

# Communication Primitives for Distributed Training Algorithms



ECE 826

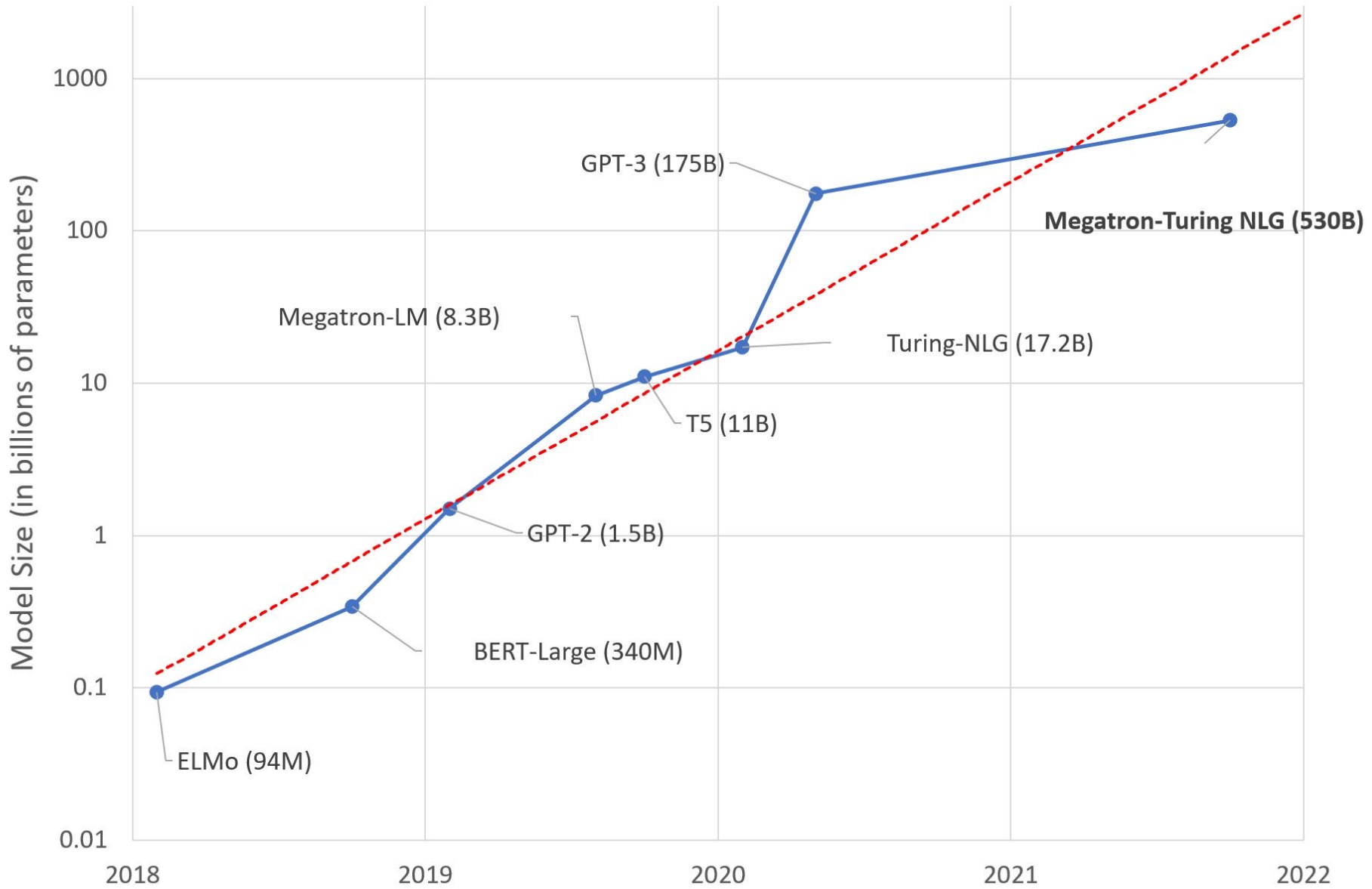


**Carnegie  
Mellon  
University**

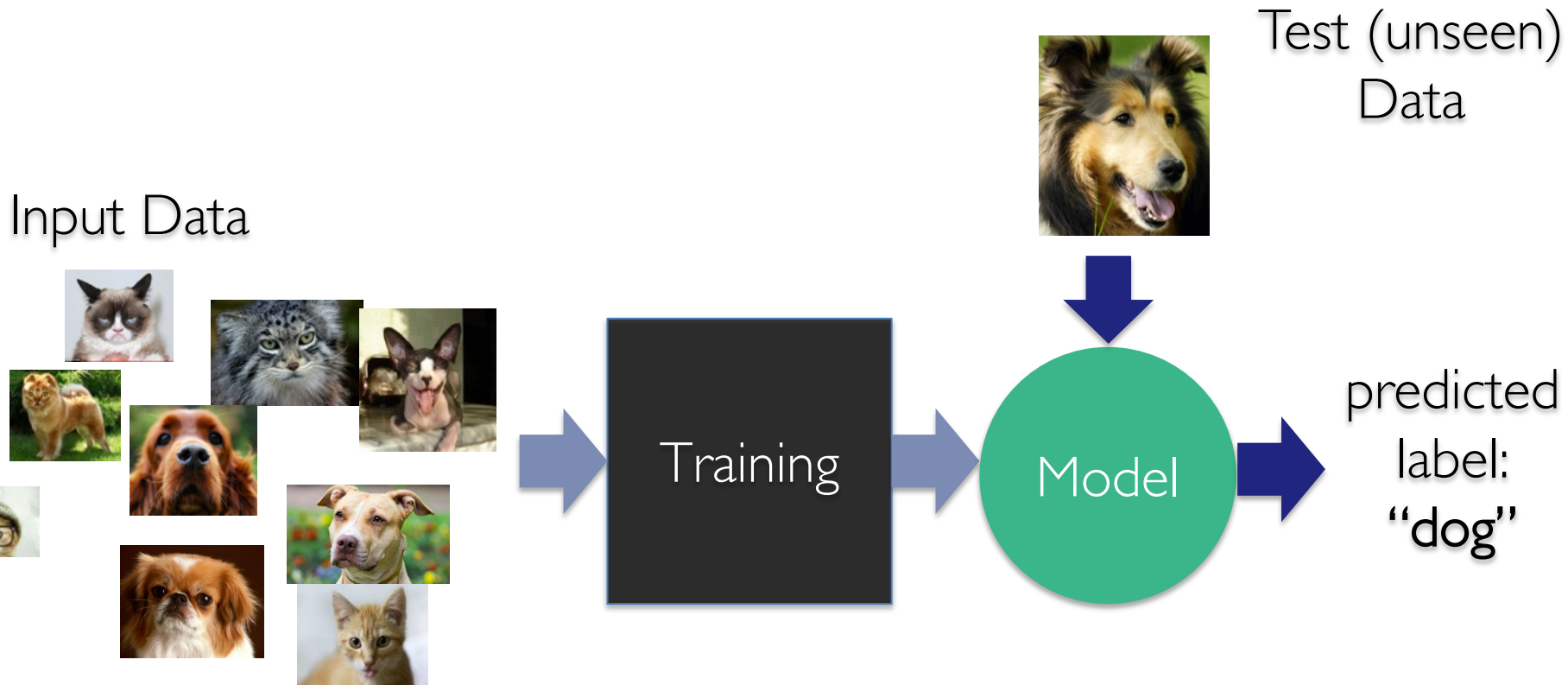
**Hongyi Wang**  
Machine Learning Department  
Carnegie Mellon University  
April 14th

# Overview

- Motivations of Distributed ML
- Parallelism & Comm. Topologies
  - ❑ Data parallelism (PS, all-reduce)
  - ❑ Model parallelism
  - ❑ Pipeline parallelism
  - ❑ Hybrid parallelism
- Communication Bottlenecks & Solutions



# A Simplified ML Pipeline



*However,  
model can contain 1.75 trillion params [1].  
+ data set can be of 400 million of data  
points [2]*

[1] *Wu Dao 2.0*

[2] *OpenAI CLIP*

---

Training a ResNet-50 on ImageNet takes 2.5 days  
on a single GPU

*[EC2 P3.2xlarge (Tesla V100) + PyTorch]*

# Distributed Model Training

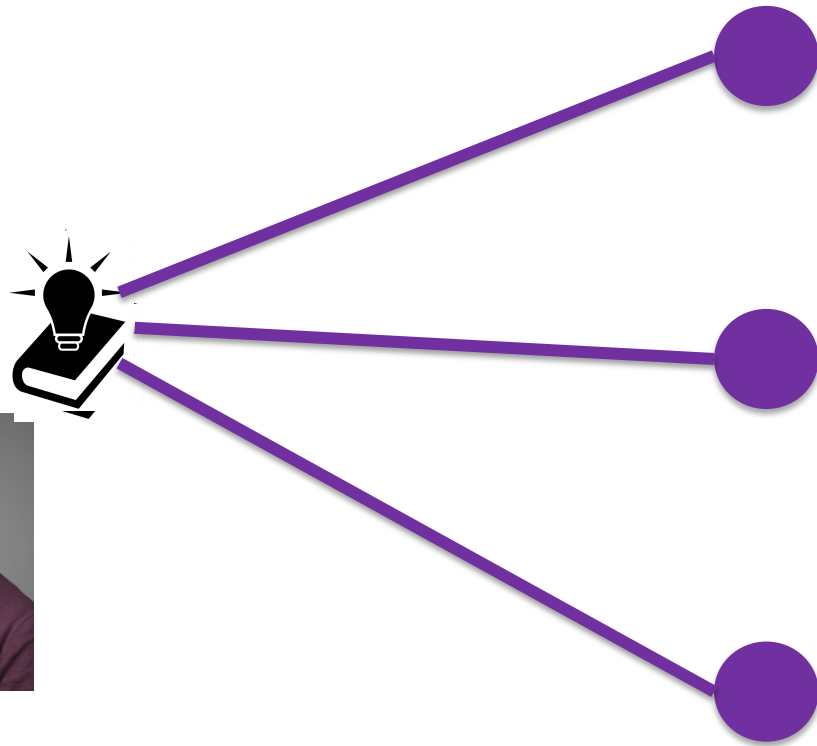


# Overview

- Motivations of Distributed ML
- Parallelism & Comm. Topologies
  - ❑ Data parallelism (PS, all-reduce)
  - ❑ Model parallelism
  - ❑ Pipeline parallelism
  - ❑ Hybrid parallelism
- Communication Bottlenecks & Solutions

# Intro to Comm. Operations

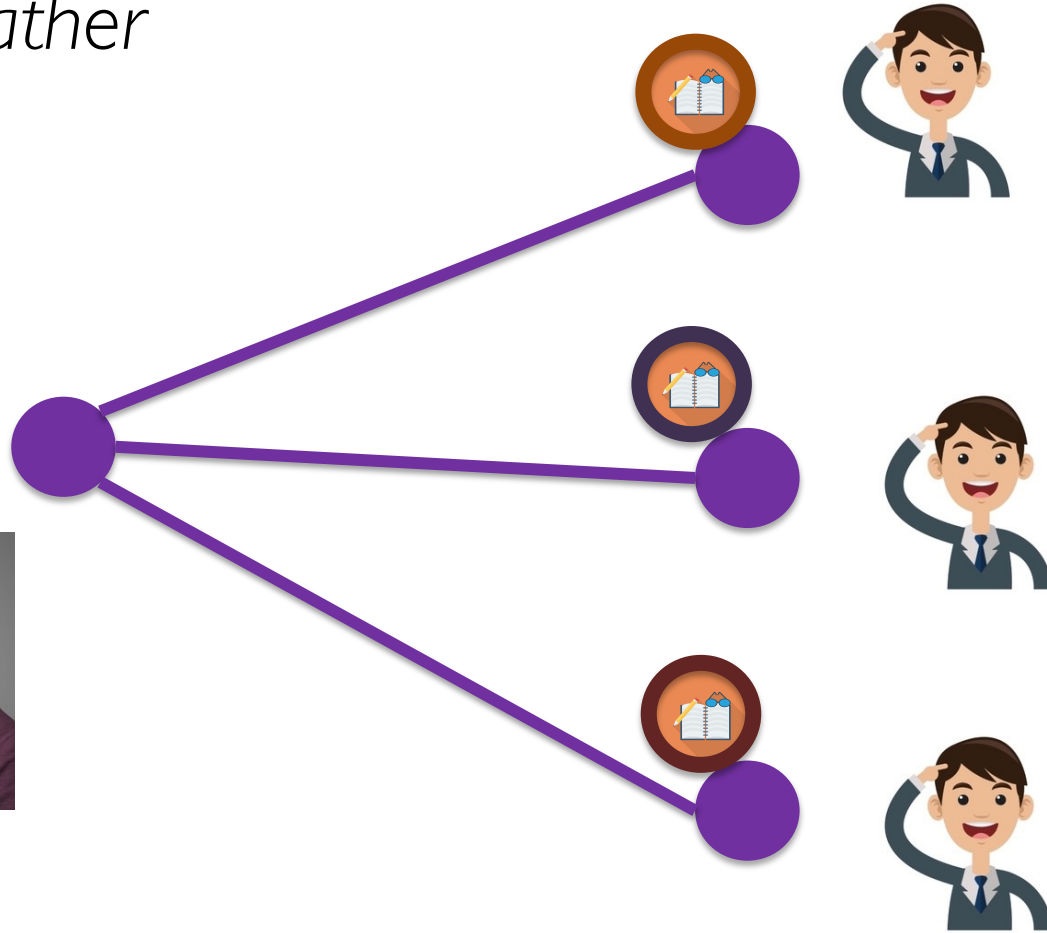
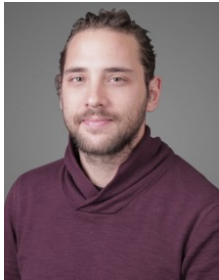
Ops I: *Broadcast*





# Intro to Comm. Operations

Ops2: *Gather*



# Data-parallel mini-batch SGD

$$\min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{w}; \mathbf{z}_i) \quad \mathcal{X} = \{z_1, \dots, z_n\}$$

For iteration  $t$ :

- Sample a mini-batch:  $\mathcal{B}_t \subset \mathcal{X}$

- Compute gradient:  $\nabla \frac{1}{|\mathcal{B}_t|} \sum_{j=1}^{|\mathcal{B}_t|} \ell(w_t; z_j)$

■  $\nabla \ell(w_t; z_1)$

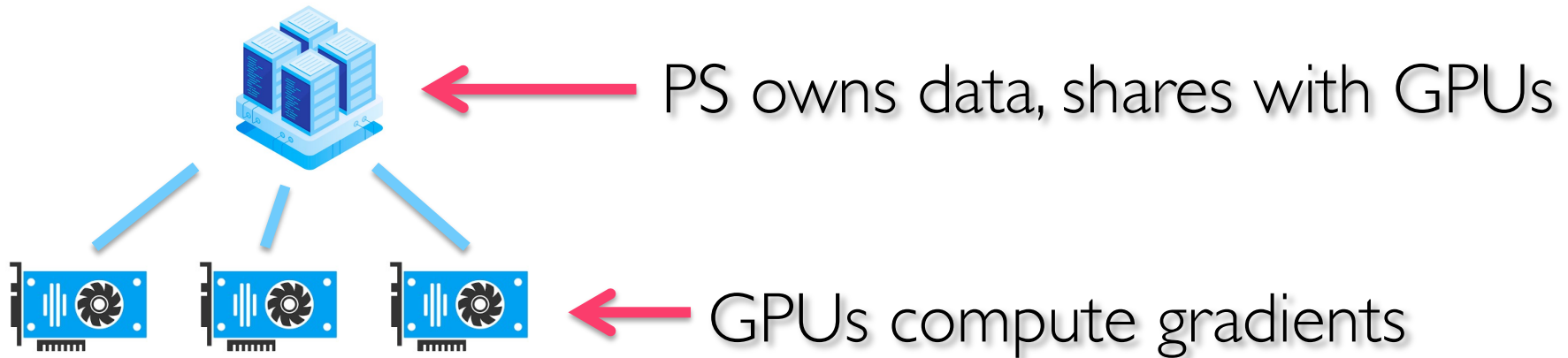
▲  $\nabla \ell(w_t; z_2)$

●  $\nabla \ell(w_t; z_3)$

$$= \frac{1}{|\mathcal{B}_t|} \sum_{j=1}^{|\mathcal{B}_t|} \nabla \ell(w_t; z_j)$$

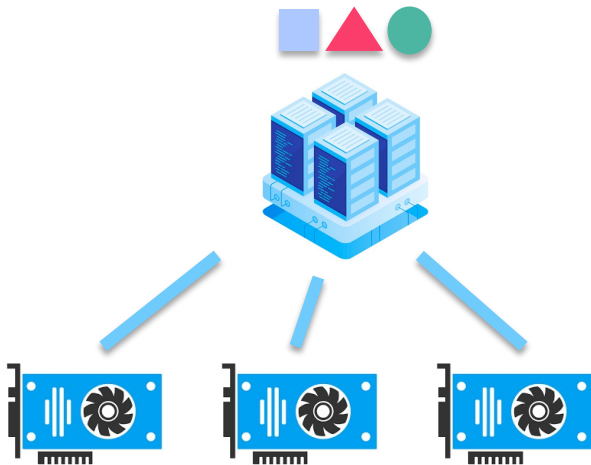
# Data-parallel mini-batch SGD

parameter server (PS)



# Data-parallel mini-batch SGD

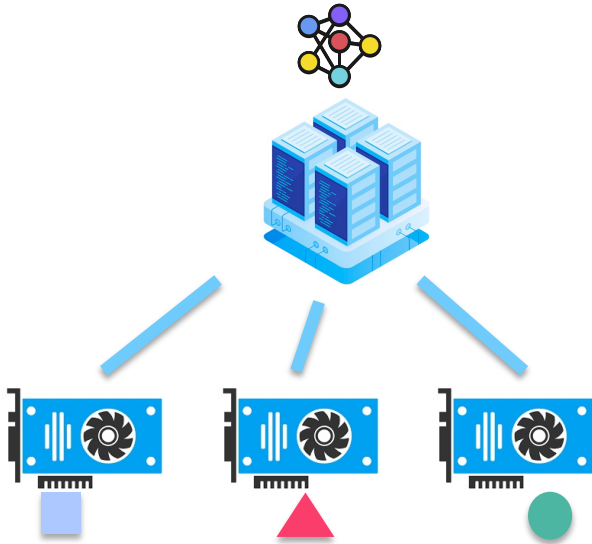
parameter server (PS)



Data Batches Assignment

# Data-parallel mini-batch SGD

parameter server (PS)

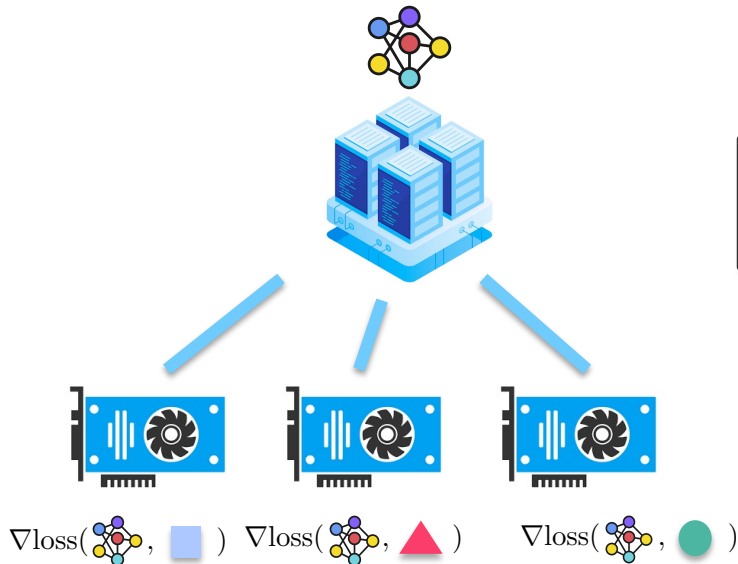


Model Broadcast

# Data-parallel mini-batch SGD

parameter server (PS)

Gradient Gathering



# Comm. Cost of PS

parameter server (PS)



$$\nabla_{\text{loss}}(\text{model}, \square) \quad \nabla_{\text{loss}}(\text{model}, \triangle) \quad \nabla_{\text{loss}}(\text{model}, \bullet)$$

1. Model weights and gradients are  $d$ -dimensional vectors.

2. There are  $p$  GPUs in the cluster for computing.

Both broadcast and gather require to communicate  $\mathcal{O}(pd)$  information

PS requires to incur TWO communication operations

# A Generalized Comm. Cost Model

$\alpha - \beta$  Cost model

Comm. cost =  $\beta \times$  Bandwidth cost

Total amount of info  
communicated.

$pd$

+  $\alpha \times$  Latency cost

The entire number of  
communication rounds.

2



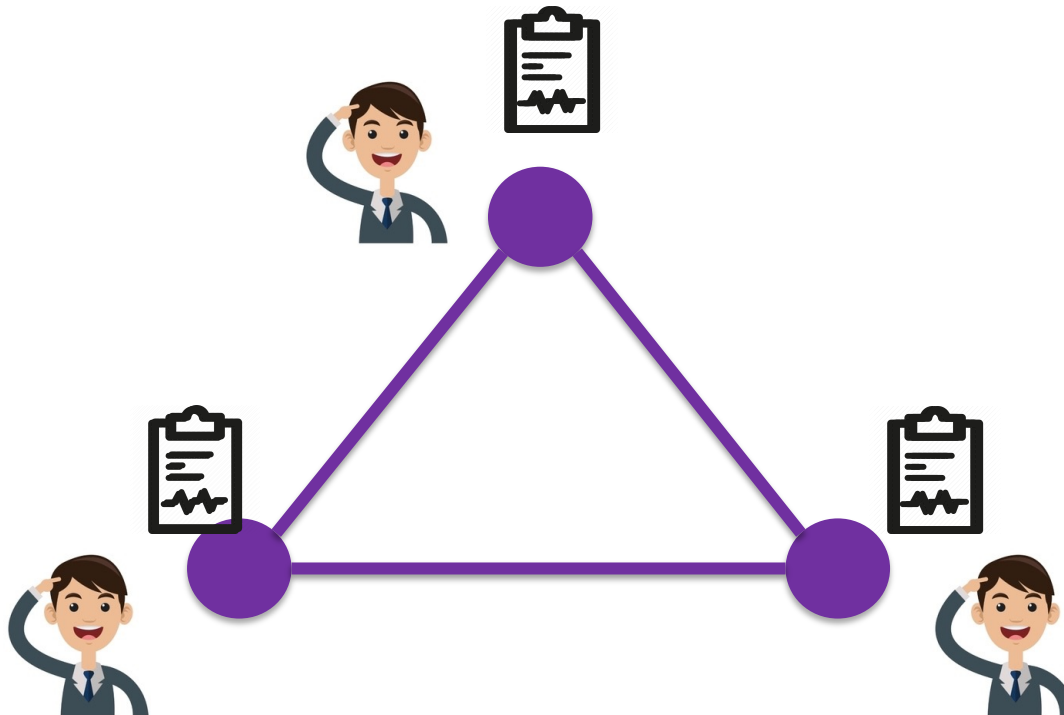
# A Generalized Comm. Cost Model

$\alpha - \beta$  Cost model

Comm. Topology	Bandwidth	Latency
Param Server	$\beta \times 2pd$	$\alpha \times 2$

# Intro to Comm. Operations

Ops3: *All-reduce*



# Data-parallel mini-batch SGD

$$\min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{w}; \mathbf{z}_i)$$

For iteration  $t$  :

- Sample a mini-batch:  $\mathcal{B}_t \subset \mathcal{X}$

- Compute gradient:  $\nabla \frac{1}{|\mathcal{B}_t|} \sum_{j=1}^{|\mathcal{B}_t|} \ell(w_t; z_j)$

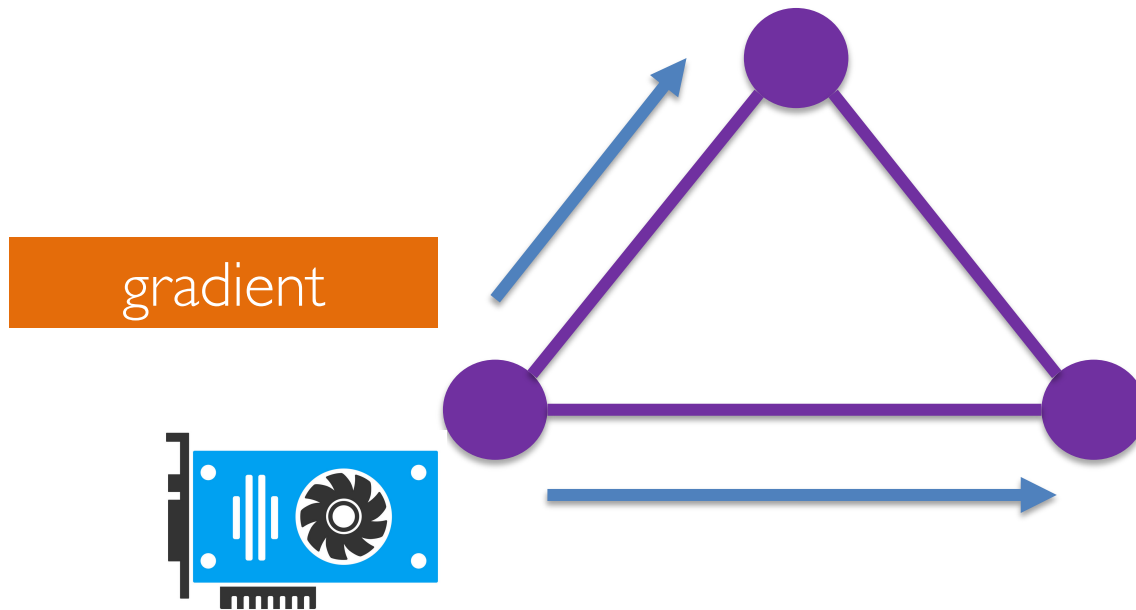
■  $\nabla \ell(w_t; z_1)$

▲  $\nabla \ell(w_t; z_2)$

●  $\nabla \ell(w_t; z_3)$

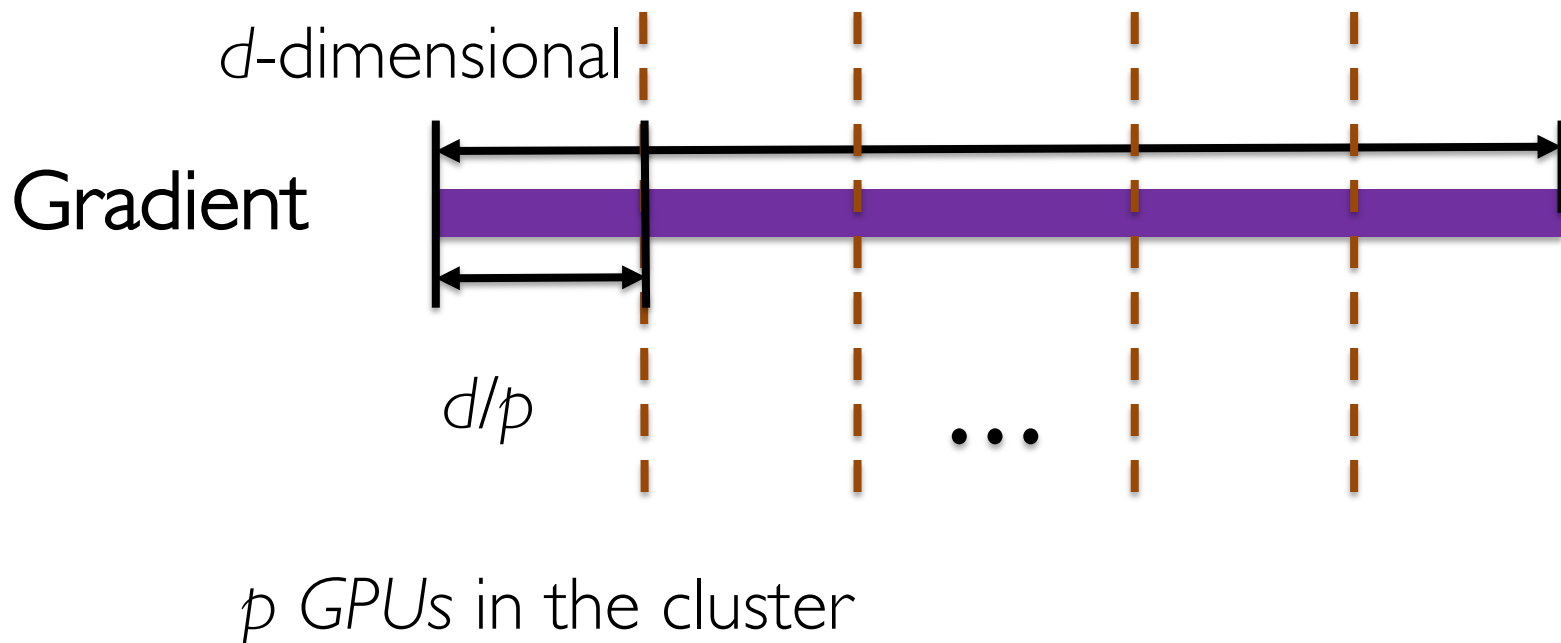
$$= \frac{1}{|\mathcal{B}_t|} \sum_{j=1}^{|\mathcal{B}_t|} \nabla \ell(w_t; z_j)$$

# All-reduce SGD



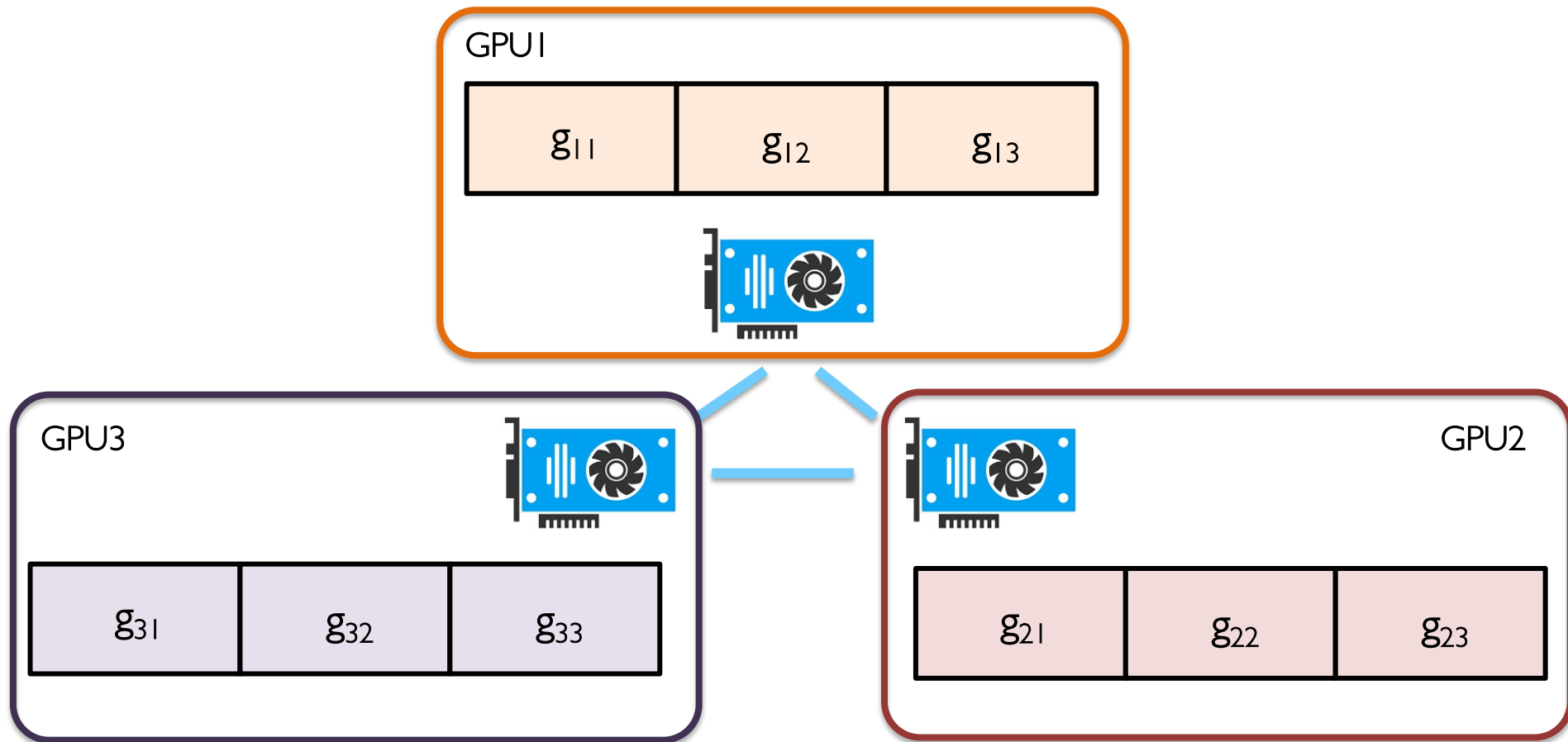
# Data-parallel SGD

## Ring-reduce SGD



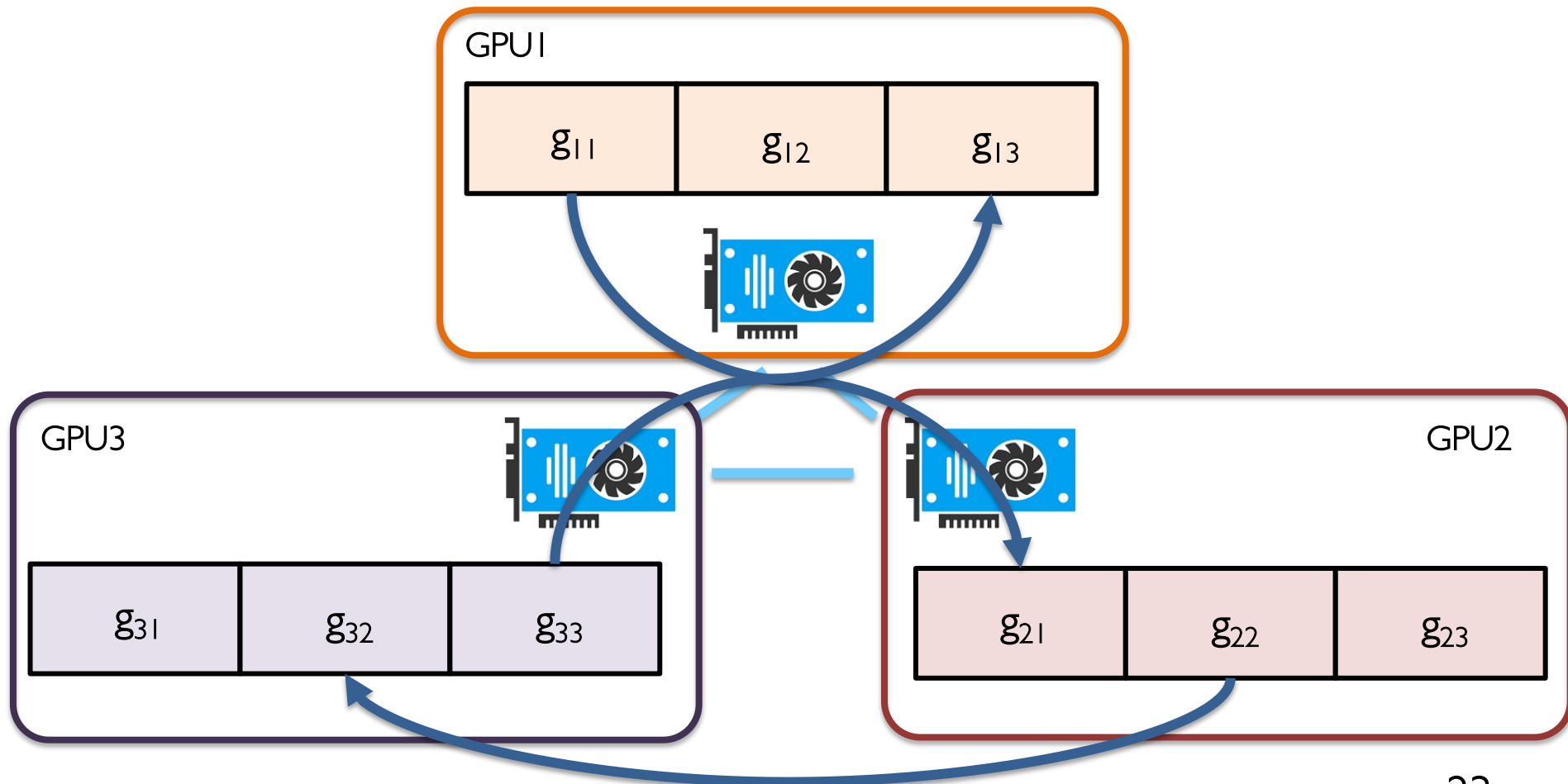
# Ring-reduce SGD

Initialization



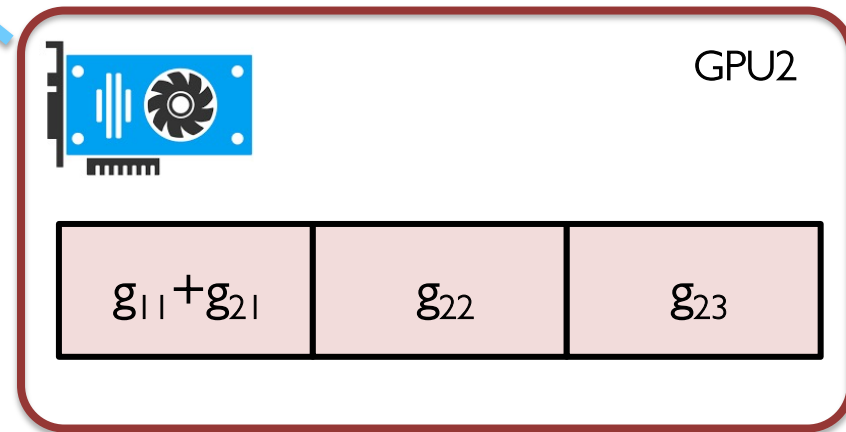
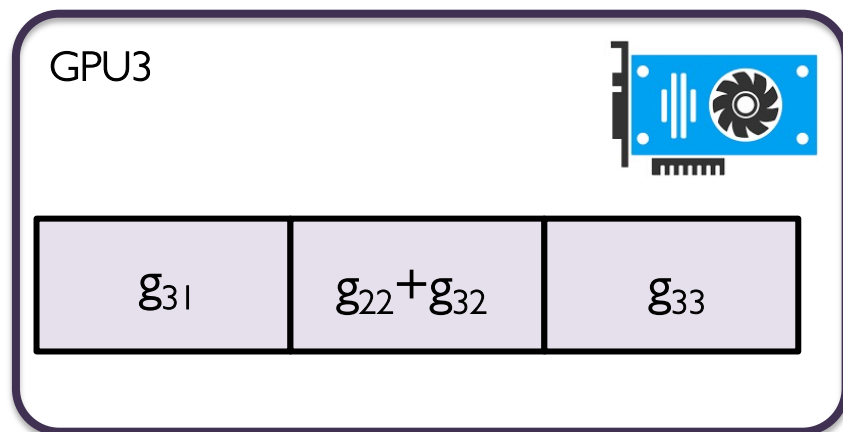
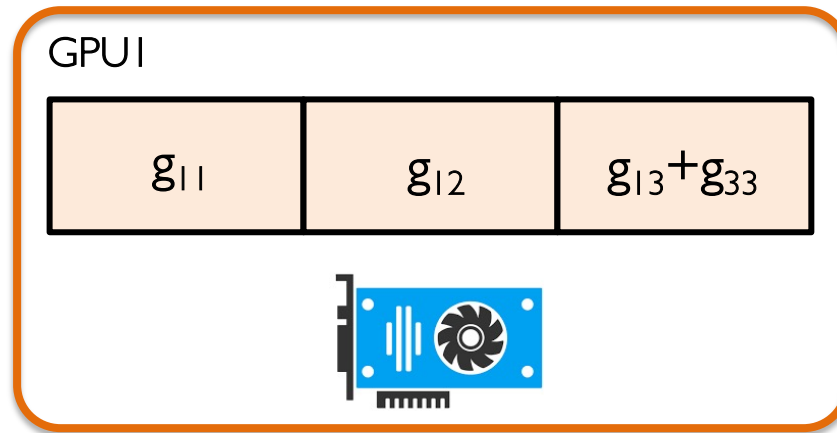
# Ring-reduce SGD

Round 1



# Ring-reduce SGD

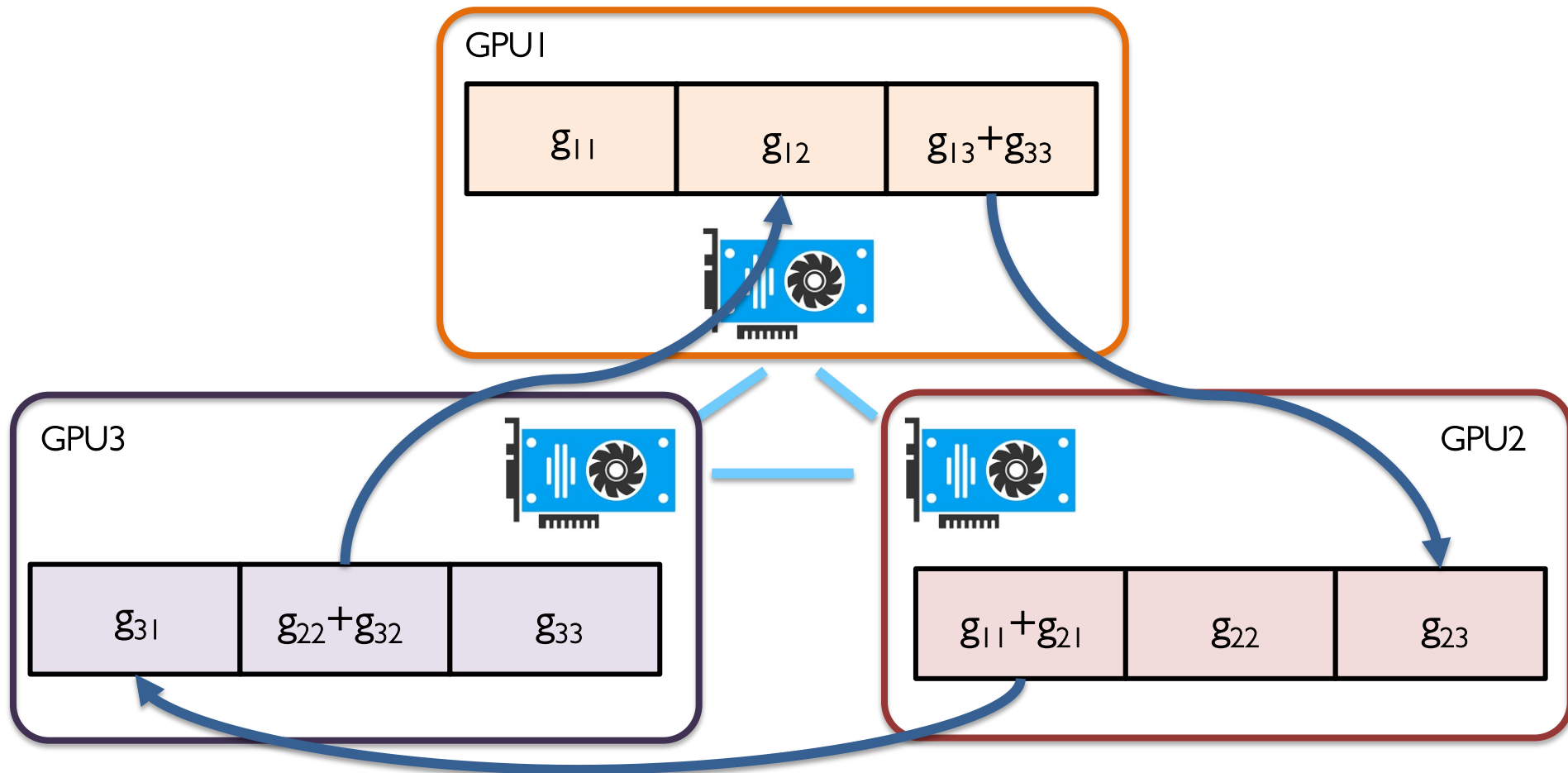
Round 1





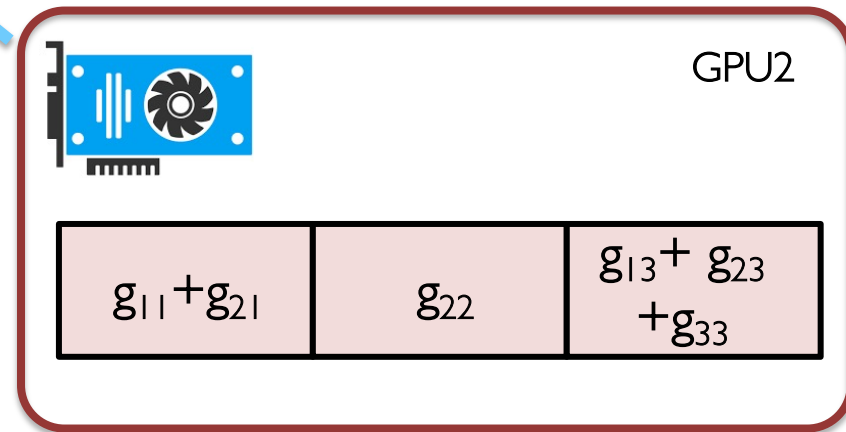
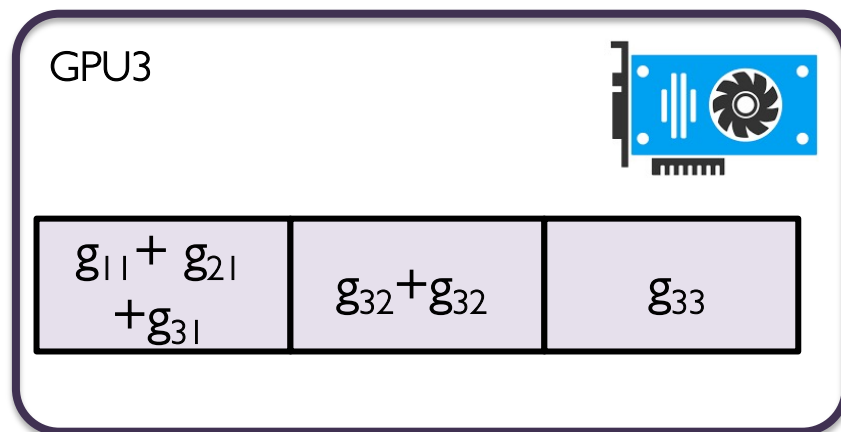
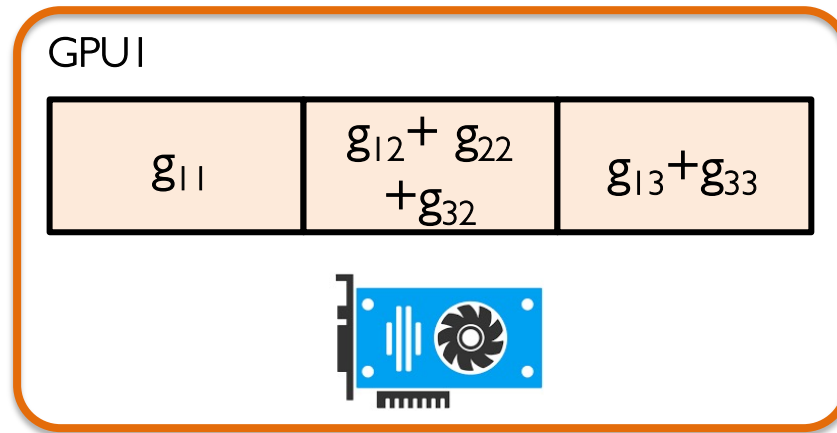
# Ring-reduce SGD

Round 2



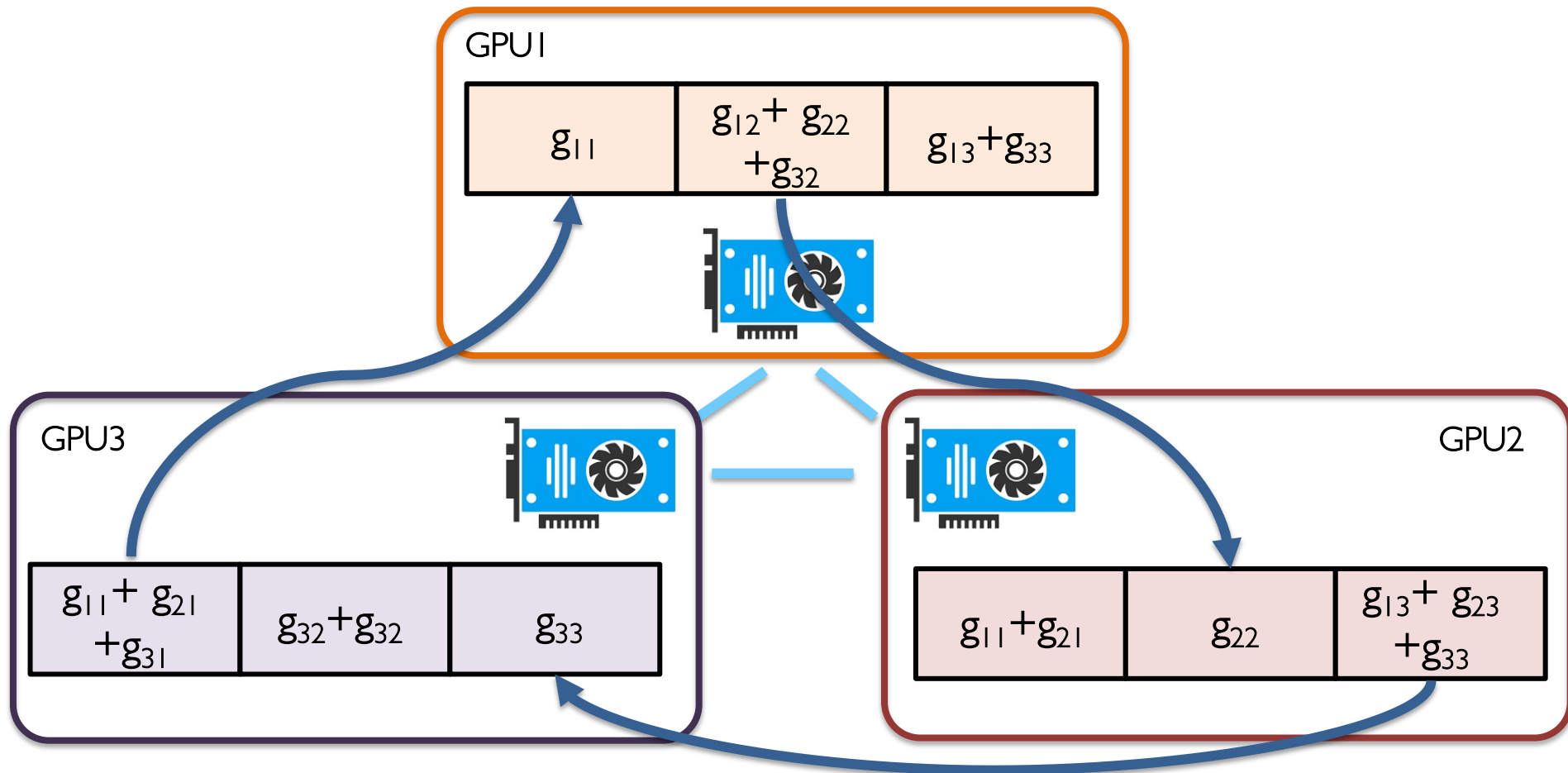
# Ring-reduce SGD

Round 2



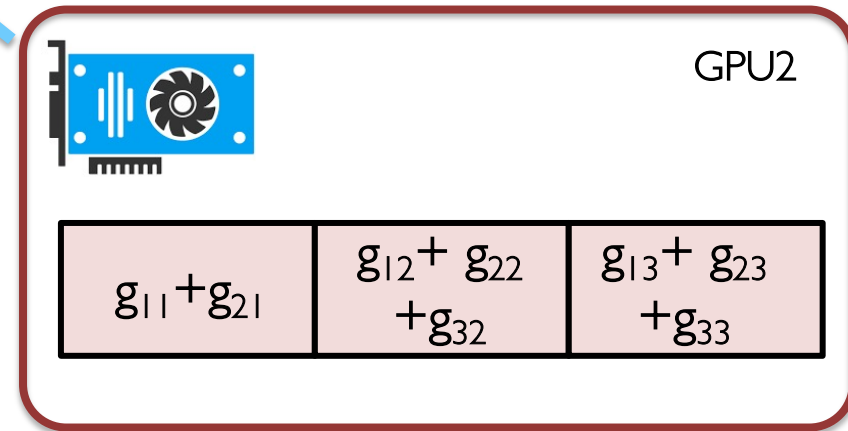
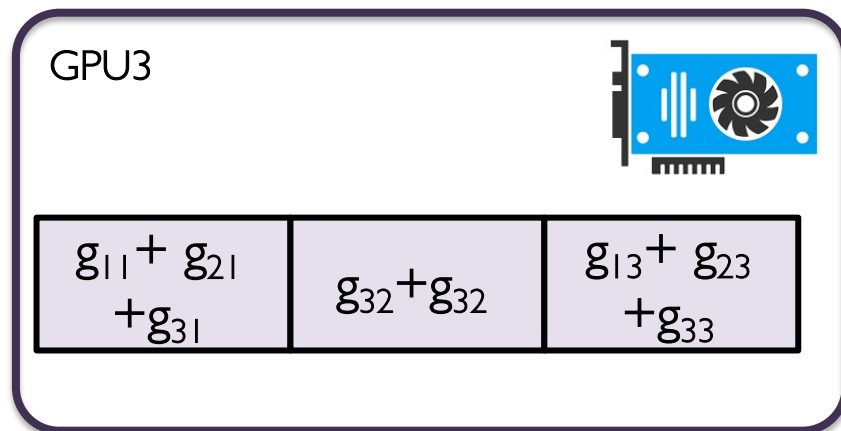
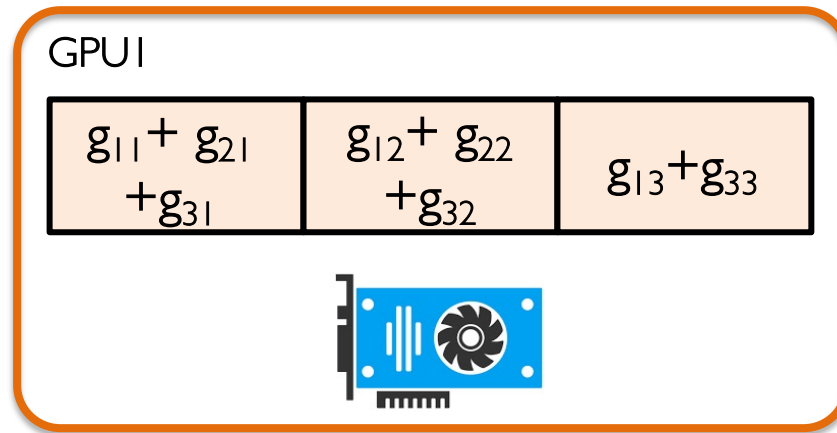
# Ring-reduce SGD

Round 3



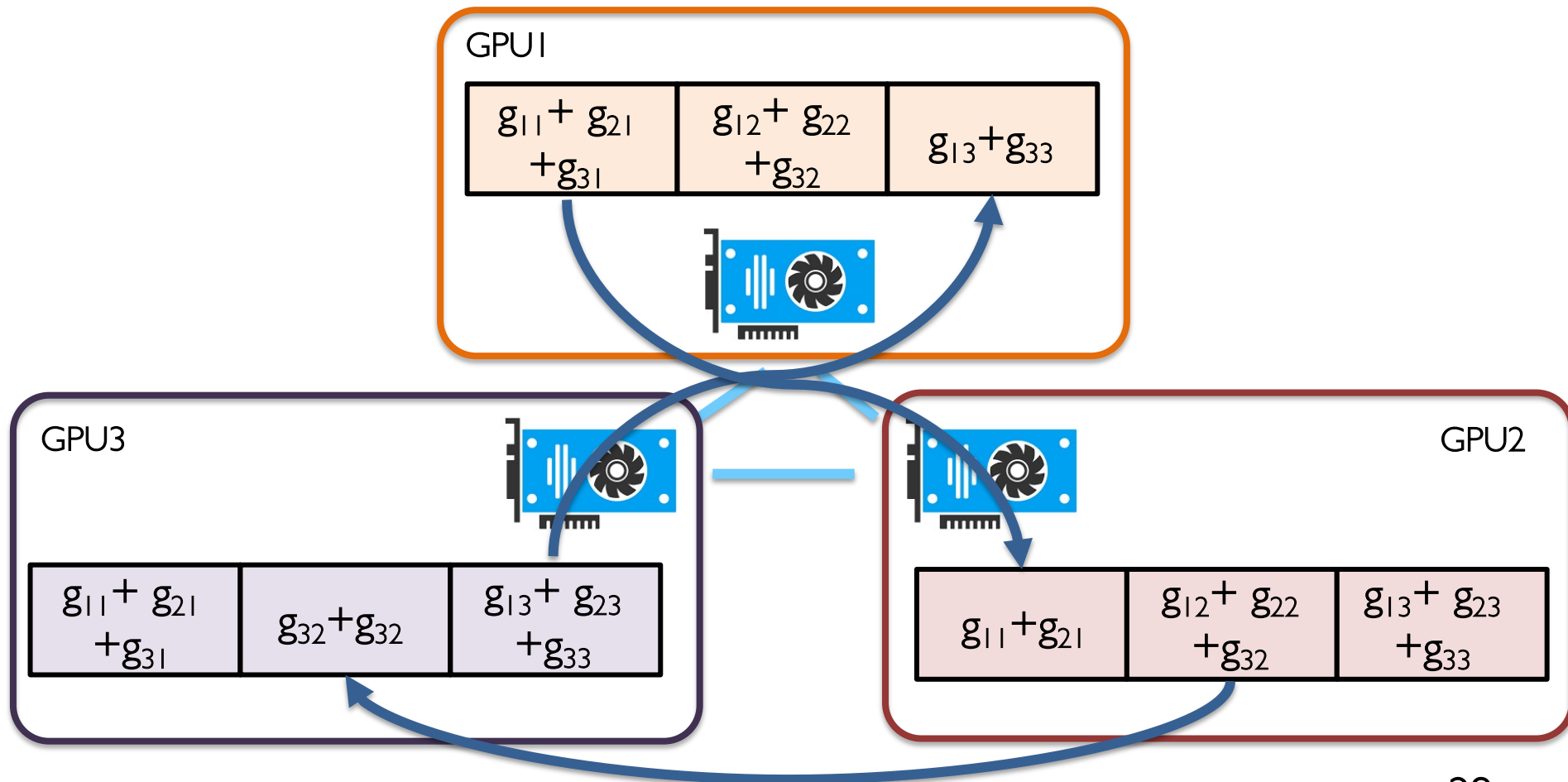
# Ring-reduce SGD

Round 3



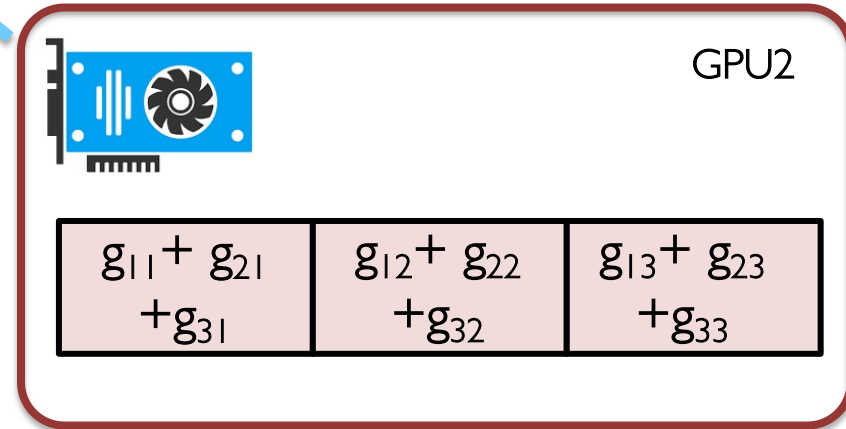
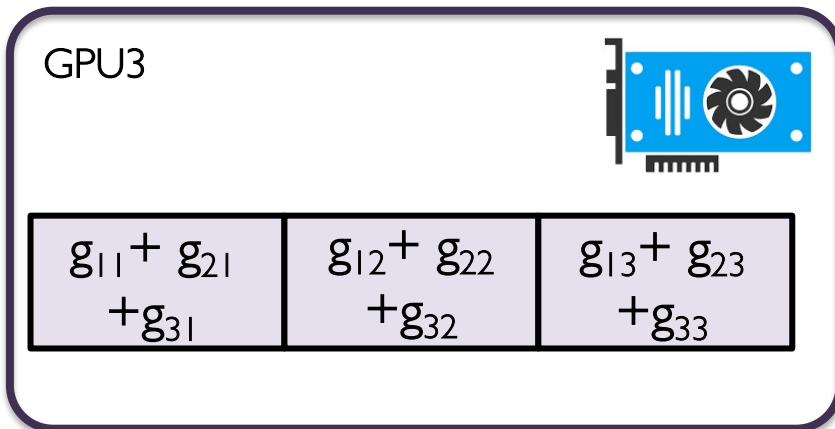
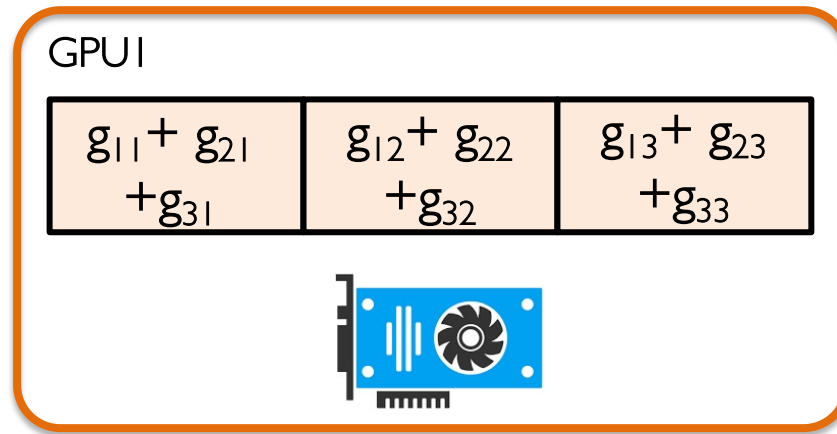
# Ring-reduce SGD

Round 4



# Ring-reduce SGD

Round 4



# Ring-reduce SGD

$\alpha - \beta$  Cost model

$$\text{Comm. cost} = \beta \times \text{Bandwidth cost} + \alpha \times \text{Latency cost}$$

Amount of info communicated per communication round  $2 \frac{d}{p} \cdot (p - 1)$

The entire number of communication rounds  $2(p - 1)$

# A Generalized Comm. Cost Model

$\alpha - \beta$  Cost model

Comm. Topology	Bandwidth	Latency
Param Server	$\beta \times 2pd$	$\alpha \times 2$
Ring-reduce	$\beta \times \frac{2(p-1)d}{p}$	$\alpha \times 2(p-1)$



# Overview

- Motivations of Distributed ML
- Parallelism & Comm. Topologies
  - ❑ Data parallelism (PS, all-reduce)
  - ❑ Model parallelism
  - ❑ Pipeline parallelism
  - ❑ Hybrid parallelism
- Communication Bottlenecks & Solutions

# An Issue of Data-parallelism

Each GPU must hold the entire copy of model parameters.

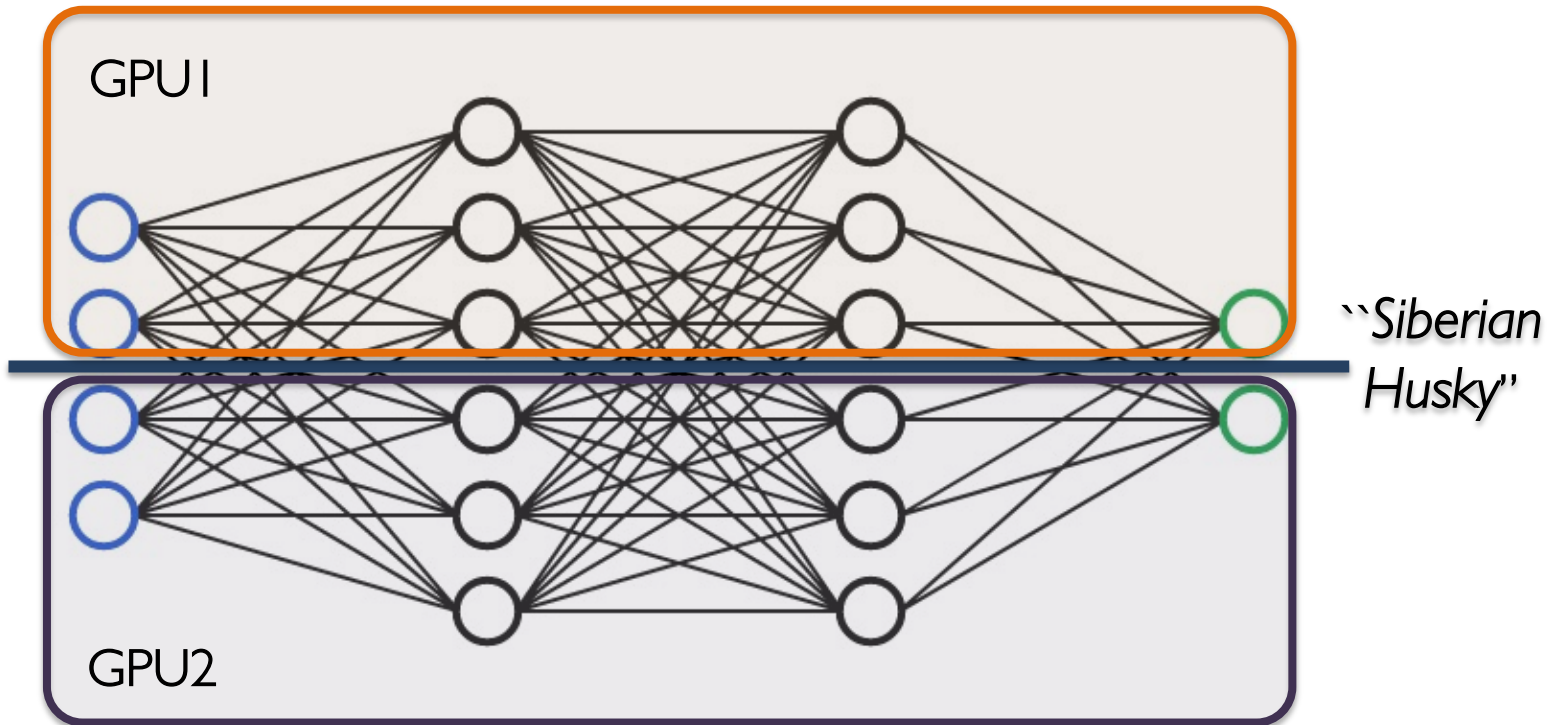
Solution:

Partition the model parameters among GPUs

- The best GPU on earth (Nvidia A100) has only 80GB of GPU memory.

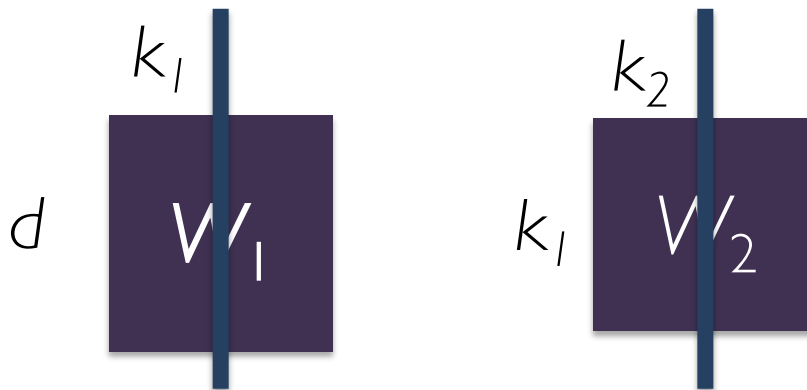
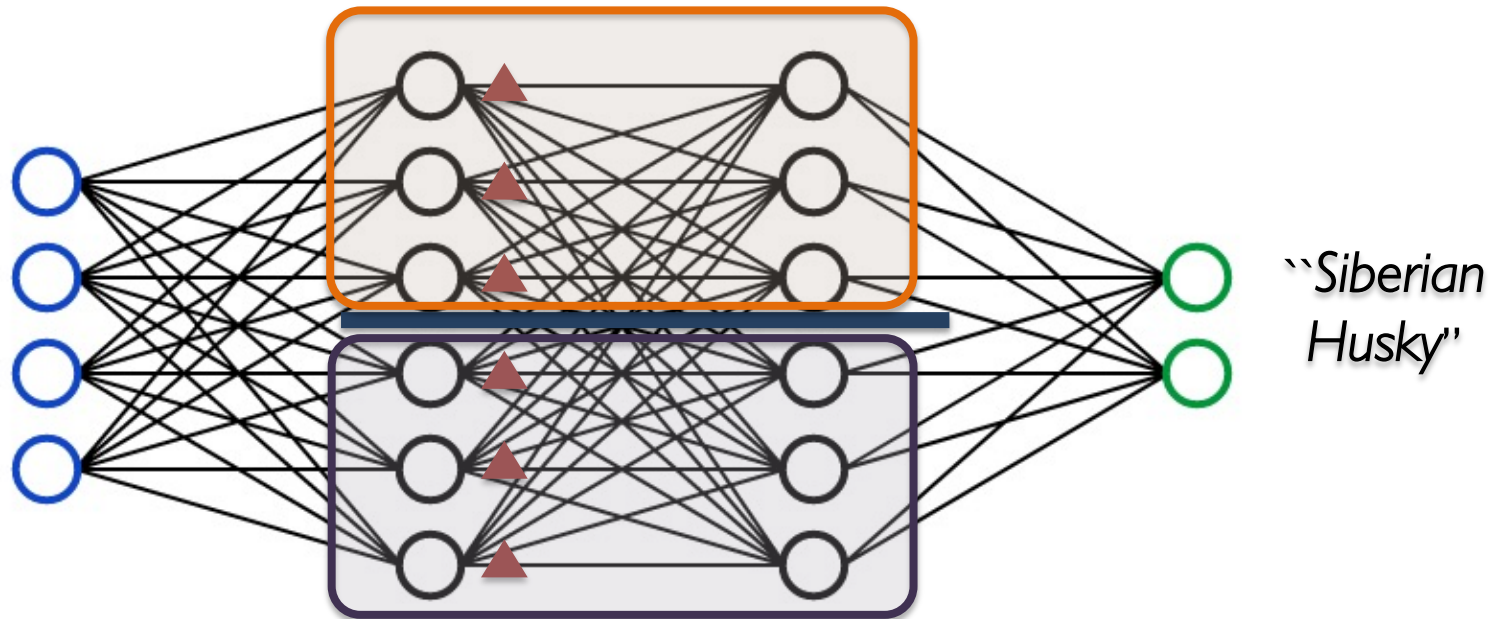
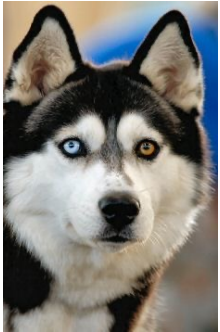
# Model Parallelism

*Also known as tensor-model parallelism*



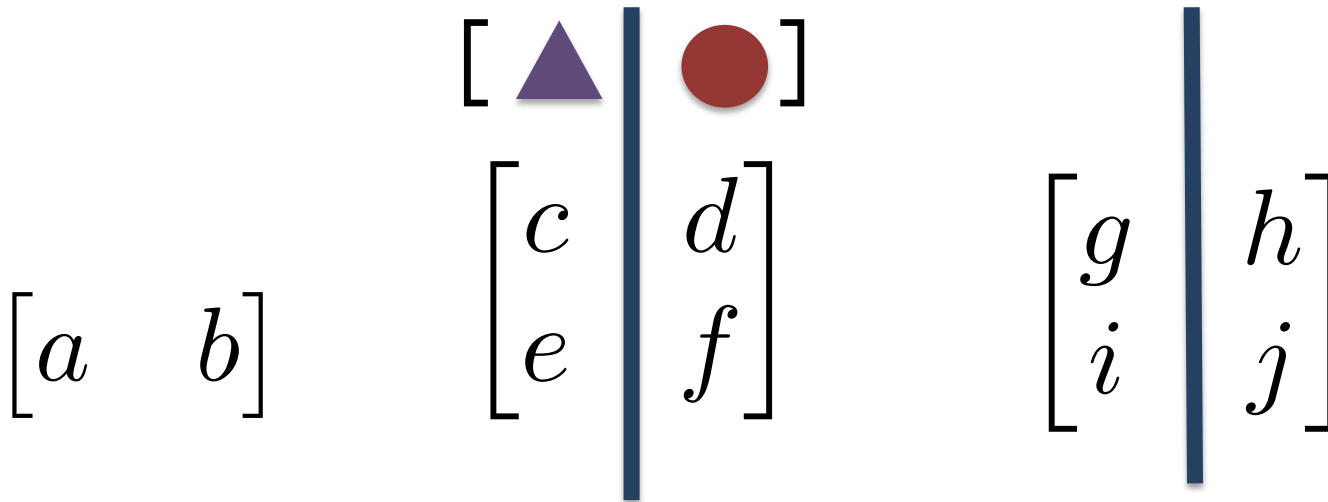
# Model Parallelism

*Also known as tensor-model parallelism*



# Model Parallelism

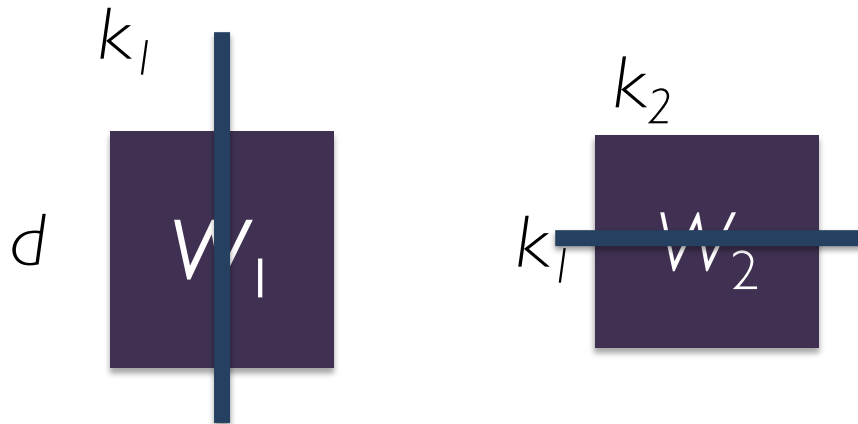
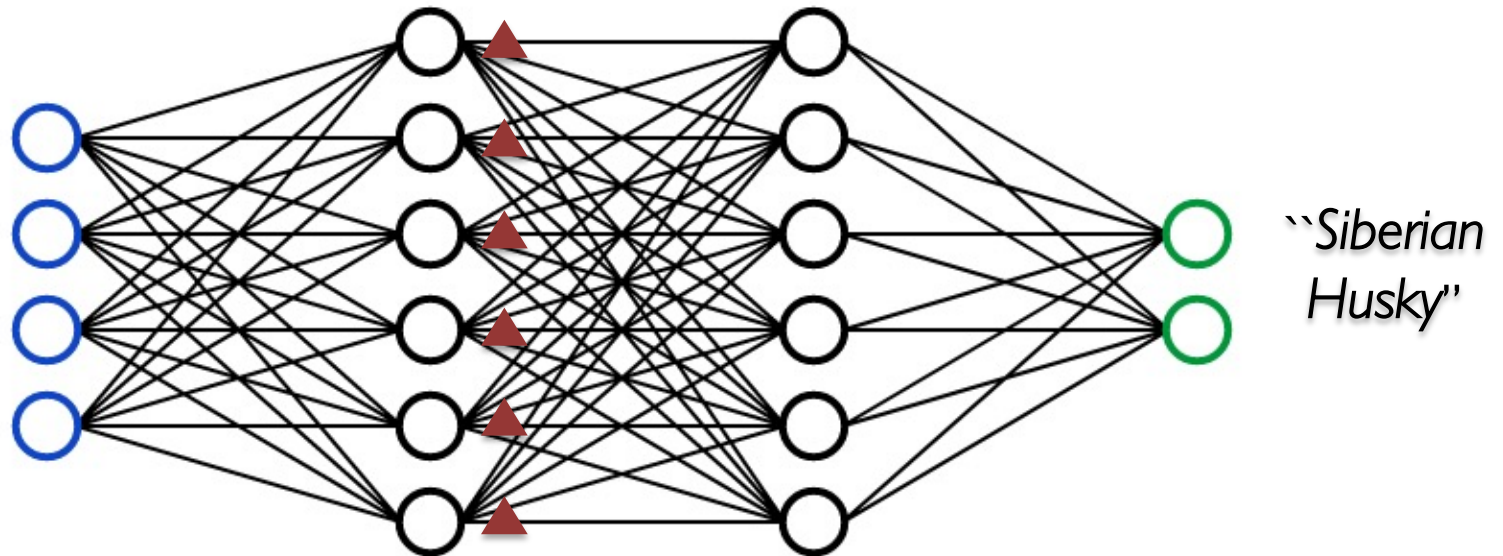
*Also known as tensor-model parallelism*



It requires to  
communicate per  
layer

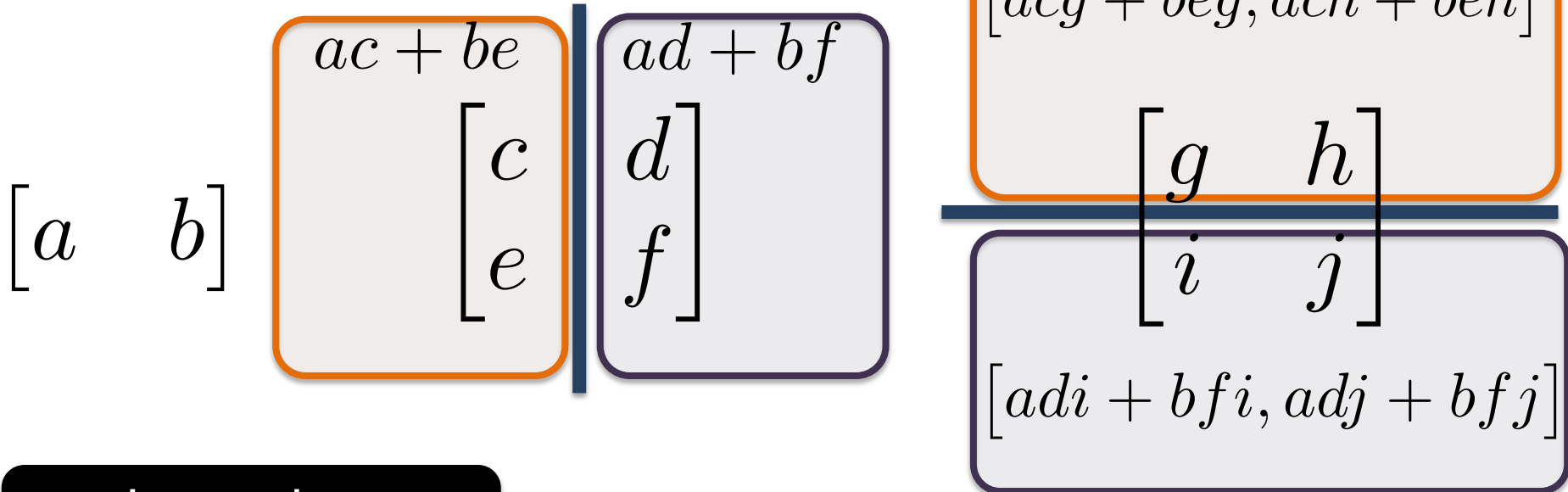
# Model Parallelism

Also known as *tensor-model parallelism*



# Model Parallelism

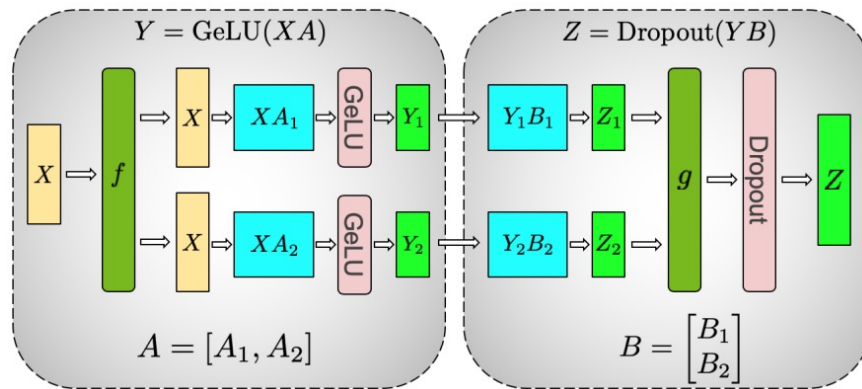
*Also known as tensor-model parallelism*



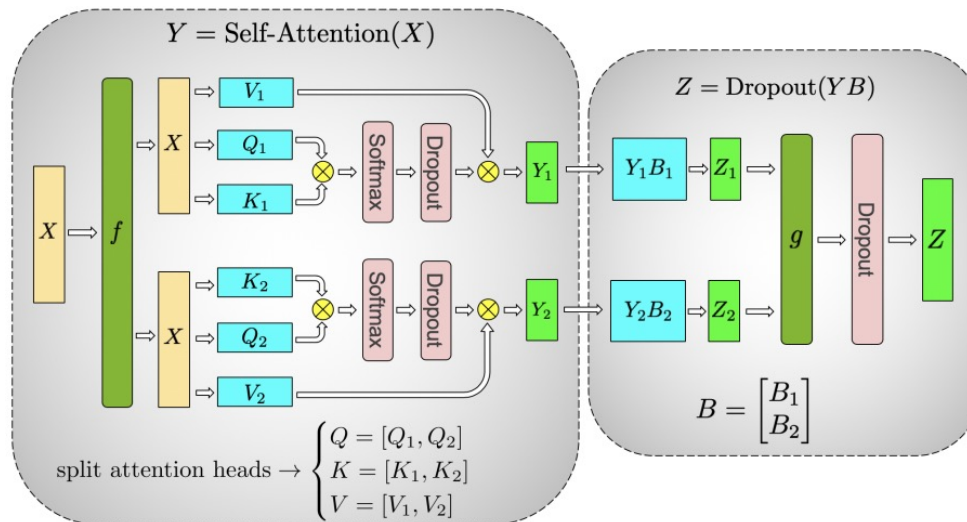
It requires  
communication  
per two layers

# Model Parallelism

Also known as *tensor-model parallelism*



MLP

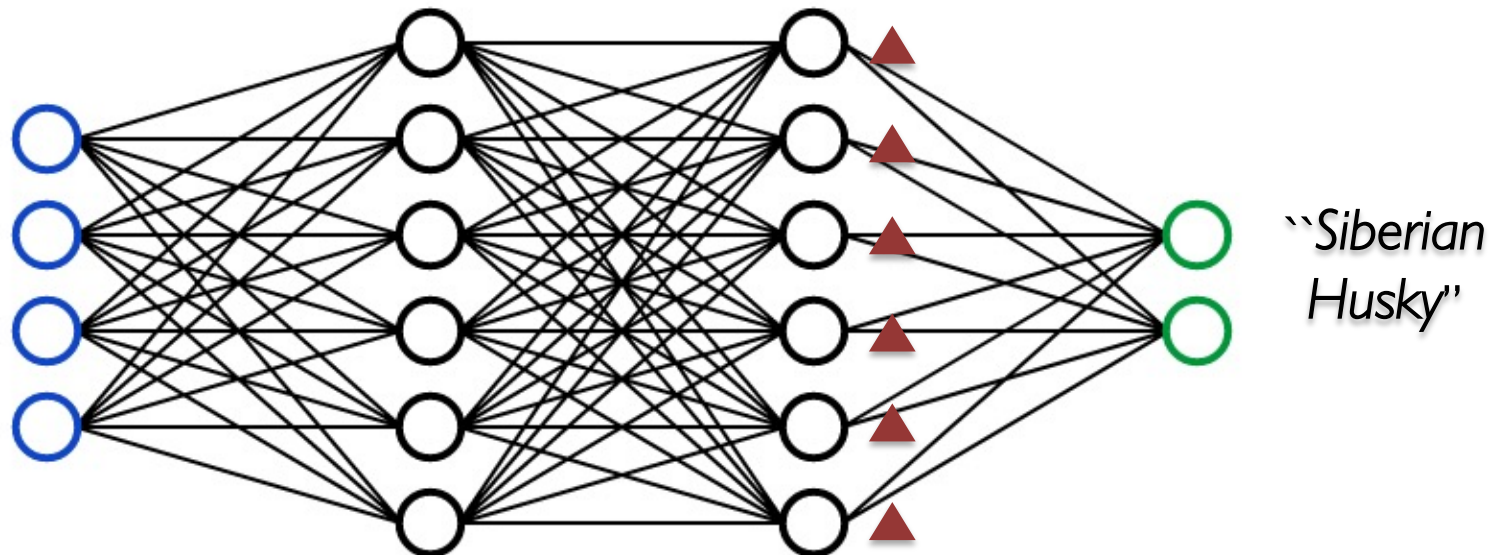


Self-Attention



# An Issue of Model Parallelism

Amount of comm. scales with # layers  
and batch size.

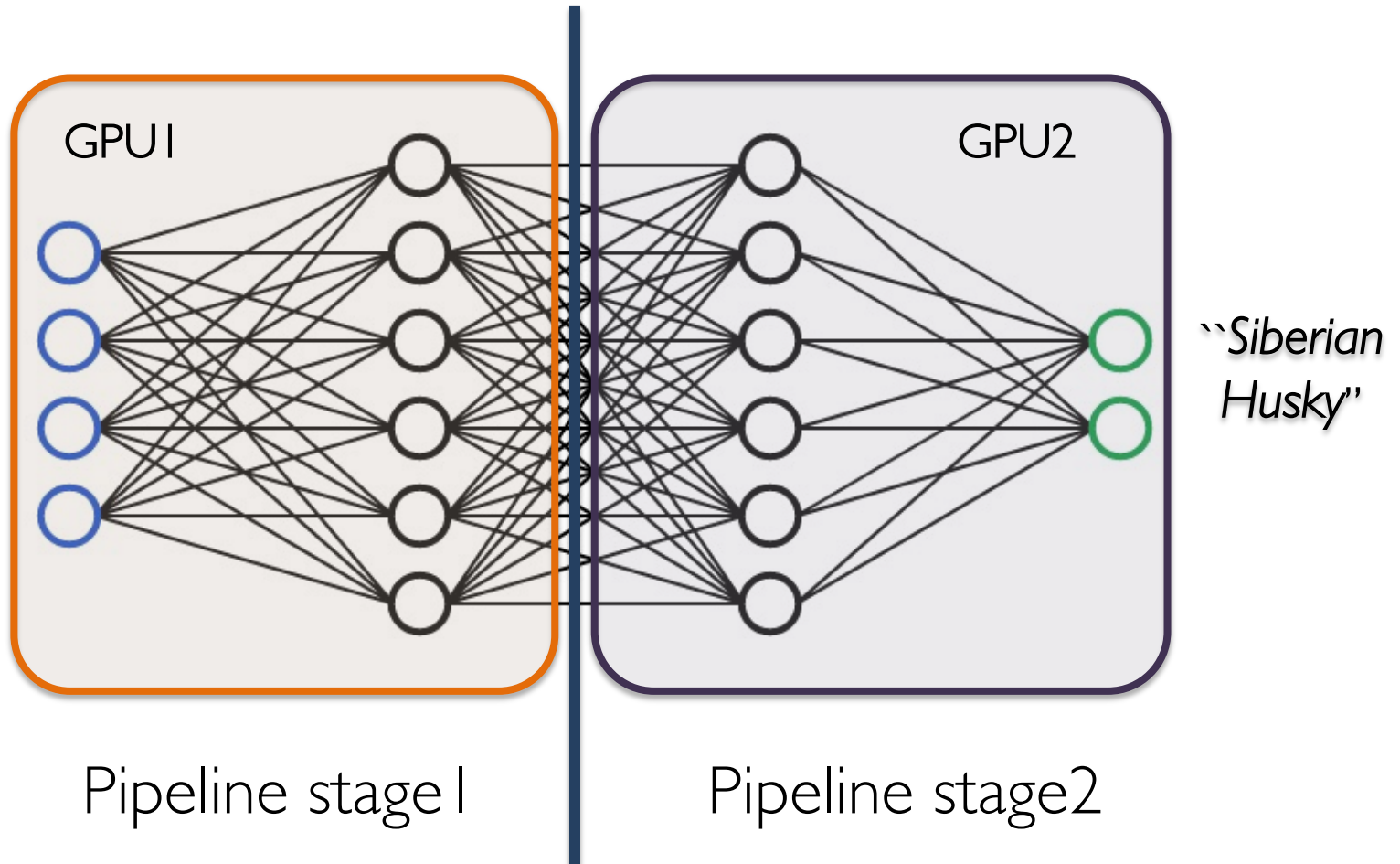
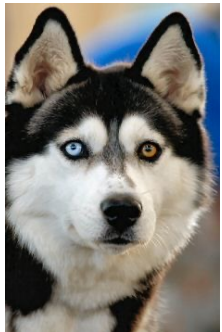


(batch size, hidden dim)

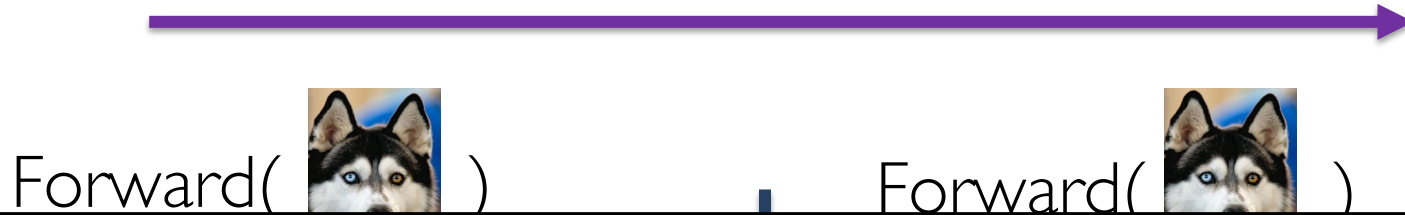
# Overview

- Motivations of Distributed ML
- Parallelism & Comm. Topologies
  - ❑ Data parallelism (PS, all-reduce)
  - ❑ Model parallelism
  - ❑ Pipeline parallelism
  - ❑ Hybrid parallelism
- Communication Bottlenecks & Solutions

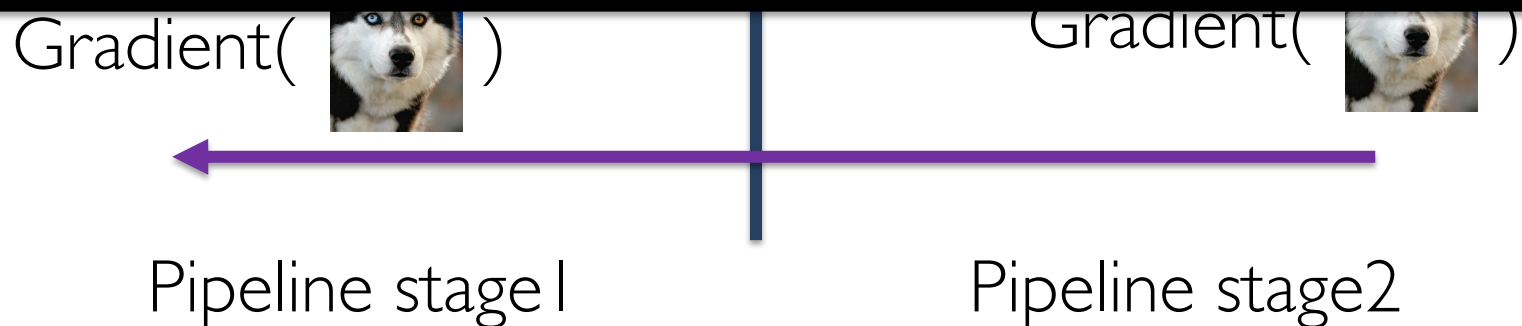
# Pipeline Parallelism



# Pipeline Parallelism

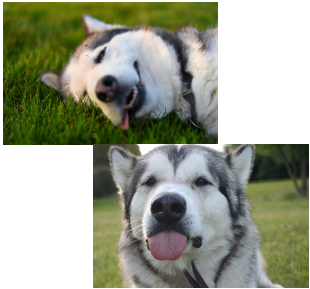
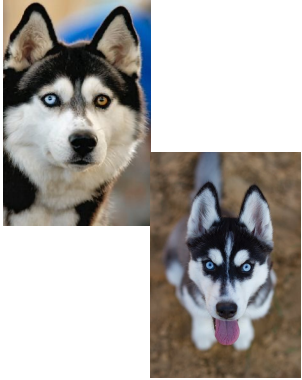


Solution:  
Split data mini-batches further

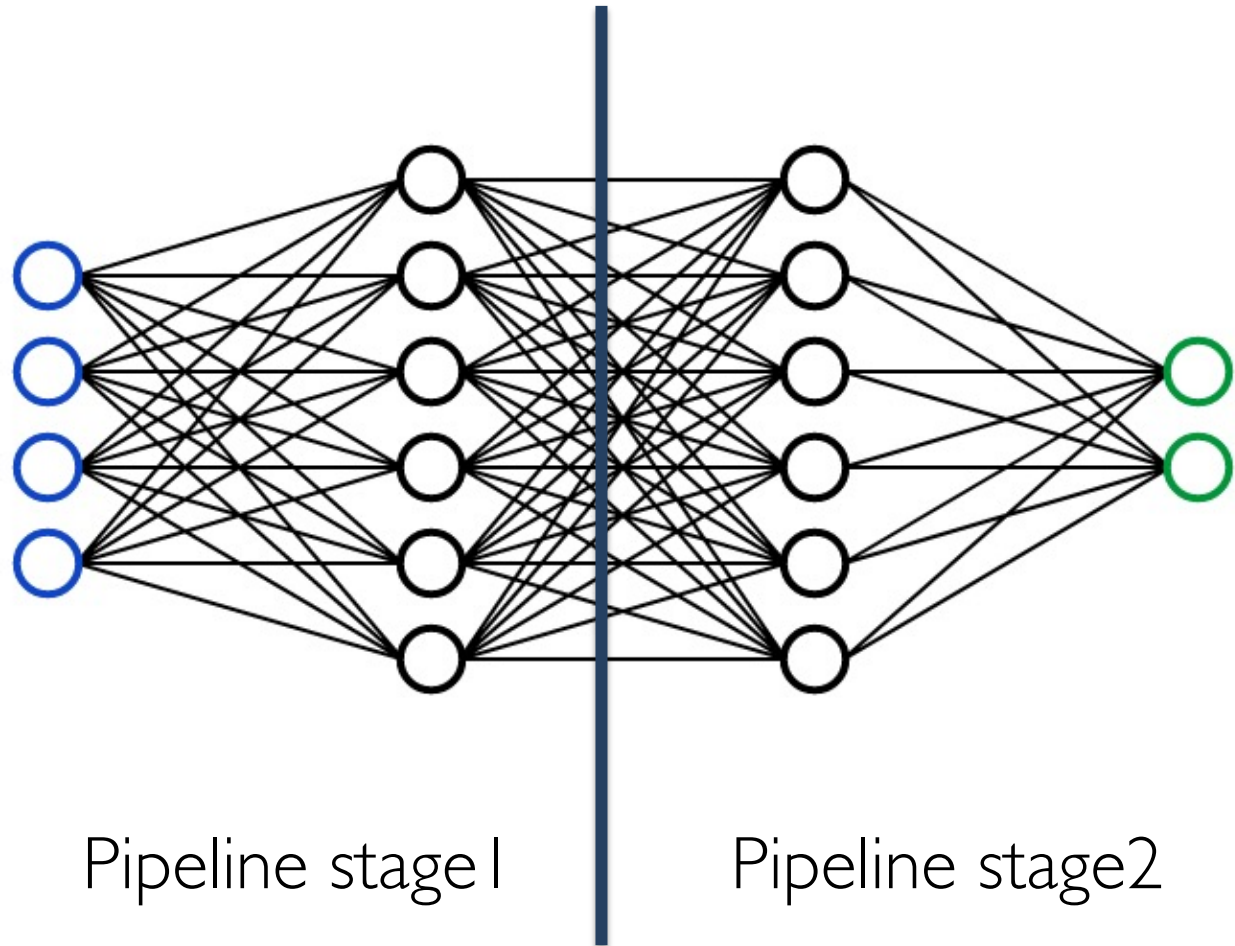


# Pipeline Parallelism

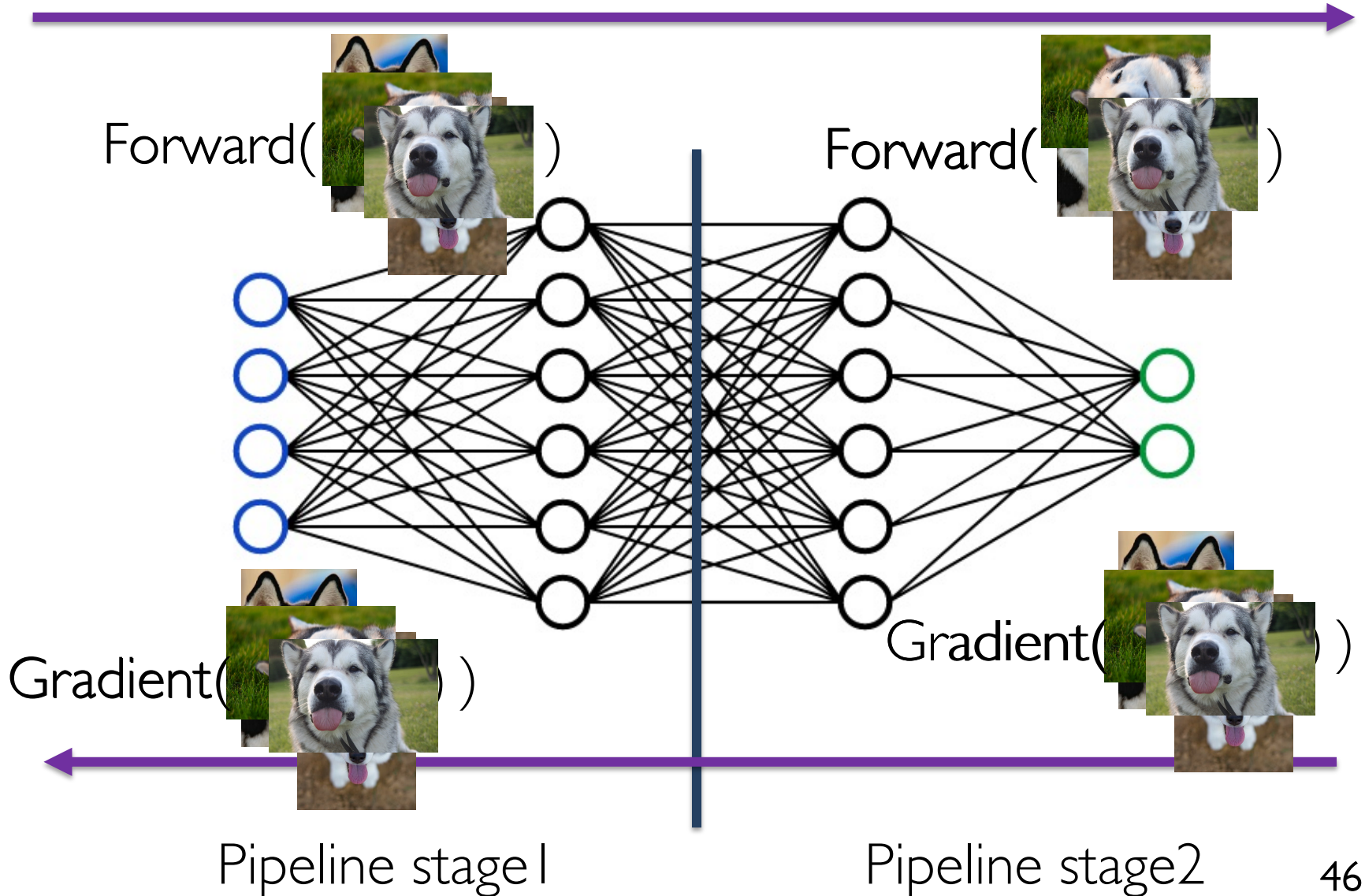
*Micro-batch 1*



*Micro-batch 2*

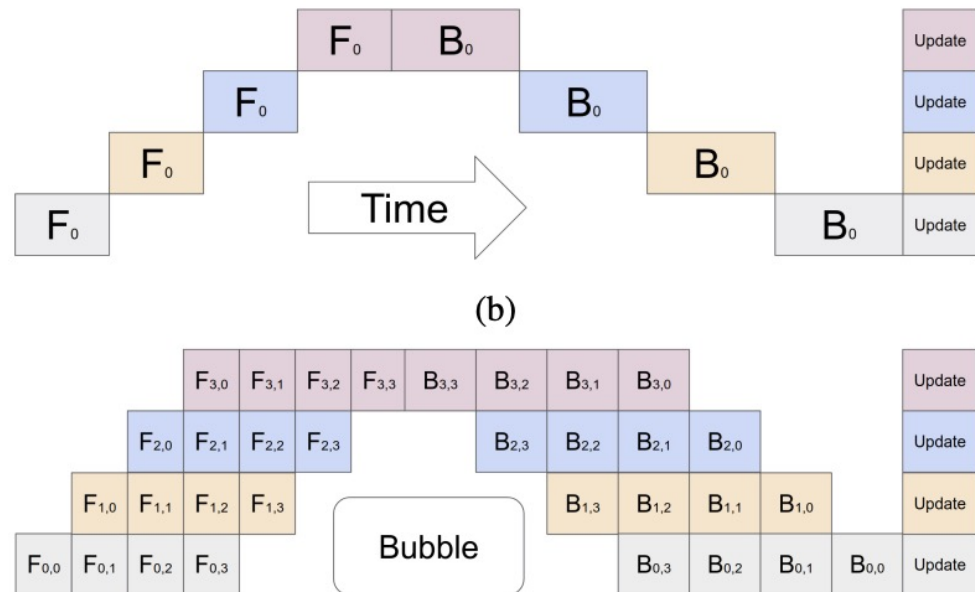
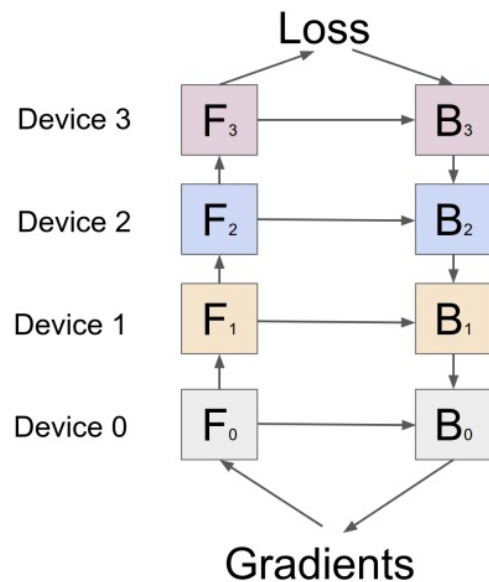


# Pipeline Parallelism



# Issue of Pipeline Parallelism

Pipeline bubble needs to be controlled carefully.

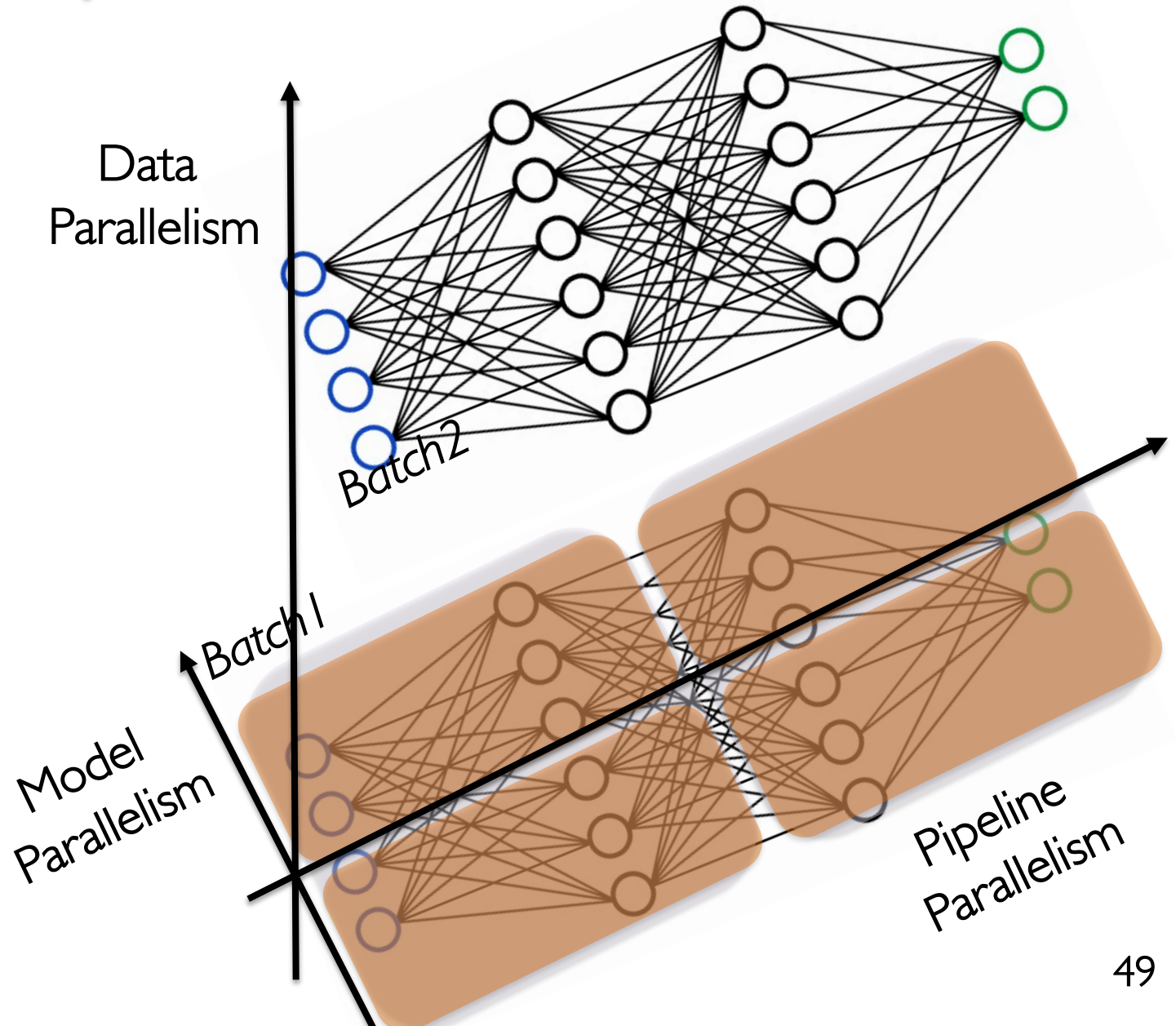


# Overview

- Motivations of Distributed ML
- Parallelism & Comm. Topologies
  - ❑ Data parallelism (PS, all-reduce)
  - ❑ Model parallelism
  - ❑ Pipeline parallelism
  - ❑ Hybrid parallelism
- Communication Bottlenecks & Solutions



# Hybrid/3D Parallelism



# Overview

- Motivations of Distributed ML
- Parallelism & Comm. Topologies
  - ❑ Data parallelism (PS, all-reduce)
  - ❑ Model parallelism
  - ❑ Pipeline parallelism
  - ❑ Hybrid parallelism
- Communication Bottlenecks & Solutions

# Data-parallel mini-batch SGD

$$\min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{w}; \mathbf{z}_i) \quad \mathcal{X} = \{z_1, \dots, z_n\}$$

For iteration  $t$  :

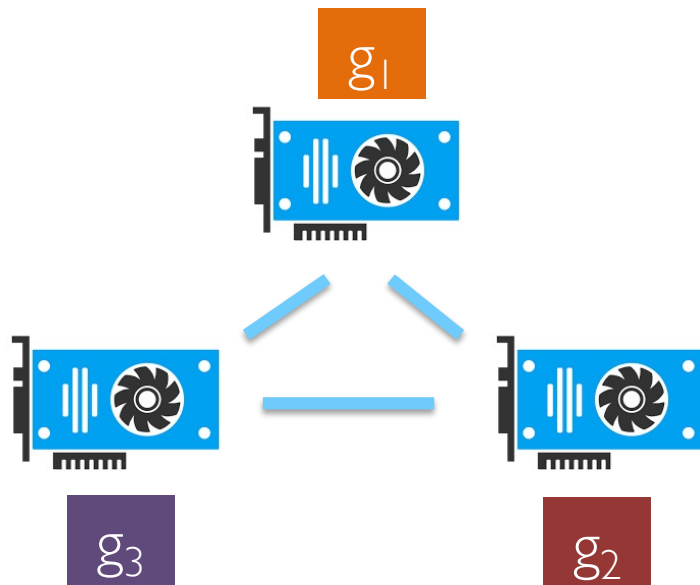
- Sample a mini-batch:  $\mathcal{B}_t \subset \mathcal{X}$

- Compute gradient:  $\nabla \frac{1}{|\mathcal{B}_t|} \sum_{j=1}^{|\mathcal{B}_t|} \ell(w_t; z_j)$

$$= \frac{1}{|\mathcal{B}_t|} \sum_{j=1}^{|\mathcal{B}_t|} \nabla \ell(w_t; z_j)$$

# Data-parallel SGD

All-reduce SGD



Repeat until convergence

[Cotter, Shamir, Srebro, Sridharan, NIPS 11]

[Dekel, Gilad-Bachrach, Shamir, Xiao, JMLR 2012]

[Friedlander and Schmidt, SIAM JSC 2012]

[Takác, Bijral, Richtárik, Srebro, ICML 2013]

[Li, Zhang, Chen, Smola, KDD 2014]

[Jain, Kakade, Kidambi, Netrapalli, Sidford, arxiv'16]

[De, Yadav, Jacobs, Goldstein, arxiv'16]

The ideal speedup should be proportional to  
#compute nodes

TL;DR:

Communication bottlenecks

# Gradient compression

Communication / worker:

$$O(\text{size of gradient}) * 32 \text{ bits}$$

Gradient Quantization

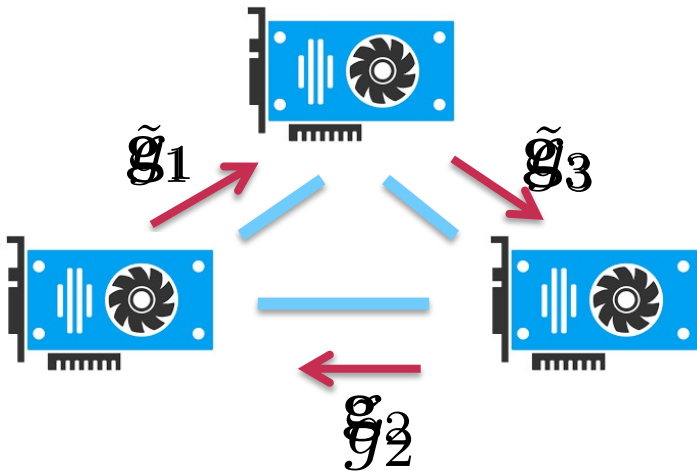
Quantize to precision:

$$O(\text{size of gradient}) * 2/4/8 \text{ bits}$$

Gradient Sparsification

Sparsified Gradients (k-sparse):

$$O(k)$$



# Top-K (Sparsification)

Examples:

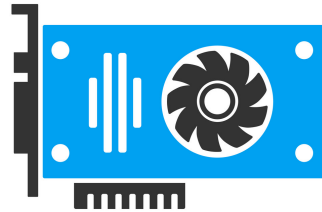
$$\text{top}_2 \left( \begin{bmatrix} 2 \\ -3 \\ 1 \\ -4 \end{bmatrix} \right) \Rightarrow \begin{bmatrix} 0 \\ -3 \\ 0 \\ -4 \end{bmatrix}$$

# *sign*SGD (Quantization)

$$\text{signsgd} \left( \begin{bmatrix} -0.2 \\ 0.1 \\ 0.3 \\ -1.8 \end{bmatrix} \right) \Rightarrow \begin{bmatrix} -1 \\ 1 \\ 1 \\ -1 \end{bmatrix}$$



# ATOMO & PowerSGD (Low-rank Factorization)



Factorization  $(U \quad V)$

# Next Week

- More about gradient compression.
- Convergence rate of gradient compressed distributed training.
- Federated learning.

# Reading List

- Sergeev, A. and Del Balso, M., 2018. Horovod: fast and easy distributed deep learning in TensorFlow. arXiv preprint arXiv:1802.05799.
- Shoeybi, M., Patwary, M., Puri, R., LeGresley, P., Casper, J. and Catanzaro, B., 2019. Megatron-1m: Training multi-billion parameter language models using model parallelism. arXiv preprint arXiv:1909.08053.
- Jia, Z., Zaharia, M. and Aiken, A., 2019. Beyond Data and Model Parallelism for Deep Neural Networks. Proceedings of Machine Learning and Systems, 1, pp.1-13.
- Narayanan, D., Harlap, A., Phanishayee, A., Seshadri, V., Devanur, N.R., Ganger, G.R., Gibbons, P.B. and Zaharia, M., 2019, October. PipeDream: generalized pipeline parallelism for DNN training. In Proceedings of the 27th ACM Symposium on Operating Systems Principles (pp. 1-15).
- Gholami, A., Azad, A., Jin, P., Keutzer, K. and Buluc, A., 2018, July. Integrated model, batch, and domain parallelism in training neural networks. In Proceedings of the 30th on Symposium on Parallelism in Algorithms and Architectures (pp. 77-86).
- Ben-Nun, T. and Hoefler, T., 2019. Demystifying parallel and distributed deep learning: An in-depth concurrency analysis. ACM Computing Surveys (CSUR), 52(4), pp.1-43.