# AI data center

**Omer Shabtai, SW Arch**

# Agenda

Distributed AI training scale and patterns
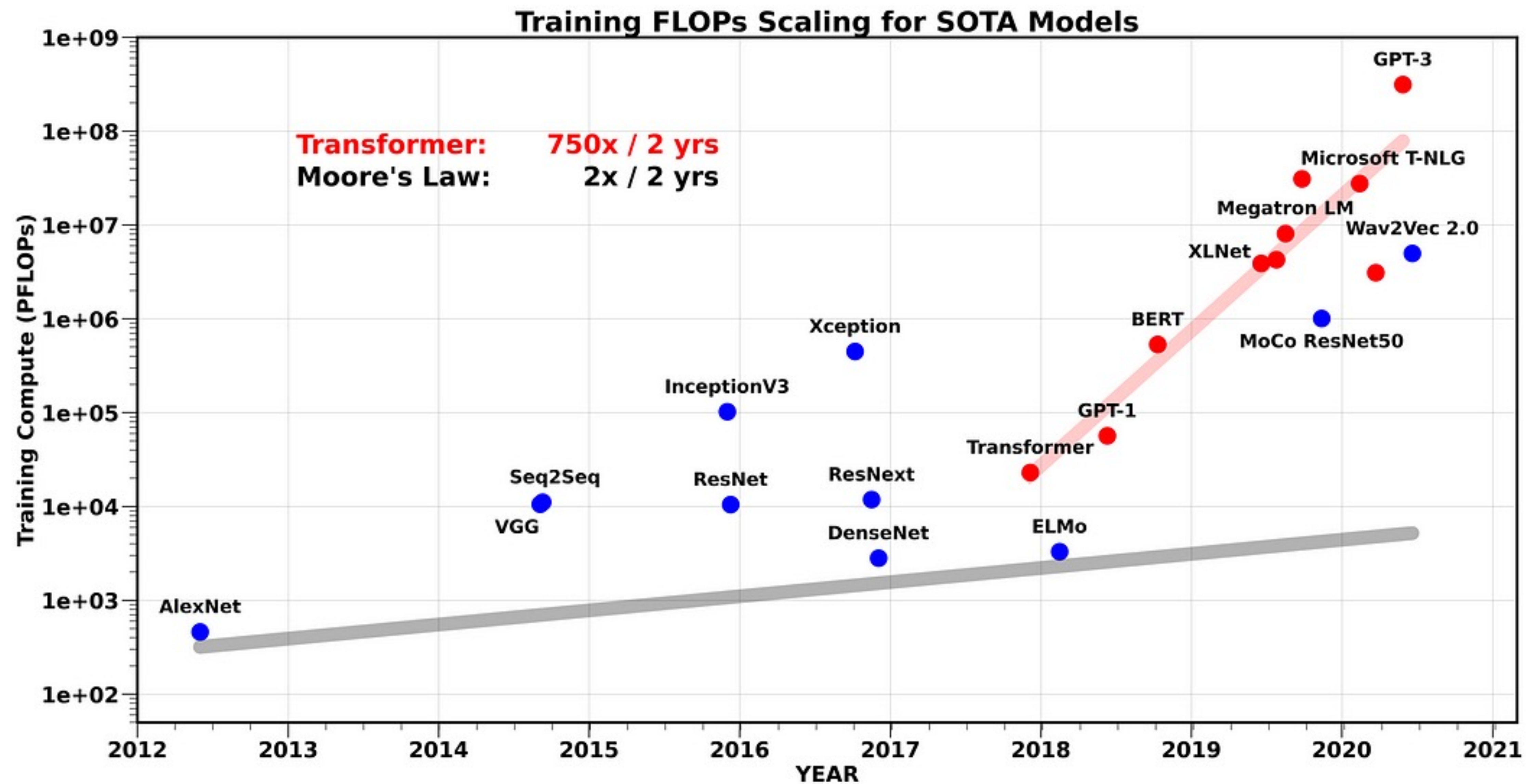
SO vs SU

AI Collectives

GPU direct RDMA

Goals

Evaluation focus

# Single GPU vs workload Flops



Training FLOPs Scaling for SOTA Models

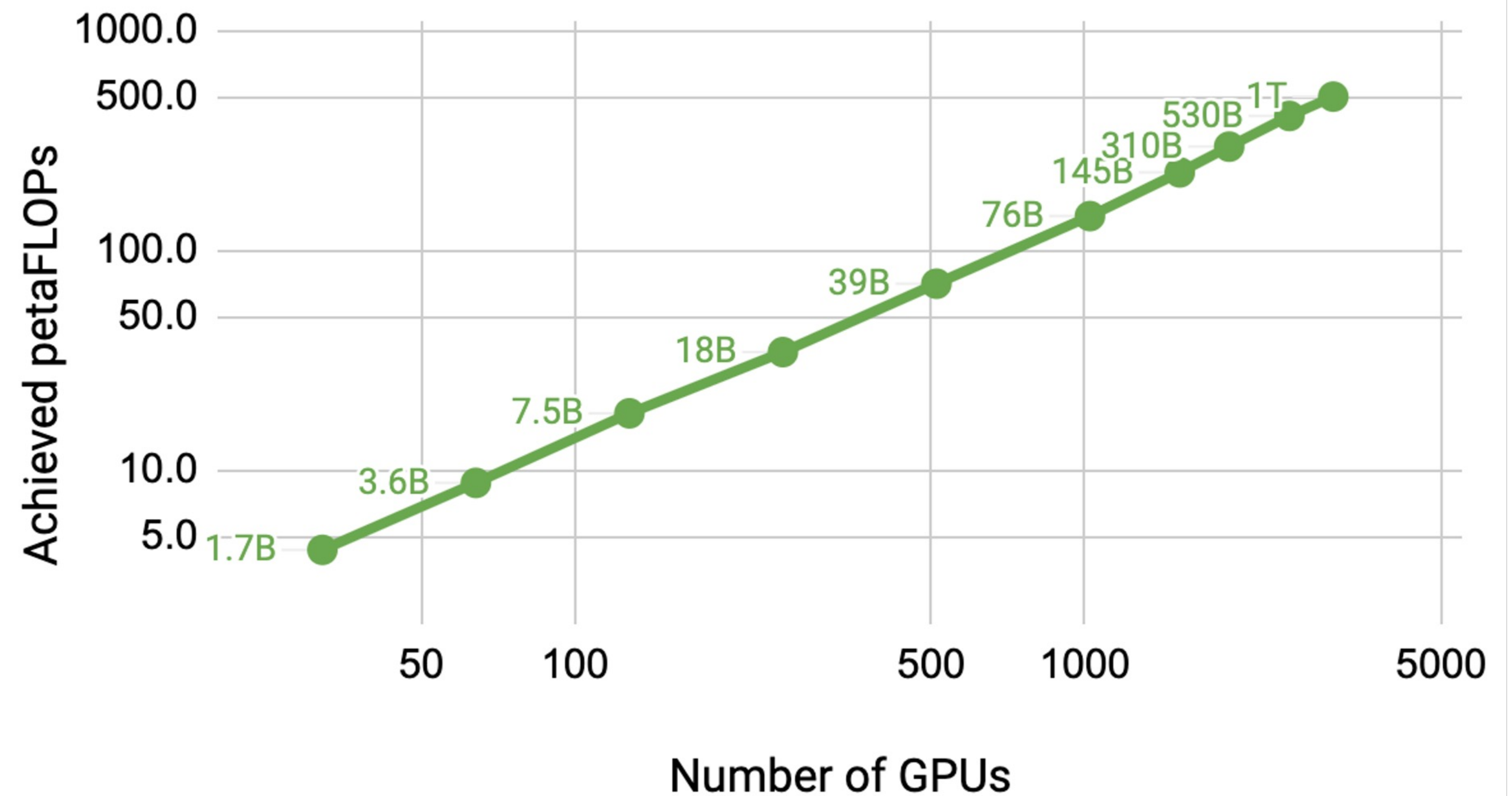Transformer: 750x / 2 yrs
Moore's Law: 2x / 2 yrs

GPU Flops ~2X every 2 years
>> 32X over same period

## Need for massive scale-up
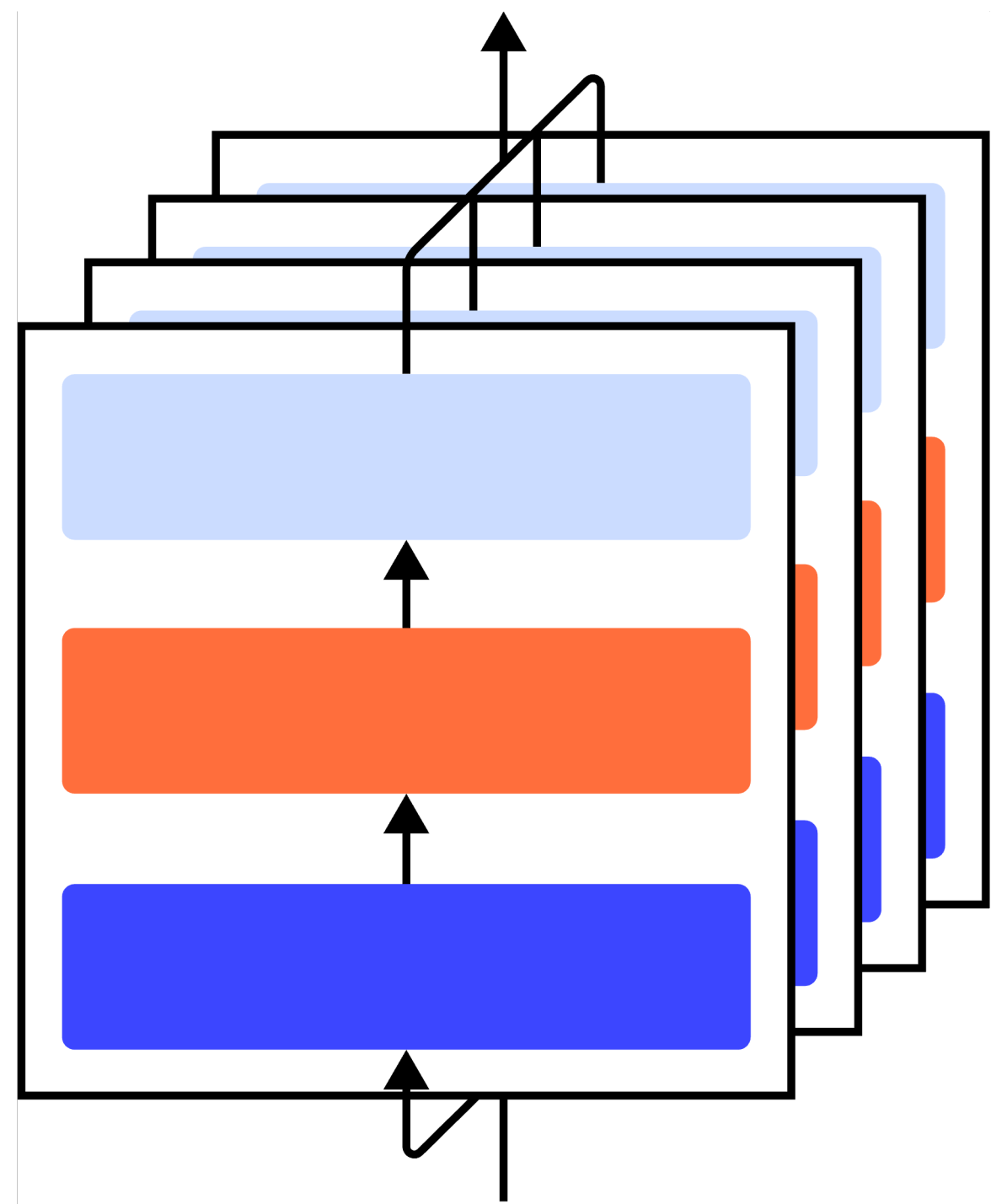
# AI training at scale

- Efficiency at scale
  - GPUs are strong
  - Goal – achieve linear scaling of training time to compute

- Synchronous training
  - Tail sensitive
  - Tightly coupled
  - Low Entropy

- computation / communication overlap
  - Induce complexity to framework and training
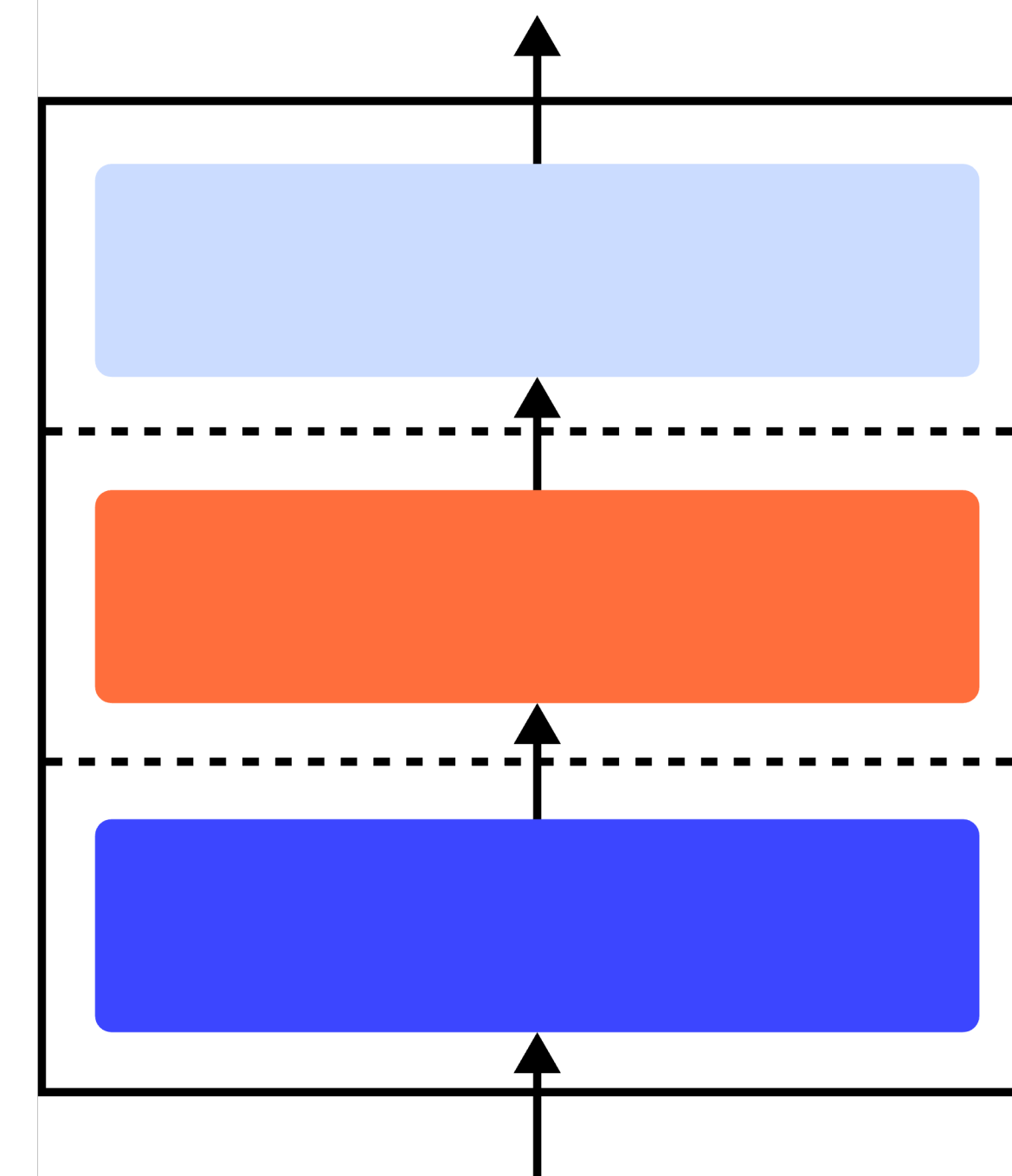
# Parallelism strategies

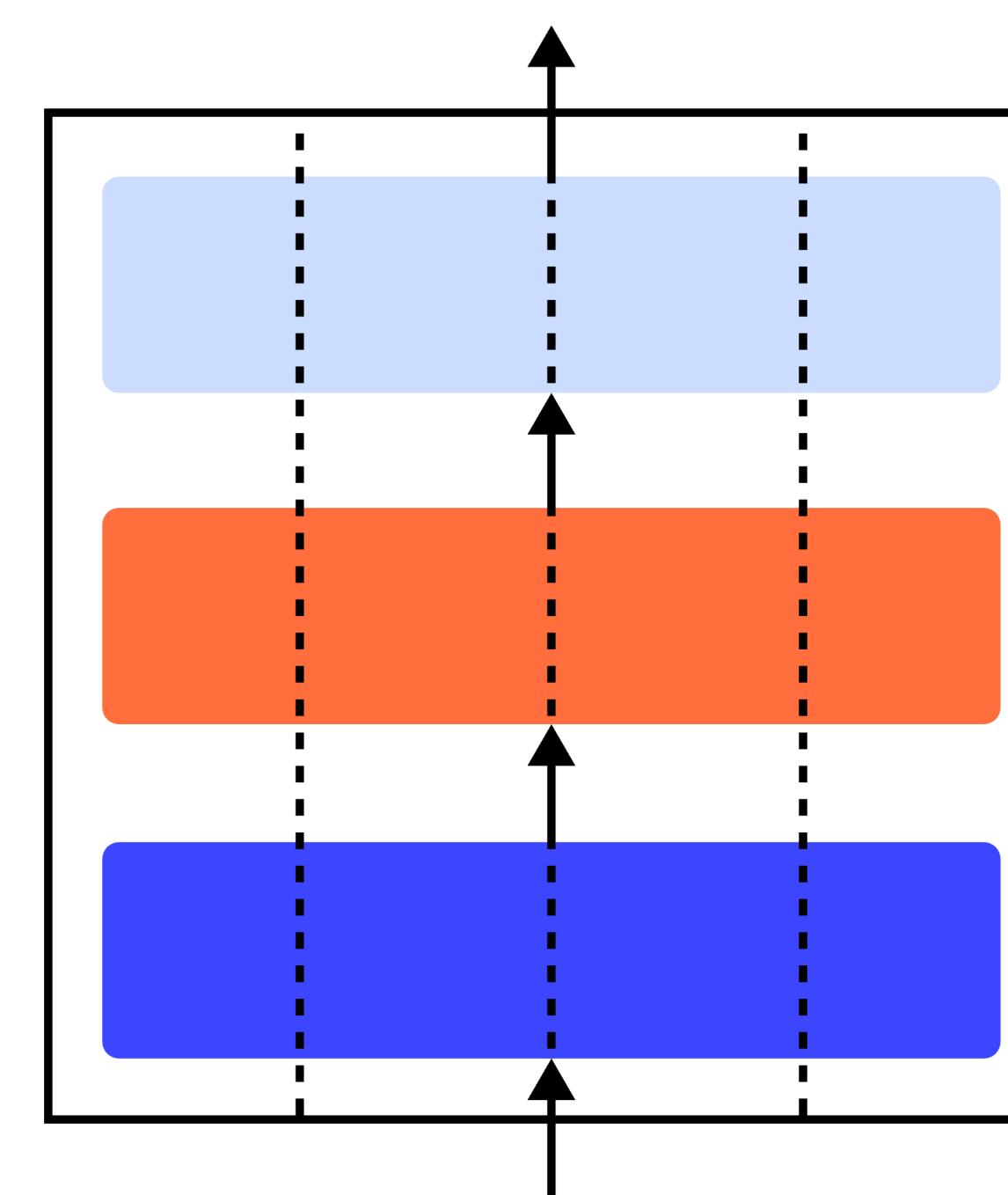## Model and data parallelism

**Data Parallelism**

**Pipeline Parallelism**

**Tensor Parallelism**

**Expert Parallelism**

**Simple to implement**

**Communication cheap**

**Communication expensive**

**Communication expensive**

Good performance at larger batch sizes (pipeline stall amortized)

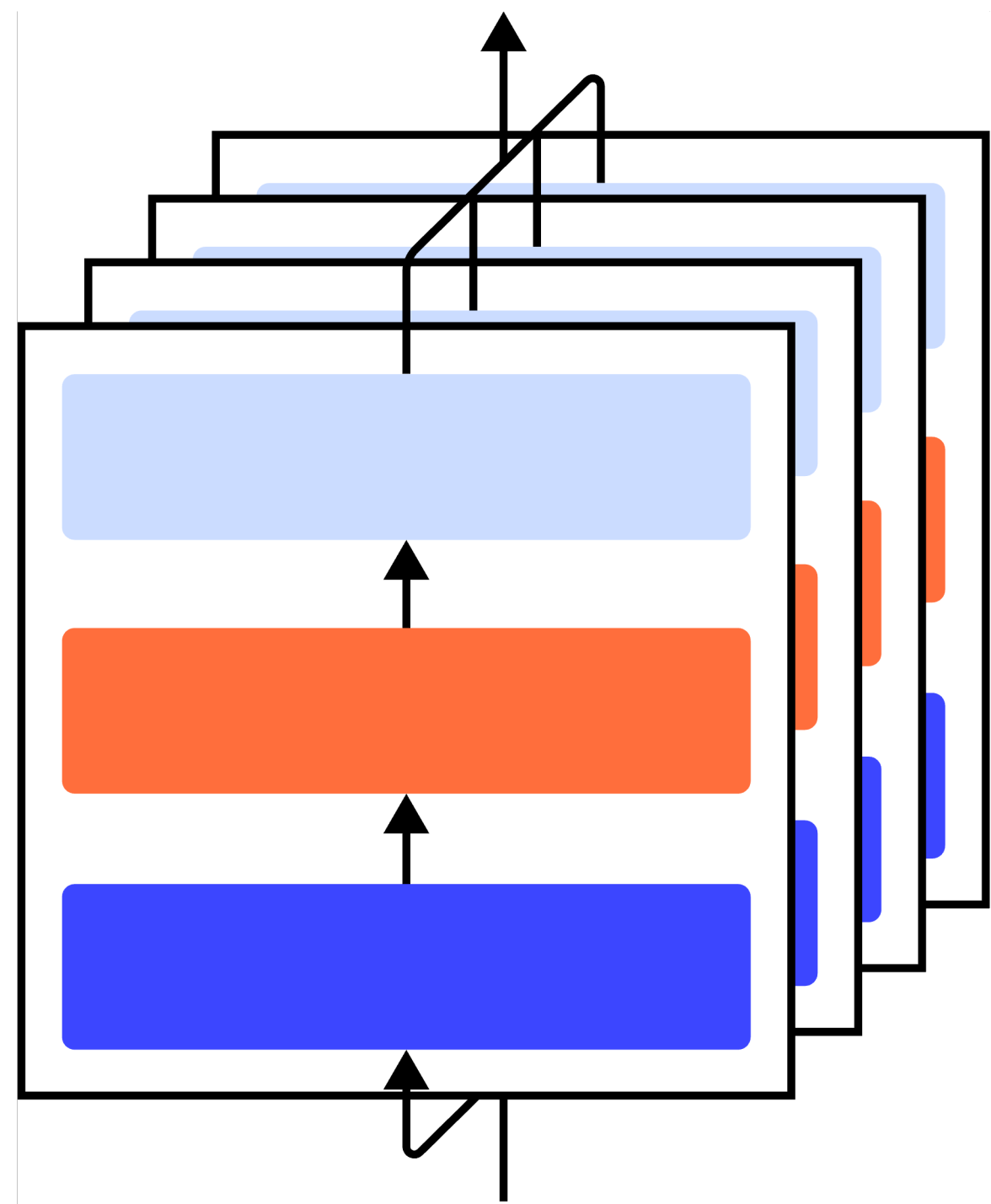**Good performance across batch sizes**

**Scalable sparsity**

Large batch size → better GPU efficiency and speedup
**Modern workloads - hybrid parallelism**

# Parallelism strategies

## Typical Collectives

**Data Parallelism**

**Pipeline Parallelism**

**Tensor Parallelism**

**Expert Parallelism**

**All-reduce
Or
All-gather + reduce scatter**

**Point to point**

**All to All, reduce scatter**

**All to all**

Large (~GBs)
Tightly coupled
Low entropy

Small (~MB)
Latency sensitive

Small, but frequent
Extremely latency sensitive
Scale-up only

Med (<100MB)
congested

6

# Popular parallelism Combinations

| | Data parallel | Tensor parallel | Pipeline parallel | Distributed Models |
|---|---|---|---|---|
| Light models, CNNs, For edge | X | | | |
| DLRM | X | | | X |
| FDSP (LLM) | X | | | X |
| 3D, GPT-3 like (LLM) | X | X | X | |
| MoE (LLM) | X | X | X | X |

AI training paradigm is rapidly evolving

New approaches pop fast

Increasing perf by significant multipliers

Over same HW generation

Optimization usually done for best perf

Network resiliency and jitter are mostly neglected

**Bad networking is usually exposed**

**but hardly detected**

# AI traffic patters

- **GPU comms**
  - Heterogeneous - Intra-node + inter-node
  - Single flow can saturate wire speed
  - Low flow count - minimize GPU's management resources (SMs)
  - Limited msg size – minimize GPU's expensive buffers

- **Collectives**
  - Synchronous
  - Reliable
  - Some incast guaranties

  Tightly coupled, both BW and tail latency sensitive

- **Accelerated transport - RDMA**
  - Zero copy – mem BW
  - Low Latency – GPU direct

# HGX/DGX - The AI building block

- Nvlink for scale-up
  - 7.2 Tbps GPU-GPU, full mesh
  - in-network aggregation

- 1:1 DPU/GPU for scale-out
  - 400gbps per GPU

- Dedicated interfaces for storage/mgmt



BlueField-3    Connect-X 7 / BlueField-3    NVMe    PCIe Switches    NVSwitch    —— PCIe    50 GbE    —— CPU communication

# GPU Direct RDMA
## Quick overview

- GPU Direct RDMA
    - CPU submit a job to the GPU
    - GPU executes the job
    - CPU notifies NIC that the data is ready
    - NIC reads the data from the GPU's memory
    - NIC sends the data to the network

# Scale up vs scale out

- Nvlink provides ~9X BW at lowest latency
  - All-reduce, all-gather speed-up of up to (8X in Hopper)
  - Tensor parallel – mandatory
  - Message aggregation for all to all

- Implications
  - Inter-node GPUs communicate only with their adjacent counter-parts
  - Rail optimized topology for better latency and scale
    - Increase 1-hop reach by 8X
    - Better isolation and latency guarantees

# NCCL Bandwidth
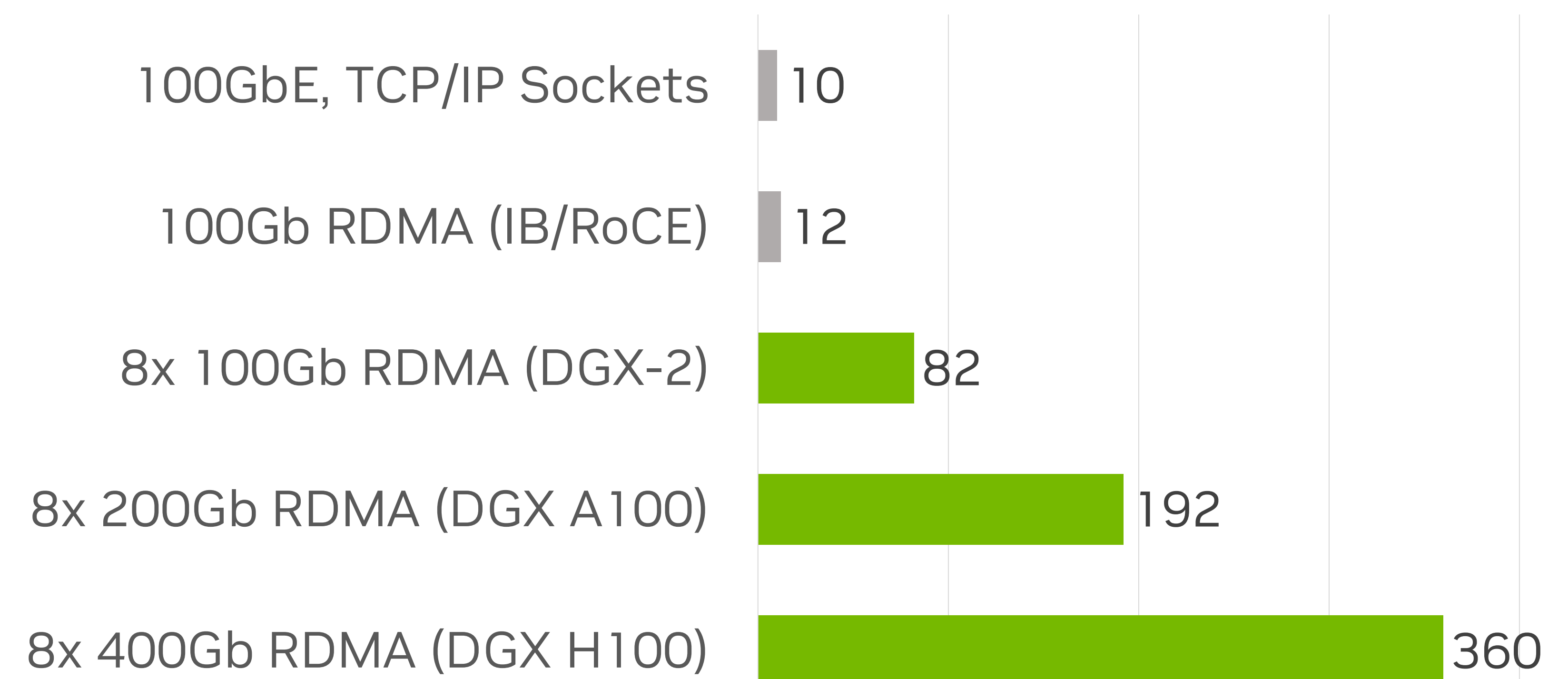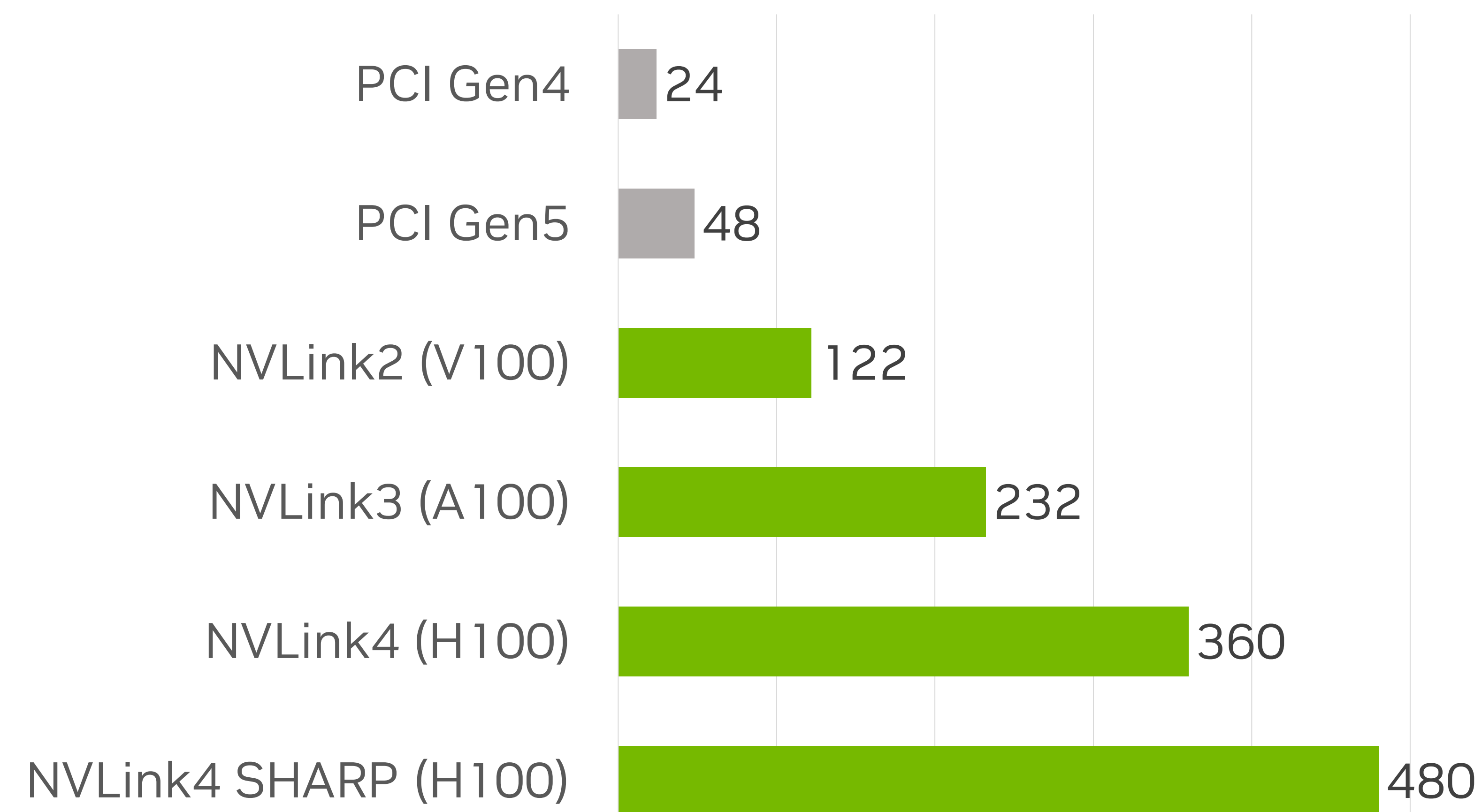


| Multi-GPU | |
|-----------|--|
| PCI | |
| NVLink | |

| Multi-node | |
|------------|--|
| TCP/IP | |
| IB/RoCE | |

| | |
|---|---|
| PCI Gen4 | 24 |
| PCI Gen5 | 48 |
| NVLink2 (V100) | 122 |
| NVLink3 (A100) | 232 |
| NVLink4 (H100) | 360 |
| NVLink4 SHARP (H100) | 480 |

| | |
|---|---|
| 100GbE, TCP/IP Sockets | 10 |
| 100Gb RDMA (IB/RoCE) | 12 |
| 8x 100Gb RDMA (DGX-2) | 82 |
| 8x 200Gb RDMA (DGX A100) | 192 |
| 8x 400Gb RDMA (DGX H100) | 360 |

*NCCL Tests Allreduce Bus Bandwidth in GB/s*

# AI datacenters fabric requirements

Any work load, any location

**Full cluster utilization**

**Full Bisection BW,
Latency/jitter sensitive**

Allocation agnostic,
Workload agnostic

Expensive
Virtualized

A-symmetry resiliency
Global awareness

Tightly coupled
collective operations

Performance isolation
(chatty neighbor)

# Adaptive Routing and Congestion Control
## A generic solution

- End to end solution, Per packet load balancing for RDMA
  - Up to 95% bisection BW utilization, at low tail latency
  - AR Convergence time << Congestion control convergence time
  - Global Fabric asymmetry and failiure support
  - Zero-copy OOO delivery – reduce dependency on tail events.

- Topology agnostic
  - Asymmetric topology
  - fragmented job allocation

- Workload agnostic
  - Doesn't relay on defined patterns like rings/trees

- Congestion control focus:
  - Incast scenarios
    - Collective first
  - Blocking fabric (mainly failures)
    - Slow forming

# AI Training networking infrastructure

## Evaluation

**E2E Training**
- DDP
- LLM (3D, FSDP, MoE)
- DLRM
- ...

**Collectives**
- NCCL
  - P2P,A2A,AR,AG,RS

**RDMA**
- Bisection BW and latency under full load

**Resiliency**
- Link failiure
- Slow receiver
- BER

**Workload Isolation**
- Congestion-free noise
- Congested-prune noise

**Single workload**
- Consistency
- Allocation agnostic