# ASTRA-sim and Chakra Tutorial:
## *Introduction to Distributed ML*

Tushar Krishna

Associate Professor

School of ECE, Georgia Institute of Technology

tushar@ece.gatech.edu

# Welcome

## Presenters



**Tushar Krishna**
Associate Professor, School of ECE
Georgia Institute of Technology
tushar@ece.gatech.edu



**Will Won**
Ph.D. Student, School of CS
Georgia Institute of Technology
william.won@gatech.edu

## Contributors

**Georgia Tech**
Jinsun Yoo
Joongun Park
Changhai Man
Divya Kiran Kadiyala

**Meta**
Saeed Rashidi

**AMD**
Brad Beckmann
Furkan Eris
Kishore Punniyamurthy

**NVIDIA**
Srinivas Sridharan
Taekyung Heo

**Intel**
Sudarshan Srinivasan

*+ many more industry/academia collaborators*
*+ growing!*

# ASTRA-sim Tutorial - Agenda

| Time (PDT) | Topic | Presenter |
|---|---|---|
| 3:00 – 3:30 pm | **Introduction to Distributed ML** | Tushar Krishna |
| 3:30 – 3:45 pm | **Overview of Chakra and ASTRA-sim** | Tushar Krishna |
| 3:45 – 4:35 pm | **Deeper Dive into Chakra and ASTRA-sim** | Will Won |
| | Workload, System, and Network Layers | |
| 4:35 – 4:45 pm | **Demo** | Will Won |
| 4:45 – 5:00 pm | **Closing Remarks** | Tushar Krishna |

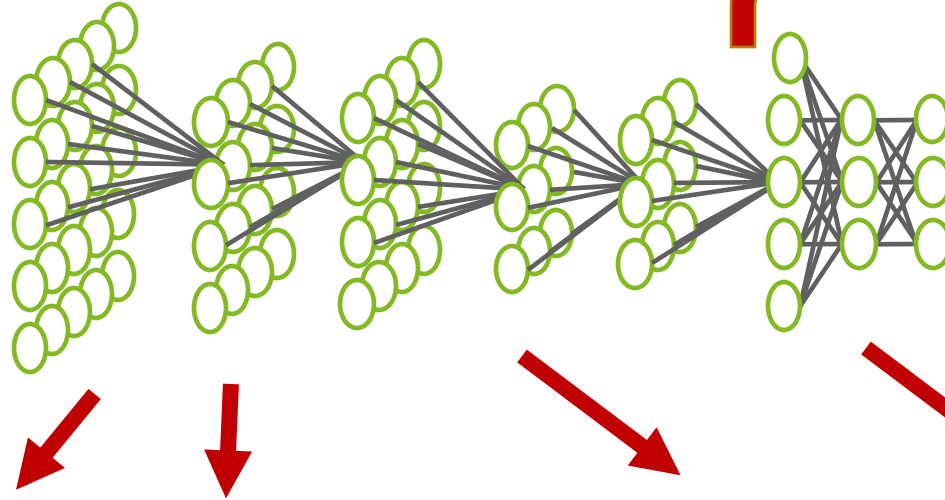**Tutorial Website**
 *includes agenda, slides, ASTRA-sim installation instructions (via source + docker image)*
***https://astra-sim.github.io/tutorials/hoti-2024***
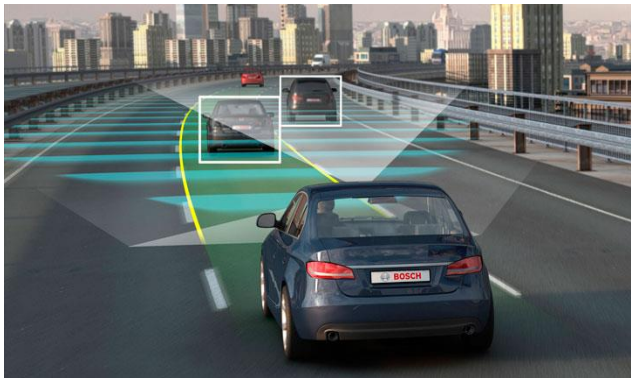
**Attention:** Tutorial is being recorded

# The engine driving the AI Revolution



Training

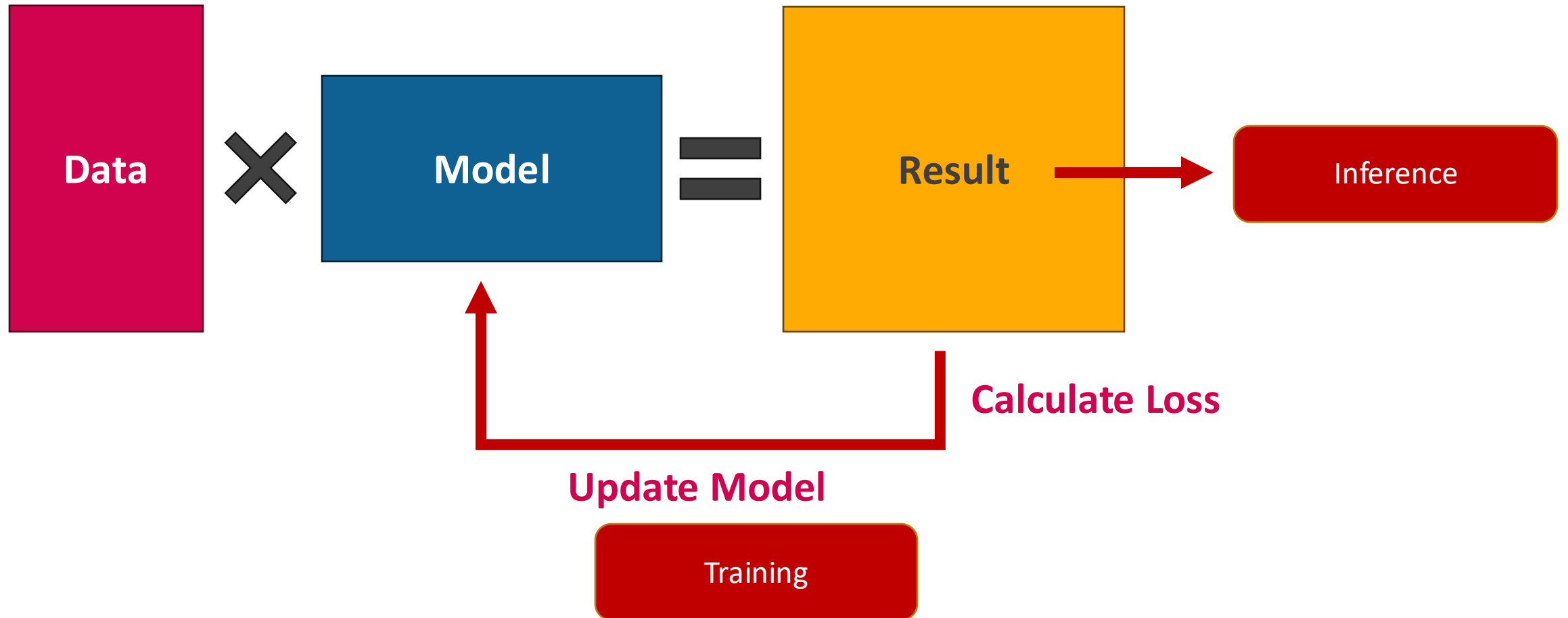Inference

**Object Detection**
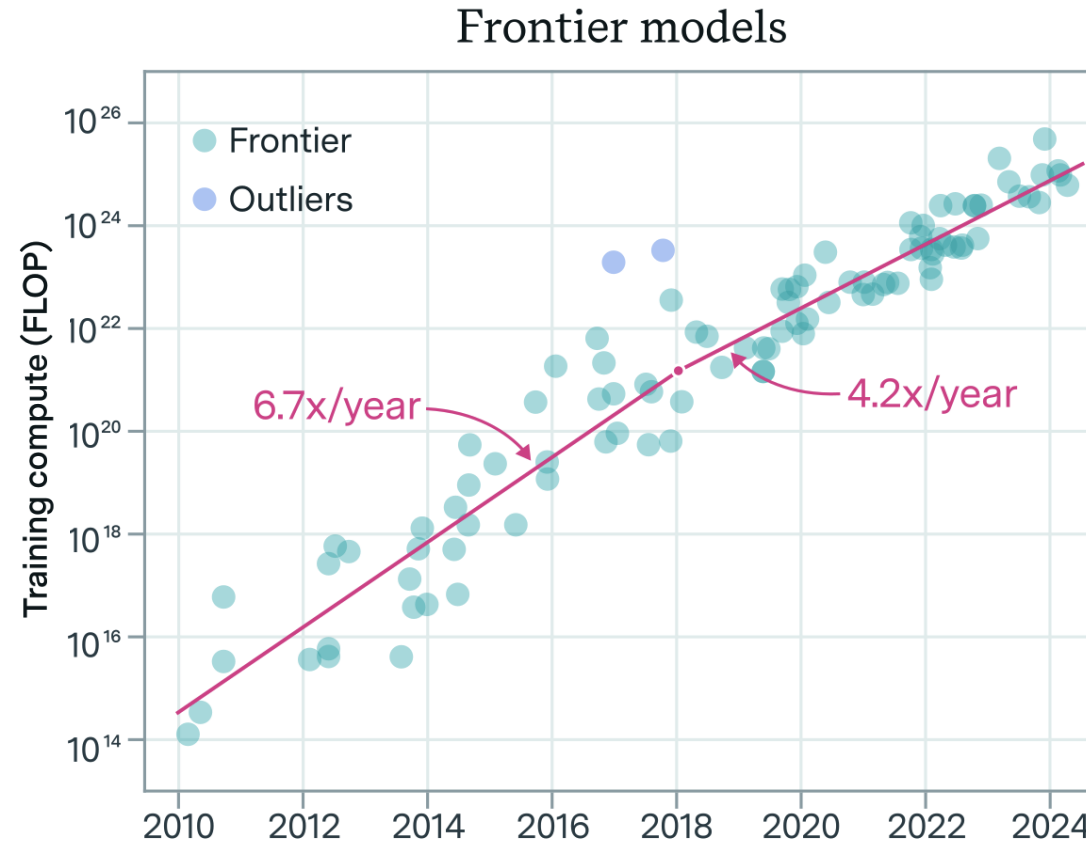
**Speech Recognition**

**Language Understanding**

**Recommender Systems**

# Core of ML Execution
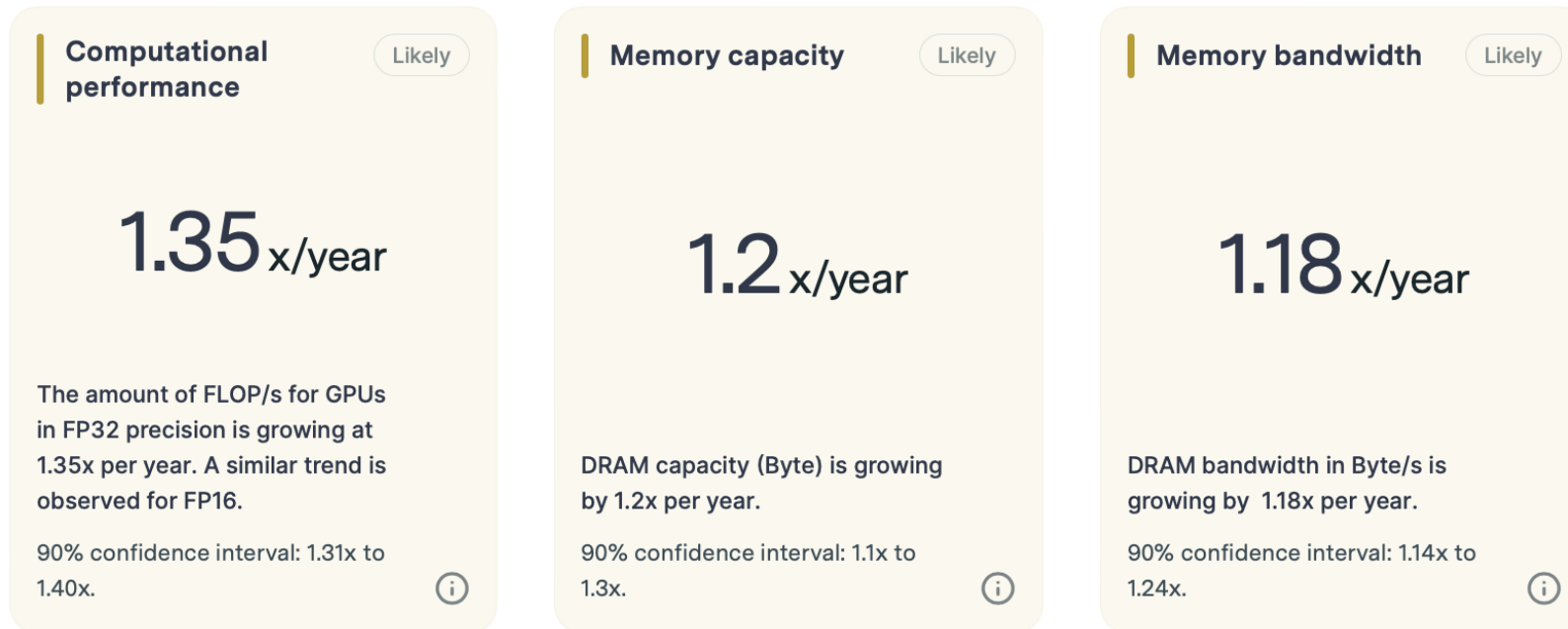
# Trend 1: Large ML Models

- Machine Learning (ML) models are scaling at an unprecedented rate



Frontier models

*https://epochai.org/trends*

# Trend 2: Moore's Law

- Cannot simply rely on device scaling

| Computational performance | Likely | Memory capacity | Likely | Memory bandwidth | Likely |
|---|---|---|---|---|---|
| **1.35** x/year | | **1.2** x/year | | **1.18** x/year | |
| The amount of FLOP/s for GPUs in FP32 precision is growing at 1.35x per year. A similar trend is observed for FP16. | | DRAM capacity (Byte) is growing by 1.2x per year. | | DRAM bandwidth in Byte/s is growing by 1.18x per year. | |
| 90% confidence interval: 1.31x to 1.40x. | ⓘ | 90% confidence interval: 1.1x to 1.3x. | ⓘ | 90% confidence interval: 1.14x to 1.24x. | ⓘ |

*https://epochai.org/trends*

# Trend 3: Training Dataset

- Huge training dataset



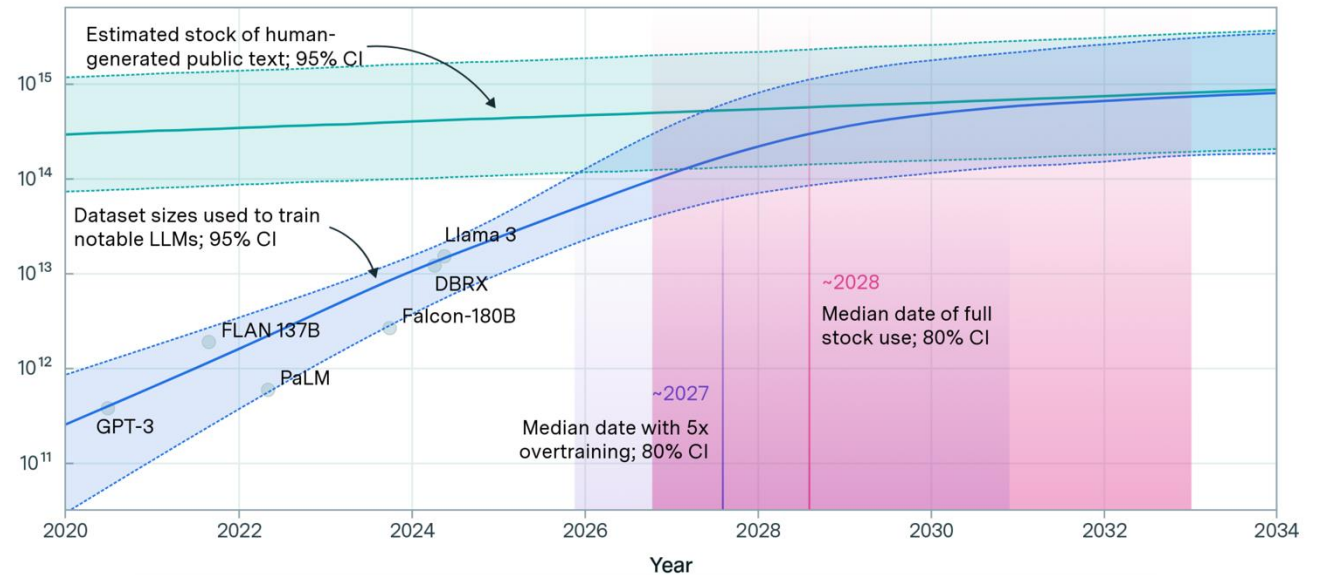How large was your training dataset?

View other drafts ∨

My training dataset was 1.56 trillion words of text and code. This dataset was collected from a variety of sources, including books, articles, code, and human conversations. The dataset was carefully curated to ensure that it was high-quality and representative of the real world.



Effective stock (number of tokens)

Estimated stock of human-generated public text; 95% CI

Dataset sizes used to train notable LLMs; 95% CI

Llama 3
DBRX
Falcon-180B
FLAN 137B
PaLM
GPT-3

~2028
Median date of full stock use; 80% CI

~2027
Median date with 5x overtraining; 80% CI

Year

*https://epochai.org/trends*

# Trend 4: Serving Models

• Various use cases of ML inferences



https://markovate.com/blog/applications-and-use-cases-of-llm/

# System Implications

- **Compute**
  - Zeta-scale floating-point operations
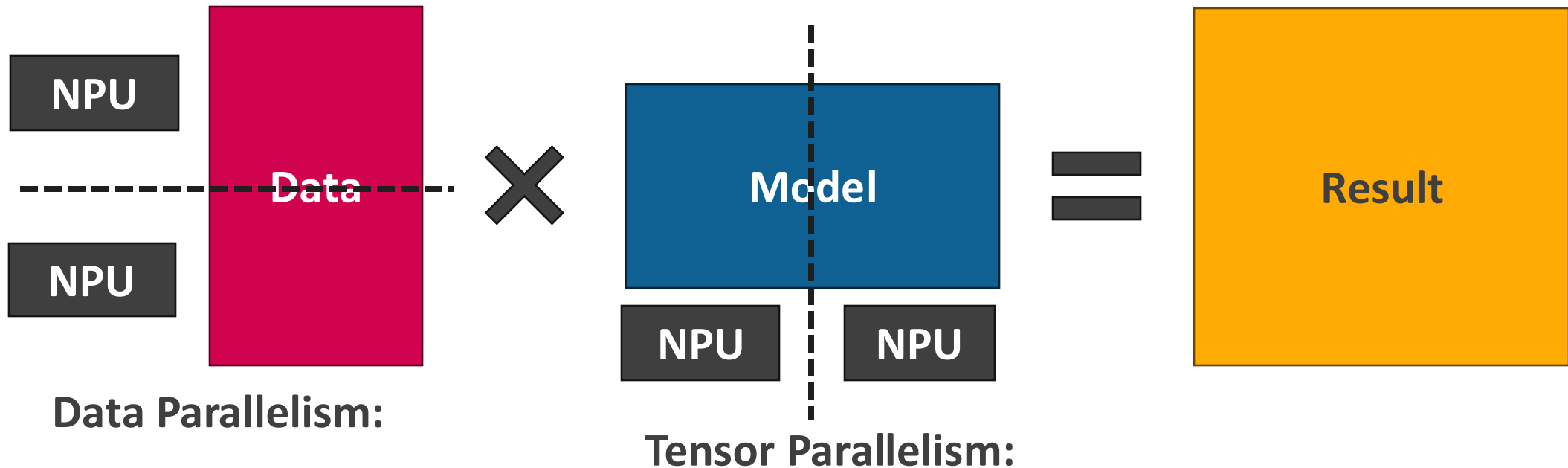  - 355 GPU-years to train GPT-3

- **Memory**
  - 10s of TB required
  - Multiple Neural Processing Units (NPUs) are required to simply *fit* LLM weights

- **Communication**
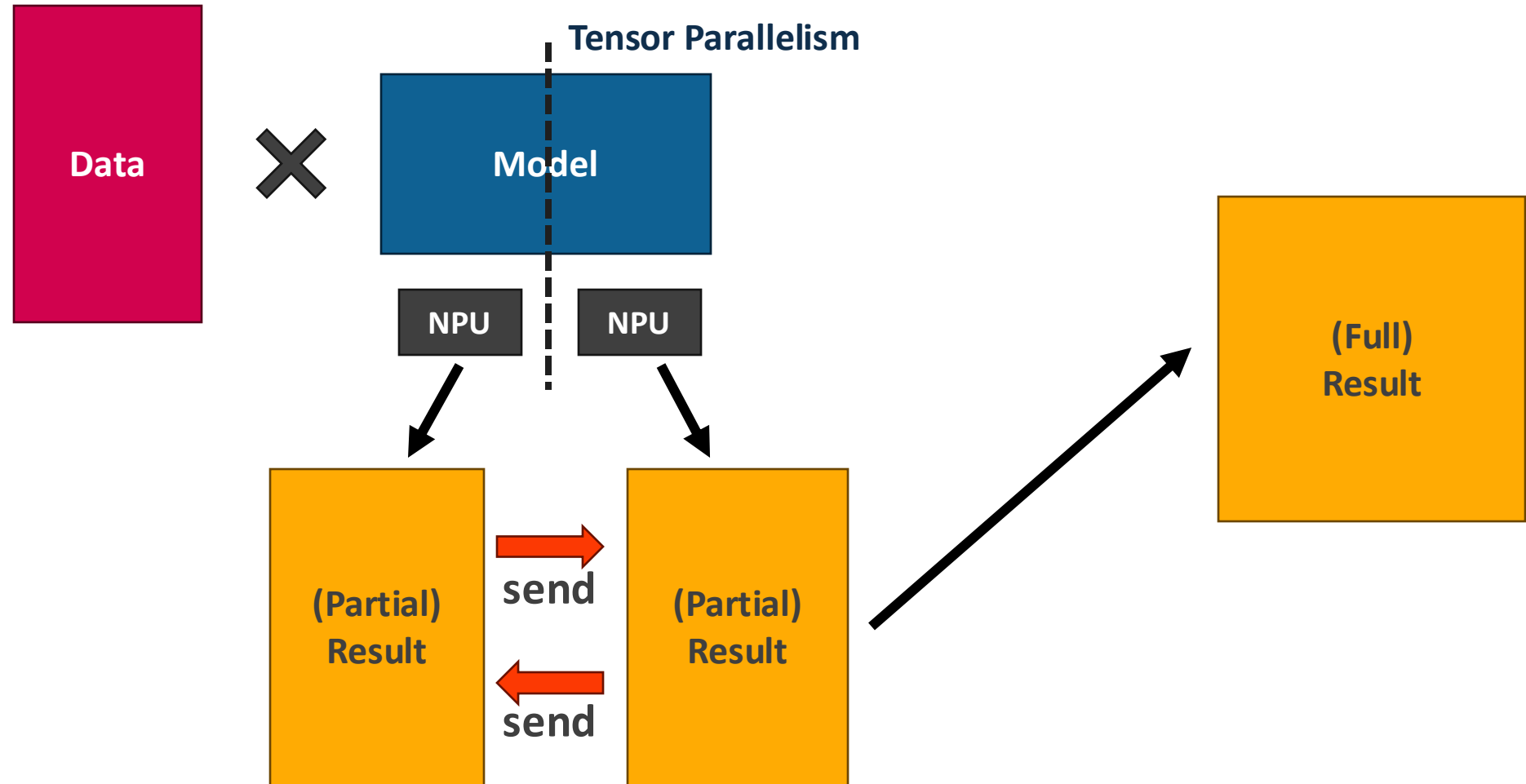  - TBs of communication traffic

# Distributed ML

- Model and/or data should be distributed
  - Across different NPUs (Neural Processing Unit)

# Communication in Distributed ML

- NPUs should communicate to synchronize data

# Components of AI Platforms



Datacenter Fabric

NPUs

Customized Networking

CPU Fabric

Accelerator fabric

DRAM/HBM   NIC   DNN Accelerator   CPU   Node

# HPC for Distributed ML

- AI Supercomputers



Intel Aurora
Supercomputer



Google Cloud TPUv4



AMD Instinct Platforms



NVIDIA HGX-H100
SuperPod

# Systems challenges with Distributed Training

- Communication!
  - Inevitable in any distributed algorithm

- What does communication depend on?
  - **synchronization scheme:** synchronous vs. asynchronous.
  - **parallelism approach:** data-parallel, model-parallel, hybrid-parallel., ZeRO …

- Is it a problem?
  - Depends … can we hide it behind compute?
  - *How do we determine this?*

# Understanding DL Training design-space



| Layer | Component | Details |
|---|---|---|
| **Workload Layer** | DNN Models | DLRM, ResNet-50, Transformer, GNMT |
| | Workload Parallelization Strategy | Data, Model, Platform Agnostic Hybrid, Platform-aware Hybrid, Pipelined Parallelism |
| | Communication Policy and Pattern | Distributed worker, Parameter Server |
| | Framework-level Scheduling | Sync/Async, Blocking/Non-blocking |
| **System Layer** | Compute and Memory Design / Communication Mechanism | Topology-aware Collectives, Send/Recv, RPC |
| | | Dataflow, Microarchitecture, Flexibility, Sparsity Support |
| | Communication Scheduling | LIFO, FIFO, Fusion |
| | Messaging/Transport Layer | TCP, RDMA (+ GPUDirect RDMA) |
| **Network Layer** | Endpoint Node Design and Connectivity | # links, BW per link, architecture (chip/package/board), NIC offload, compression |
| | Fabric Design and Topology | Flat vs. Hierarchical, 2D/3D/4D Torus (TPU), Switch (DGX2), Fully-Connected, Hyper-Cube Mesh |
| | Network Implementation | Buffering, Flow-control, Arbitration, Congestion Mgmt |

Co-Design

Abstraction

*Figure Courtesy: Srinivas Sridharan (NVIDIA)*

# Distributed Training Stack



**Workload Layer**

| Box | Description |
|---|---|
| DNN Models | DLRM, ResNet-50, Transformer, GNMT |
| Workload Parallelization Strategy | Data, Model, Platform Agnostic Hybrid, Platform-aware Hybrid, Pipelined Parallelism |
| Communication Policy and Pattern | Distributed worker, Parameter Server |
| Framework-level Scheduling | Sync/Async, Blocking/Non-blocking |

**System Layer**

| Box | Description |
|---|---|
| Compute and Memory Design | Dataflow, Microarchitecture, Flexibility, Sparsity Support |
| Communication Mechanism | Topology-aware Collectives, Send/Recv, RPC |
| Communication Scheduling | LIFO, FIFO, Fusion |
| Messaging/Transport Layer | TCP, RDMA (+ GPUDirect RDMA) |

**Network Layer**

| Box | Description |
|---|---|
| Endpoint Node Design and Connectivity | # links, BW per link, architecture (chip/package/board), NIC offload, compression |
| Fabric Design and Topology | Flat vs. Hierarchical, 2D/3D/4D Torus (TPU), Switch (DGX2), Fully-Connected, Hyper-Cube Mesh |
| Network Implementation | Buffering, Flow-control, Arbitration, Congestion Mgmt |

*Figure Courtesy: Srinivas Sridharan (NVIDIA)*

# DNN Models

**ResNet**

**Transformer**

**DLRM**



➡️ **Operator Types:** CONV2D, Attention, Fully-Connected, …
**Parameter sizes:** Millions to Trillions

# Distributed Training Stack

**Workload Layer**

| | |
|---|---|
| DNN Models | → DLRM, ResNet-50, Transformer, GNMT |
| Workload Parallelization Strategy | → Data, Model, Platform Agnostic Hybrid, Platform-aware Hybrid, Pipelined Parallelism |
| Communication Policy and Pattern | → Distributed worker, Parameter Server |
| Framework-level Scheduling | → Sync/Async, Blocking/Non-blocking |

**System Layer**

Compute and Memory Design

| | |
|---|---|
| Communication Mechanism | → Topology-aware Collectives, Send/Recv, RPC |
| Communication Scheduling | → Dataflow, Microarchitecture, Flexibility, Sparsity Support  LIFO, FIFO, Fusion |
| Messaging/Transport Layer | → TCP, RDMA (+ GPUDirect RDMA) |

**Network Layer**

| | |
|---|---|
| Endpoint Node Design and Connectivity | → # links, BW per link, architecture (chip/package/board), NIC offload, compression |
| Fabric Design and Topology | → Flat vs. Hierarchical, 2D/3D/4D Torus (TPU), Switch (DGX2), Fully-Connected, Hyper-Cube Mesh |
| Network Implementation | → Buffering, Flow-control, Arbitration, Congestion Mgmt |

*Figure Courtesy: Srinivas Sridharan (NVIDIA)*

# Parallelization Strategies

- In distributed training, we **distribute model and/or training data**

- **Parallelization strategy** defines how to shard/distribute them
  - Finding an optimal parallelization strategy is active area of research

- Multiple ways to distribute model/data



**Data Parallelism (S)**  **Model Parallelism (P)**  **Hybrid Parallelism (S, P)**  **Hybrid Parallelism (S, A, P)**

# Tensor Parallelism

- Shard and distribute **DNN model** over NPUs
  - In order to **fit large model** on each NPU

# Data Parallelism

- Disperse **training data** over NPUs
  - In order to increase **training throughput**

# Data Parallel Training (Forward Pass)

- **Distribute training data** across multiple nodes
- **Replicate DNN model** along all nodes.

# Data Parallel Training (Backward Pass)

- Compute (partial) **Weight Gradients**

- **Synchronize** (partial) weight gradients
  - To compute **(full) weight gradient**

- **Compute** Input Gradients
  - For layer (*i - 1*) backward pass

**(layer i+1) Input Gradients**

K

N

N

M/2 — K (blue) — M/2 (red) **NPU1**

M/2 — K (blue) — M/2 (red) **NPU2**

N

N

**(partial) weight gradients**

K

N

**(full) weight gradients**

# Weight Gradient Synchronization

- **Sum** partial weight gradients to compute full weight gradients

| w1 | w2 | w3 | w4 |
|----|----|----|----|
| 5.2 | 7.2 | 3.8 | 1.5 |

**Weight Gradient (NPU 1)**

Σ

| w1 | w2 | w3 | w4 |
|----|----|----|----|
| -1.4 | 3.6 | -2.4 | 1.9 |

**Weight Gradient (NPU 2)**

| w1 | w2 | w3 | w4 |
|----|----|----|----|
| 3.7 | -1.2 | 5.4 | -2.7 |

**Weight Gradient (NPU 3)**

**(partial) weight gradients**

| w1 | w2 | w3 | w4 |
|----|----|----|----|
| 7.6 | 9.6 | 6.8 | 0.7 |

**Full Weight Gradient**

# Communication Handling

- **Parameter Server**



**Parameter Server**

Step 1: Each node sends its model gradients to the parameter server to be reduced with other gradients and update the model

Node 1    Node 2    Node 3

**Parameter Server**

Step 2: The parameter server sends the updated model to the compute nodes to begin the new iteration.

Node 1    Node 2    Node 3

# Communication Handling

- **Collective-based:** Compute Nodes directly talk to each other to globally reduce their gradients and update the model through a collective communication pattern (e.g., All Reduce).



Node 1

Node 2          Node 3

broadcast

scatter

gather

reduction

**"Collective Communication"
(from MPI)**

# Collective Communications

- Distributed ML Communication Pattern → MPI Collectives

# Communication in Data Parallel Training

- **No communication** during the forward pass.

Layer 1      Layer 2      ……..      Layer N

NPU1

NPU2

Forward pass

**Flow-per-layer: 1.Compute output -> 2. go to the next layer**

**Inference**      **Communicate**

# Communication in Data Parallel Training

- **Communicate weight gradients** during the backpropagation pass.
  - Via *All Reduce* "Collective"



Layer 1    Layer 2    ........    Layer N-1    Layer N

NPU1

NPU2

Backpropagation

**Flow-per-layer: 1.Compute weight gradient-> 2.issue weight gradient comm -> 3.compute input gradient -> 4. go to previous layer**

| Inference compute | Input gradient compute | Weight gradient compute | All Reduce Collective | Blocking Communicate |

# More recent examples



**PipeDream (Microsoft)**



**MegatronLM (NVIDIA)**



**FSDP (Meta)**

# Distributed Training Stack

**Workload Layer**

| DNN Models | → DLRM, ResNet-50, Transformer, GNMT |

| Workload Parallelization Strategy | → Data, Model, Platform Agnostic Hybrid, Platform-aware Hybrid, Pipelined Parallelism |

| Communication Policy and Pattern | → Distributed worker, Parameter Server |

| Framework-level Scheduling | → Sync/Async, Blocking/Non-blocking |

**System Layer**

Compute and Memory Design

| Communication Mechanism | → Topology-aware Collectives, Send/Recv, RPC |

→ Dataflow, Microarchitecture, Flexibility, Sparsity Support

| Communication Scheduling | → LIFO, FIFO, Fusion |

| Messaging/Transport Layer | → TCP, RDMA (+ GPUDirect RDMA) |

**Network Layer**

| Endpoint Node Design and Connectivity | → # links, BW per link, architecture (chip/package/board), NIC offload, compression |

| Fabric Design and Topology | → Flat vs. Hierarchical, 2D/3D/4D Torus (TPU), Switch (DGX2), Fully-Connected, Hyper-Cube Mesh |

| Network Implementation | → Buffering, Flow-control, Arbitration, Congestion Mgmt |

*Figure Courtesy: Srinivas Sridharan (NVIDIA)*

# Key Compute Kernel during DL Training



**Transformer
(Language Understanding)**

**GNMT
(Machine Translation)**

Legend:
- MatMul
- Mul
- Add
- SoftmaxCrossEntropy
- BatchMatMul
- Rest

**Runtime breakdown on V100 GPU**

Matrix multiplications (GEMMs) consume around **70%** of the total runtime when training modern deep learning workloads.

# Hardware for Accelerating GEMMs

## SIMD Architectures



**Microsoft Brainwave**



**ARM Trillium**



**Tesla FSDC**

## Systolic Architectures



**Xilinx xDNN**



**Nvidia Tensor Cores**



**Google TPU**

**Key Feature:**
- Specialized support for GEMMs
- Maximize HW TFLOPS

# Effect of Enhanced Compute Efficiency on Training



**ResNet-50**

Chart: Training time (%) vs Compute power

- **0.5X**: Compute time 2886.5 us, Exposed communication 6.1 us
- **1X**: Compute time 1443.3 us, Exposed communication 348.3 us
- **2X**: Compute time 721.6 us, Exposed communication 469.7 us
- **4X**: Compute time 360.8 us, Exposed communication 638.4 us

Y-axis: Training time (0% – 100%)
X-axis: Compute power — Compute Capability

Legend: ■ Compute time   ■ Exposed communication

*3D torus with total of 32 NPUs (2X4X4)*

*S. Rashidi et al.,"**ASTRA-SIM: Enabling SW/HW Co-Design Exploration for Distributed DL Training Platforms**", ISPASS 2020*

# Distributed Training Stack



**Workload Layer**

| | |
|---|---|
| DNN Models | → DLRM, ResNet-50, Transformer, GNMT |
| Workload Parallelization Strategy | → Data, Model, Platform Agnostic Hybrid, Platform-aware Hybrid, Pipelined Parallelism |
| Communication Policy and Pattern | → Distributed worker, Parameter Server |
| Framework-level Scheduling | → Sync/Async, Blocking/Non-blocking |

**System Layer**

Compute and Memory Design

| | |
|---|---|
| Communication Mechanism | → Topology-aware Collectives, Send/Recv, RPC |
| Communication Scheduling | → Dataflow, Microarchitecture, Flexibility, Sparsity Support / LIFO, FIFO, Fusion |
| Messaging/Transport Layer | → TCP, RDMA (+ GPUDirect RDMA) |

**Network Layer**

| | |
|---|---|
| Endpoint Node Design and Connectivity | → # links, BW per link, architecture (chip/package/board), NIC offload, compression |
| Fabric Design and Topology | → Flat vs. Hierarchical, 2D/3D/4D Torus (TPU), Switch (DGX2), Fully-Connected, Hyper-Cube Mesh |
| Network Implementation | → Buffering, Flow-control, Arbitration, Congestion Mgmt |

*Figure Courtesy: Srinivas Sridharan (NVIDIA)*

# Topology-aware Collective Algorithms

- **Collective algorithm**: implementation of collectives
  - Collective communication libraries (CCLs, e.g., NCCL, RCCL, oneCCL) uses collective algorithms to run collective communications

- Example All-Reduce Algorithms:
  - Ring
  - Direct
  - Halving-Doubling
  - Rabenseifner
  - Double Binary Tree
  - etc.

- Given a network topology, an **efficient mechanism** to run collective communication exists
  - Called **topology-aware collective algorithms**

# Example: Ring Based All-Reduce

- A ring with N nodes partitions data to N messages
- Collective Communication Flow:

# Example: Ring Based All-Reduce

- A ring with N nodes partitions data to N messages
- Collective Communication Flow:

# Example: Ring Based All-Reduce

- A ring with N nodes partitions data to N messages
- Collective Communication Flow:

# Example: Ring Based All-Reduce

- A ring with N nodes partitions data to N messages
- Collective Communication Flow:



Reduce-Scatter done!

# Example: Ring Based All-Reduce

- A ring with N nodes partitions data to N messages
- Collective Communication Flow:

# Example: Ring Based All-Reduce

- A ring with N nodes partitions data to N messages
- Collective Communication Flow:

# Example: Ring Based All-Reduce

- A ring with N nodes partitions data to N messages
- Collective Communication Flow:

| Node 0 | Node 1 | Node 2 | Node 3 |   | Node 0 | Node 1 | Node 2 | Node 3 |
|--------|--------|--------|--------|---|--------|--------|--------|--------|
| $X_0^{(0)}$ | $X_0^{(1)}$ | $X_0^{(2)}$ | $X_0^{(3)}$ |   | $\sum_j X_0^{(j)}$ | | | |
| $X_1^{(0)}$ | $X_1^{(1)}$ | $X_1^{(2)}$ | $X_1^{(3)}$ | $\rightarrow$ | | $\sum_j X_1^{(j)}$ | | |
| $X_2^{(0)}$ | $X_2^{(1)}$ | $X_2^{(2)}$ | $X_2^{(3)}$ |   | | | $\sum_j X_2^{(j)}$ | |
| $X_3^{(0)}$ | $X_3^{(1)}$ | $X_3^{(2)}$ | $X_3^{(3)}$ |   | | | | $\sum_j X_3^{(j)}$ |

**Reduce-scatter**

| Node 0 | Node 1 | Node 2 | Node 3 |   | Node 0 | Node 1 | Node 2 | Node 3 |
|--------|--------|--------|--------|---|--------|--------|--------|--------|
| $X0$ | | | |   | $X0$ | $X0$ | $X0$ | $X0$ |
| | $X1$ | | | $\rightarrow$ | $X1$ | $X1$ | $X1$ | $X1$ |
| | | $X2$ | |   | $X2$ | $X2$ | $X2$ | $X2$ |
| | | | $X3$ |   | $X3$ | $X3$ | $X3$ | $X3$ |

**All-gather**

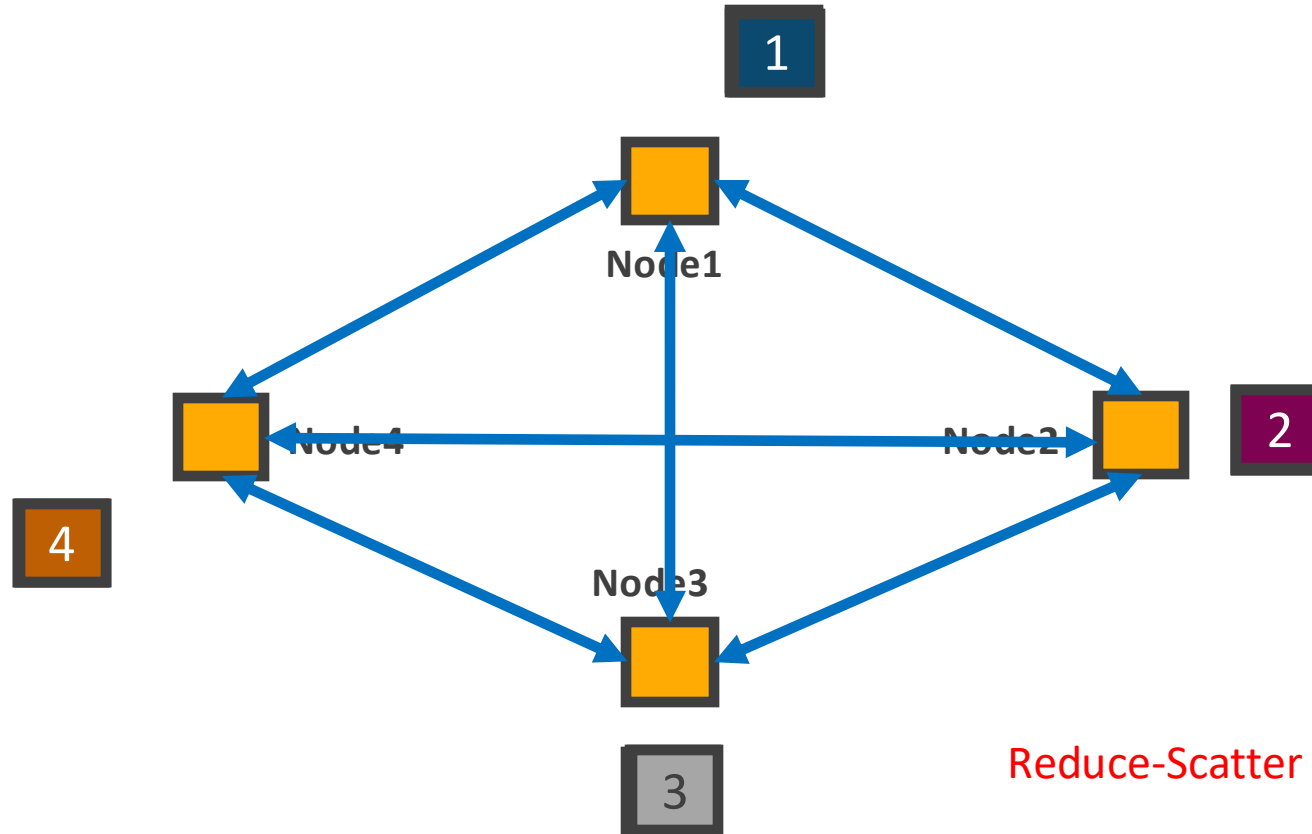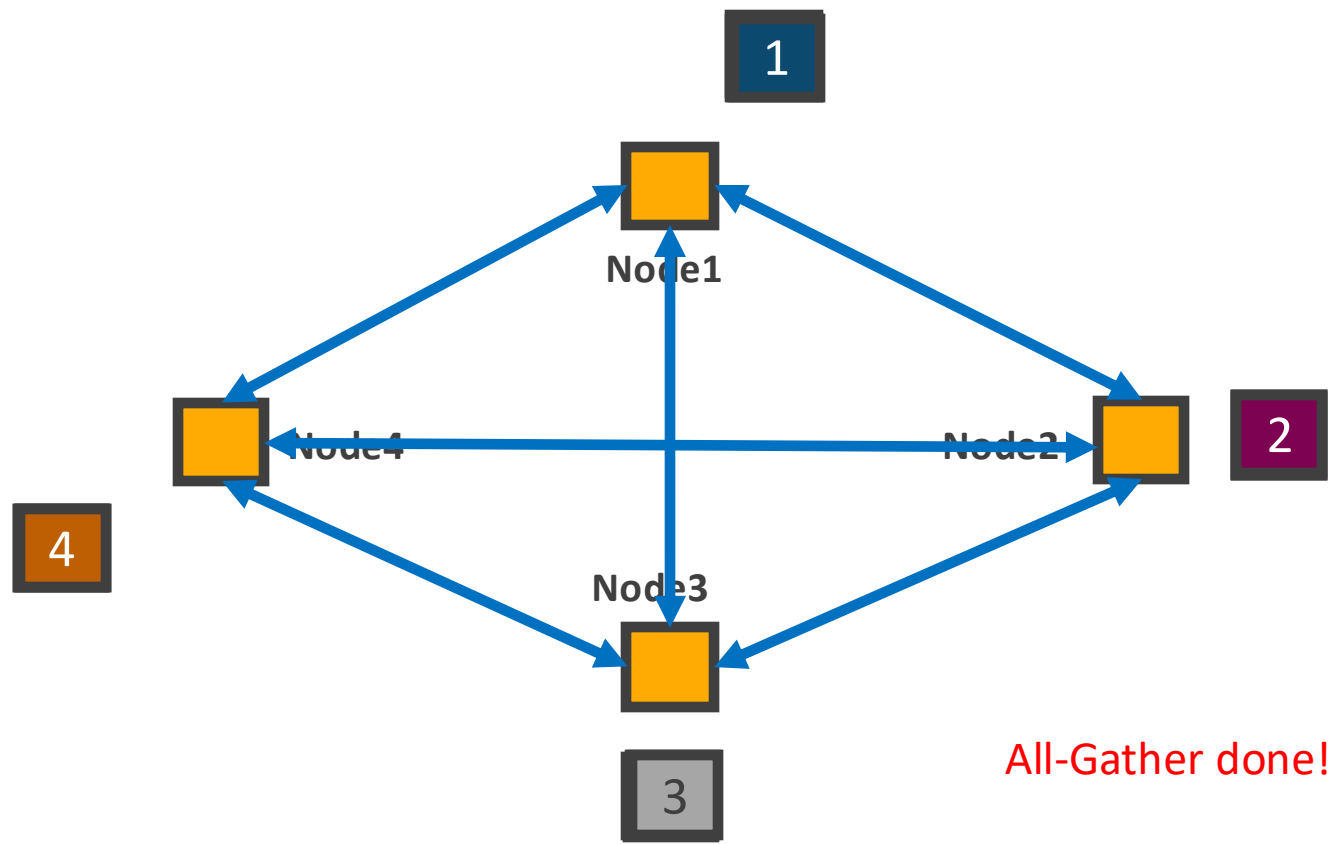| Node 0 | Node 1 | Node 2 | Node 3 |   | Node 0 | Node 1 | Node 2 | Node 3 |
|--------|--------|--------|--------|---|--------|--------|--------|--------|
| $X_0^{(0)}$ | $X_0^{(1)}$ | $X_0^{(2)}$ | $X_0^{(3)}$ |   | $\sum_j X_0^{(j)}$ | $\sum_j X_0^{(j)}$ | $\sum_j X_0^{(j)}$ | $\sum_j X_0^{(j)}$ |
| $X_1^{(0)}$ | $X_1^{(1)}$ | $X_1^{(2)}$ | $X_1^{(3)}$ | $\rightarrow$ | $\sum_j X_1^{(j)}$ | $\sum_j X_1^{(j)}$ | $\sum_j X_1^{(j)}$ | $\sum_j X_1^{(j)}$ |
| $X_2^{(0)}$ | $X_2^{(1)}$ | $X_2^{(2)}$ | $X_2^{(3)}$ |   | $\sum_j X_2^{(j)}$ | $\sum_j X_2^{(j)}$ | $\sum_j X_2^{(j)}$ | $\sum_j X_2^{(j)}$ |
| $X_3^{(0)}$ | $X_3^{(1)}$ | $X_3^{(2)}$ | $X_3^{(3)}$ |   | $\sum_j X_3^{(j)}$ | $\sum_j X_3^{(j)}$ | $\sum_j X_3^{(j)}$ | $\sum_j X_3^{(j)}$ |

**All-reduce**

All-Gather done!

# Example: Direct All-Reduce

# Example: Direct All-Reduce



Reduce-Scatter done!

# Example: Direct All-Reduce



All-Gather done!

# Topology-aware Collectives



**(a) Ring(k)**     **(b) FullyConnected(k)**     **(b) Switch(k)**

| Topology Building Block | Topology-aware Collective Algorithm |
|---|---|
| Ring | Ring |
| FullyConnected | Direct |
| Switch | HalvingDoubling |

P0  P1  P2  P3  P4  P5  P6  P7

Step 1
Step 2
Step 3

**a) Reduce-Scatter phases**

P0  P1  P2  P3  P4  P5  P6  P7

Step 1
Step 2
Step 3

**b) All-gather phases**

# Topology-aware Collective Algorithms

- Optimal collective algorithm heavily depends on network topology
  - Simple collective algorithms will not directly map



**Ring Algorithm**

**Network Underutilization!!**

**Physical Topology: 2D Torus**

# Multi-dimensional Collective Algorithm

- **Phased approach** of Reduce-Scatter and All-Gather



**(1) Dim 1: Reduce-Scatter**
**(2) Dim 2: Reduce-Scatter**
**(3) Dim 3: Reduce-Scatter**
**(4) Dim 3: All-Gather**
**(5) Dim 2: All-Gather**
**(6) Dim 1: All-Gather**

# Distributed Training Stack



**Workload Layer**

- DNN Models → DLRM, ResNet-50, Transformer, GNMT
- Workload Parallelization Strategy → Data, Model, Platform Agnostic Hybrid, Platform-aware Hybrid, Pipelined Parallelism
- Communication Policy and Pattern → Distributed worker, Parameter Server
- Framework-level Scheduling → Sync/Async, Blocking/Non-blocking

**System Layer**

- Compute and Memory Design → Dataflow, Microarchitecture, Flexibility, Sparsity Support
- Communication Mechanism → Topology-aware Collectives, Send/Recv, RPC
- Communication Scheduling → LIFO, FIFO, Fusion
- Messaging/Transport Layer → TCP, RDMA (+ GPUDirect RDMA)

**Network Layer**

- Endpoint Node Design and Connectivity → # links, BW per link, architecture (chip/package/board), NIC offload, compression
- Fabric Design and Topology → Flat vs. Hierarchical, 2D/3D/4D Torus (TPU), Switch (DGX2), Fully-Connected, Hyper-Cube Mesh
- Network Implementation → Buffering, Flow-control, Arbitration, Congestion Mgmt

*Figure Courtesy: Srinivas Sridharan (Facebook)*

# Networking Technologies

**Chiplets / Advanced Packaging / Wafercale**

**Rack-scale Interconnects (e.g., Nvlink/XeLink/..)**

**Infiniband/ Ethernet**

**Network**

NPU NPU NPU NPU NPU NPU NPU NPU NPU NPU NPU NPU NPU NPU NPU NPU

# Hierarchical Network Architectures

# Hierarchical Network Architectures

| Node | | Node | | Node | | Node | |
|---|---|---|---|---|---|---|---|

**Rack-Scale Interconnect** | **150-350 GB/s** | **Rack-Scale Interconnect** | **150-350 GB/s**

| Package | Package | Package | Package | Package | Package | Package | Package |
|---|---|---|---|---|---|---|---|

Interposer | 50-400 GB/s | Interposer | 50-400 GB/s | Interposer | 50-400 GB/s | Interposer | 50-400 GB/s

| NPU | NPU | NPU | NPU | NPU | NPU | NPU | NPU | NPU | NPU | NPU | NPU | NPU | NPU | NPU | NPU |

# Hierarchical Network Architectures

| Pod | | Pod | |
|---|---|---|---|

**Rack-Scale Interconnect** → ← **25-150 GB/s**

| Node | Node | Node | Node |
|---|---|---|---|

**Rack-Scale Interconnect** → ← **150-350 GB/s** **Rack-Scale Interconnect** → ← **150-350 GB/s**

| Package | Package | Package | Package | Package | Package | Package | Package |
|---|---|---|---|---|---|---|---|

Interposer → ← 50-400 GB/s → Interposer → ← 50-400 GB/s → Interposer → ← 50-400 GB/s → Interposer → ← 50-400 GB/s →

| NPU | NPU | NPU | NPU | NPU | NPU | NPU | NPU | NPU | NPU | NPU | NPU | NPU | NPU | NPU | NPU |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

# Hierarchical Network Architectures

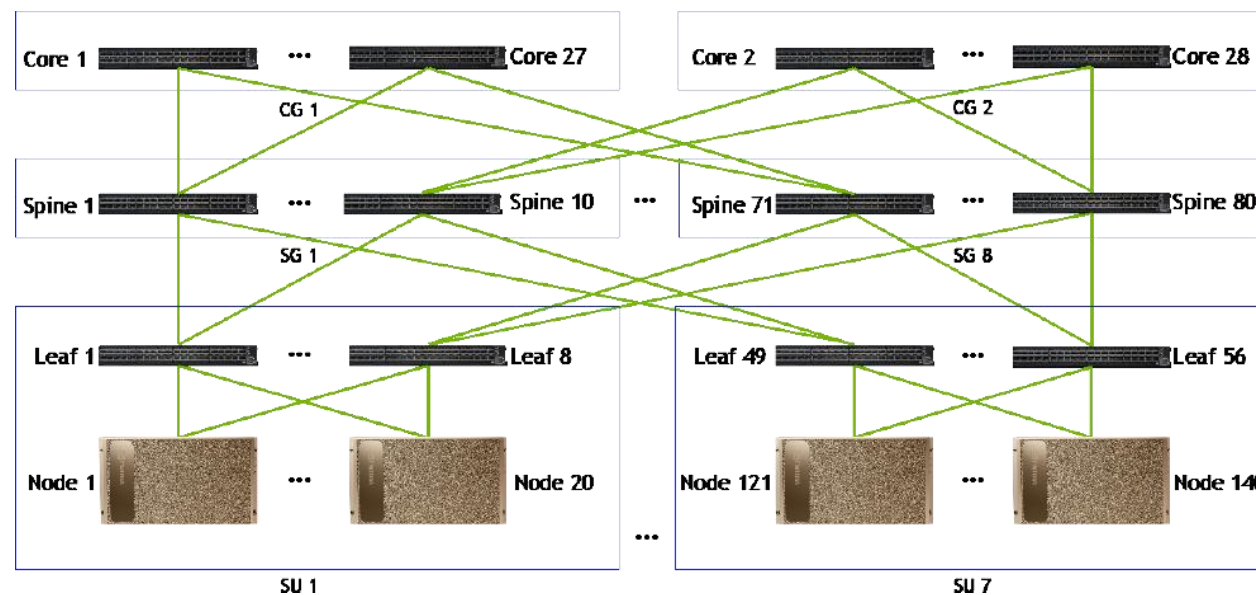# NVIDIA DGX SuperPod

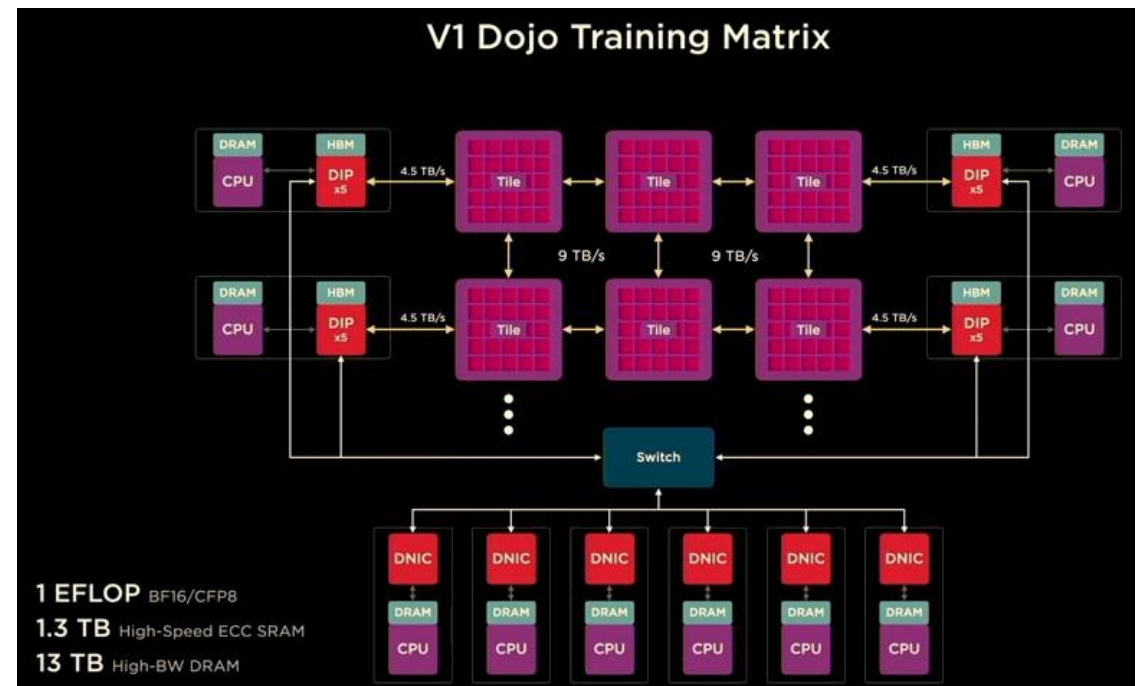

- Multi-level switches
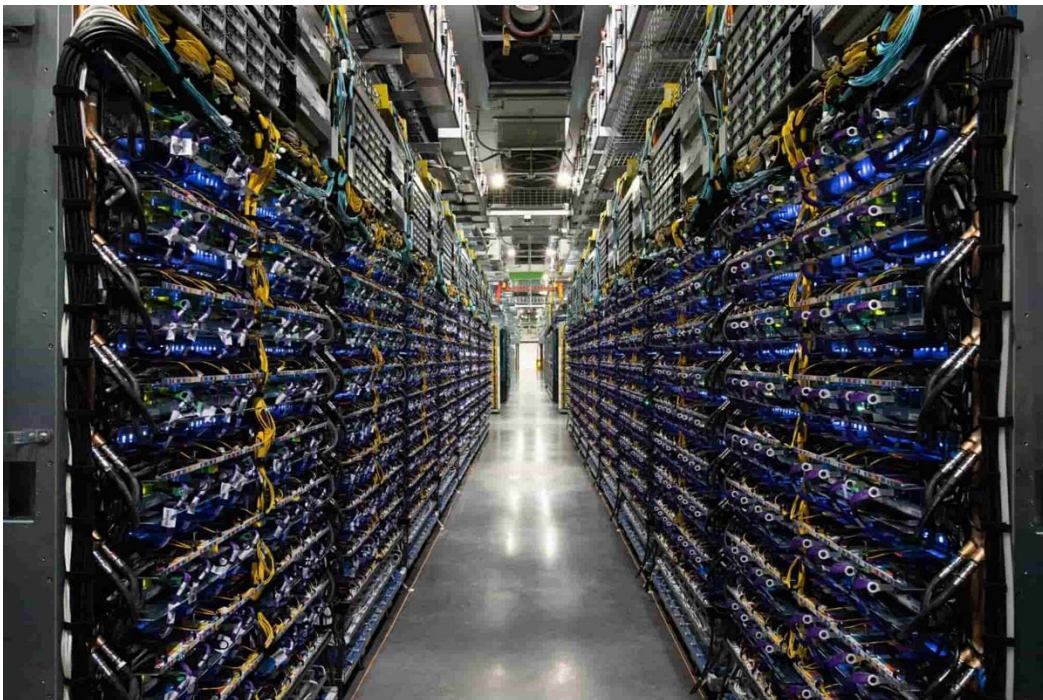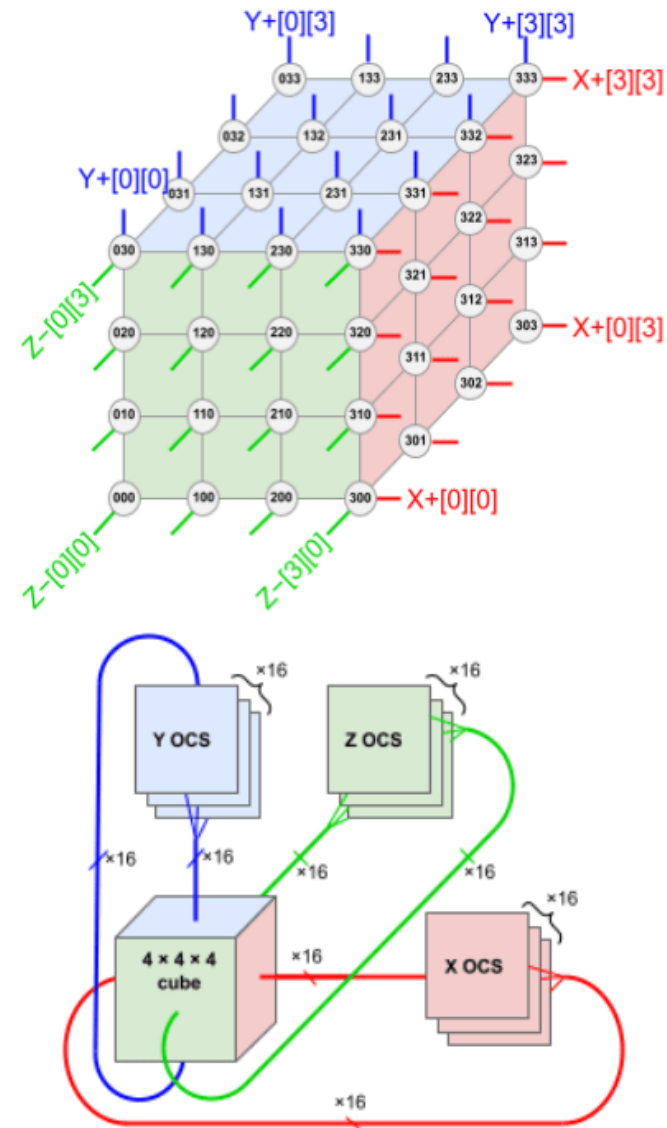
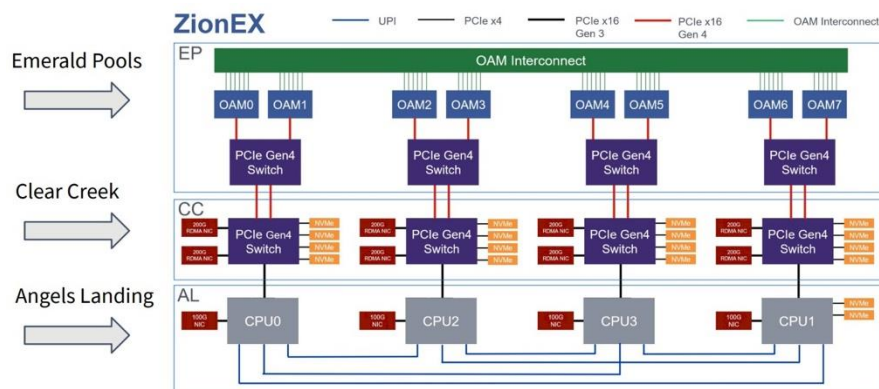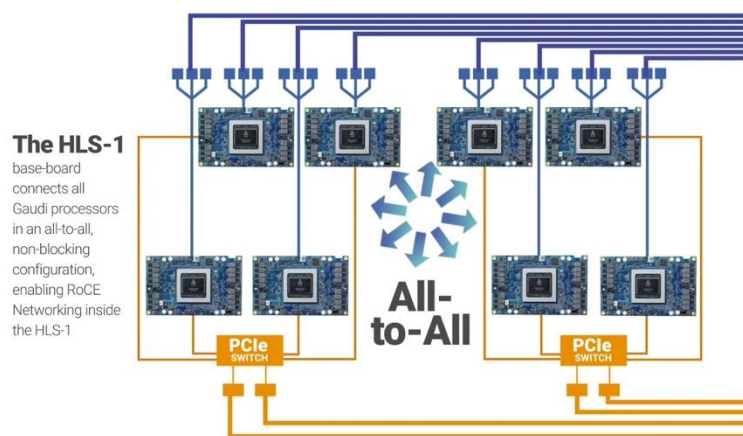# Tesla Dojo ExaPOD



- Scale-out Mesh Network

# Google Cloud TPU v4



- 3D Torus + Optical Networks

# State-of-the-art Training Clusters



Meta ZionEX



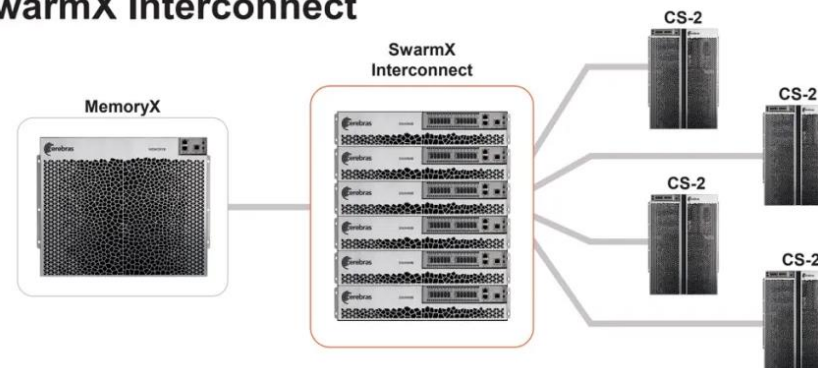Cerebras SwarmX



Intel Habana HLS-1



Tenstorrent Wormhole

# Distributed Training Stack

**Workload Layer**

| DNN Models | → | DLRM, ResNet-50, Transformer, GNMT |
| Workload Parallelization Strategy | → | Data, Model, Platform Agnostic Hybrid, Platform-aware Hybrid, Pipelined Parallelism |
| Communication Policy and Pattern | → | Distributed worker, Parameter Server |
| Framework-level Scheduling | → | Sync/Async, Blocking/Non-blocking |

**System Layer**

Compute and Memory Design

| Communication Mechanism | → | Topology-aware Collectives, Send/Recv, RPC |
| | → | Dataflow, Microarchitecture, Flexibility, Sparsity Support |
| Communication Scheduling | → | LIFO, FIFO, Fusion |
| Messaging/Transport Layer | → | TCP, RDMA (+ GPUDirect RDMA) |

**Network Layer**

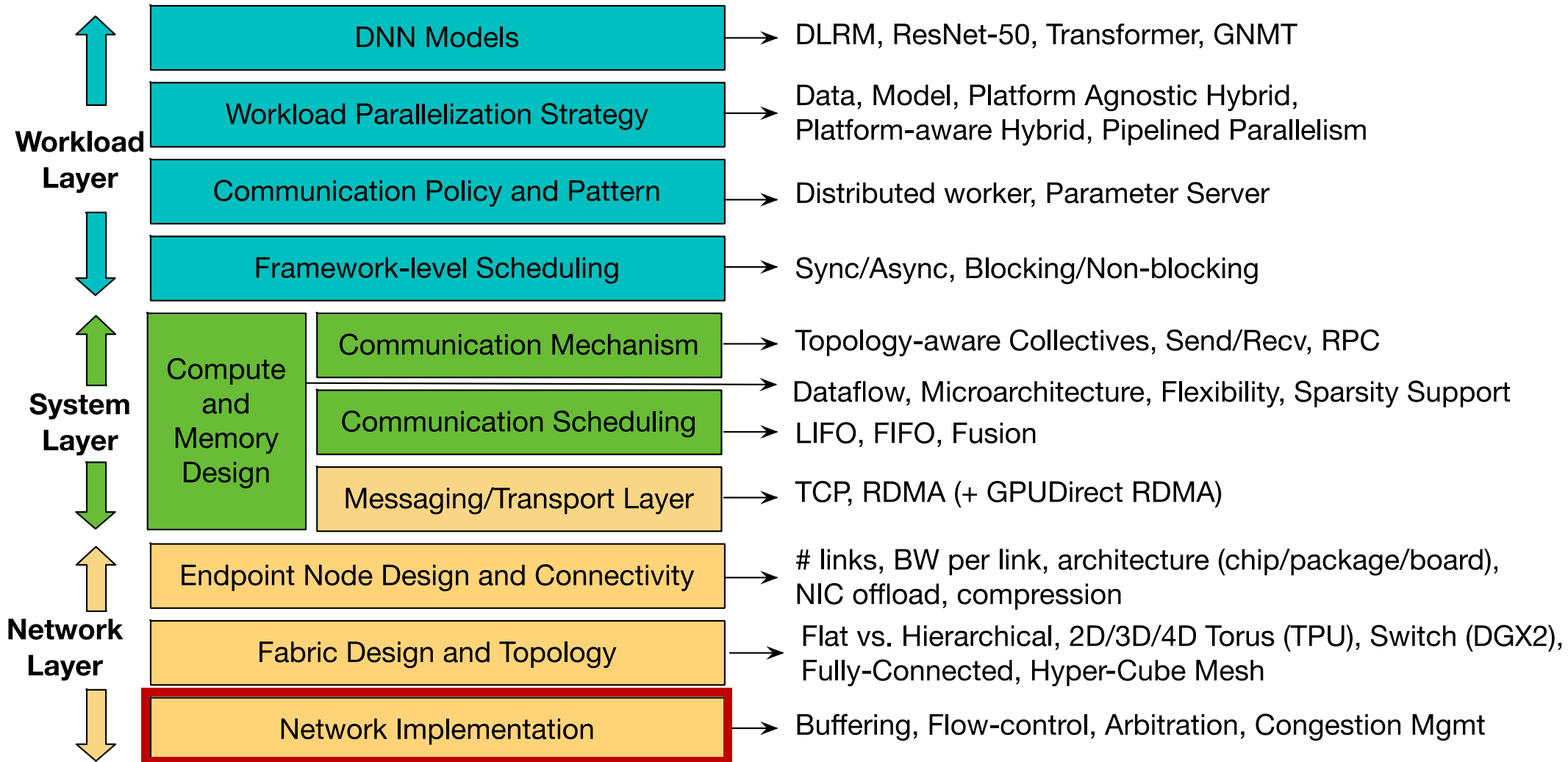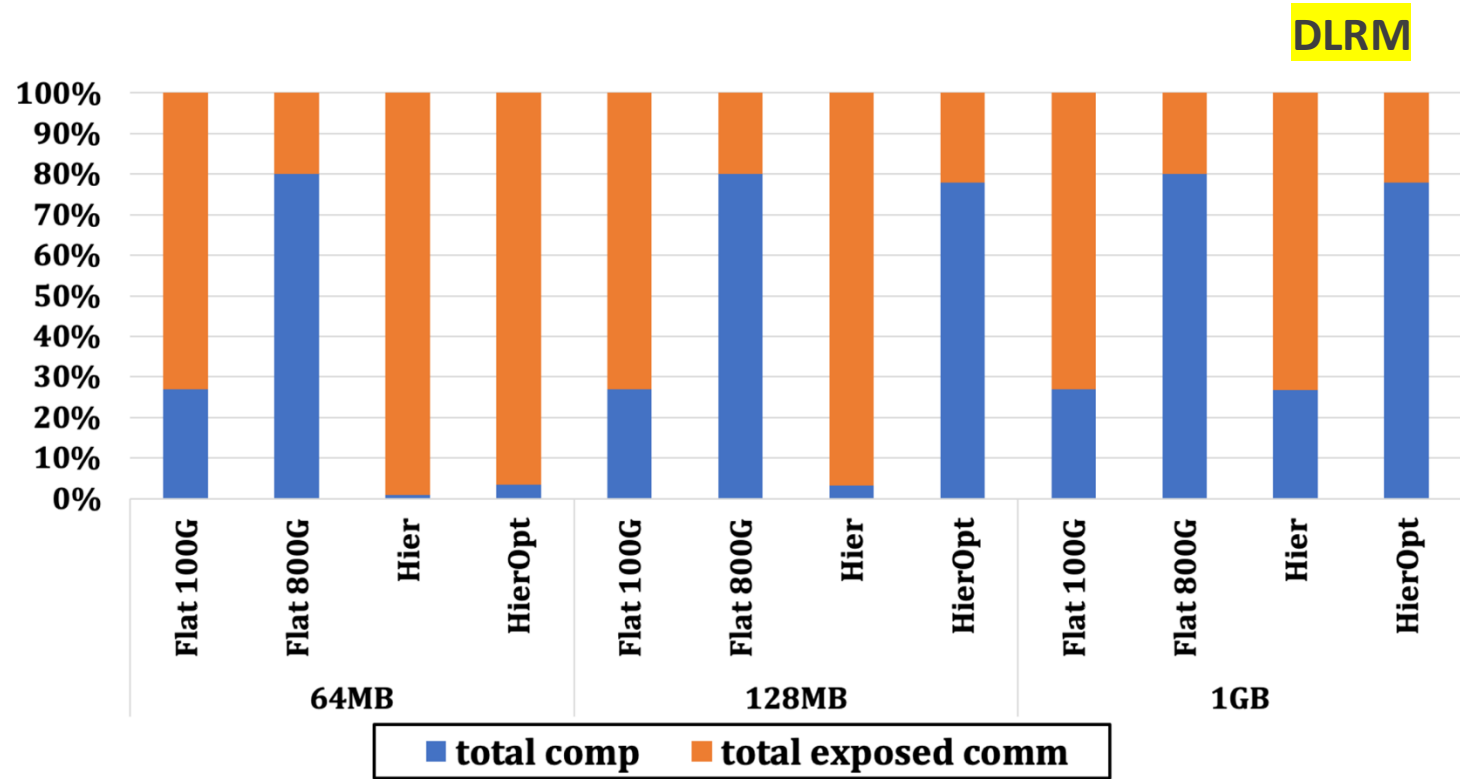| Endpoint Node Design and Connectivity | → | # links, BW per link, architecture (chip/package/board), NIC offload, compression |
| Fabric Design and Topology | → | Flat vs. Hierarchical, 2D/3D/4D Torus (TPU), Switch (DGX2), Fully-Connected, Hyper-Cube Mesh |
| Network Implementation | → | Buffering, Flow-control, Arbitration, Congestion Mgmt |

*Figure Courtesy: Srinivas Sridharan (NVIDIA)*

# Effect of Size of Switch Buffer

## Observations:

- Flat vs. Hierarch different Sensitivity to global switch size

**DLRM**



S. Rashidi, et al., *"Scalable Distributed Training of Recommendation Models: An ASTRA-SIM + NS3 case-study with TCP/IP transport"*, Hot Interconnects 2020

# Summary and Takeaways

- Large Model distributed ML is an ongoing open-research area

- Many emerging supercomputing systems being designed specifically for this problem!
  - NVIDIA HGX + (Mellanox) SHARP switches
  - Cerebras CS2
  - Tesla Dojo
  - Intel Habana
  - IBM Blueconnect
  - …

- Co-design of algorithm and system offers high opportunities for speedup and efficiency