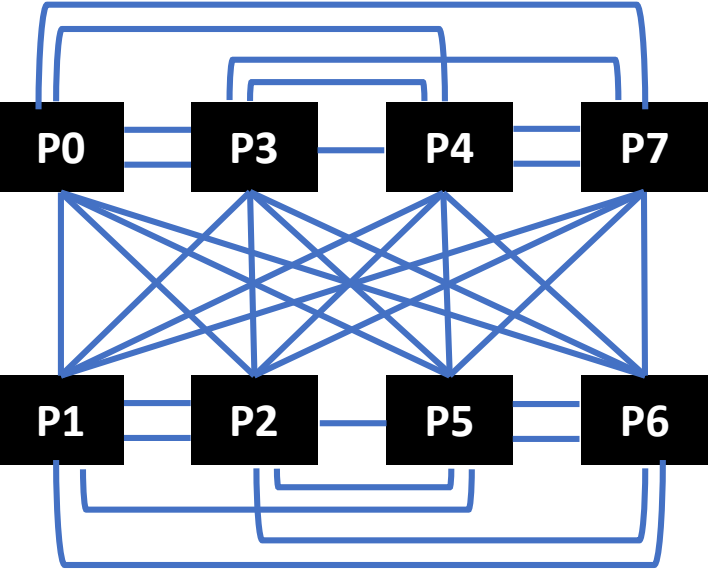    - Define vendor-agnostic common form factor

| AMD | Habana | Intel |
| NVIDIA | ... | |

**facebook**
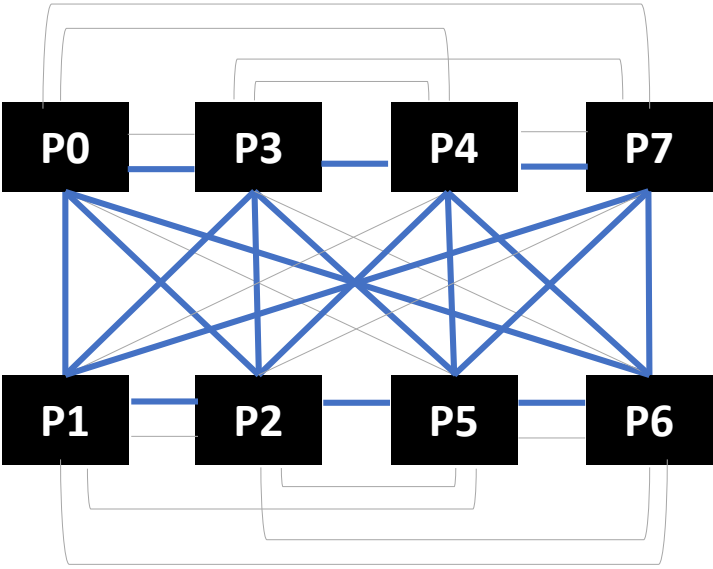
# Accelerator Interconnect Topology

- Challenge: vendors support different topologies: FC, AFC, HCM, …
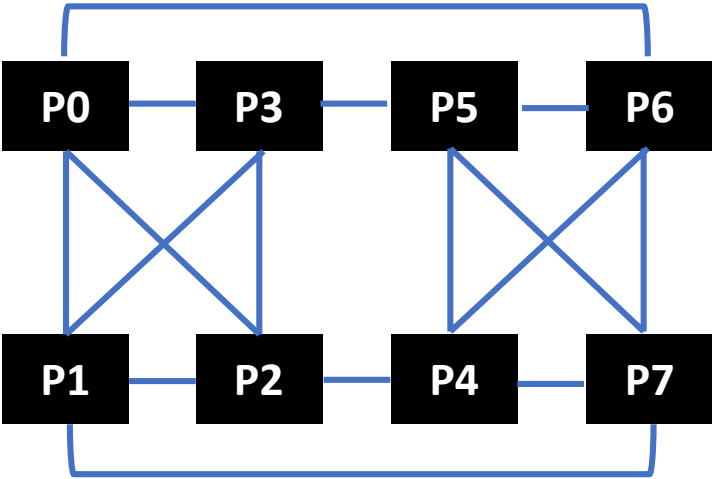- Solution: superset physical topology



facebook

# Example: Embedding Hypercube Mesh



Superset topology
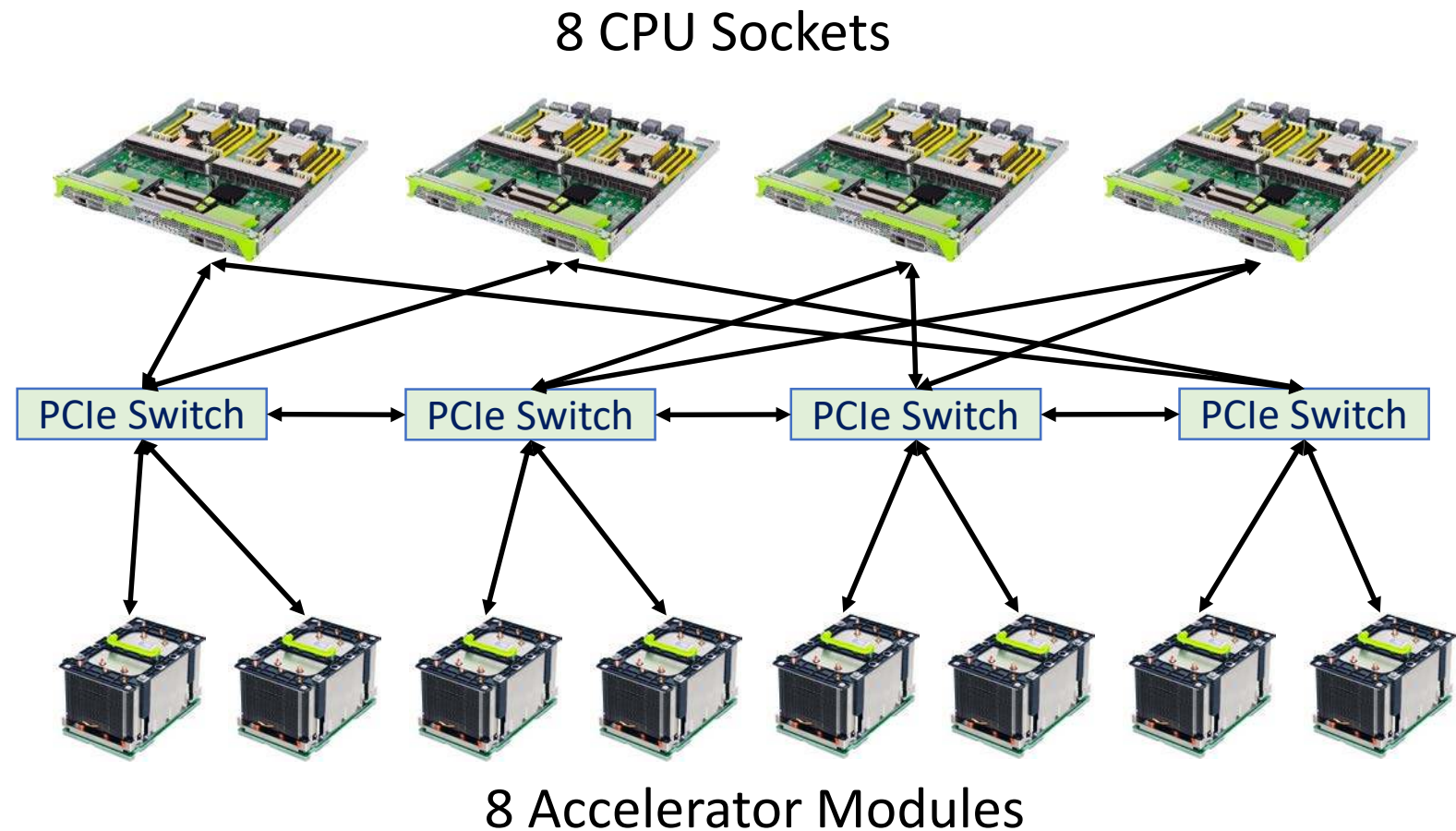
Remove unused links

Rotate 4,7,5,6 by 180° ➜ HCM

facebook

# Software Flexibility
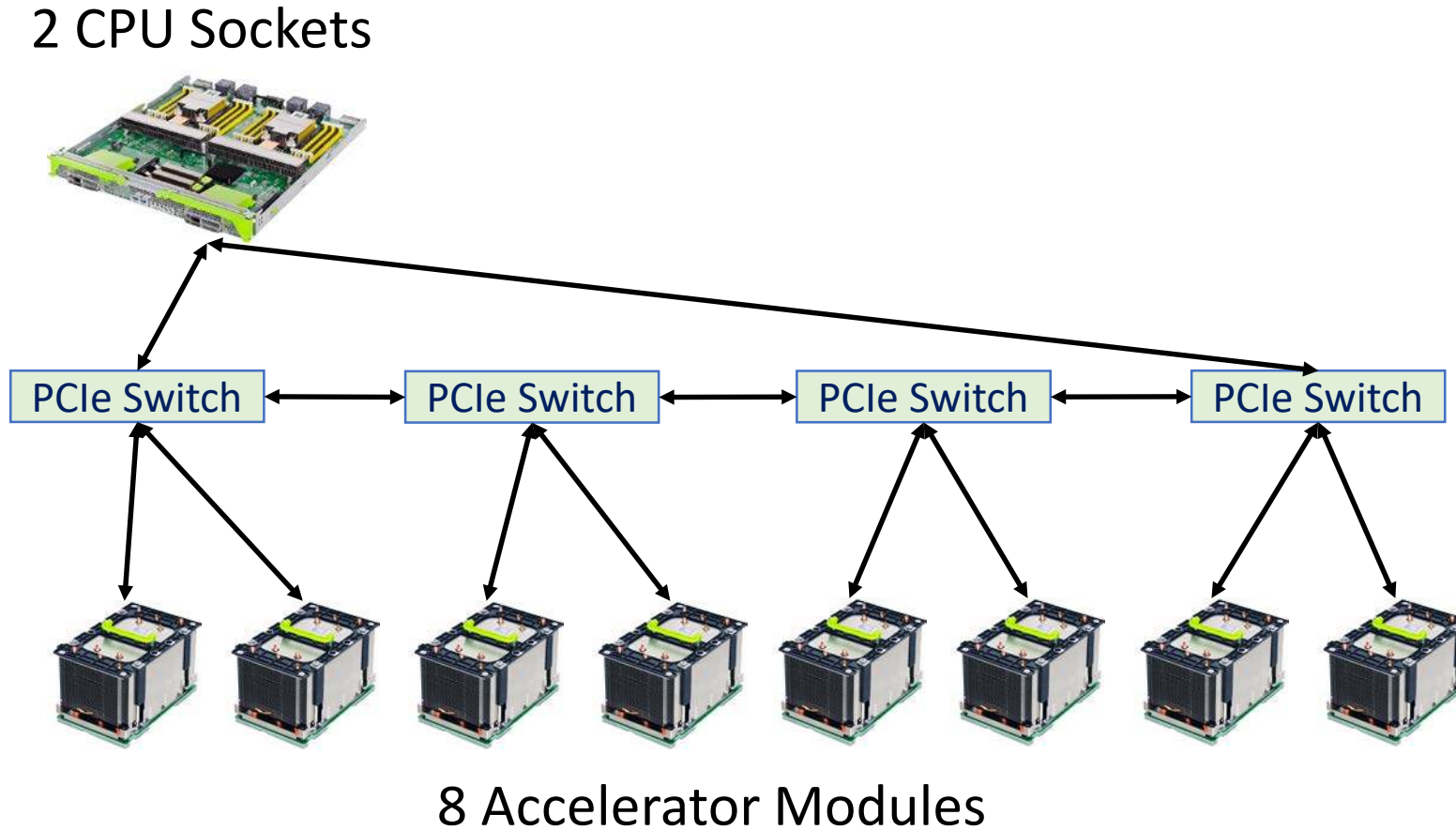
- User can gradually increase SW complexity (and performance)

1. CPU-only

2. CPU for embeddings + Accelerators for MLP

3. Use Accelerator HBM for embeddings as well
   - Challenge: table accesses have different frequencies
   - Benefits from run-time profile driven table partitioning

4. Distributed training

- Creates continuum of dev efficiency vs performance tradeoffs
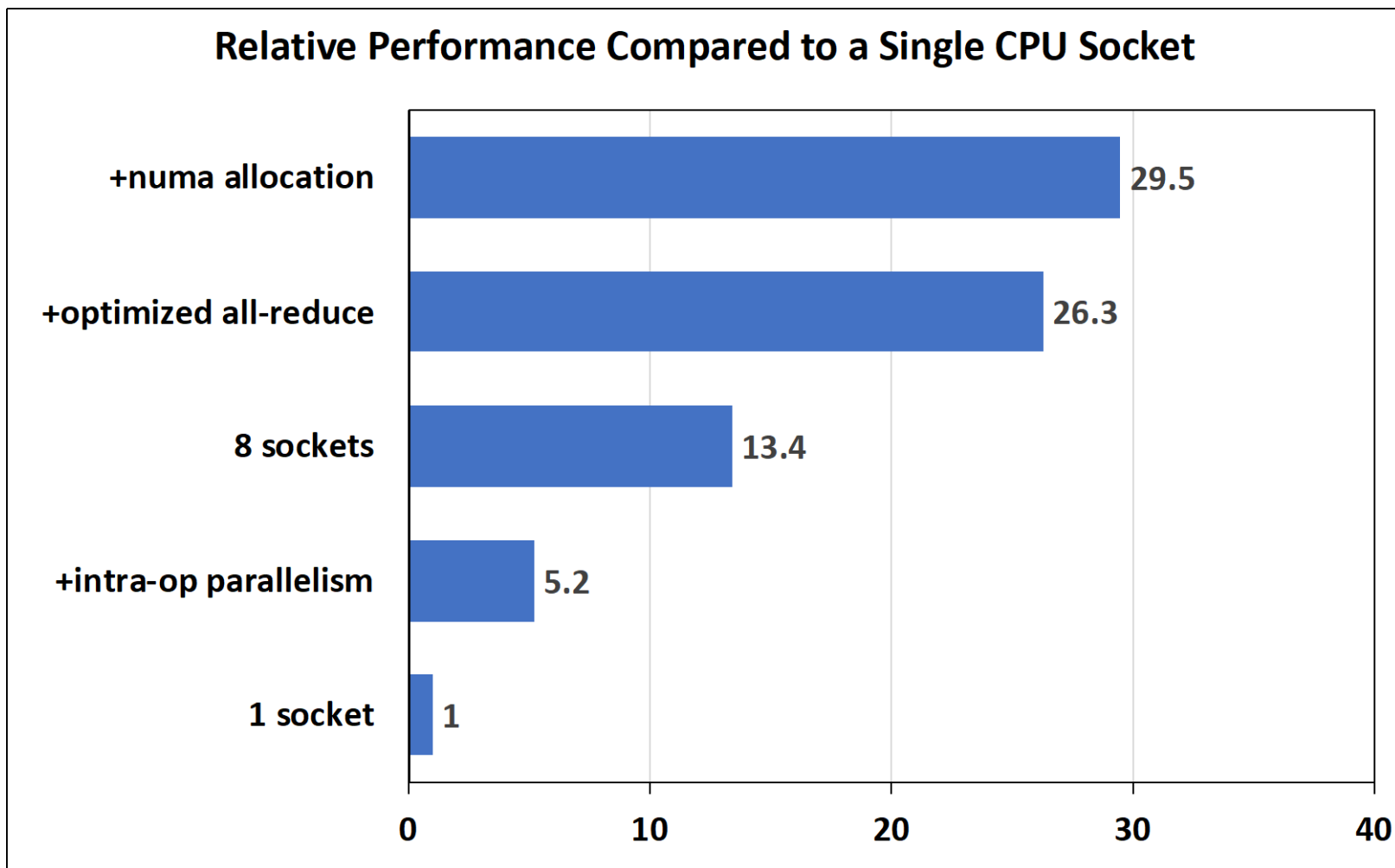
facebook

# Hardware Flexibility

- Four 2S modules are identical

- Configured based on workloads needs as
  - Up to four 2S systems
  - One or two 4S systems
  - One 8S

- SW -> BMC -> CM -> Configure board IDs to be 2S, 4S or 8S to power on

8 CPU Sockets

| PCIe Switch | PCIe Switch | PCIe Switch | PCIe Switch |

8 Accelerator Modules

facebook

18

# One 2S Re-configuration Example

2 CPU Sockets



| PCIe Switch | PCIe Switch | PCIe Switch | PCIe Switch |

8 Accelerator Modules

**facebook**

# Production Performance Results (CPU Only)



**Relative Performance Compared to a Single CPU Socket**

| Category | Value |
|---|---|
| +numa allocation | 29.5 |
| +optimized all-reduce | 26.3 |
| 8 sockets | 13.4 |
| +intra-op parallelism | 5.2 |
| 1 socket | 1 |

**facebook**

# Comparison with GPU-based Platform



| | Big Basin | Zion |
|---|---|---|
| **Accelerator** | NVIDIA GPU Only | Different accelerators |
| **Interconnect** | Hypercube mesh via NVLINK | Richer set of topologies |
| **Memory Capacity** | O(100) GB | O(1000) GB |
| **Number of CPUs** | Single headnode | Reconfigurable |

**facebook**

# Conclusions

- Zion is FB next generation flexible training platform

- Co-designed to target demanding recommendation models

- Adopts new vendor-agnostic OCP accelerator module

- Building block that can scale out to a bigger system

facebook