# Hardware-Software Co-Design

Abhinav Bhatele, Daniel Nichols

UNIVERSITY OF
MARYLAND

# Announcements

- For those that haven't presented, submit videos by May 1

- Extra credit due May 7

- Exam grades out; submit regrade requests by Friday 4/25

DEPARTMENT OF
COMPUTER SCIENCE

# What is HW/SW Co-Design?

- So far we have been changing our algorithms to optimally match hardware

- But what if we changed both?

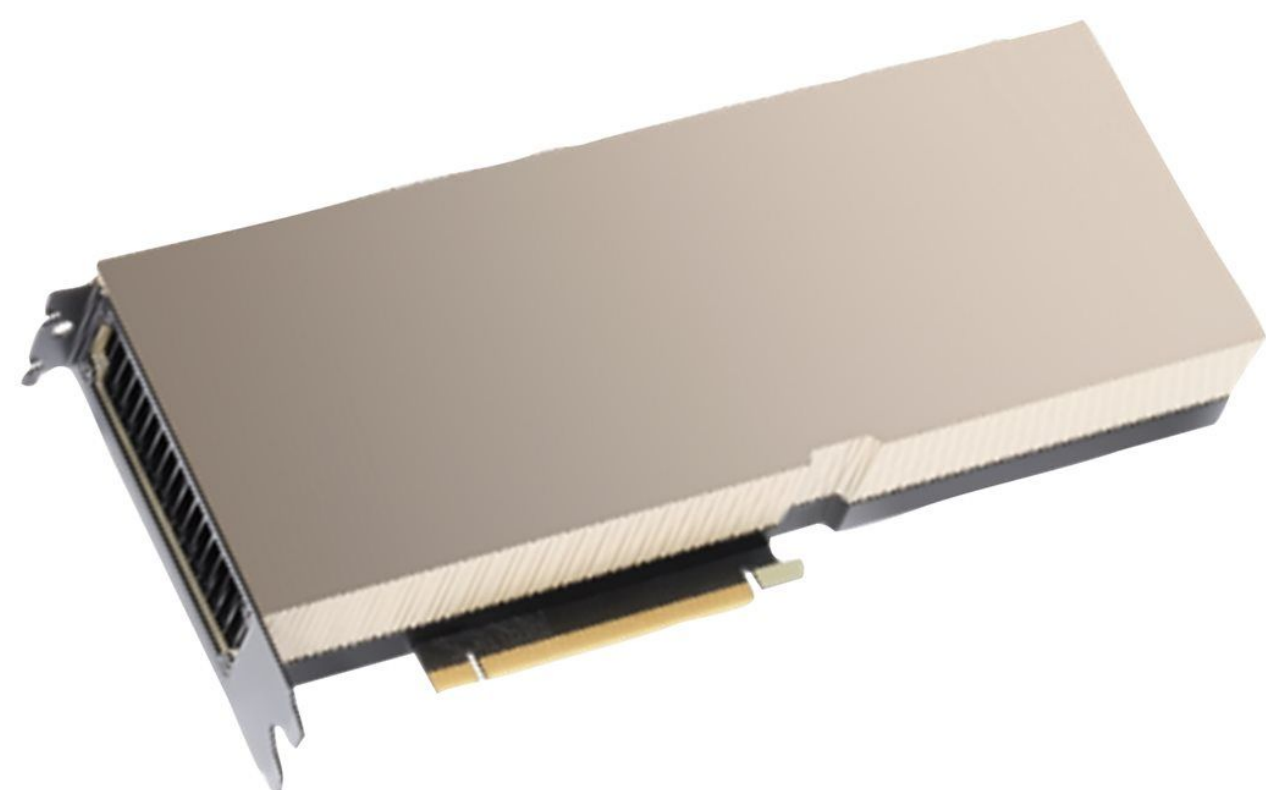What are some HW inefficiencies we've seen often in this class?

# Types of HW/SW Co-Design

- Standalone accelerators for specific domains

  - great efficiency but not very general

- Extend existing hardware with task specific components

  - great efficiency but can mess with original performance
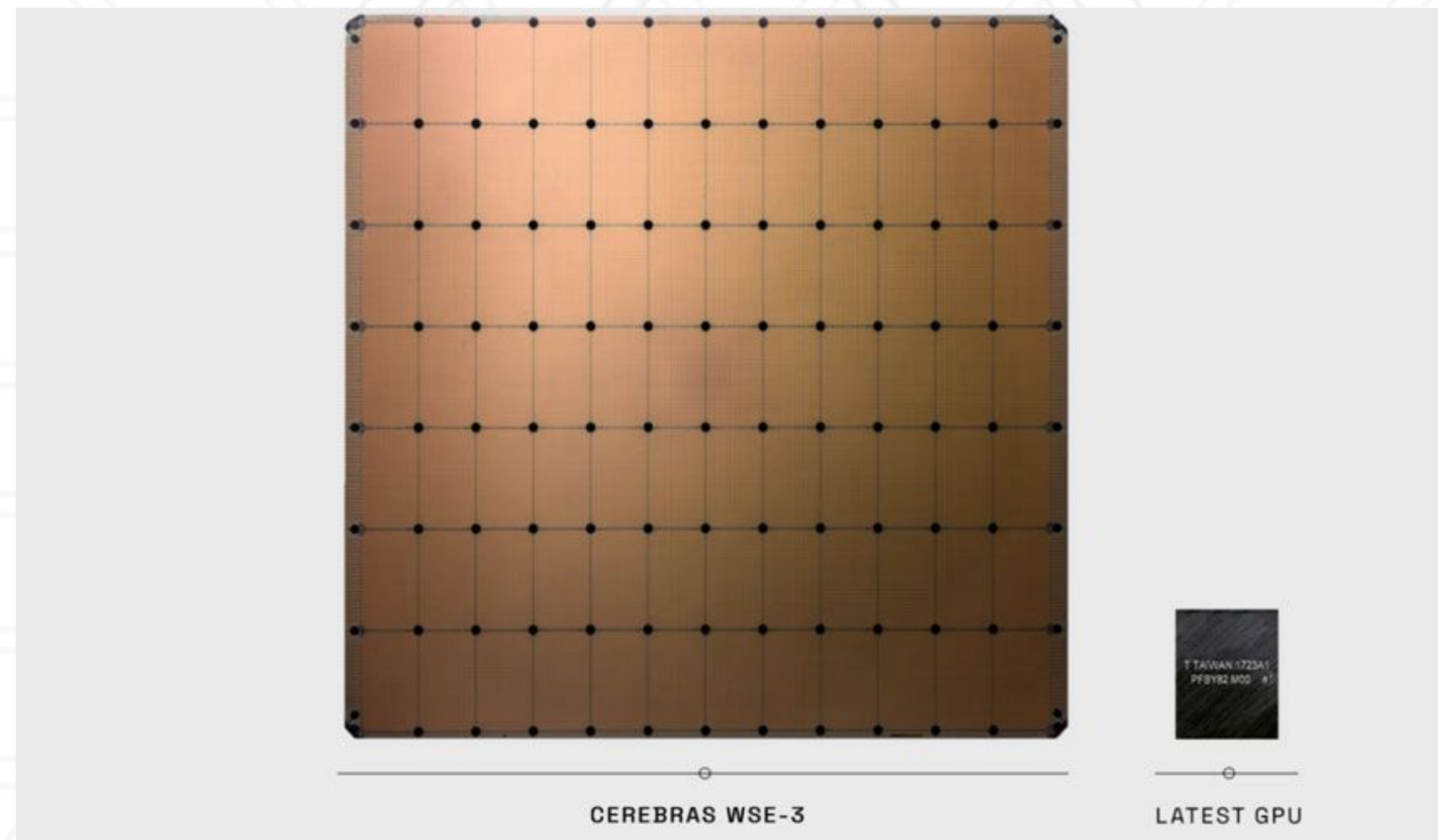
- Improve existing hardware

# Examples



TENSOR CORE 4X4X4 MATRIX-MULTIPLY ACC

$$D = \begin{pmatrix} A_{0,0} & A_{0,1} & A_{0,2} & A_{0,3} \\ A_{1,0} & A_{1,1} & A_{1,2} & A_{1,3} \\ A_{2,0} & A_{2,1} & A_{2,2} & A_{2,3} \\ A_{3,0} & A_{3,1} & A_{3,2} & A_{3,3} \end{pmatrix} \begin{pmatrix} B_{0,0} & B_{0,1} & B_{0,2} & B_{0,3} \\ B_{1,0} & B_{1,1} & B_{1,2} & B_{1,3} \\ B_{2,0} & B_{2,1} & B_{2,2} & B_{2,3} \\ B_{3,0} & B_{3,1} & B_{3,2} & B_{3,3} \end{pmatrix} + \begin{pmatrix} C_{0,0} & C_{0,1} & C_{0,2} & C_{0,3} \\ C_{1,0} & C_{1,1} & C_{1,2} & C_{1,3} \\ C_{2,0} & C_{2,1} & C_{2,2} & C_{2,3} \\ C_{3,0} & C_{3,1} & C_{3,2} & C_{3,3} \end{pmatrix}$$

FP16 or FP32          FP16                    FP16                    FP16 or FP32

8  NVIDIA

Abhinav Bhatele, Daniel Nichols (CMSC828G)

DEPARTMENT OF
COMPUTER SCIENCE

# Examples



CEREBRAS WSE-3                LATEST GPU
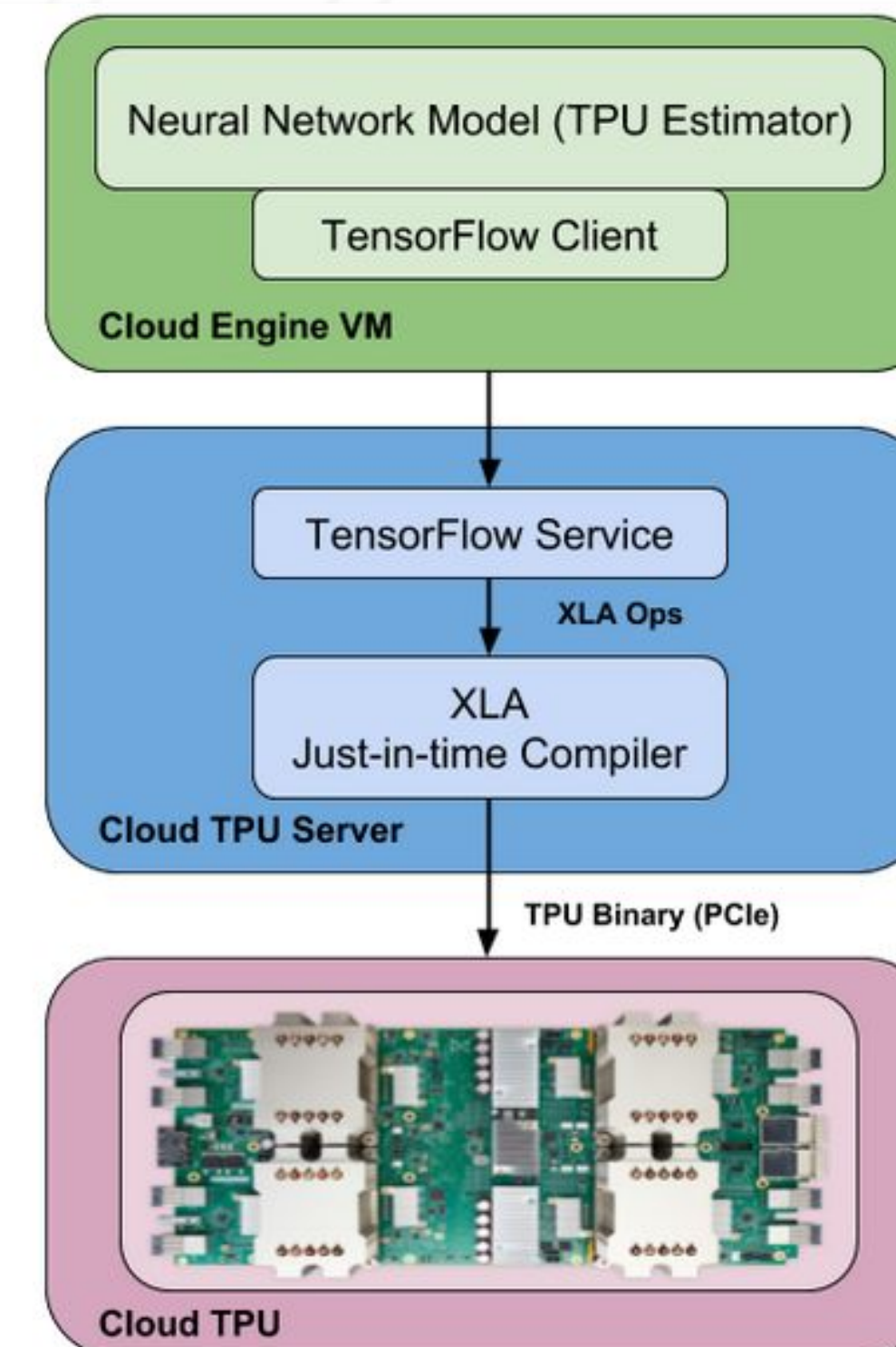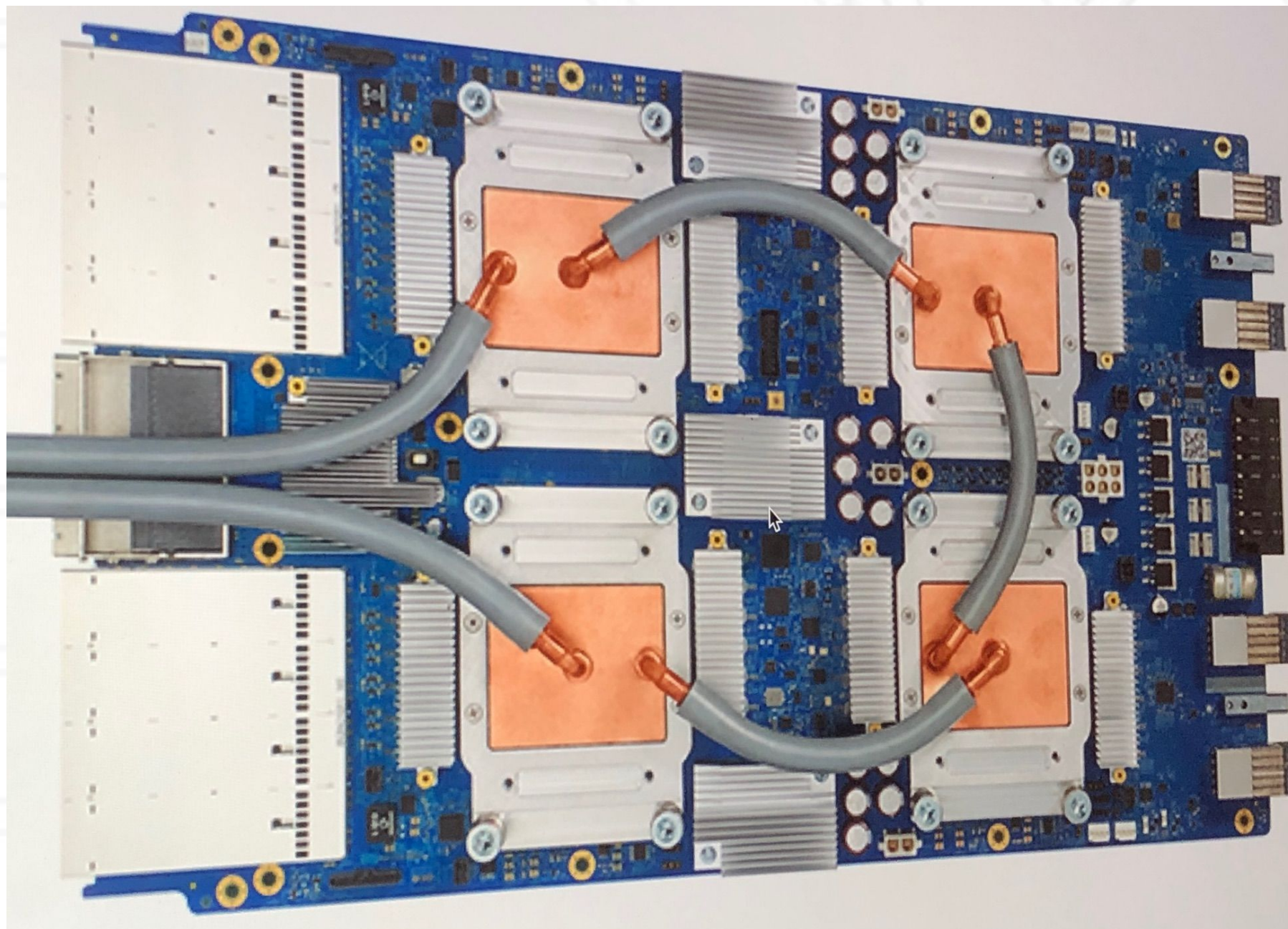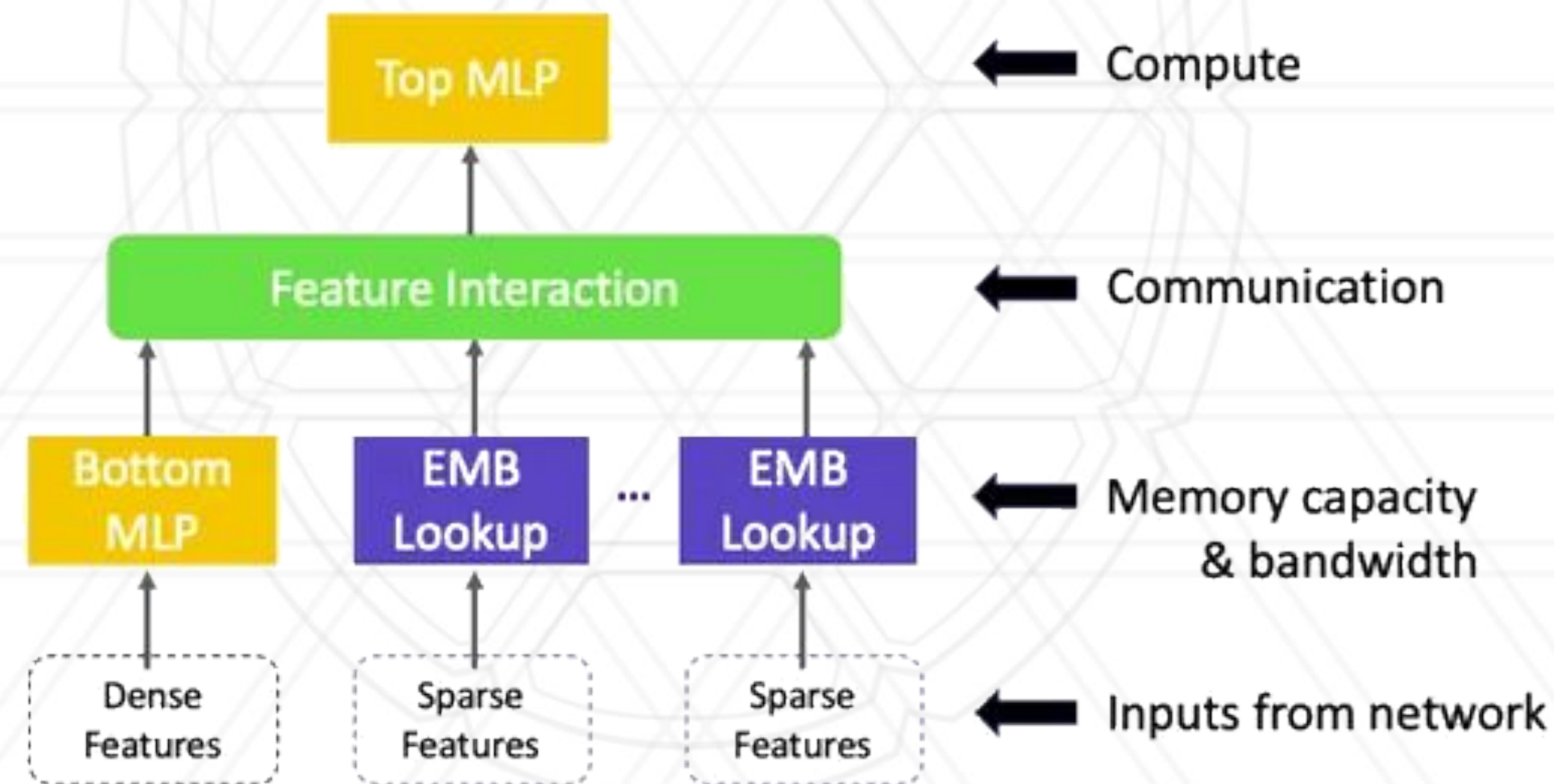
# Examples

# Goals of Co-Design

- Data movement and locality optimizations

- Specialized computation components

    - higher throughput, lower latency

- Reduced power consumption

- Software development ease
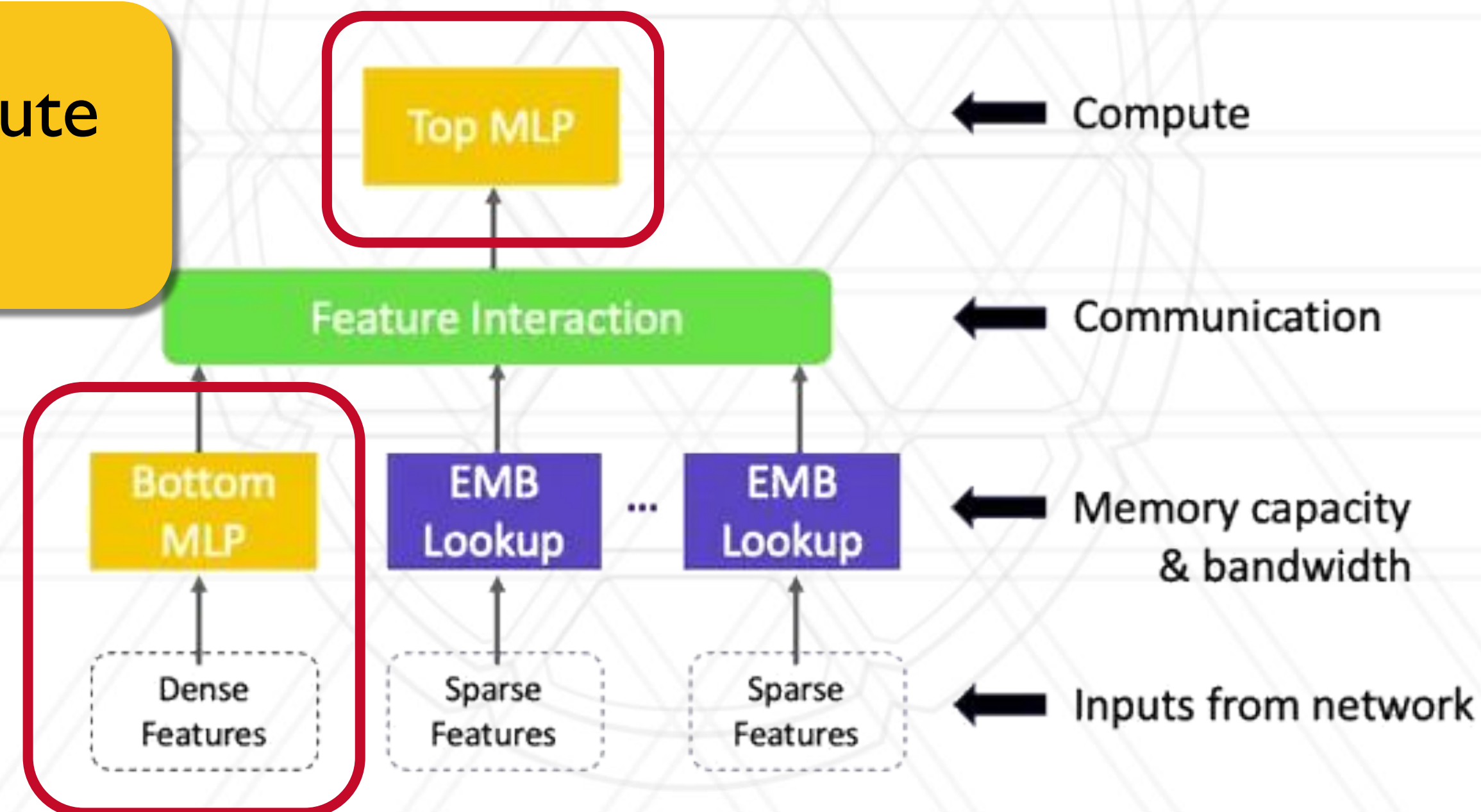
- Reduce costs

# DLRMs Overview

- Deep Learning Recommendation Models
- "Deep Learning Recommendation Model for Personalization and Recommendation Systems", M. Naumov et al
- Online and offline training

# DLRMs Overview

- Deep Learning Recommendation Models
- "Deep Learning Recommendation Model for Personalization and Recommendation Systems", M. Naumov et al
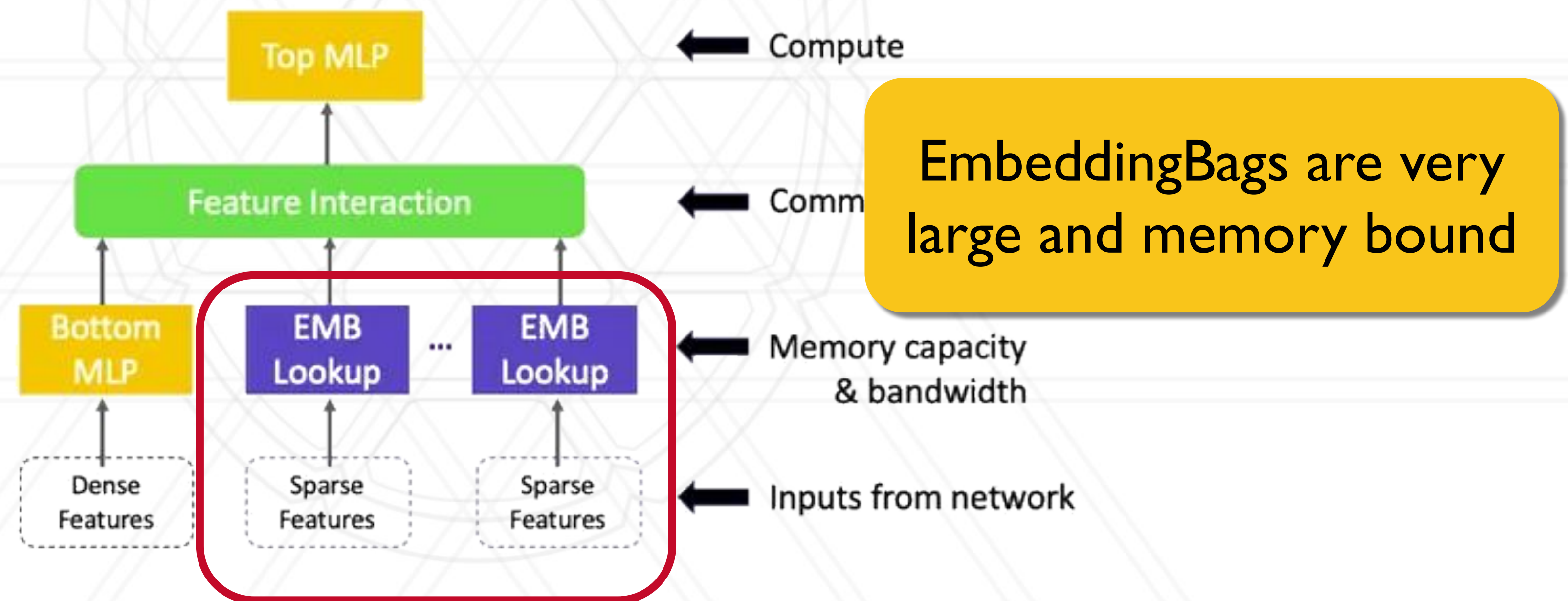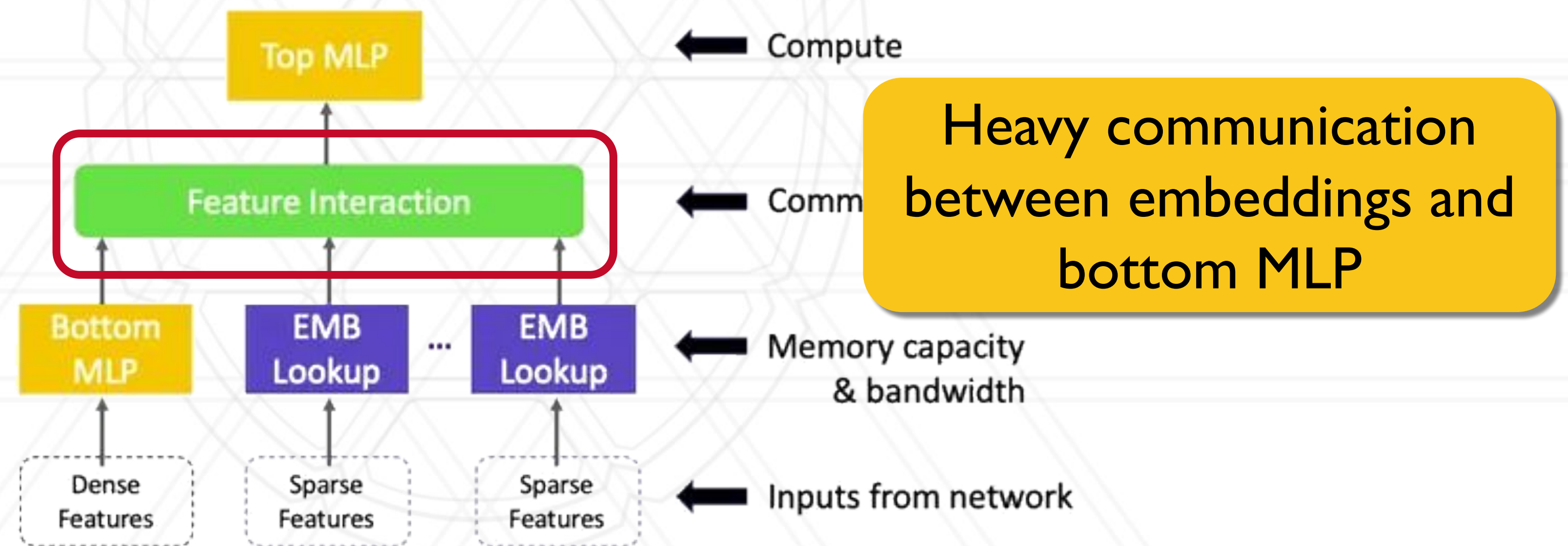- Online and offline training



**MLP layers are compute intensive**

Top MLP ← Compute

Feature Interaction ← Communication

Bottom MLP | EMB Lookup | ... | EMB Lookup ← Memory capacity & bandwidth

Dense Features | Sparse Features | Sparse Features ← Inputs from network

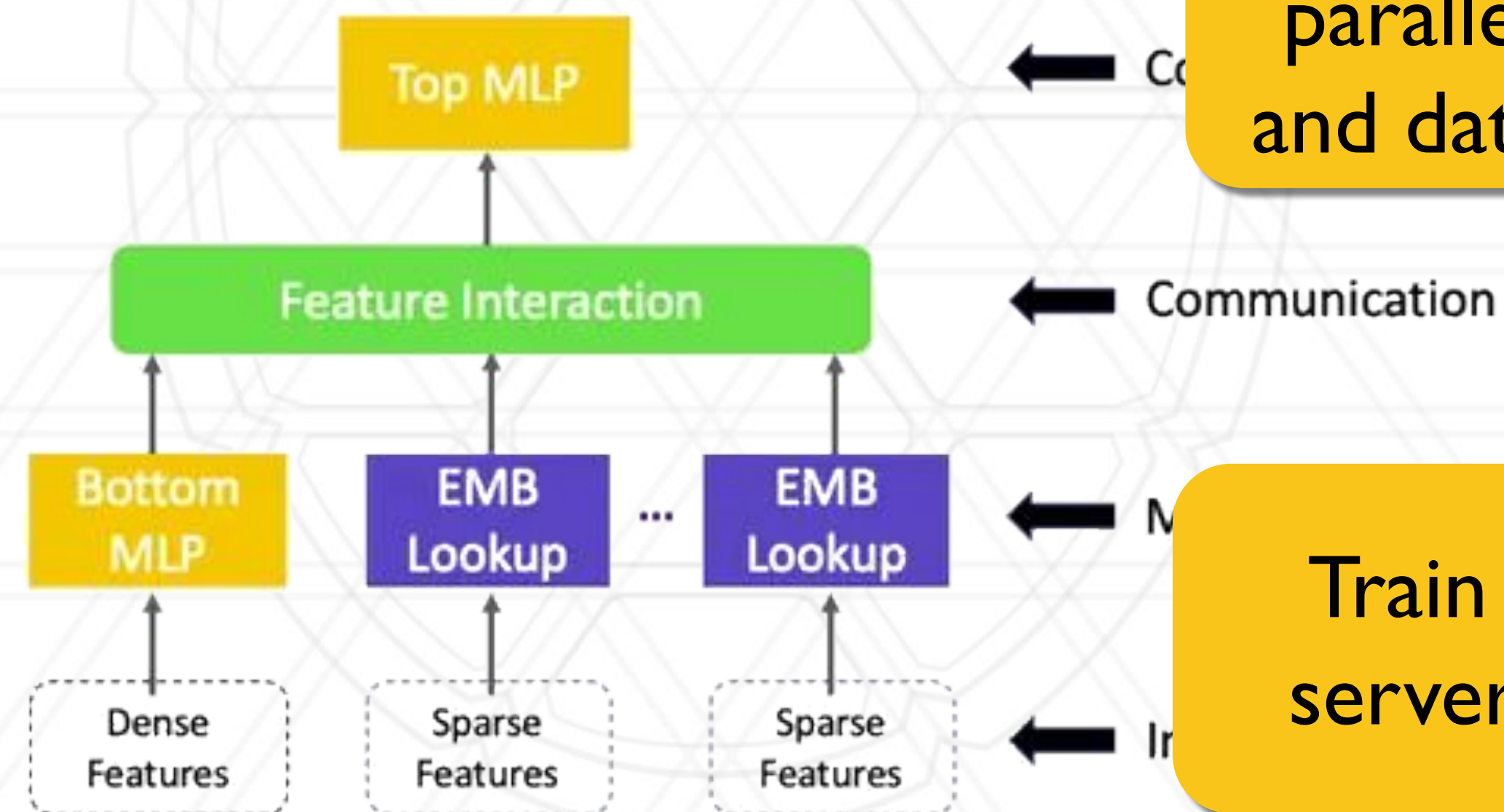DEPARTMENT OF COMPUTER SCIENCE

# DLRMs Overview

- Deep Learning Recommendation Models
- "Deep Learning Recommendation Model for Personalization and Recommendation Systems", M. Naumov et al
- Online and offline training



EmbeddingBags are very large and memory bound

# DLRMs Overview

- Deep Learning Recommendation Models
- "Deep Learning Recommendation Model for Personalization and Recommendation Systems", M. Naumov et al
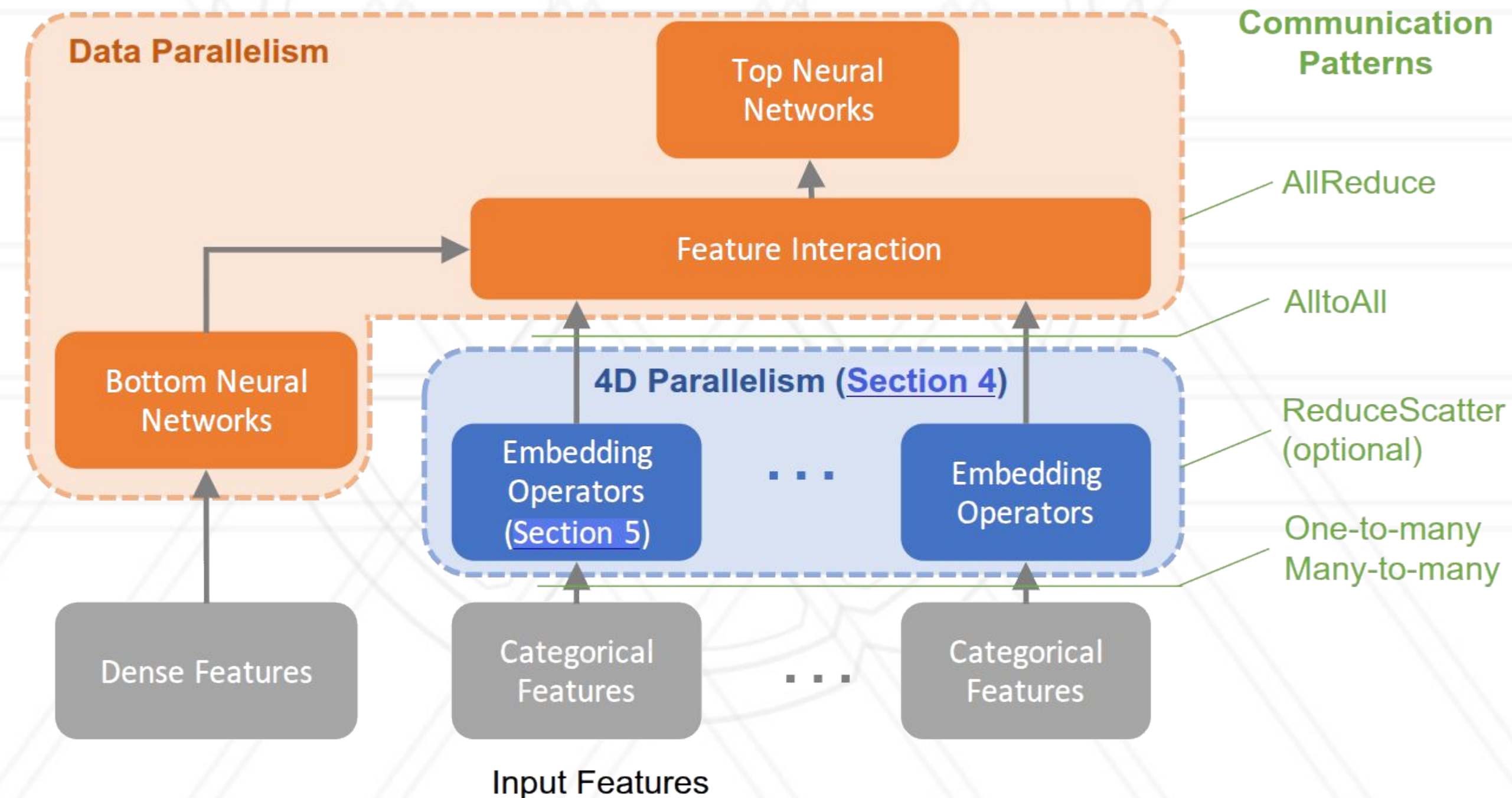- Online and offline training

# DLRMs Overview

- Deep Learning Recommendation Models
- "Deep Learning Recommendation Model for Personalization and Recommendation Systems", M. Naumov et al
- Online and offline training



Traditionally use model parallel for embeddings and data parallel for MLP
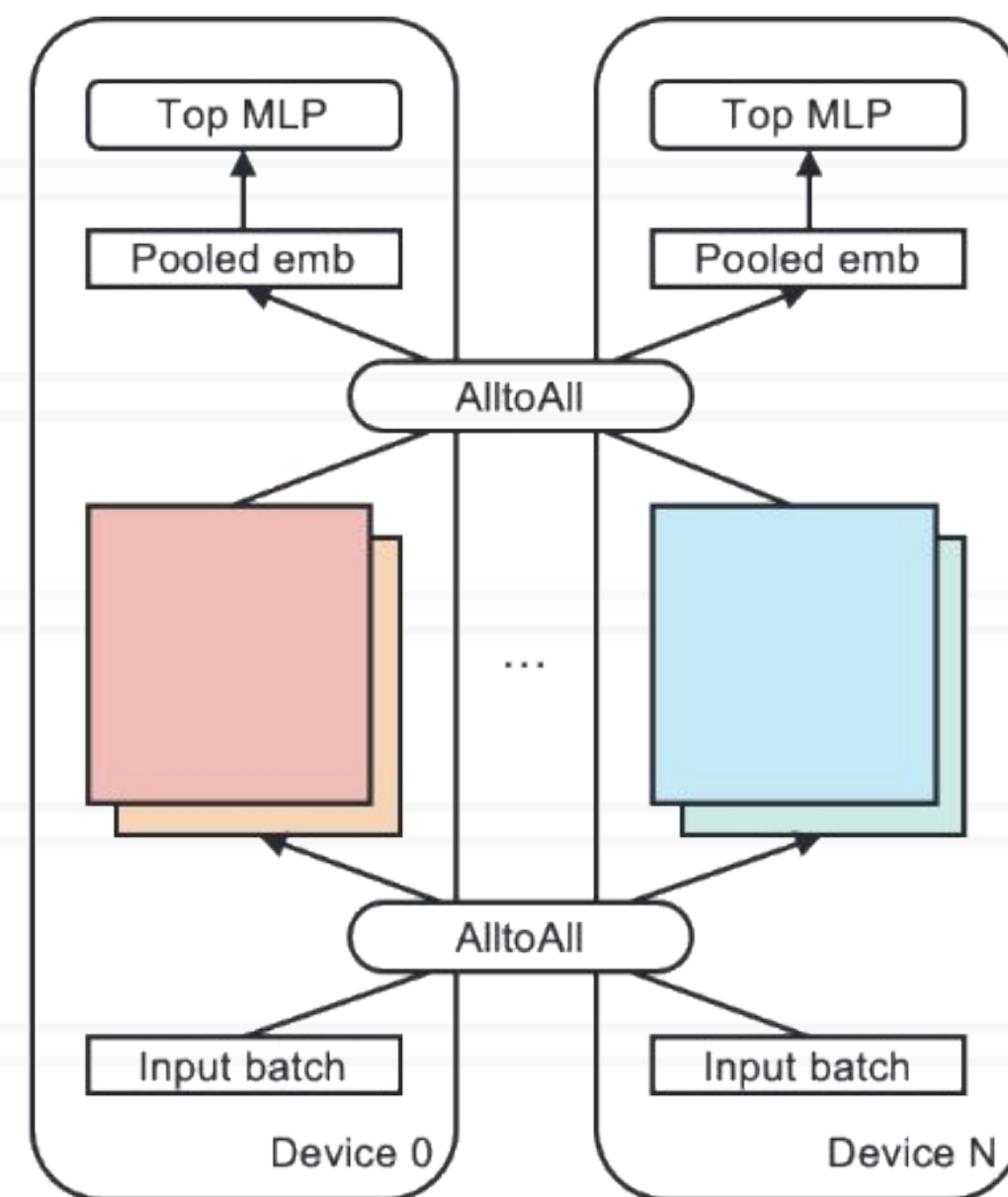
Train with parameter server and async SGD

DEPARTMENT OF
COMPUTER SCIENCE

# Neo Overview

- A DLRM training system
  - Neo software with 4D parallelism for embedding operators
  - Optimized sequential embedding implementations
  - ZionEX: a hardware system designed to efficiently run Neo
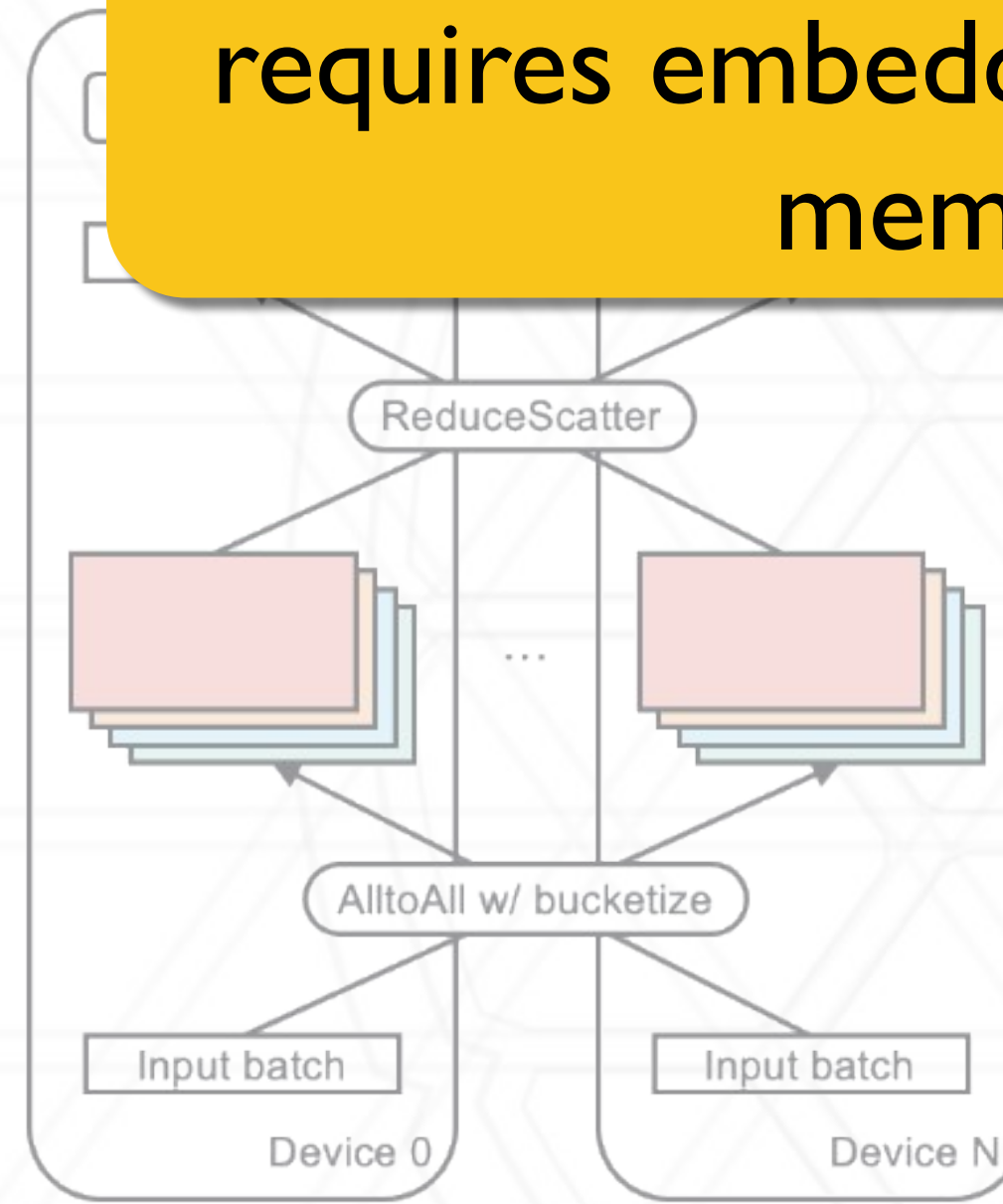
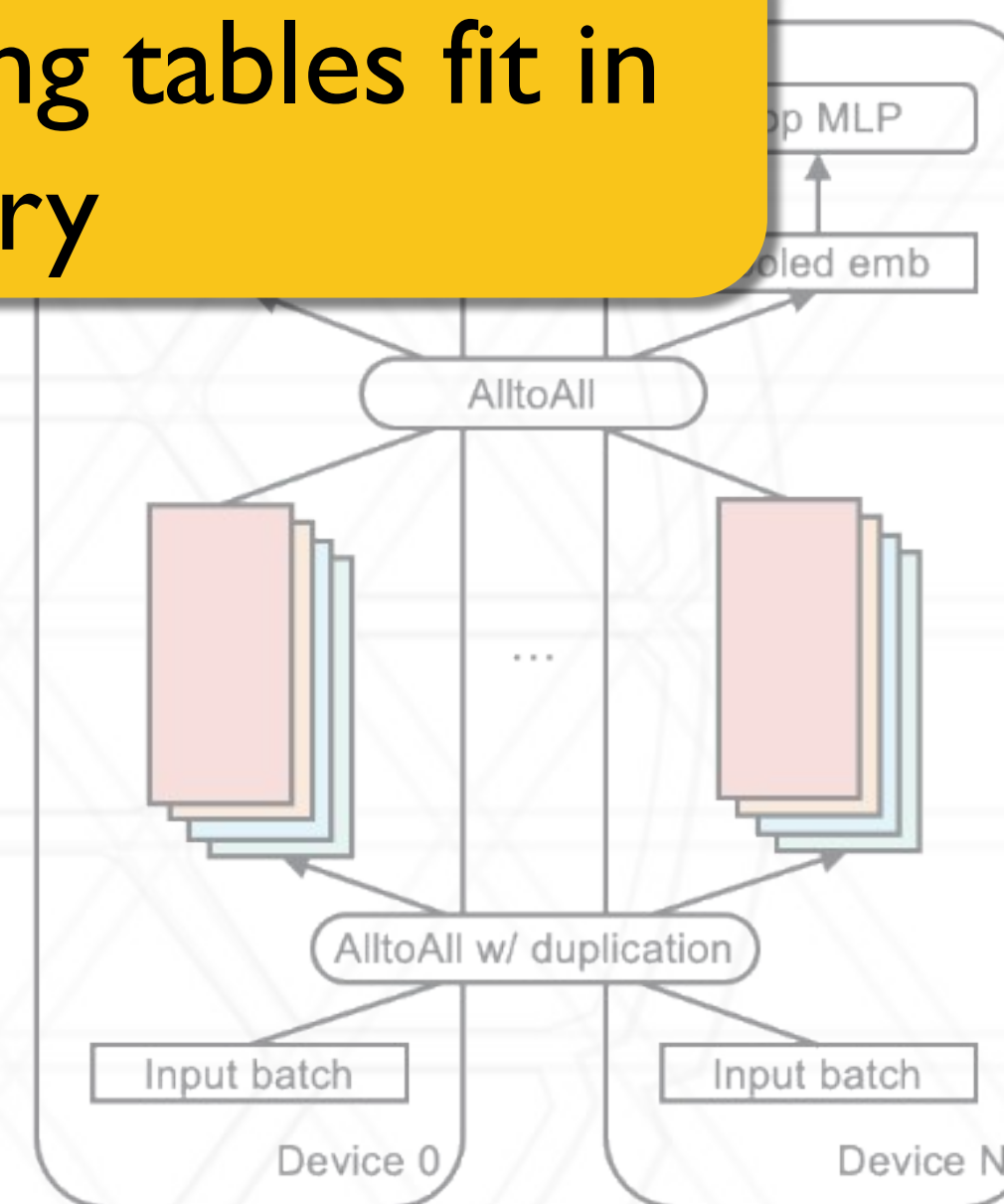# 4D Embedding Parallelism

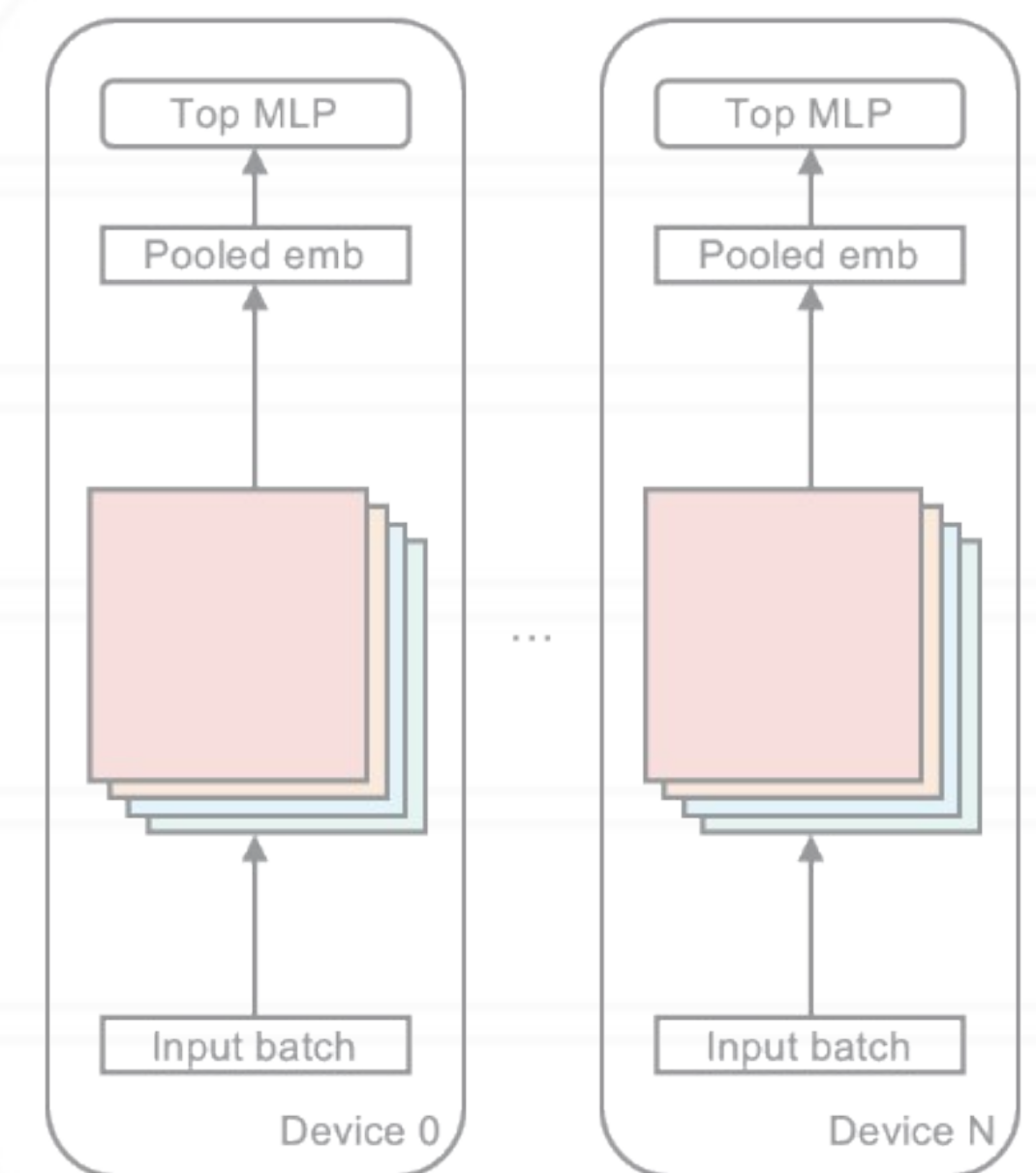Most communication optimal, but requires embedding tables fit in memory



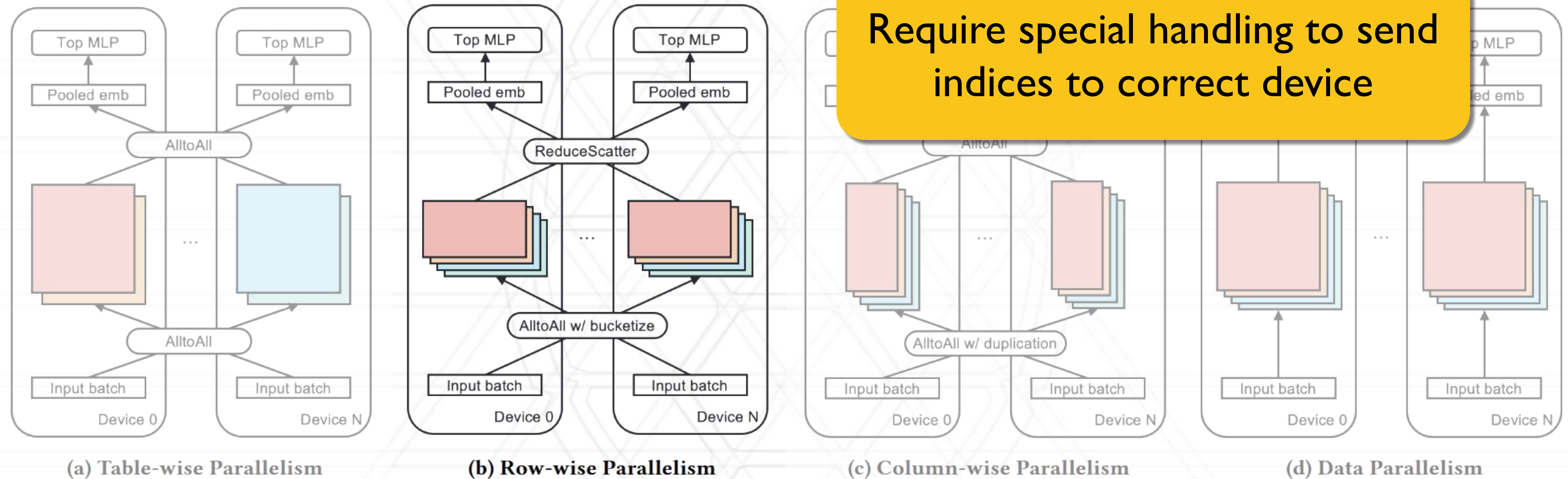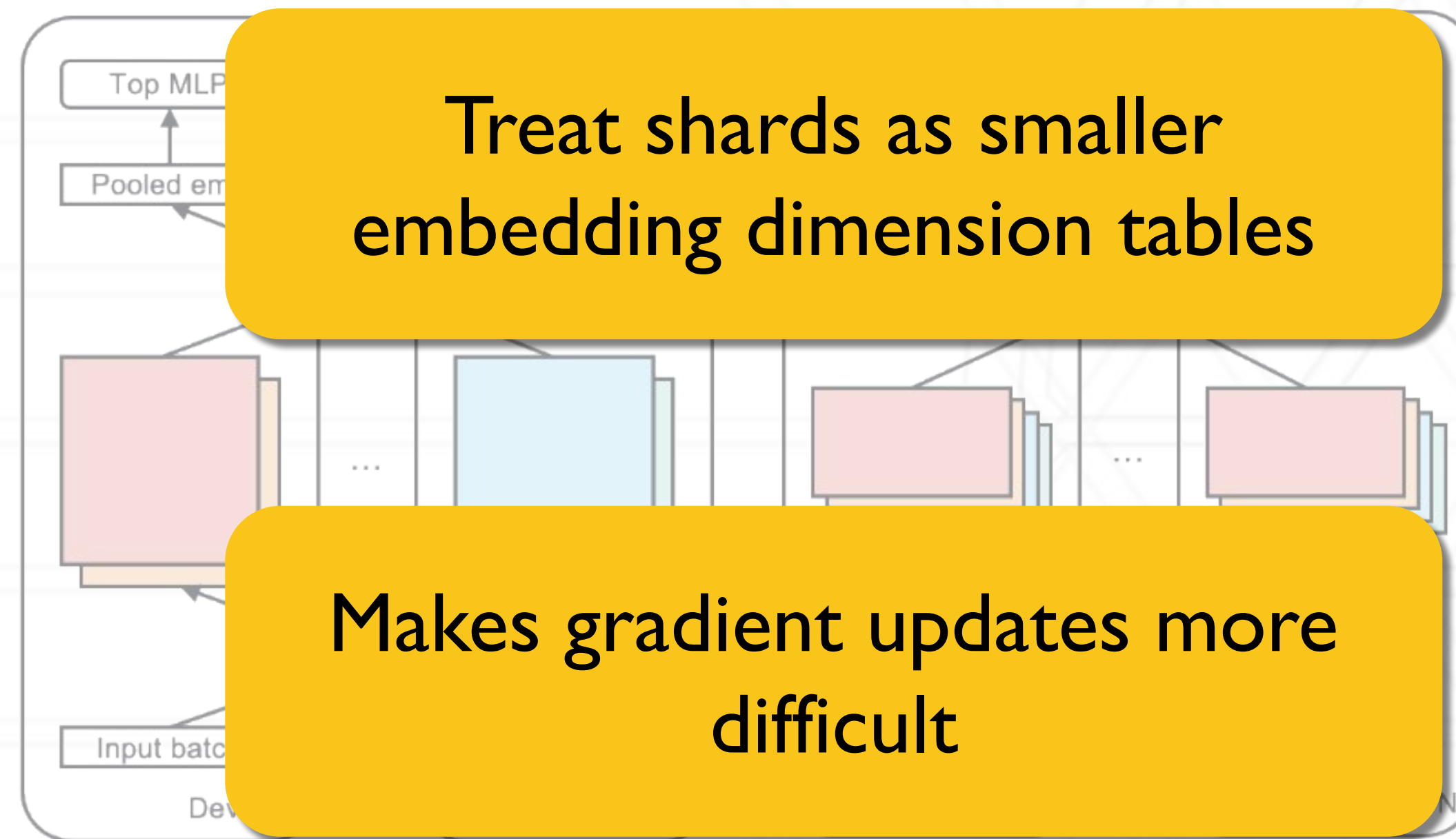(a) Table-wise Parallelism    (b) Row-wise Parallelism    (c) Column-wise Parallelism    (d) Data Parallelism
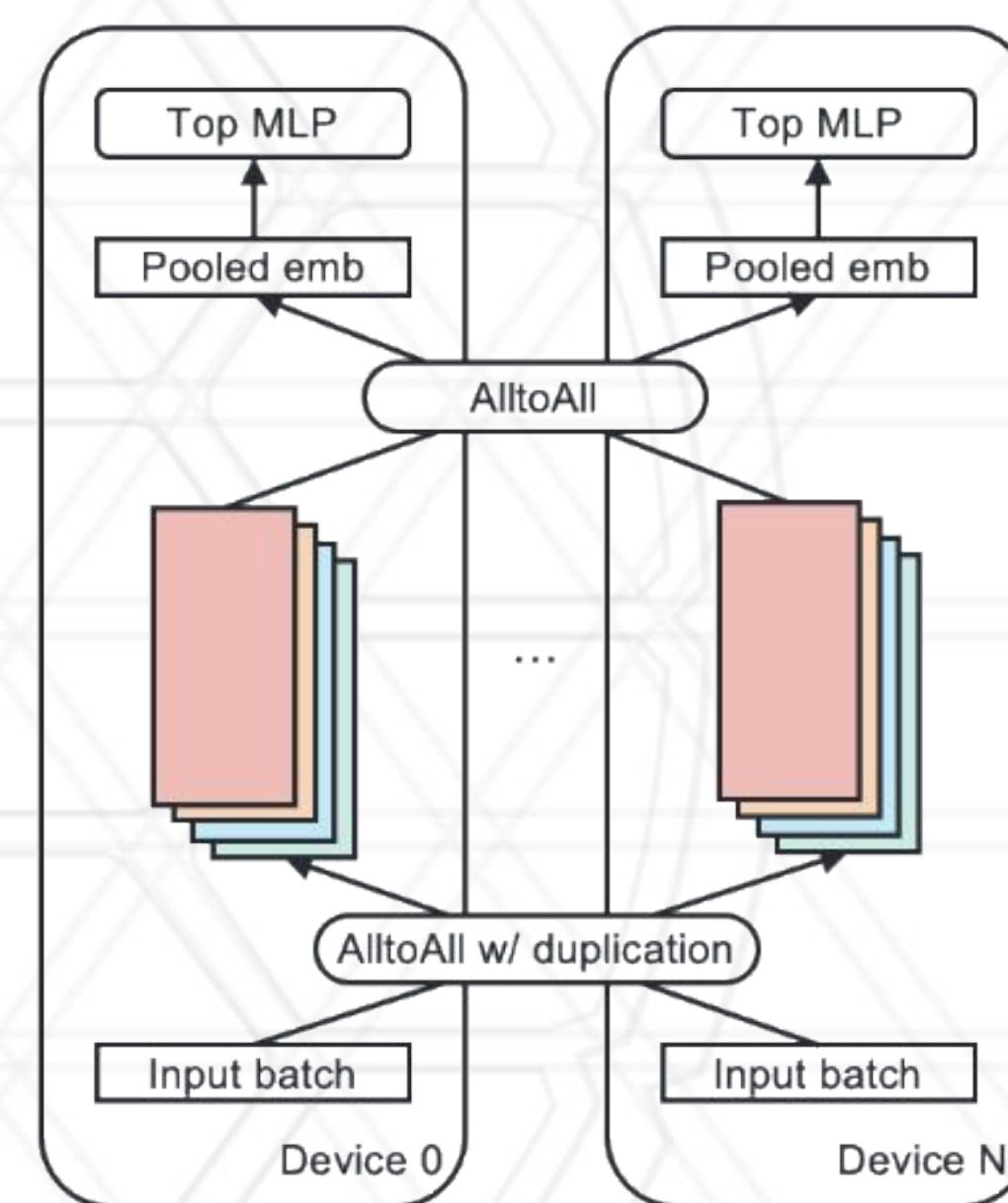
DEPARTMENT OF COMPUTER SCIENCE

# 4D Embedding Parallelism



(a) Table-wise Parallelism

(b) Row-wise Parallelism

(c) Column-wise Parallelism

(d) Data Parallelism

Require special handling to send indices to correct device

DEPARTMENT OF COMPUTER SCIENCE

# 4D Embedding Parallelism



Treat shards as smaller embedding dimension tables

Makes gradient updates more difficult

(a) Table-wise Parallelism

(b) Row-wise Parallelism

(c) Column-wise Parallelism

(d) Data Parallelism
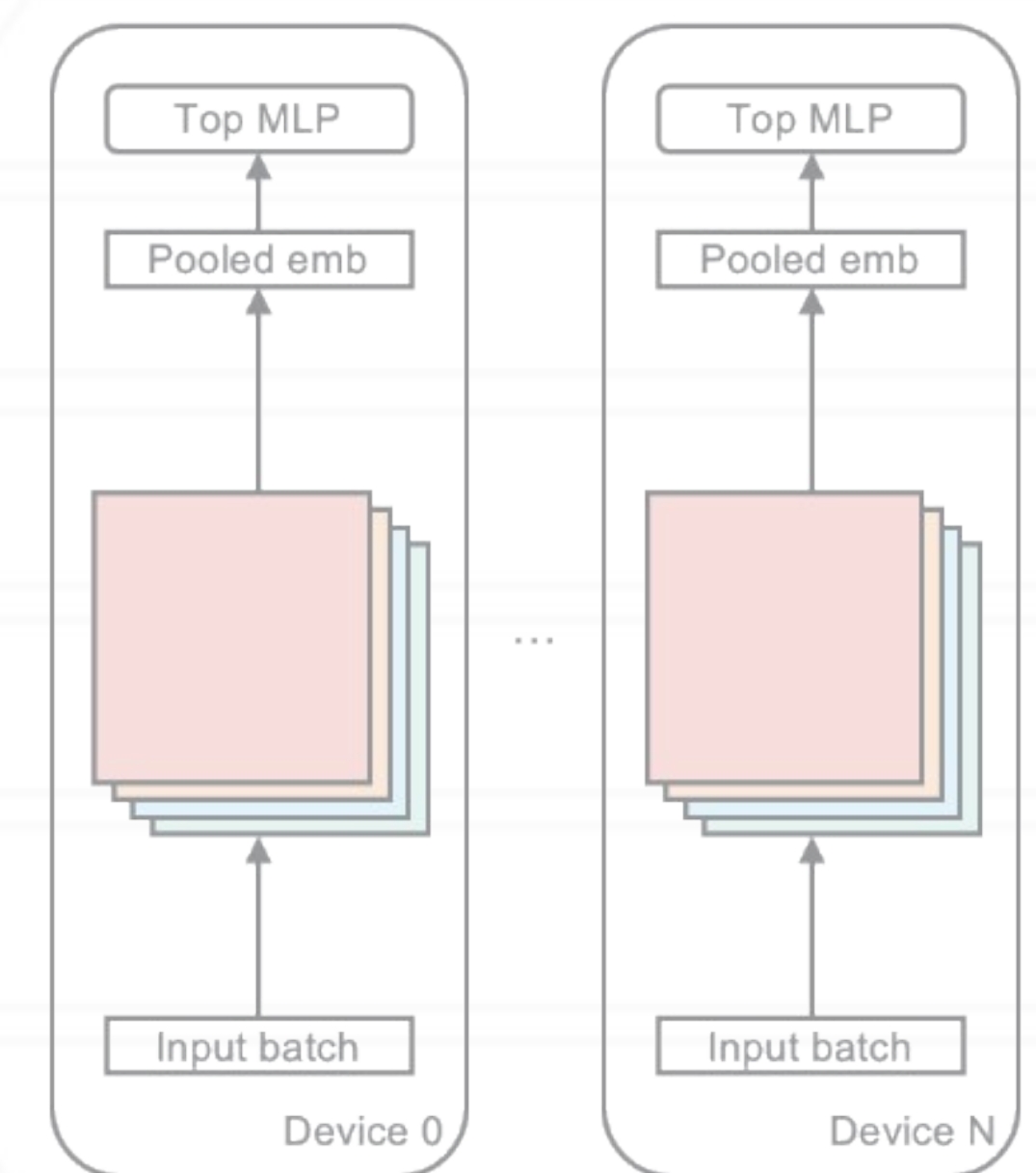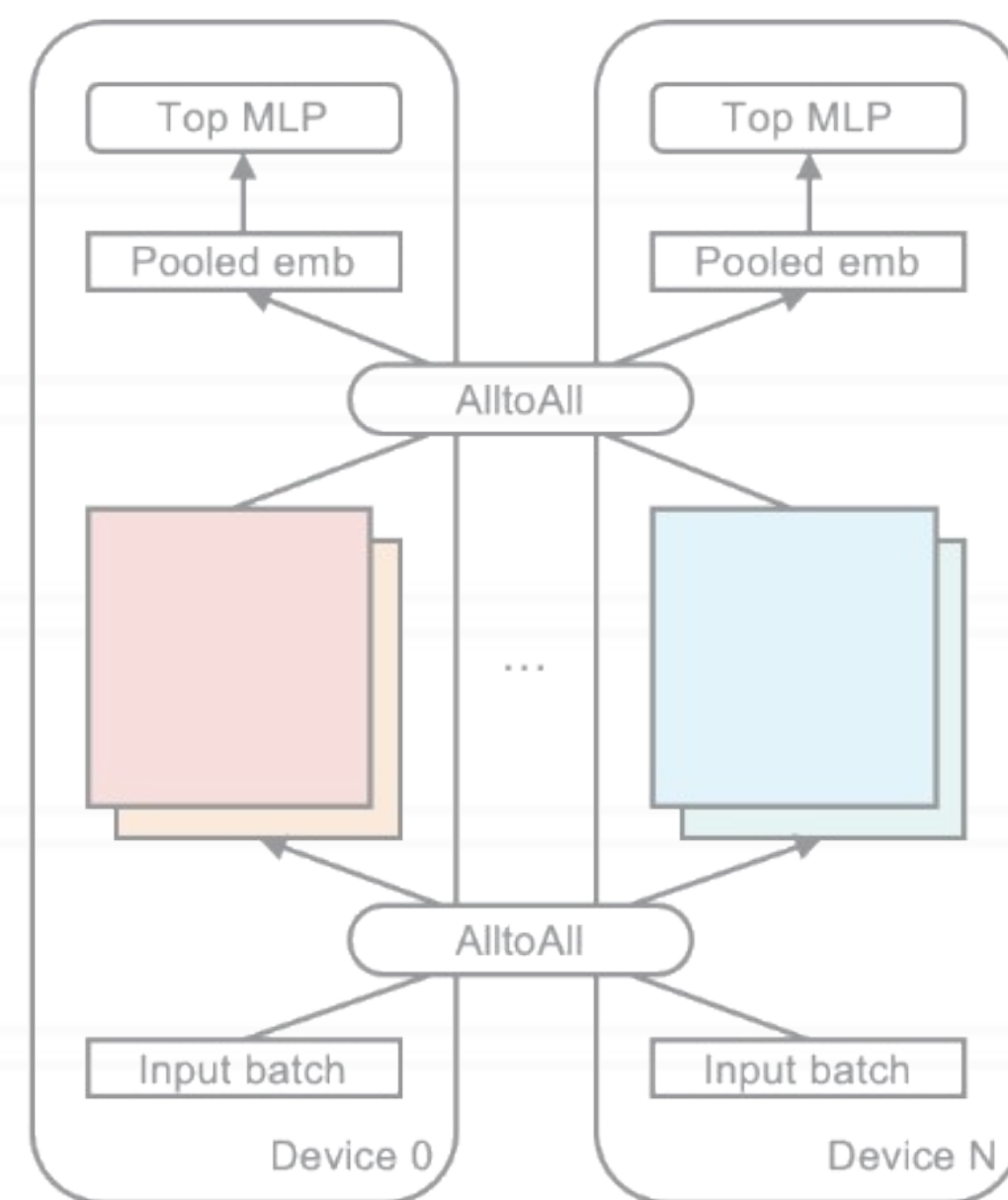
# 4D Embedding Parallelism



(a) Table-wise Parallelism    (b) Row-wise Parallelism    (c) Column-wise Parallelism    (d) Data Parallelism

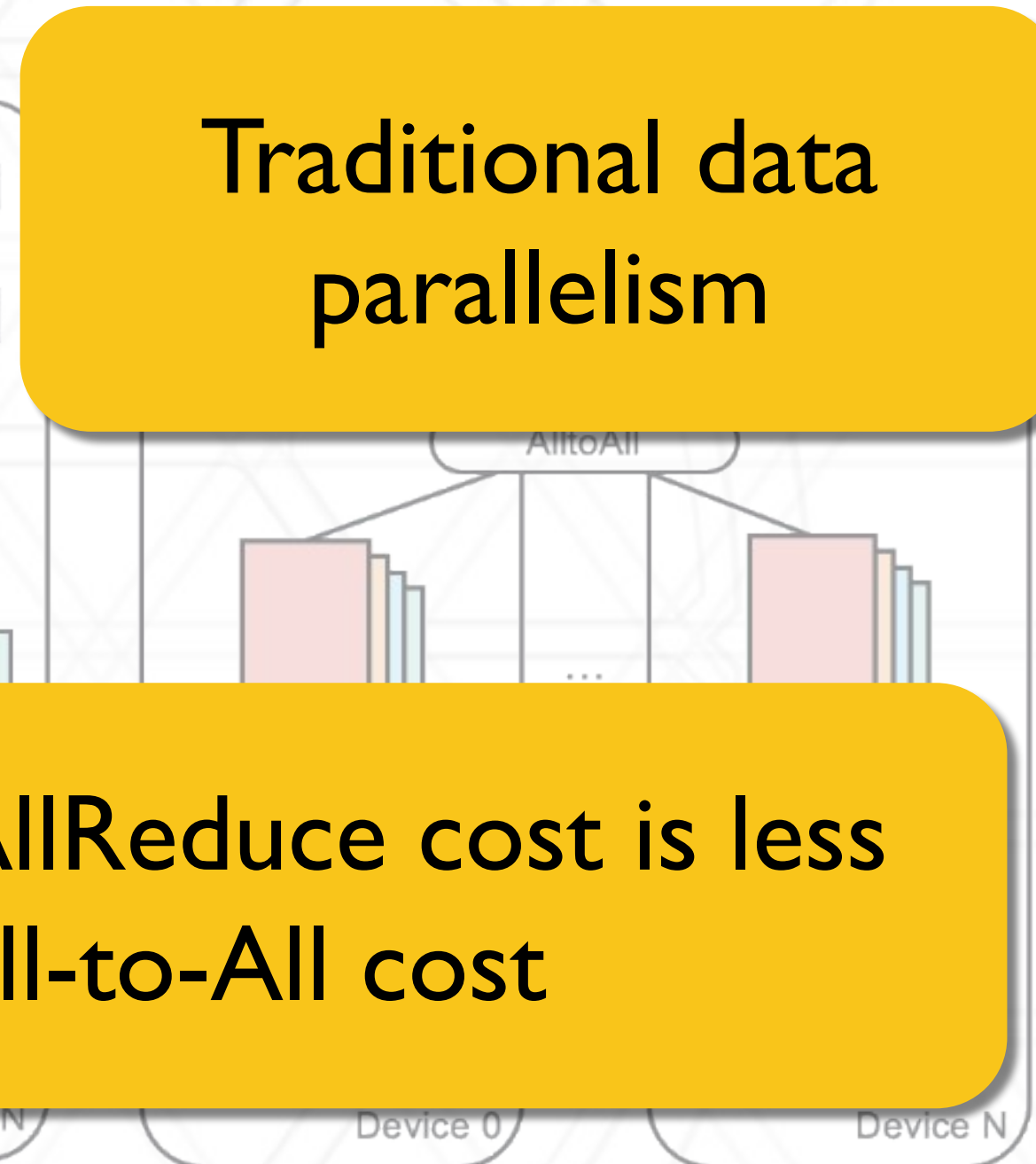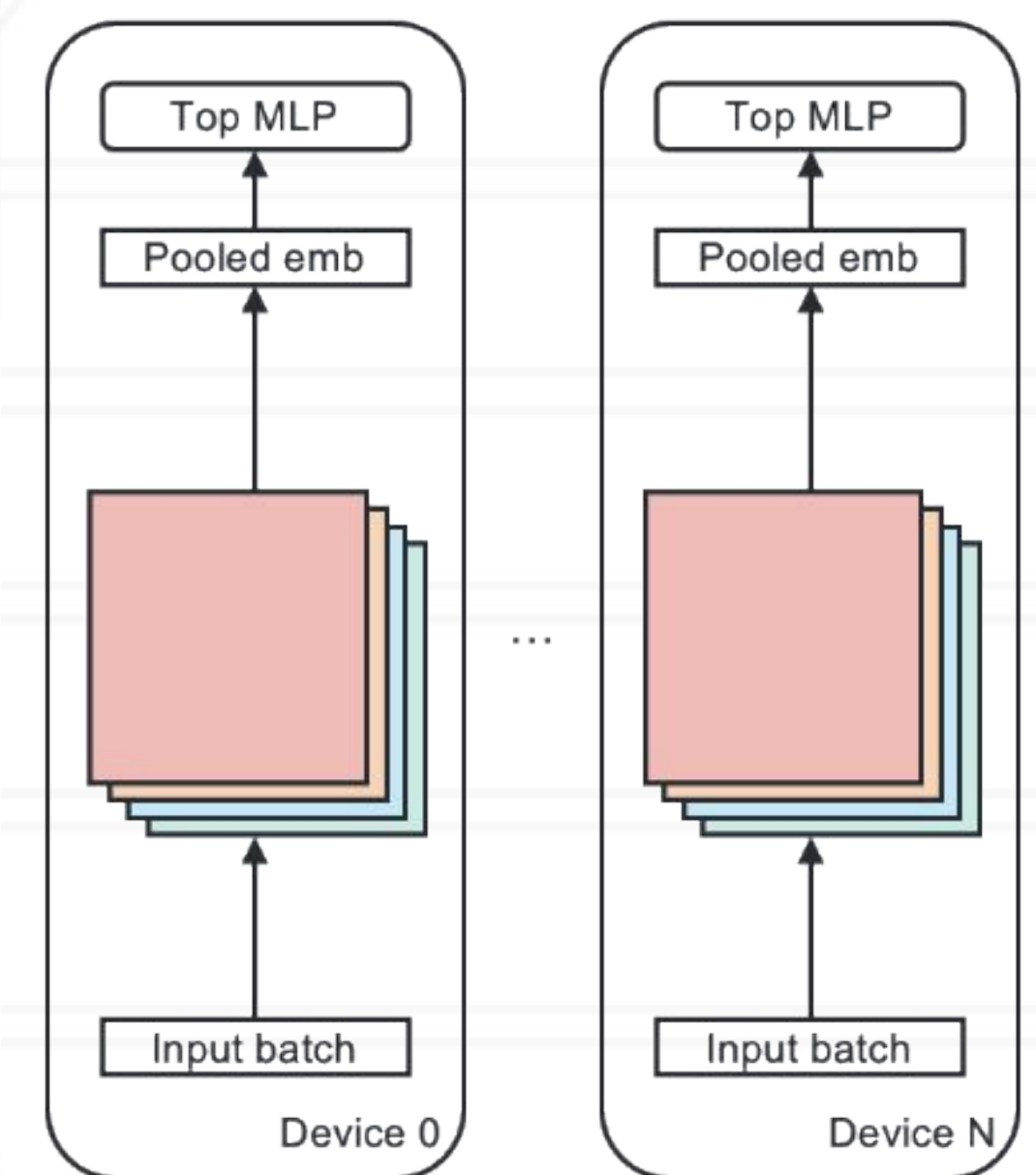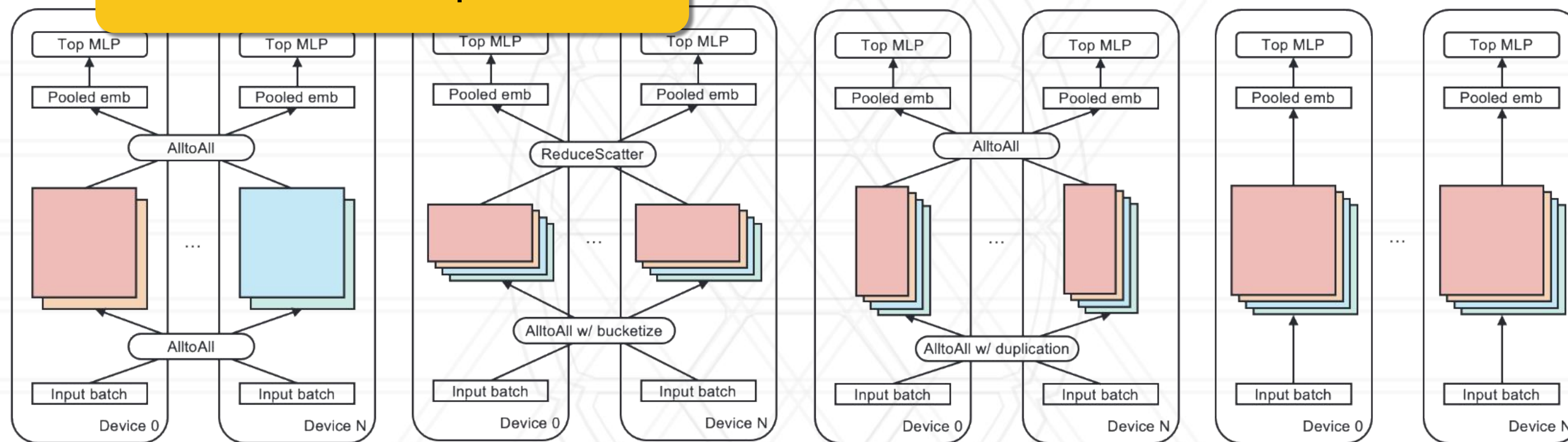Traditional data parallelism

Better when AllReduce cost is less than All-to-All cost

# 4D Embedding Parallelism

Use all 4 for ideal parallelism!



(a) Table-wise Parallelism

(b) Row-wise Parallelism

Use performance models and pipelining to find optimal configurations
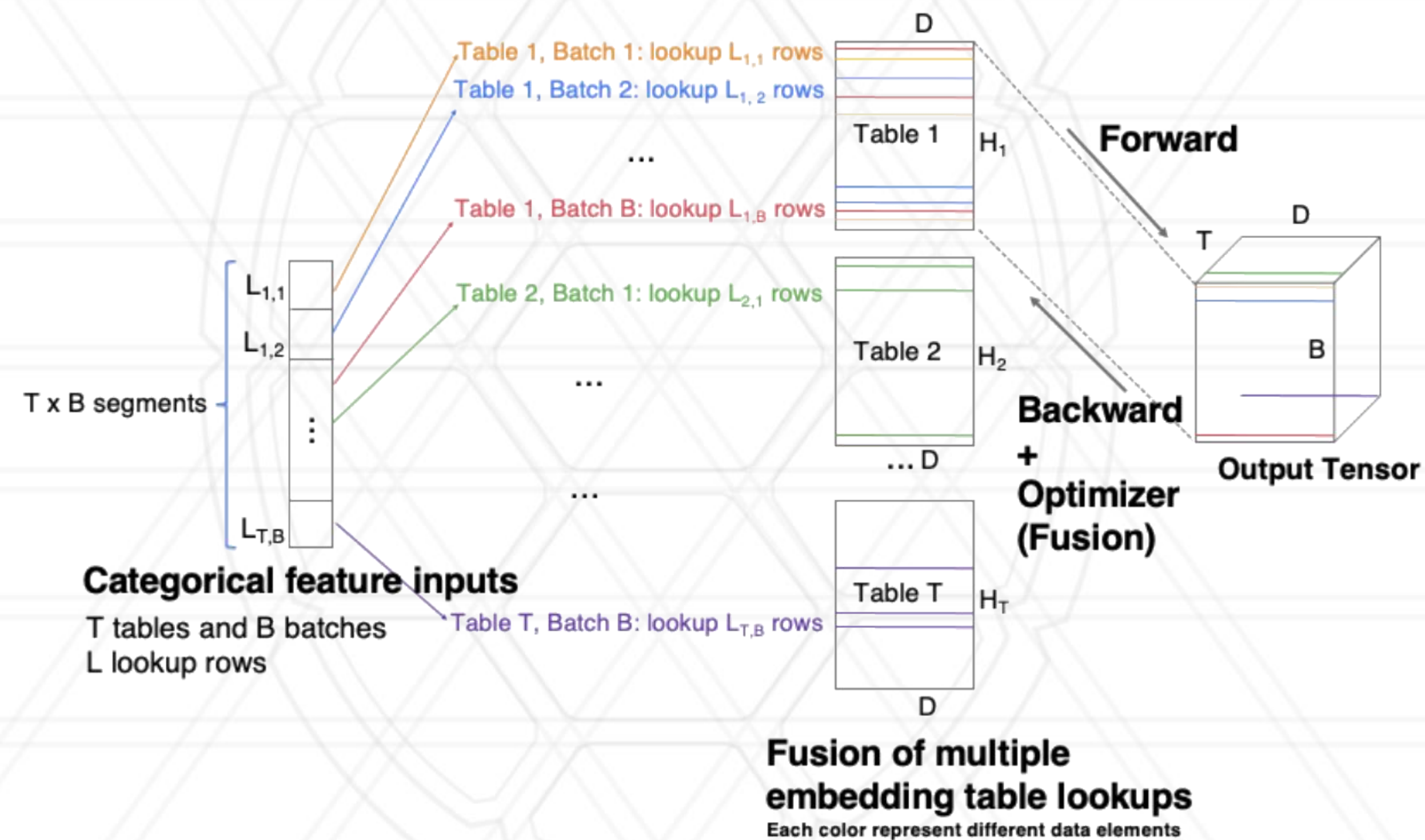
DEPARTMENT OF COMPUTER SCIENCE
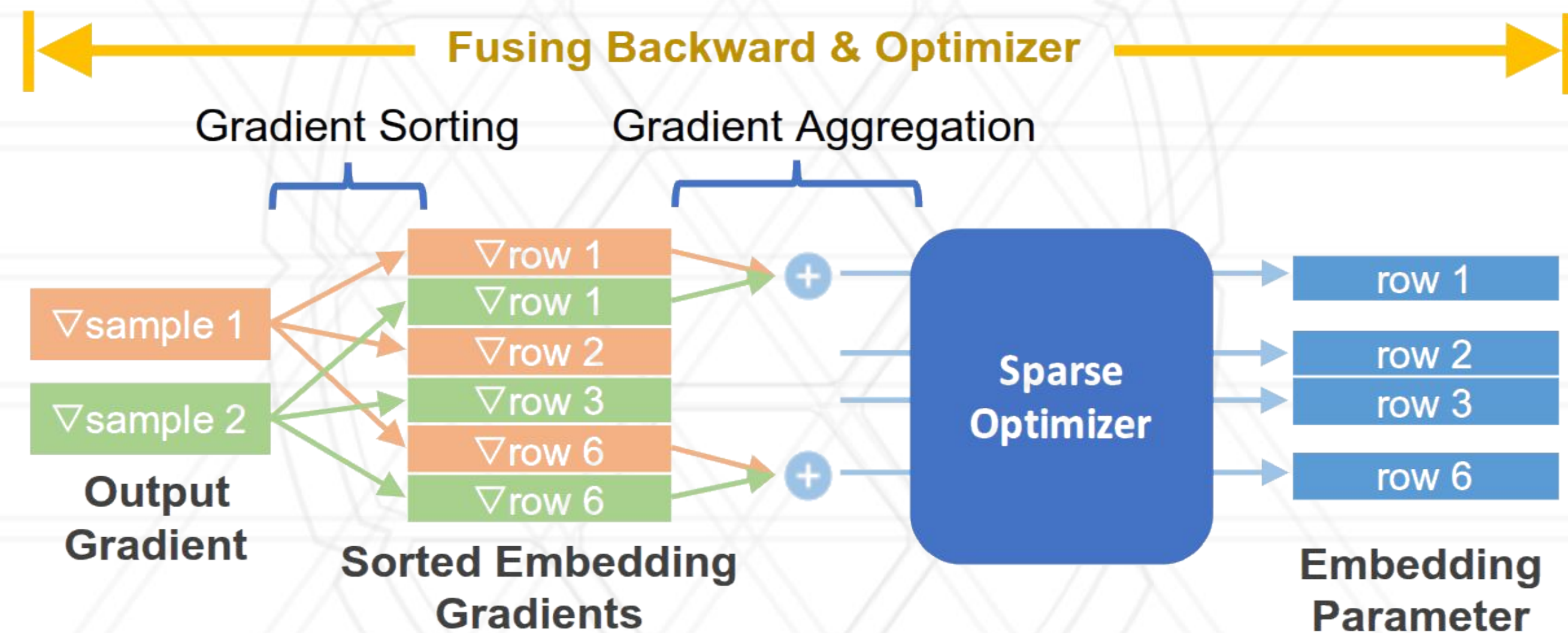
# Optimizing the Embedding Operators

- Two key inefficiencies

    - Each lookup is a single GPU kernel; incurs high overheads

    - Large tables need multi-GPU implementations

- Operator optimizations

    - Fused embeddings into single kernels; sort embedding gradients

    - Multi-GPU implementations

    - Co-Design with ZionEX to save memory

# Fused Embedding Operators

# Sort Gradients of Embedding Operators
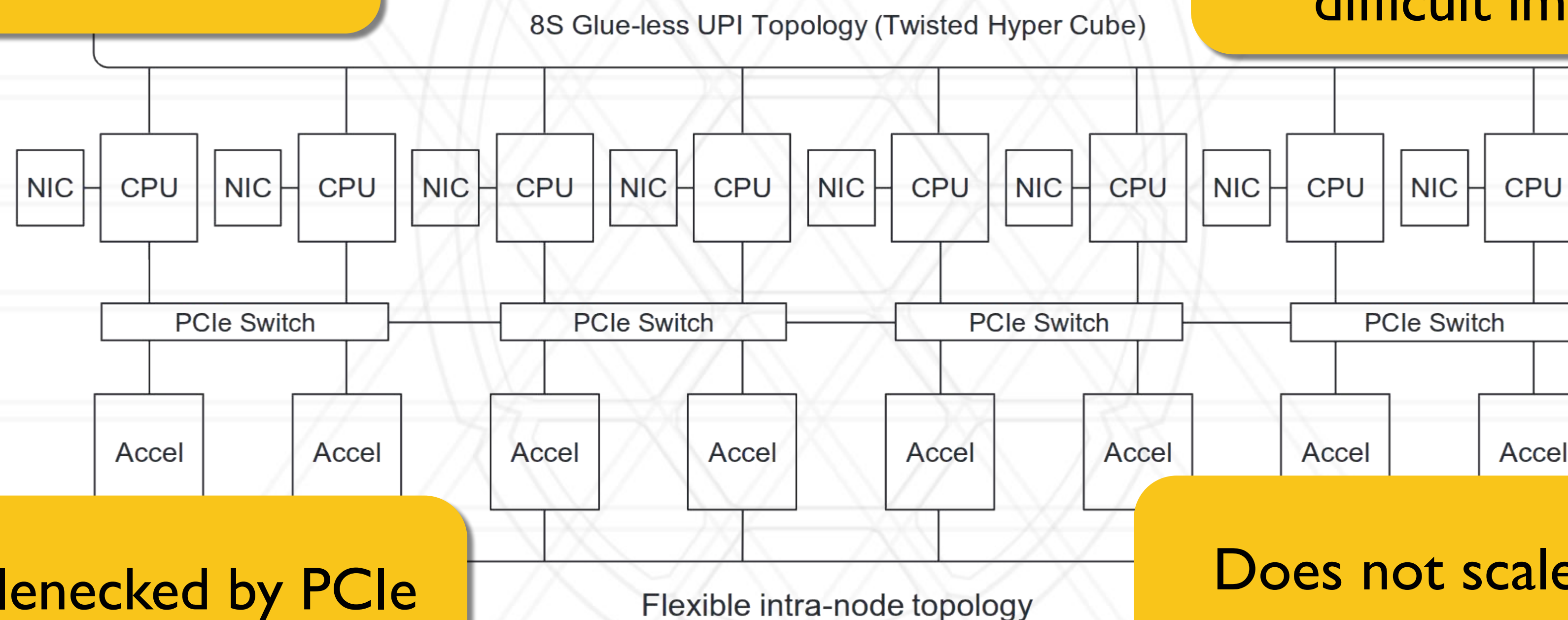
# Memory Optimizations

- Make use of device memory, host memory, and disk

- Access behavior is irregular

  - Default managed memory (like CUDA unified) experiences very poor performance

  - Implement custom cache in software

- Embedding compression

  - low precision (high precision cache and low precision embeddings)

  - sparse optimizers

# ZionEX: Co-Designing with Neo

- Previous system, Zion, had many limitations

Embedding on CPUs and MLP on GPUs

Balancing work requires careful design considerations and difficult implementation

8S Glue-less UPI Topology (Twisted Hyper Cube)

| NIC | CPU | NIC | CPU | NIC | CPU | NIC | CPU | NIC | CPU | NIC | CPU | NIC | CPU | NIC | CPU |

PCIe Switch  PCIe Switch  PCIe Switch  PCIe Switch

| Accel | Accel | Accel | Accel | Accel | Accel | Accel | Accel |

Flexible intra-node topology

GPUs are bottlenecked by PCIe for internode communication

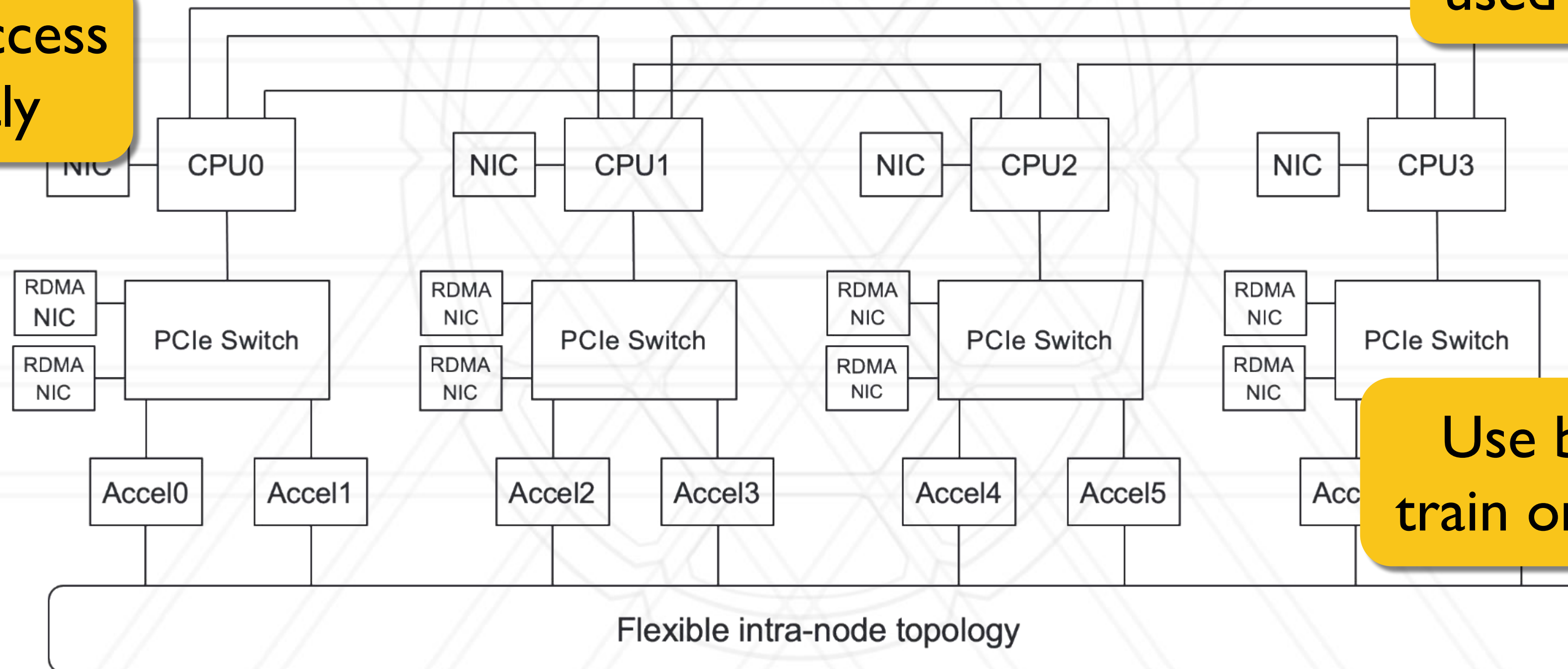Does not scale to multi-node very well!

DEPARTMENT OF
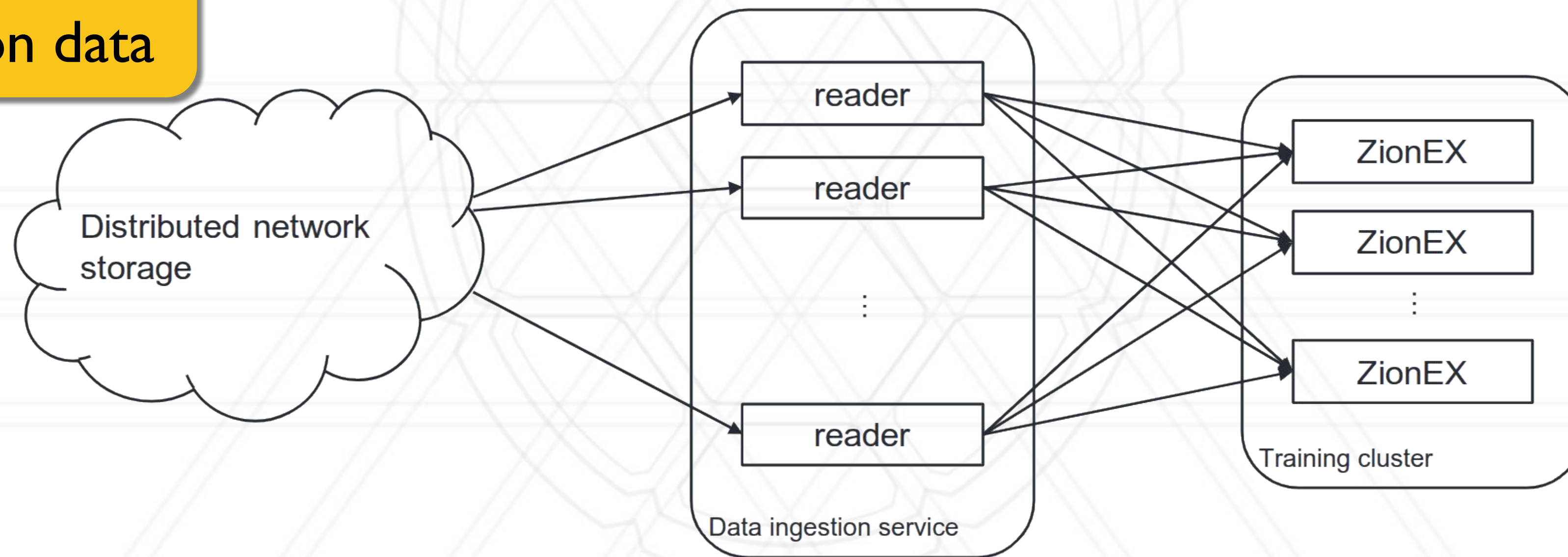COMPUTER SCIENCE

# ZionEX: Co-Designing with Neo

- ZionEx addresses these shortcomings
- Designed to support 4D parallelism and data requirements of DLRM

Separate NICs can be used for data ingestion

Allow GPUs to access network directly

Use bigger GPUs and train on them exclusively



Flexible intra-node topology

# ZionEX: Co-Designing with Neo

- ZionEx addresses these shortcomings
- Designed to support 4D parallelism and data requirements of DLRM
- Custom data ingestion servers to overcome latencies

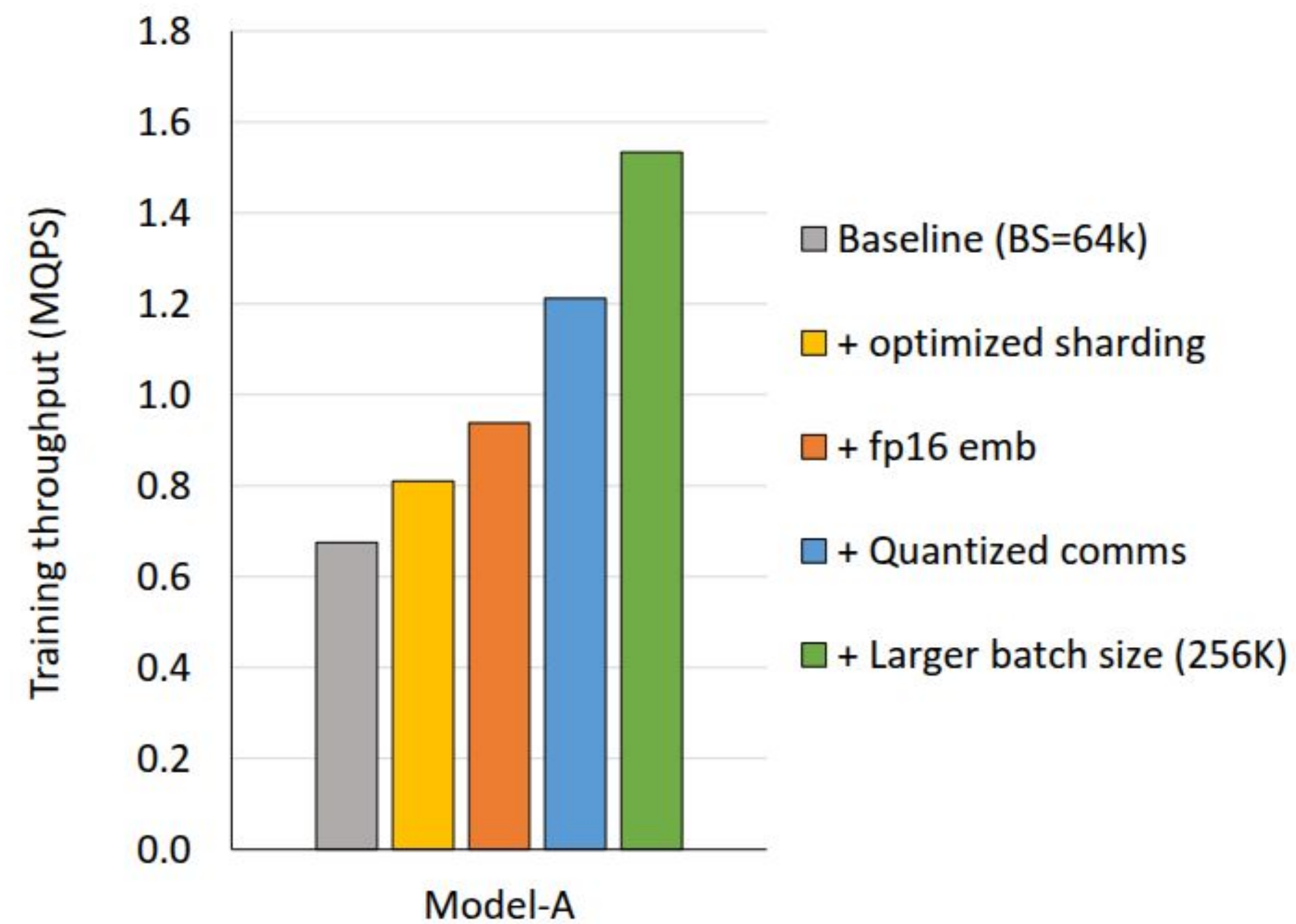Prevent ZionEx nodes from waiting on data
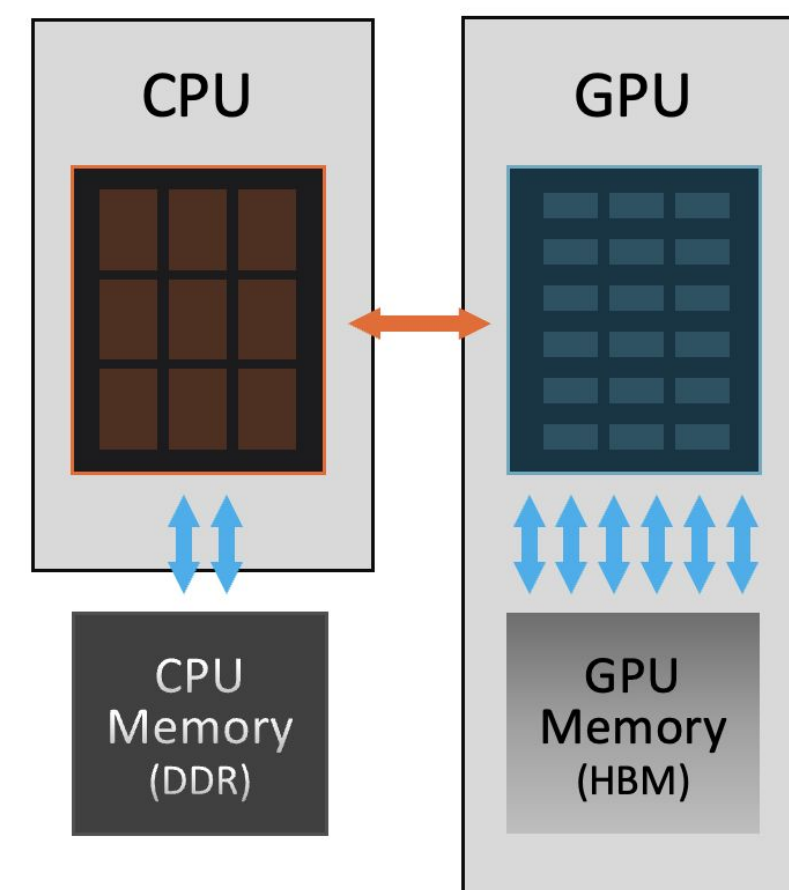
# Training Performance

# Other Training Optimizations

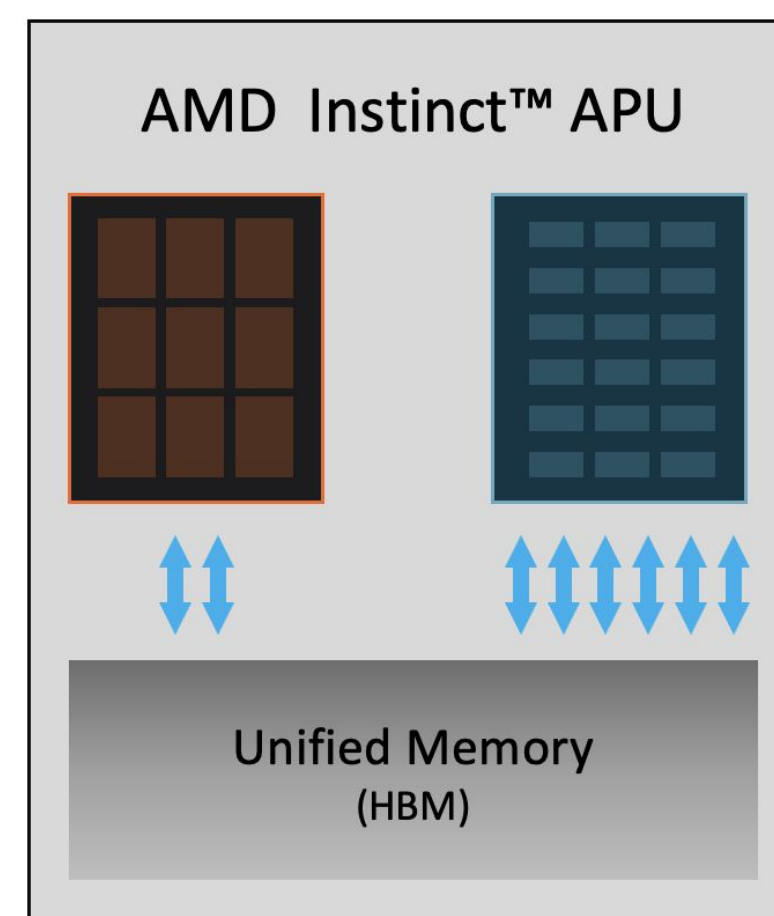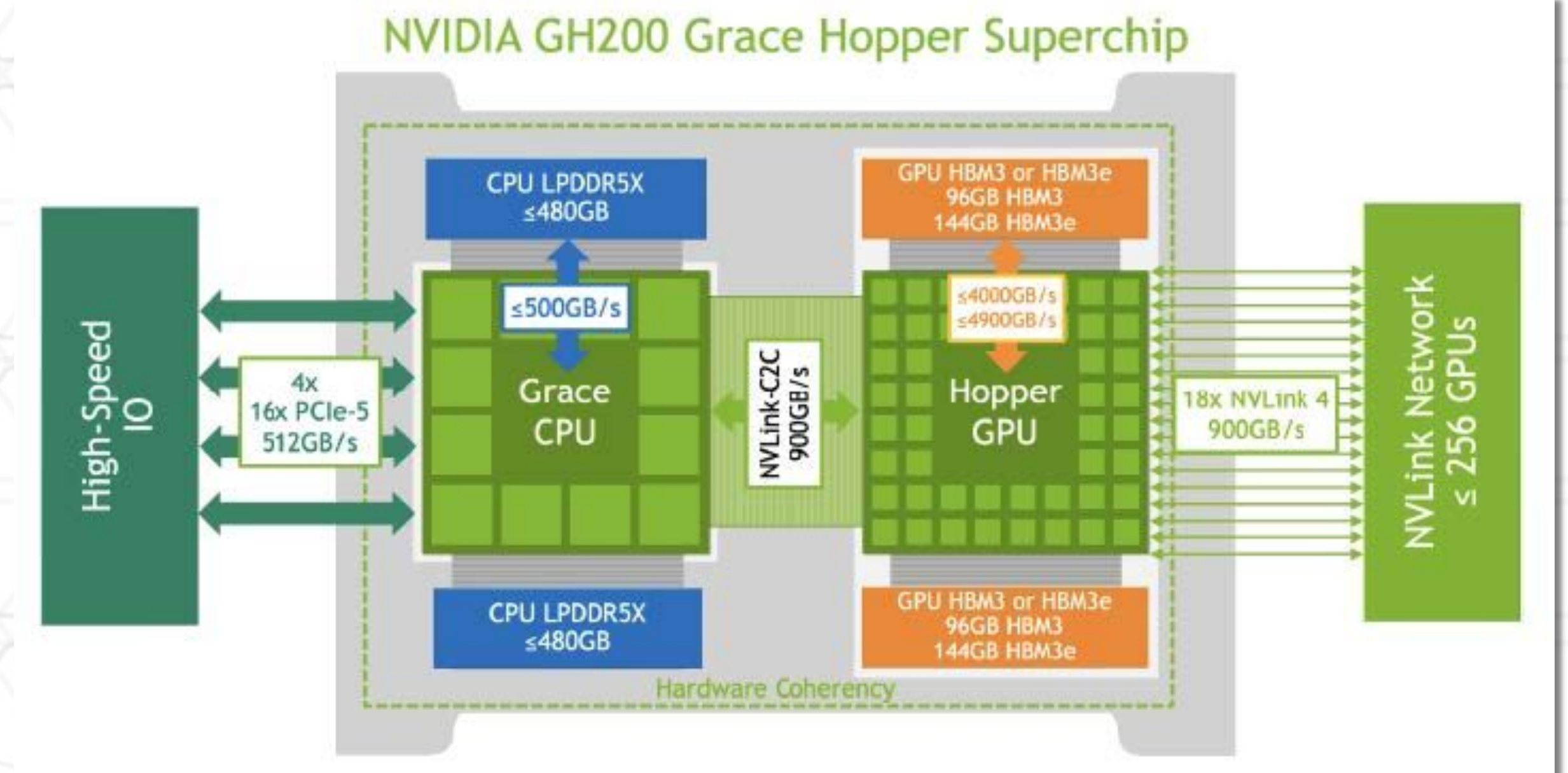# Shared Memory Nodes

MI300A and GH200 are combined CPU-GPU nodes produced by AMD and NVIDIA



(a) Discrete CPU and GPU.

(b) APU.