



# Optimization for Deep Learning (Part II)

Kun Yuan

Center for Machine Learning Research @ Peking University

Fudan University

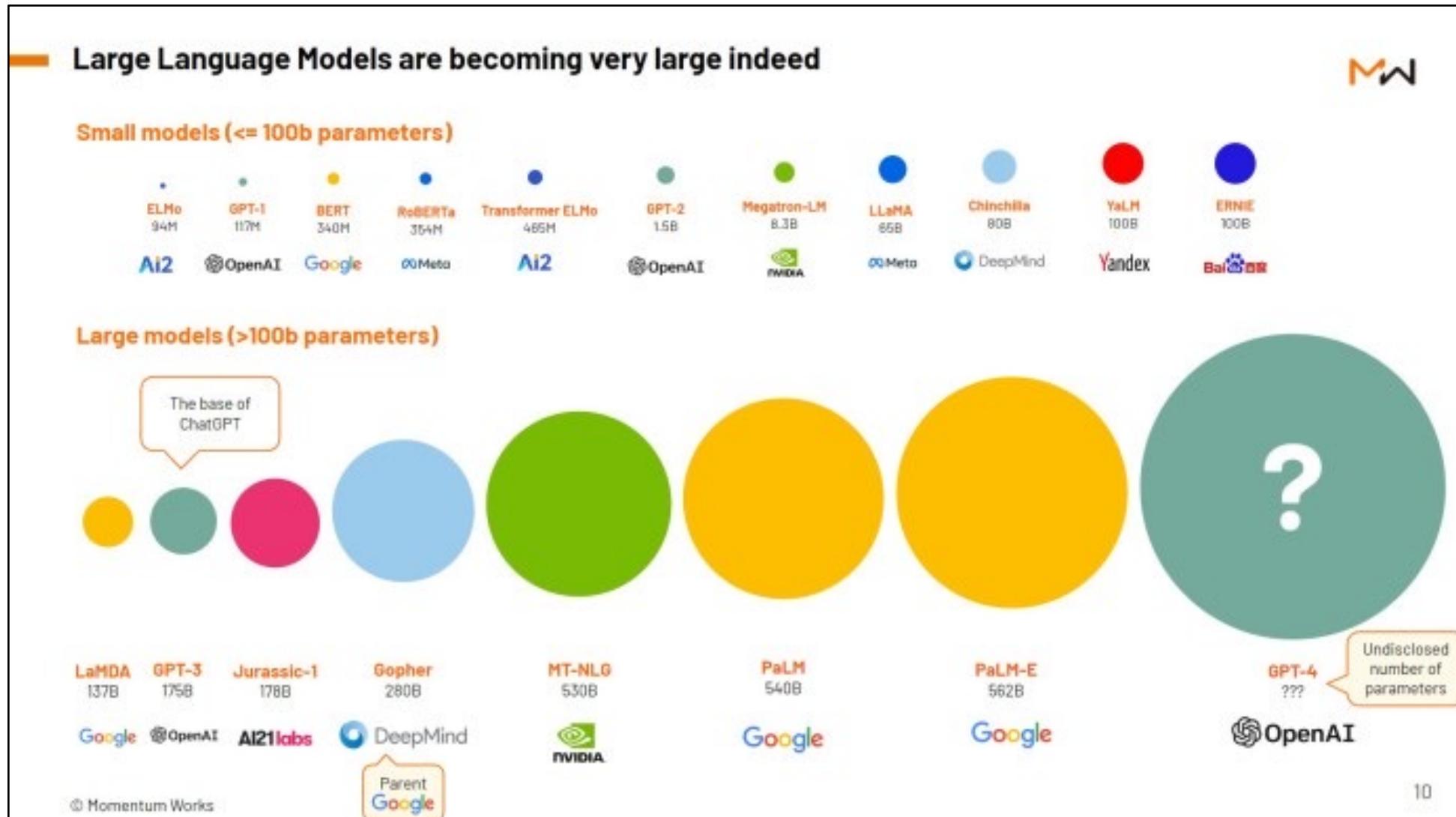
June 8, 2024

## PART 01

---

### LLM parameters and memories

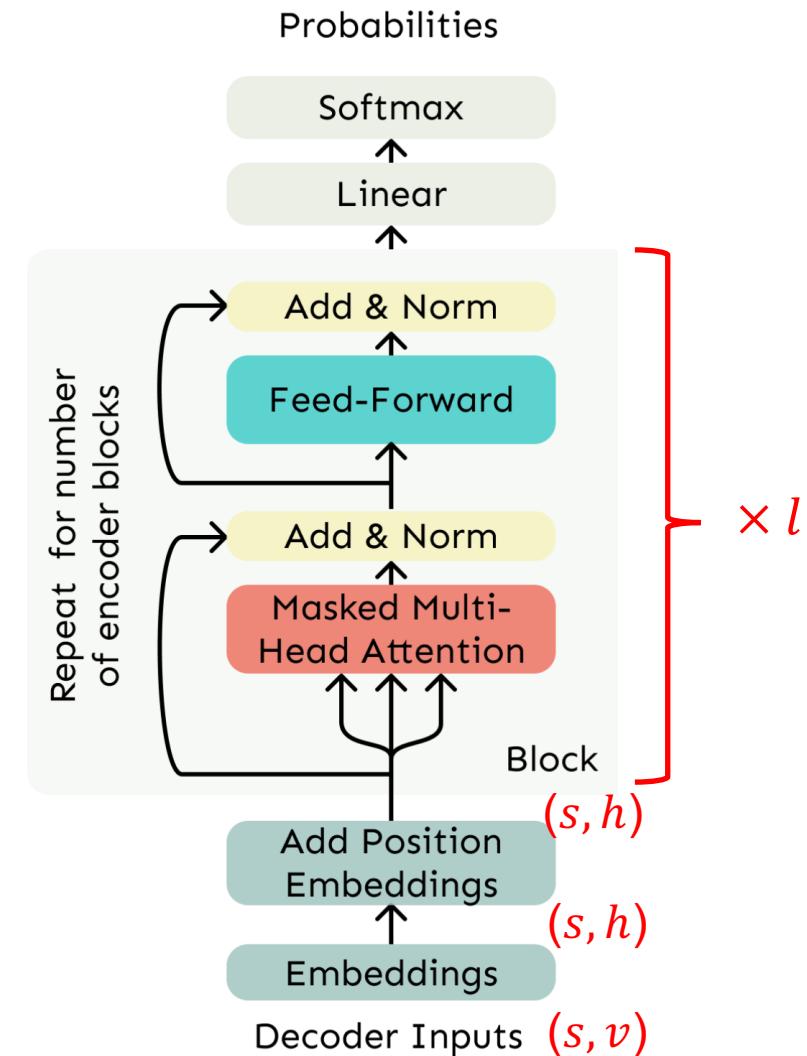
# Large language models



[The emergence of Large Language Models (LLMs)]

**GPTs are based on decoder-only transformers**

- Number of the transformer layers:  $l$
  - Sequence length:  $s$
  - Vocabulary size:  $\nu$
  - Embedding representation dims:  $h$



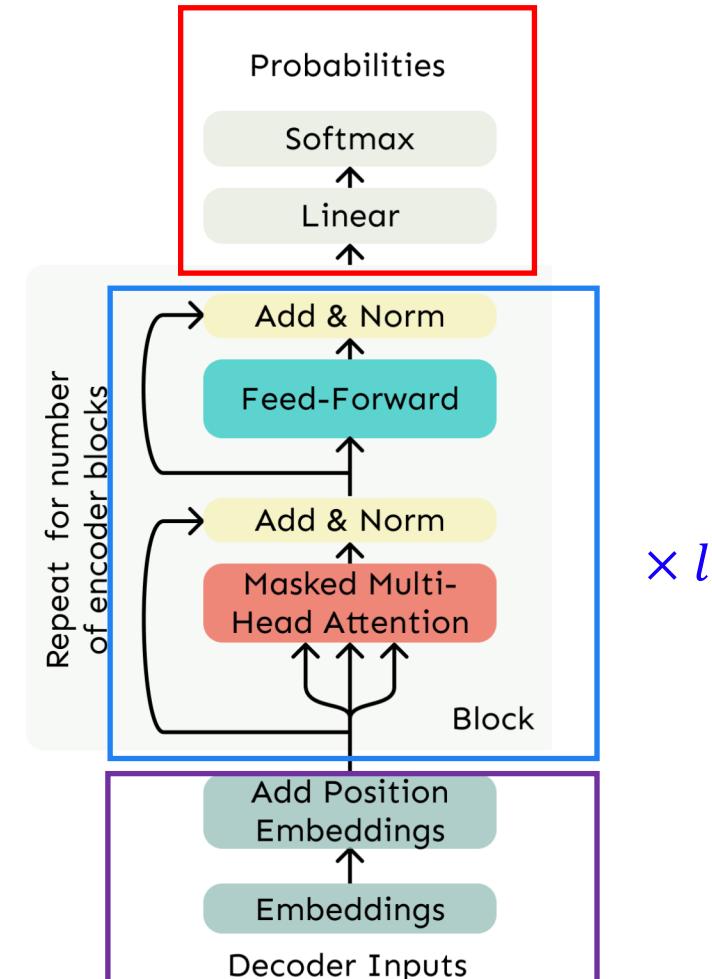
# Total parameters

---

- Embeddings:  $vh$
- Attention blocks:  $12lh^2$
- Probability predictions:  $vh$

Total parameters:

$$12lh^2 + 2vh$$



## Example: LLaMA parameters

---

- Now we compare our theoretical evaluations with LLaMA model
- $12\ell h^2 + 2vh$  is a very accurate estimation

实际参数量	Embedding h	Attention层数l	Vocab大小v	预估参数量
6.7B	4096	32	32000	6,704,594,944
13.0B	5120	40	32000	12,910,592,000
32.5B	6656	60	32000	32,323,665,920
65.2B	8192	80	32000	64,948,797,440

Memory = Model + Gradient + Optimizer states + Activations

- Recall the Adam algorithm

$$g_k = \nabla F(x_k; \xi_k)$$

$$m_k = \beta_1 m_{k-1} + (1 - \beta_1) g_k$$

$$s_k = \beta_2 s_{k-1} + (1 - \beta_2) g_k \odot g_k$$

$$x_{k+1} = x_k - \frac{\gamma}{\sqrt{s_k} + \epsilon} \odot m_k$$

- The optimizer states are twice the number of model parameters

Memory = Model + Gradient + Optimizer states + Activations

- Given a model with  $P$  parameters, gradient will consume  $P$  parameters, and Optimizer states will consume  $2P$  parameters; **4P parameters in total.**
- When using FP32 to store parameters, each parameter takes **4** Bytes
- When using FP16 or BF16 to store parameters, each parameter takes **2** Bytes

Memory = **Model + Gradient + Optimizer states + Activations**

$$(48l h^2 + bl(2s^2a + 14sh)) \times 4 \text{ Bytes}$$

- When hidden state **h** is large, the model parameters dominate the memory
- When batch-size **b** or sequence length **s** is large, the activation dominates the memory

# Memory examples

---

Model Name	$n_{\text{params}}$	$n_{\text{layers}}$	$d_{\text{model}}$	$n_{\text{heads}}$	$d_{\text{head}}$	Batch Size	Learning Rate
GPT-3 Small	125M	12	768	12	64	0.5M	$6.0 \times 10^{-4}$
GPT-3 Medium	350M	24	1024	16	64	0.5M	$3.0 \times 10^{-4}$
GPT-3 Large	760M	24	1536	16	96	0.5M	$2.5 \times 10^{-4}$
GPT-3 XL	1.3B	24	2048	24	128	1M	$2.0 \times 10^{-4}$
GPT-3 2.7B	2.7B	32	2560	32	80	1M	$1.6 \times 10^{-4}$
GPT-3 6.7B	6.7B	32	4096	32	128	2M	$1.2 \times 10^{-4}$
GPT-3 13B	13.0B	40	5140	40	128	2M	$1.0 \times 10^{-4}$
GPT-3 175B or “GPT-3”	175.0B	96	12288	96	128	3.2M	$0.6 \times 10^{-4}$

GPT3 has 175B parameters; its model consumes  $4 \times 175 \times 10^9$  Bytes = 700 GB

Its gradient takes 700 GB parameters; Optimizer states take 1.4 TB

GPT has sequence length  $s = 2048$ . When  $b=1$ , its activation takes 444 GB, 63% of the model

When  $b=128$ , its activation is 81 times of the model size

The massive parameters and memory costs in LLM  
calls for distributed optimizers

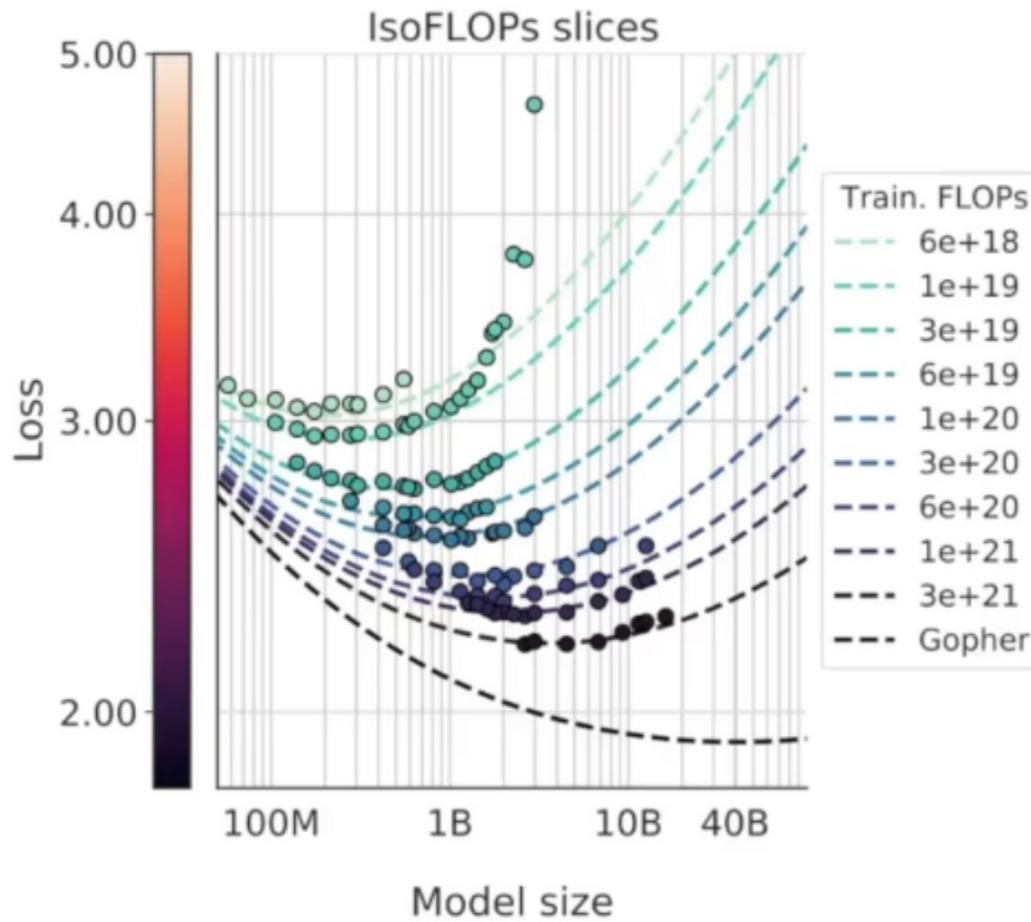
## PART 02

---

### **LLM training systems and their scalability**

# Chinchilla law

[Training Compute-Optimal Large Language Models, 2022]



The “RGB” elements in LLM:

**Larger dataset**  
+

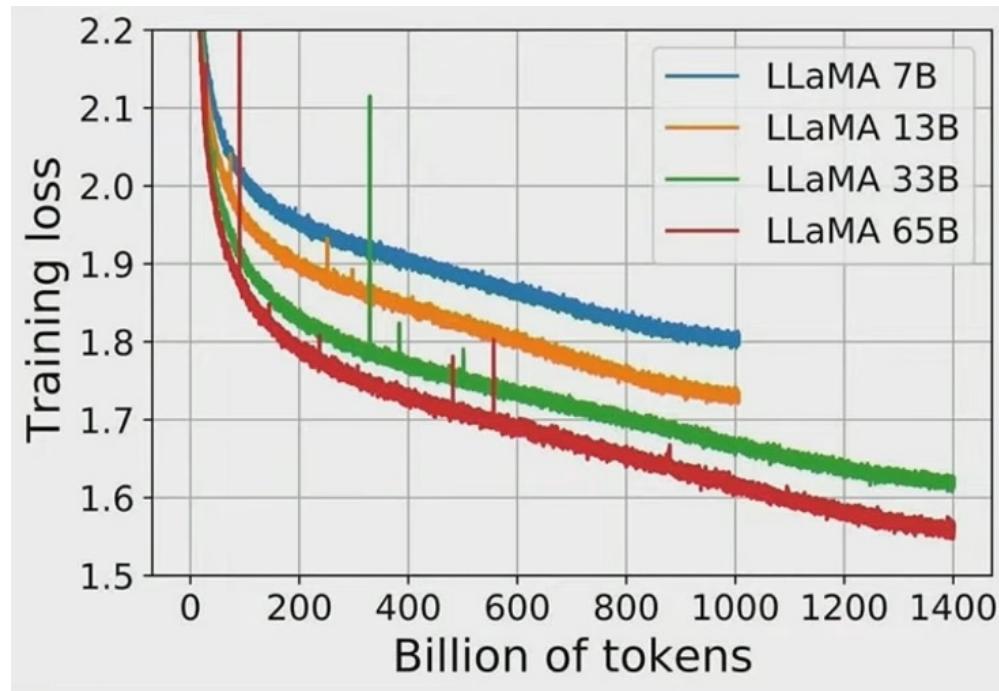
**Bigger model**  
+

**Longer training**

= **Better LLM**

# LLaMA follows Chinchilla law

[LLaMA: Open and Efficient Foundation Language Models, 2023]



According to Chinchilla law:

- LLM model gets increasingly bigger
- Dataset gets increasingly larger
- Distributed training gets increasingly important

# Thousands of GPUs are needed to train LLM



[State of GPT, 2023]

## 2 example models

### GPT-3 (2020)

50,257 vocabulary size  
2048 context length  
175B parameters  
Trained on 300B tokens

Model Name	$n_{\text{params}}$	$n_{\text{layers}}$	$d_{\text{model}}$	$n_{\text{heads}}$	$d_{\text{head}}$	Batch Size	Learning Rate
GPT-3 Small	125M	12	768	12	64	0.5M	$6.0 \times 10^{-4}$
GPT-3 Medium	350M	24	1024	16	64	0.5M	$3.0 \times 10^{-4}$
GPT-3 Large	760M	24	1536	16	96	0.5M	$2.5 \times 10^{-4}$
GPT-3 XL	1.3B	24	2048	24	128	1M	$2.0 \times 10^{-4}$
GPT-3 2.7B	2.7B	32	2560	32	80	1M	$1.6 \times 10^{-4}$
GPT-3 6.7B	6.7B	32	4096	32	128	2M	$1.2 \times 10^{-4}$
GPT-3 13B	13.0B	40	5140	40	128	2M	$1.0 \times 10^{-4}$
GPT-3 175B or "GPT-3"	175.0B	96	12288	96	128	3.2M	$0.6 \times 10^{-4}$

Table 2.1: Sizes, architectures, and learning hyper-parameters (batch size in tokens and learning rate) of the models which we trained. All models were trained for a total of 300 billion tokens.

### Training: (rough order of magnitude to have in mind)

- O(1,000 - 10,000) V100 GPUs
- O(1) month of training
- O(1-10) \$M

### LLaMA (2023)

32,000 vocabulary size  
2048 context length  
65B parameters  
Trained on 1-1.4T tokens

params	dimension	$n_{\text{heads}}$	$n_{\text{layers}}$	learning rate	batch size	$n_{\text{tokens}}$
6.7B	4096	32	32	$3.0e^{-4}$	4M	1.0T
13.0B	5120	40	40	$3.0e^{-4}$	4M	1.0T
32.5B	6656	52	60	$1.5e^{-4}$	4M	1.4T
65.2B	8192	64	80	$1.5e^{-4}$	4M	1.4T

Table 2: Model sizes, architectures, and optimization hyper-parameters.

### Training for 65B model:

- 2,048 A100 GPUs
- 21 days of training
- \$5M

Feb. 23, 2024

## MegaScale: Scaling Large Language Model Training to More Than **10,000 GPUs**

Ziheng Jiang<sup>1,\*</sup> Haibin Lin<sup>1,\*</sup> Yinmin Zhong<sup>2,\*</sup> Qi Huang<sup>1</sup> Yangrui Chen<sup>1</sup> Zhi Zhang<sup>1</sup>  
Yanghua Peng<sup>1</sup> Xiang Li<sup>1</sup> Cong Xie<sup>1</sup> Shibiao Nong<sup>1</sup> Yulu Jia<sup>1</sup> Sun He<sup>1</sup> Hongmin Chen<sup>1</sup>  
Zhihao Bai<sup>1</sup> Qi Hou<sup>1</sup> Shipeng Yan<sup>1</sup> Ding Zhou<sup>1</sup> Yiyao Sheng<sup>1</sup> Zhuo Jiang<sup>1</sup>  
Haohan Xu<sup>1</sup> Haoran Wei<sup>1</sup> Zhang Zhang<sup>1</sup> Pengfei Nie<sup>1</sup> Leqi Zou<sup>1</sup> Sida Zhao<sup>1</sup>  
Liang Xiang<sup>1</sup> Zherui Liu<sup>1</sup> Zhe Li<sup>1</sup> Xiaoying Jia<sup>1</sup> Jianxi Ye<sup>1</sup> Xin Jin<sup>2,†</sup> Xin Liu<sup>1,†</sup>

<sup>1</sup>*ByteDance* <sup>2</sup>*Peking University*

### Abstract

We present the design, implementation and engineering experience in building and deploying MegaScale, a production system for training large language models (LLMs) at the scale of more than 10,000 GPUs. Training LLMs at this scale brings unprecedented challenges to training efficiency and stability. We take a full-stack approach that co-designs the algorithmic and system components across model block and optimizer design, computation and communication overlapping, oper-

serving billions of users, we have been aggressively integrating AI into our products, and we are putting LLMs as a high priority to shape the future of our products.

Training LLMs is a daunting task that requires enormous computation resources. The scaling law [3] dictates that the model size and the training data size are critical factors that determine the model capability. To achieve state-of-the-art model capability, many efforts have been devoted to train large models with hundreds of billions or even trillions of parameters on hundreds of billions or even trillions of tokens. For example, GPT-3 [4] has 175 billion parameters and

Apr. 18, 2024

## Build the future of AI with Meta Llama 3

To train our largest Llama 3 models, we combined three types of parallelization: data parallelization, model parallelization, and pipeline parallelization. Our most efficient implementation achieves a compute utilization of over 400 TFLOPS per GPU when trained on 16K GPUs simultaneously. We performed training runs on two custom-built [24K GPU clusters](#). To maximize GPU uptime, we developed an advanced new training stack that automates error detection, handling, and maintenance. We also greatly improved our

[Introducing Meta Llama 3: The most capable openly available LLM to date]

# Distributed training over massive GPUs is extremely challenging

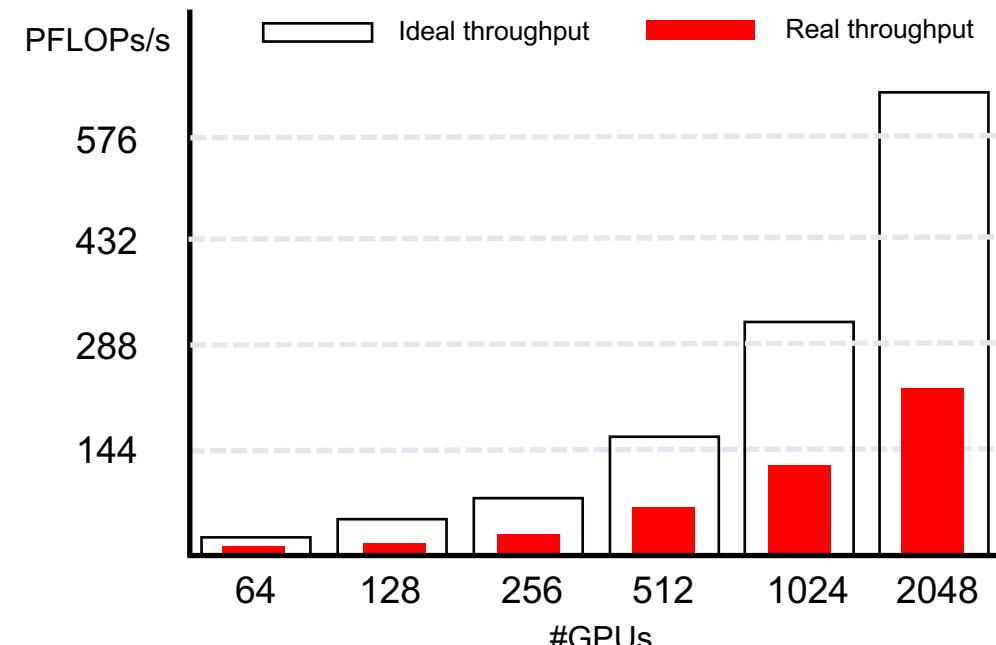
---

- Distributed training over thousands of GPUs is extremely challenging
- Two major challenges: **stability** and **scalability**
- Stability: very common that some GPU crashes during training LLMs
  - Meta OPT-175B: 175+ job restarts caused by hardware failures in 2 months
  - 书生大模型: Waste 41700 GPU hours due to training crashes (>80% are infrastructure failures)
- Stability in LLM training is very important, but this talk will not discuss it

# Distributed training over massive GPUs is extremely challenging

- The **communication overhead** and **GPU idle time** severely hamper the scalability
- Each GPU can only achieve **30%~55%** of its peak FLOPs/s during LLM training
- When GPU achieves 30% of peak FLOP/s, we say the system achieves 30% scalability

30% of its peak FLOPs/s visualization



# Existing systems suffer from severe scalability issue

[Efficient Large-Scale Language Model Training on GPU Clusters Using Megatron-LM, 2021]

Number of parameters (billion)	Attention heads	Hidden size	Number of layers	Tensor model-parallel size	Pipeline model-parallel size	Number of GPUs	Batch size	Achieved teraFLOP/s per GPU	Percentage of theoretical peak FLOP/s	Achieved aggregate petaFLOP/s
1.7	24	2304	24	1	1	32	512	137	44%	4.4
3.6	32	3072	30	2	1	64	512	138	44%	8.8
7.5	32	4096	36	4	1	128	512	142	46%	18.2
18.4	48	6144	40	8	1	256	1024	135	43%	34.6
39.1	64	8192	48	8	2	512	1536	138	44%	70.8
76.1	80	10240	60	8	4	1024	1792	140	45%	143.8
145.6	96	12288	80	8	8	1536	2304	148	47%	227.1
310.1	128	16384	96	8	16	1920	2160	155	50%	297.4
529.6	128	20480	105	8	35	2520	2520	163	52%	410.2
1008.0	160	25600	128	8	64	3072	3072	163	52%	502.0

Nvidia Megatron: the most popular distributed training framework

# Existing systems suffer from severe scalability issue

- Open AI and Meta does not disclose its training scalability, but we can infer them
- End-to-end training time can be estimated by

$$\text{Training time} = \frac{\text{Total FLOPS to train the model}}{\text{FLOPs/s of the GPU cluster}}$$
$$\approx \frac{8 T P}{N F U}$$

T: the number of tokens

P: the number of parameters

N: the number of GPUs

F: Peak FLOP/s per GPU

U: Utilization

[Efficient Large-Scale Language Model Training on GPU Clusters Using Megatron-LM, 2021]

## Existing systems suffer from severe scalability issue

---

- **Open AI GPT3-175B:** 1024 A100 GPUs, 300B tokens, training time is 35 days

$$U \approx \frac{8 \times (300 \times 10^9) \times (175 \times 10^9)}{1024 \times (312 \times 10^{12}) \times (35 \times 24 \times 3600)} = 0.44$$

The scalability is around 0.44

- **Meta LLaMA-65B:** 2048 A100 GPUs, 1.4TB tokens, training time is 21 days

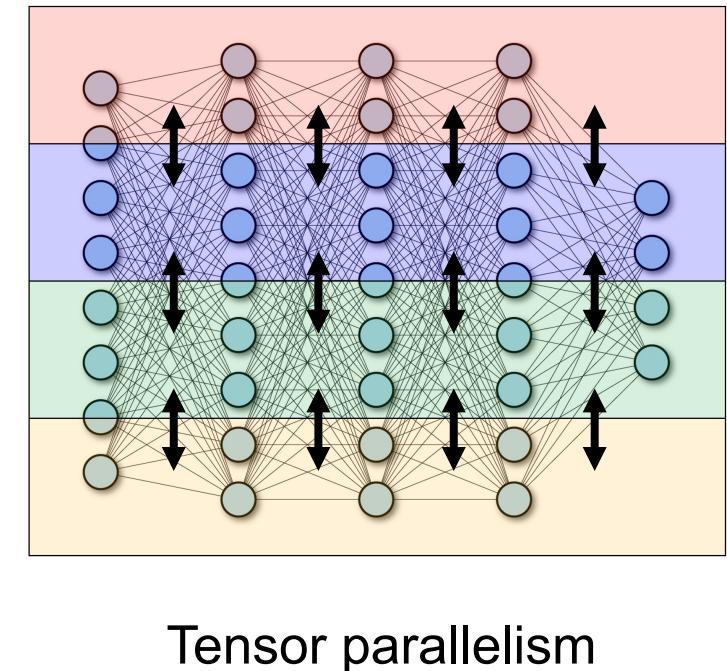
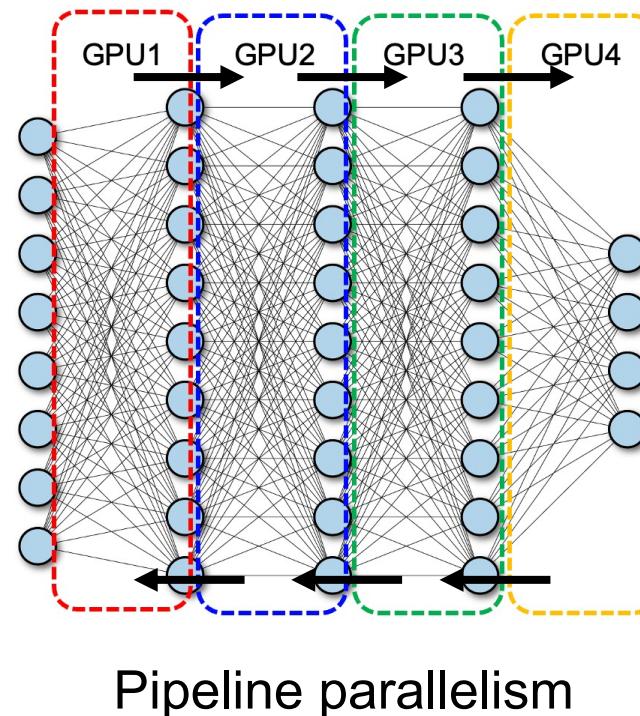
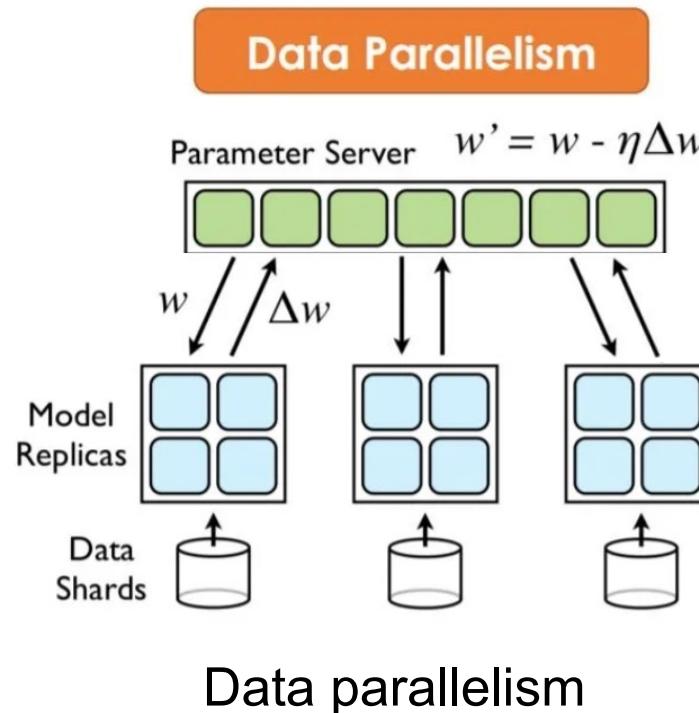
$$U \approx 0.3$$

The scalability is around 0.3

- We find even Nvidia, Open AI, and Meta **cannot** achieve strong scalability

# Communication overhead is the top factor hampering the scalability

- 3D parallelism indicates 3 orthogonal parallel techniques used to train LLM



- This talk mainly focuses on **data parallelism**.



communication

## PART 03

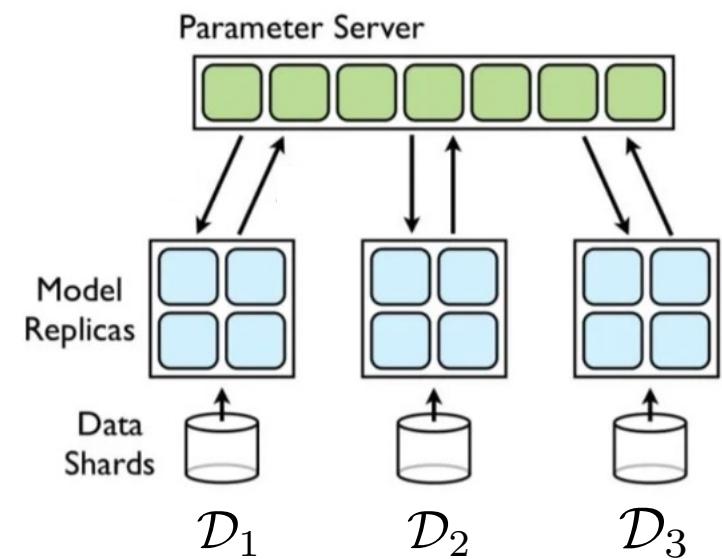
---

### Data Parallelism

A network of  $n$  nodes (devices such as GPUs) collaborate to solve the problem:

$$\min_{x \in \mathbb{R}^d} f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x), \quad \text{where } f_i(x) = \mathbb{E}_{\xi_i \sim D_i} F(x; \xi_i)$$

- Each component  $f_i : \mathbb{R}^d \rightarrow \mathbb{R}$  is local and private to node  $i$
- Random variable  $\xi_i$  denotes local data that follows distribution  $D_i$
- Each local distribution  $D_i$  is different; data heterogeneity exists



# Vanilla parallel stochastic gradient descent (PSGD)



$$\min_{x \in \mathbb{R}^d} f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x), \quad \text{where } f_i(x) = \mathbb{E}_{\xi_i \sim D_i} F(x; \xi_i).$$

**PSGD**

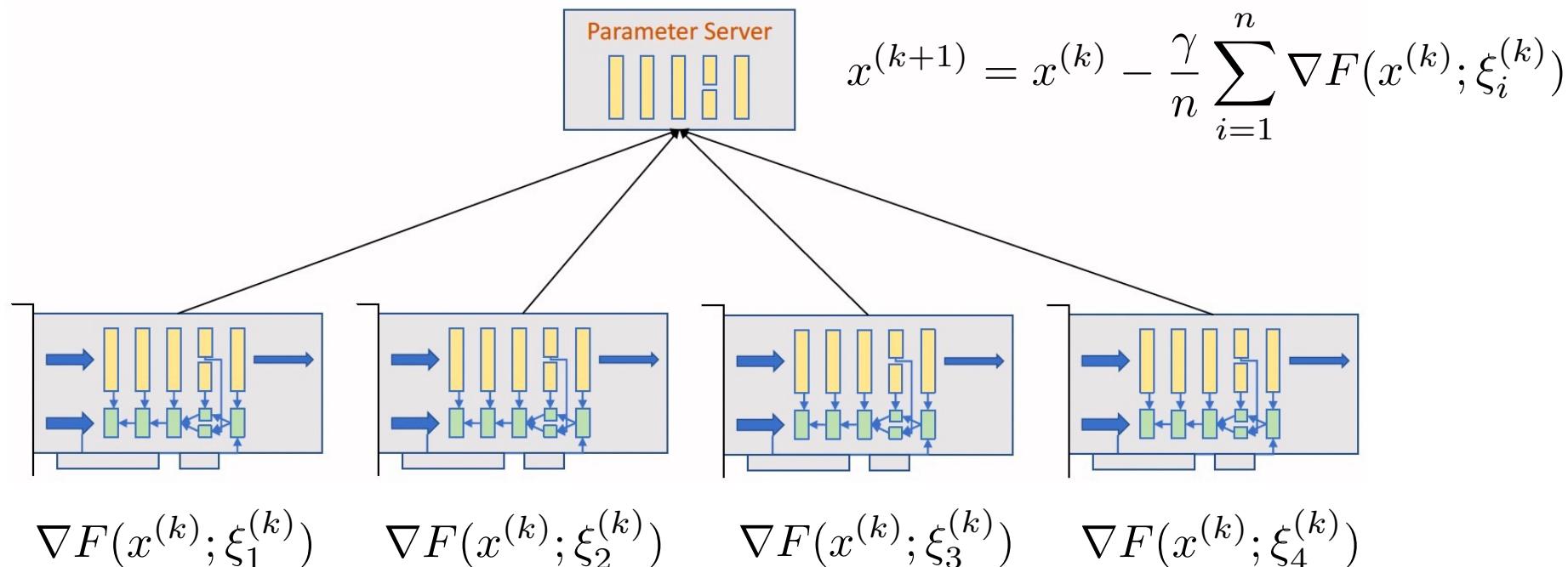
$$g_i^{(k)} = \nabla F(x^{(k)}; \xi_i^{(k)}) \quad (\text{Local compt.})$$

$$x^{(k+1)} = x^{(k)} - \frac{\gamma}{n} \sum_{i=1}^n g_i^{(k)} \quad (\text{Global comm.})$$

- Each node  $i$  samples data  $\xi_i^{(k)}$  and computes gradient  $\nabla F(x^{(k)}; \xi_i^{(k)})$
- All nodes synchronize (i.e. globally average) to update model  $x$  per iteration

# Vanilla parallel stochastic gradient descent (PSGD)

- Federated learning typically implements PSGD using parameter server



- LLM training within data-centers implements PSGD using Ring-Allreduce

## Assumption [PSGD assumption]

- (1) Each  $f_i(x)$  is  $L$ -smooth, i.e.,  $\|\nabla f_i(x) - \nabla f_i(y)\| \leq L\|x - y\|$  for any  $x, y$ ;
- (2) Each stochastic gradient is unbiased and has bounded variance, i.e.,

$$\mathbb{E}[\nabla F(x; \xi_i)] = \nabla f_i(x), \quad \mathbb{E}\|\nabla F(x; \xi_i) - \nabla f_i(x)\|^2 \leq \sigma^2$$

## Theorem [PSGD convergence]

Under the above assumptions and with proper  $\gamma$ , PSGD converges as follows

$$\frac{1}{T} \sum_{k=1}^T \mathbb{E}\|\nabla f(x^{(k)})\|^2 = \mathcal{O}\left(\frac{\sigma}{\sqrt{nT}}\right)$$

[K. Yuan, Lecture 6: stochastic gradient descent, PKU Class “Optimization for deep learning”, check Kun Yuan’s website]

## Theorem [PSGD convergence]

Under the above assumptions and with proper  $\gamma$ , PSGD converges as follows

$$\frac{1}{T} \sum_{k=1}^T \mathbb{E} \|\nabla f(x^{(k)})\|^2 = \mathcal{O}\left(\frac{\sigma}{\sqrt{nT}}\right)$$

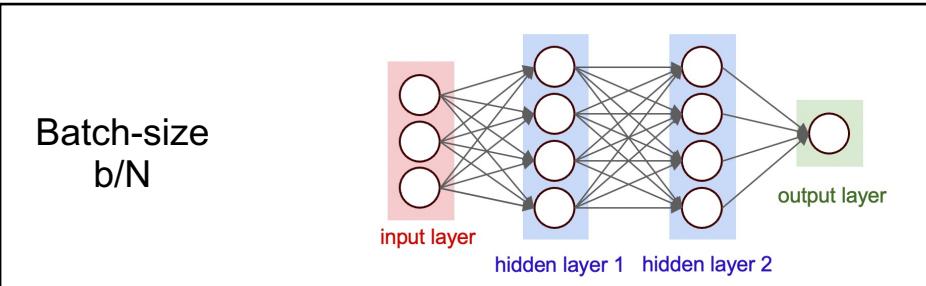
- This implies that to achieve an  $\epsilon$ -accurate solution, PSGD needs

$$\frac{\sigma}{\sqrt{nT}} \leq \epsilon \implies T \geq \frac{\sigma^2}{n\epsilon^2}$$

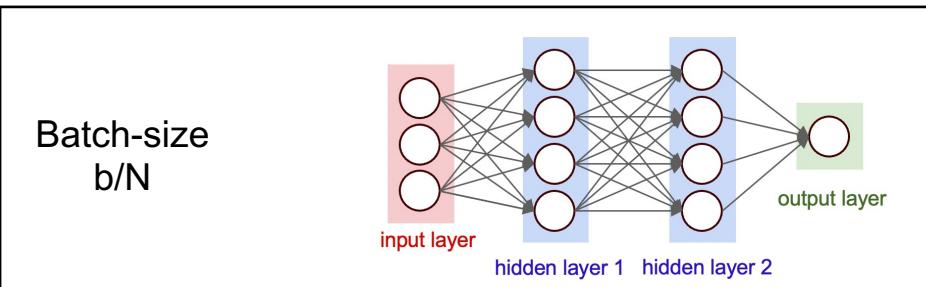
iterations, which decreases linearly with number of nodes  $n$ ; this is called **linear speedup**

# Data parallelism: memory and communication

GPU 1

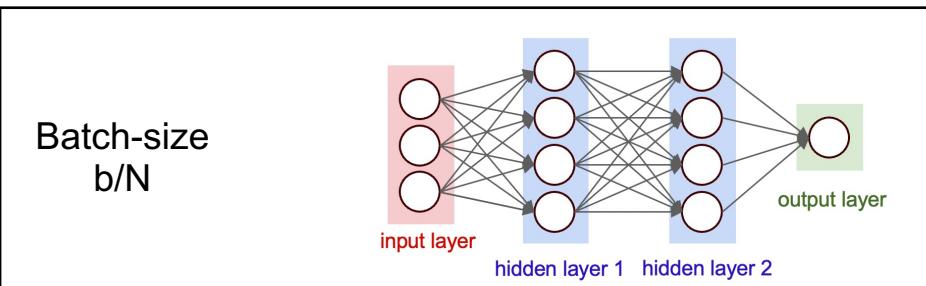


GPU 2



•      •      •  
 •      •      •  
 •      •      •

GPU N



## Pros:

- Cuts the computation time by  $N$  folds
- Reduce the activation memory by  $N$  folds

$$2[4\Phi + \frac{b}{N}\Theta)$$

**FP16 M+g+os = 4 # parameters # activations**

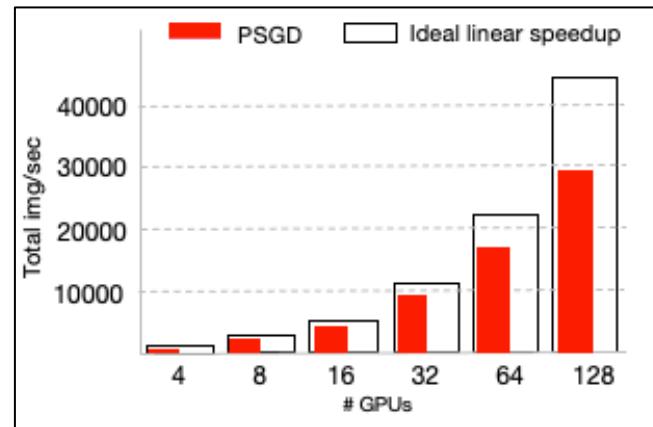
- Memory-efficient when  $\Theta \gg \Phi$

## Cons:

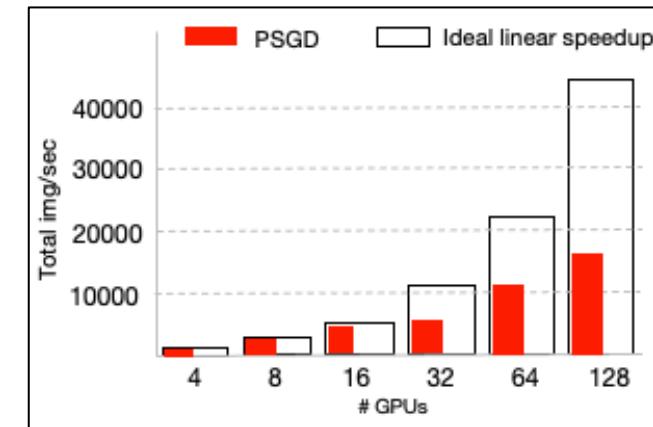
- Cannot save memory significantly if  $\Phi \gg \Theta$
- Incurs  $2\Phi$  Bytes to communicate (extra time)

# PSGD cannot achieve ideal scalability due to comm. overhead

- PSGD **cannot** achieve ideal linear speedup in throughput due to comm. overhead
- Larger comm-to-compt ratio leads to worse performance in PSGD



Small comm.-to-compt. ratio



Large comm.-to-compt. ratio

Ring-Allreduce is used in each experiment

- How can we reduce the communication overhead in PSGD?

- Global average incurs  $O(n)$  comm. overhead; proportional to network size n

## [Decentralized communication]

- Each node interacts with the server at every iteration; proportional to iteration numbers

## [Lazy communication]

- Each node sends a full model (or gradient) to the server; proportional to dimension d

## [Compressed communication]

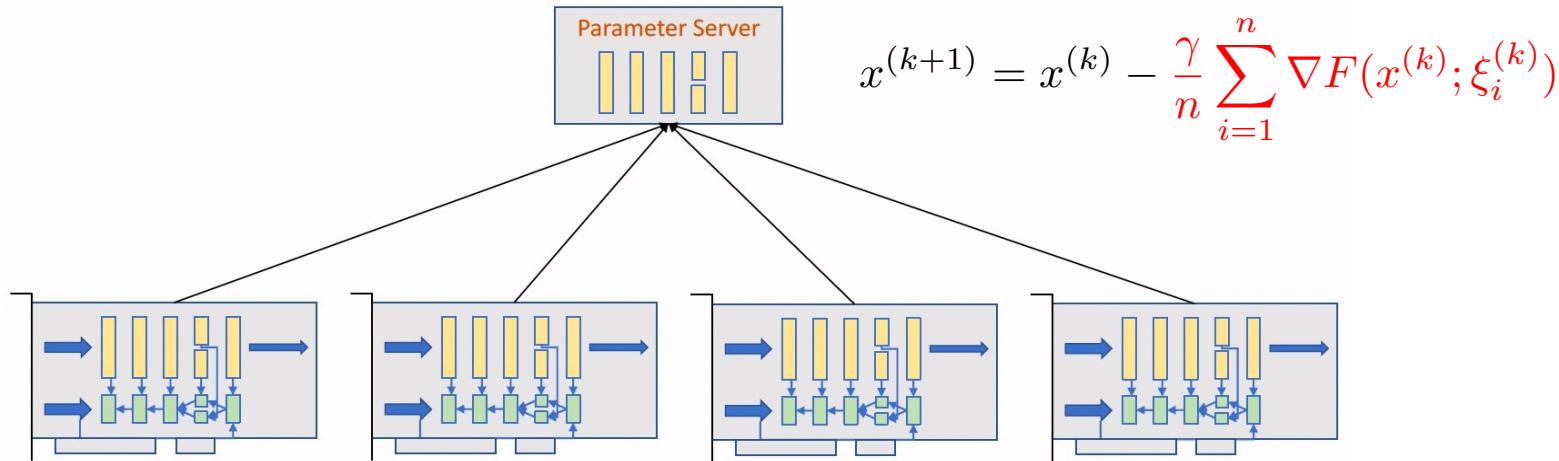
- and more (asynchronous communication; robust communication against Byzantine nodes, etc.)

## PART 04

---

# Communication-Efficient Data Parallelism

## Parameter-server architecture



- Global average using PS incurs  $O(n)$  comm. overhead; **proportional to network size n**
- Global average using Ring-Allreduce incurs  $O(n)$  latency; **proportional to network size n**
- When network size n is large (say  $n \geq 10000$ ), PSGD suffers severe communication overhead

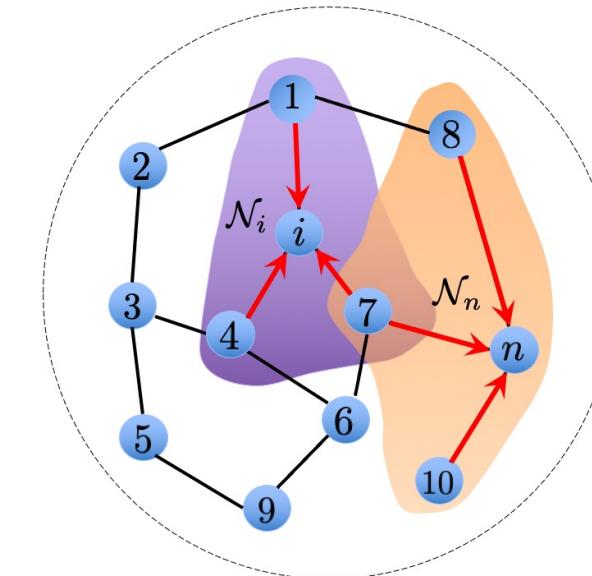
# Decentralized SGD (DSGD): reduce central bottleneck

- To break  $O(n)$  comm. overhead, we replace global average with partial average

$$x_i^{(k+\frac{1}{2})} = x_i^{(k)} - \gamma \nabla F(x_i^{(k)}; \xi_i^{(k)}) \quad (\text{Local update})$$

$$x_i^{(k+1)} = \sum_{j \in \mathcal{N}_i} w_{ij} x_j^{(k+\frac{1}{2})} \quad (\text{Partial averaging})$$

- DSGD = local SGD update + partial averaging [1, 2]
- $\mathcal{N}_i$  is the set of neighbors at node  $i$ ;  $w_{ij}$  scales information from  $j$  to  $i$
- Incurs  $O(d_{\max})$  comm. overhead per iteration where  $d_{\max} = \max_i \{|\mathcal{N}_i|\}$  is the graph maximum degree

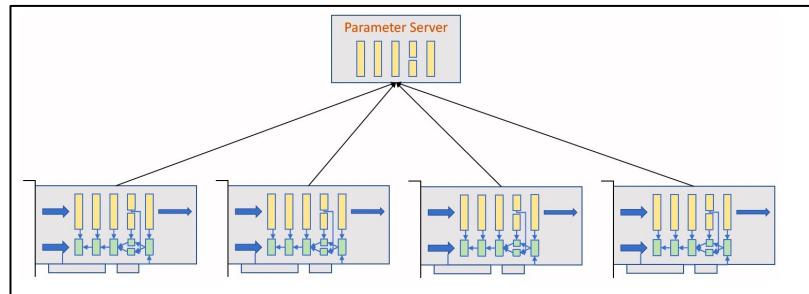


[1] C. G. Lopes and A. H. Sayed, "Diffusion least-mean-squares over adaptive networks", ICASSP, 2007

[2] A. Nedich and A. Ozdaglar, " Distributed subgradient methods for multi-agent optimization", IEEE TAC, 2009

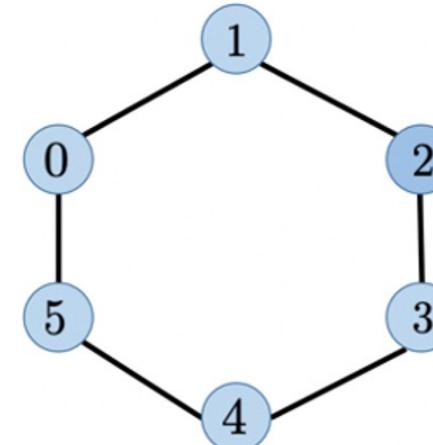
# DSGD is more communication-efficient than PSGD

- Incurs  $O(1)$  comm. overhead on **sparse** topologies; much less than global average  $O(n)$

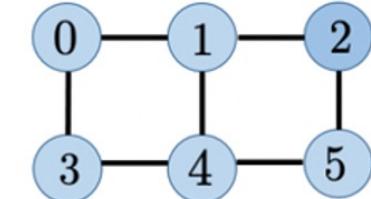


Global averaging over centralized network

comm. overhead  $O(n)$



Partial averaging over ring or grid



comm. overhead  $O(1)$

# DSGD is more communication-efficient than PSGD

---

- A real experiment on a 256-GPUs cluster [1]

Model	Ring-Allreduce	Partial average
ResNet-50 (25.5M)	278 ms	150 ms

Table. Comparison of per-iter comm. time in terms of runtime with 256 GPUs

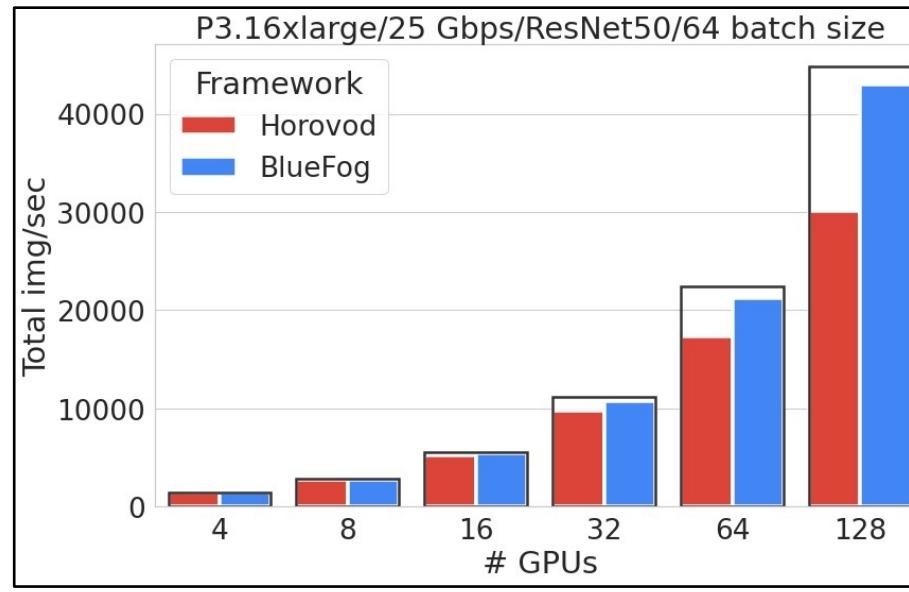
- DSGD saves more communications per iteration for larger models

---

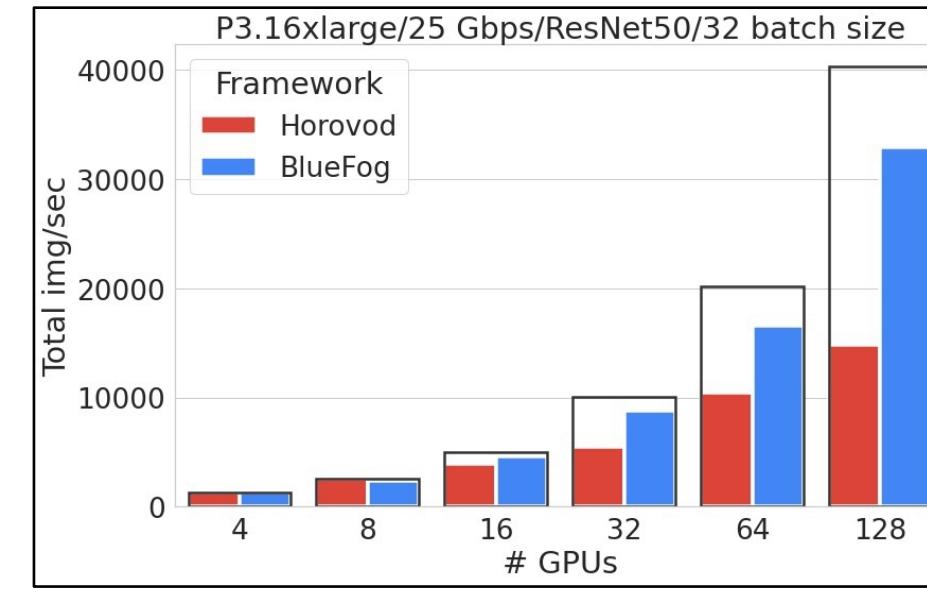
[1] Y. Chen, K. Yuan, Y. Zhang, P. Pan, Y. Xu, and W. Yin, "Accelerating Gossip SGD with Periodic Global Averaging", ICML 2021

# DSGD is more communication-efficient than PSGD

- DSGD (BlueFog) has **better scalability** than PSGD (Horovod) due to its small comm. overhead



Small comm.-to-compt. ratio



Large comm.-to-compt. ratio

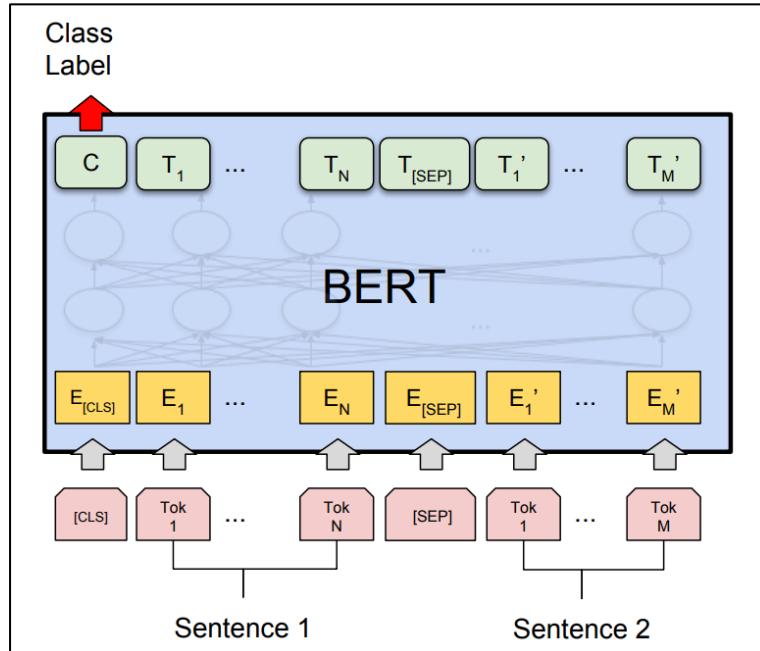
# DSGD saves more wall-clock time without severely hurting performance



nodes topology	4(4x8 GPUs) acc.	4(4x8 GPUs) time	8(8x8 GPUs) acc.	8(8x8 GPUs) time	16(16x8 GPUs) acc.	16(16x8 GPUs) time	32(32x8 GPUs) acc.	32(32x8 GPUs) time
P-SGD	76.32	11.6	76.47	6.3	76.46	3.7	76.25	2.2
Decentralized	<b>76.34</b>	<b>11.1</b>	<b>76.52</b>	<b>5.7</b>	<b>76.47</b>	<b>2.8</b>	<b>76.27</b>	<b>1.5</b>

DSGD shows very impressive linear speedup performance and saves more time than PSGD!

# DSGD saves more wall-clock time without severely hurting performance



Model: BERT-Large (330M parameters)

Dataset: Wikipedia (2500M words) and  
BookCorpus (800M words)

Hardware: 64 GPUs

Table. Comparison in loss and training time [1]

Method	Final Loss	Wall-clock Time (hrs)
P-SGD	1.75	59.02
D-SGD	1.77	30.4

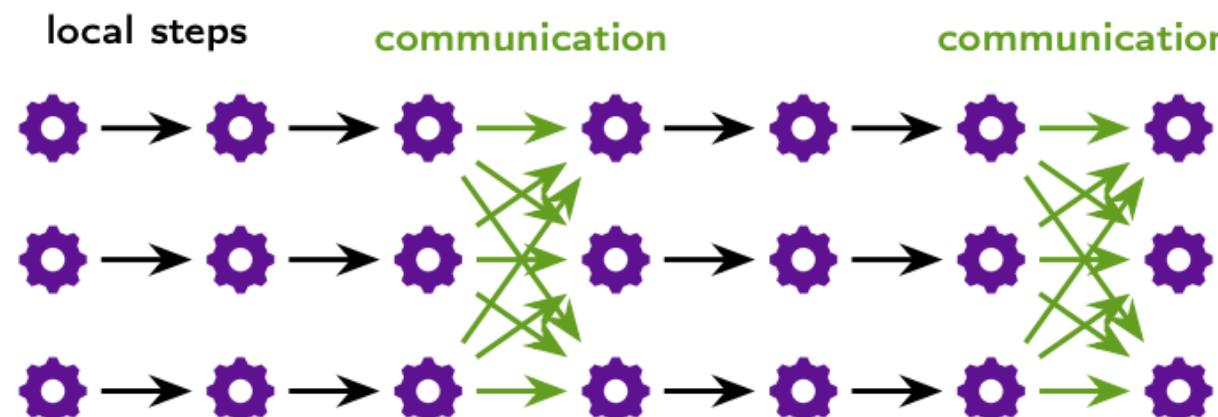
[1] Y. Chen, K. Yuan, Y. Zhang, P. Pan, Y. Xu, and W. Yin, "Accelerating Gossip SGD with Periodic Global Averaging", ICML 2021

# Local SGD: reduce communication frequency

- Local SGD saves communication by **reducing communication frequency**

$$x_i^{k+1} = \begin{cases} x_i^k - \gamma \nabla F(x_i^k; \xi_i^k) & \text{if } \text{mod}(k, \tau) \neq 0 \\ \frac{1}{n} \sum_{i=1}^n x_i^k & \text{if } \text{mod}(k, \tau) = 0 \end{cases}$$

- Communicate once every  $\tau$  iterations



# Local SGD is also known as FedAvg

- Local SGD can be rewritten in a different format

For  $r = 0, 1, \dots, R - 1$  **(R comm. rounds)**

$$x_i^{(r,0)} = x^{(r)}$$

For  $k = 0, 1, \dots, K - 1$  **(local update)**

$$x_i^{(r,k+1)} = x_i^{(r,k+1)} - \gamma \nabla F(x_i^{(r,k)}, \xi_i^{(r,k)})$$

$$x^{(r+1)} = \frac{1}{n} \sum_{i=1}^n x_i^{(r,K)} \quad \textbf{(global comm.)}$$

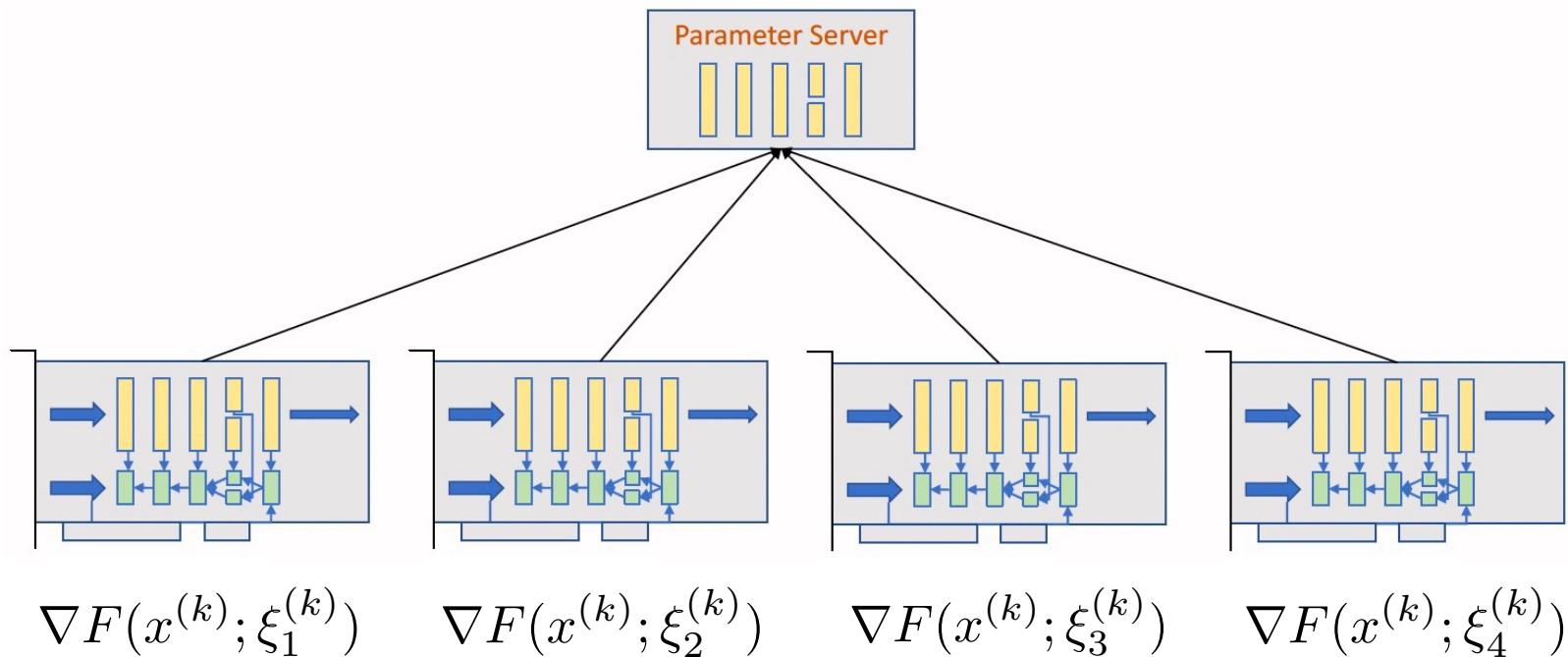
- Local SGD is also known as **FedAvg** [1,2]

[1] J. Konečný, H. B. McMahan, D. Ramage, P. Richtárik, “Federated Optimization: Distributed Machine Learning for On-Device Intelligence”, 2016

[2] B. McMahan, E. Moore, D. Ramage, S. Hampson, B. A. y Arcas, “Communication-Efficient Learning of Deep Networks from Decentralized Data”, AISTATS 2017

# Review of the vanilla parallel SGD

$$x^{(k+1)} = x^{(k)} - \frac{\gamma}{n} \sum_{i=1}^n \nabla F(x^{(k)}; \xi_i^{(k)})$$

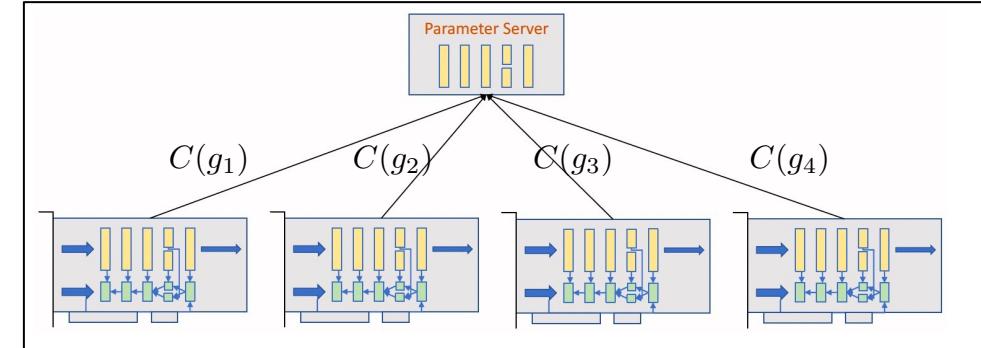


# Compressed SGD: reduce communication volume

- A basic (but not state-of-the-art) algorithm is QSGD [Alistarh et. al., 2017]

$$g_i^{(k)} = \nabla F(x_i^{(k)}; \xi_i^{(k)})$$

$$x^{(k+1)} = x^{(k)} - \frac{\gamma}{n} \sum_{j=1}^n C(g_j^{(k)})$$



- $C(\cdot)$  is a compressor. It can quantize or sparsify the full gradient



**Quantization**



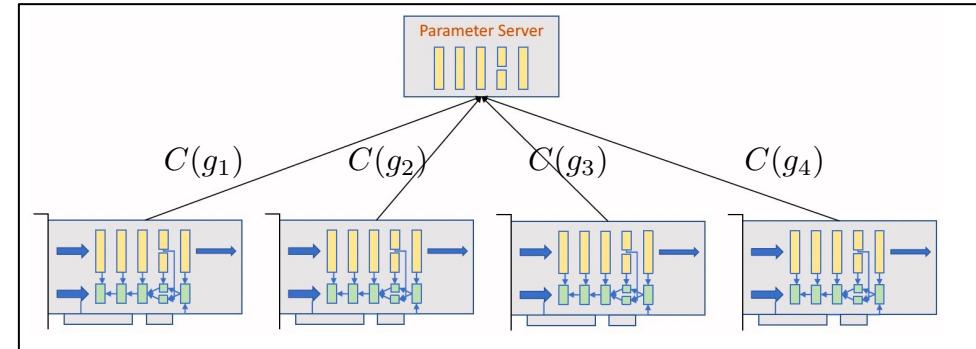
**1 bit**

# Compressed SGD

- A basic (but not state-of-the-art) algorithm is QSGD [Alistarh et. al., 2017]

$$g_i^{(k)} = \nabla F(x_i^{(k)}; \xi_i^{(k)})$$

$$x_i^{(k+1)} = x_i^{(k)} - \frac{\gamma}{n} \sum_{j=1}^n C(g_j^{(k)})$$



- $C(\cdot)$  is a compressor. It can quantize or sparsify the full gradient



**Sparsification**



# Unbiased compressor

---

- Compressor class. Most compressors in literature are either **unbiased** or **contractive**
- We let  $\mathcal{U}_\omega$  denote the set of unbiased compressors satisfying Assumption 3

**Assumption (Unbiased compressor)** The compression operator  $C : \mathbb{R}^d \rightarrow \mathbb{R}^d$  satisfies

$$\mathbb{E}[C(x)] = x, \quad \mathbb{E}[\|C(x) - x\|^2] \leq \omega \|x\|^2, \quad \forall x \in \mathbb{R}^d$$

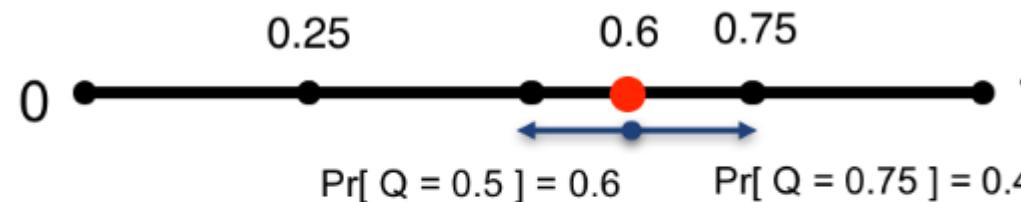
for constant  $\omega \geq 0$ , where the expectation is taken over the randomness of the compression operator  $C$ .

- Identity operator  $I$  (i.e. no compression) is an unbiased compressor with  $\omega = 0$ .

# Example: random quantization

---

- Random quantization with 5 levels:



$$\mathbb{E}[Q] = 0.6 \times 0.5 + 0.4 \times 0.75 = 0.6$$

Unbiased: 
$$\mathbb{E}[Q(x)] = \frac{Q_+(x) - x}{Q_+(x) - Q_-(x)} \cdot Q_-(x) + \frac{x - Q_-(x)}{Q_+(x) - Q_-(x)} \cdot Q_+(x) = x$$

# Contractive compressor



- We let  $\mathcal{C}_\delta$  denote the set of unbiased compressors satisfying Assumption 4

**Assumption (Contractive compressor)** The compression operator  $C : \mathbb{R}^d \rightarrow \mathbb{R}^d$  satisfies

$$\mathbb{E}[\|C(x) - x\|^2] \leq (1 - \delta)\|x\|^2, \quad \forall x \in \mathbb{R}^d$$

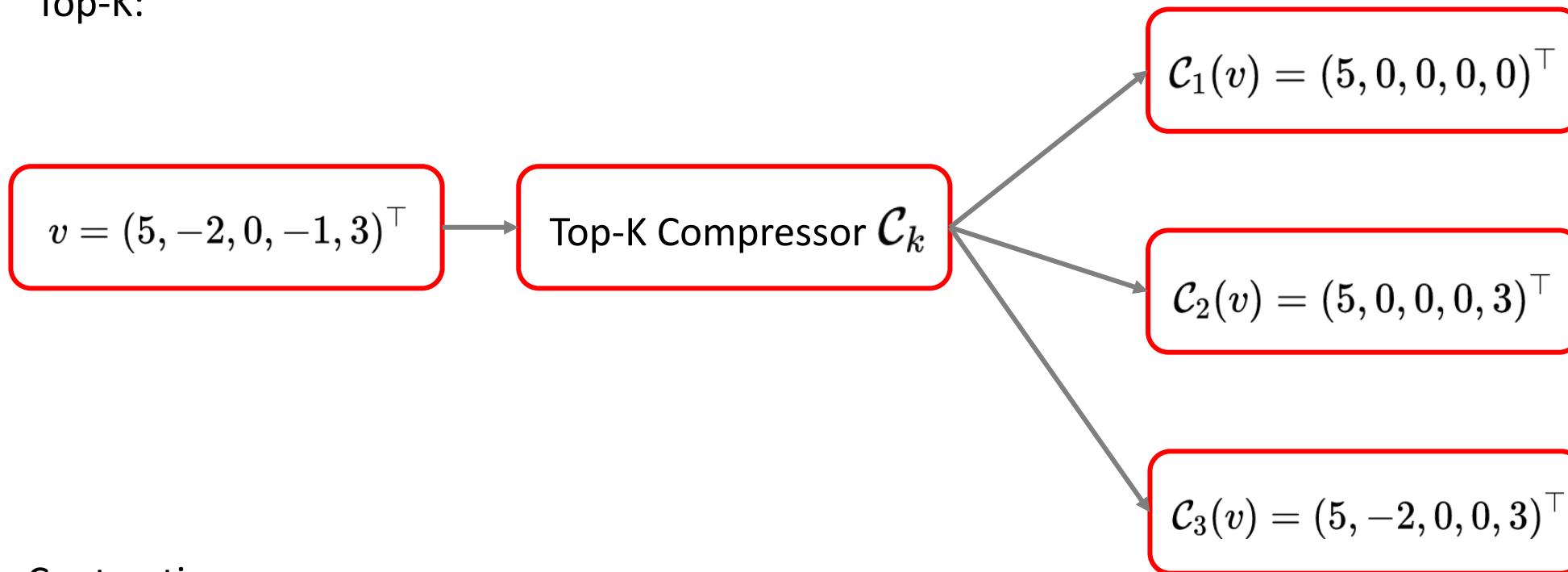
for constant  $\delta \in (0, 1]$ , where the expectation is taken over the randomness of the compression operator  $C$ .

- Identity operator  $I$  (i.e. no compression) is a contractive compressor with  $\delta = 1$ .

## Example: Top-K compressor

---

- Top-K:

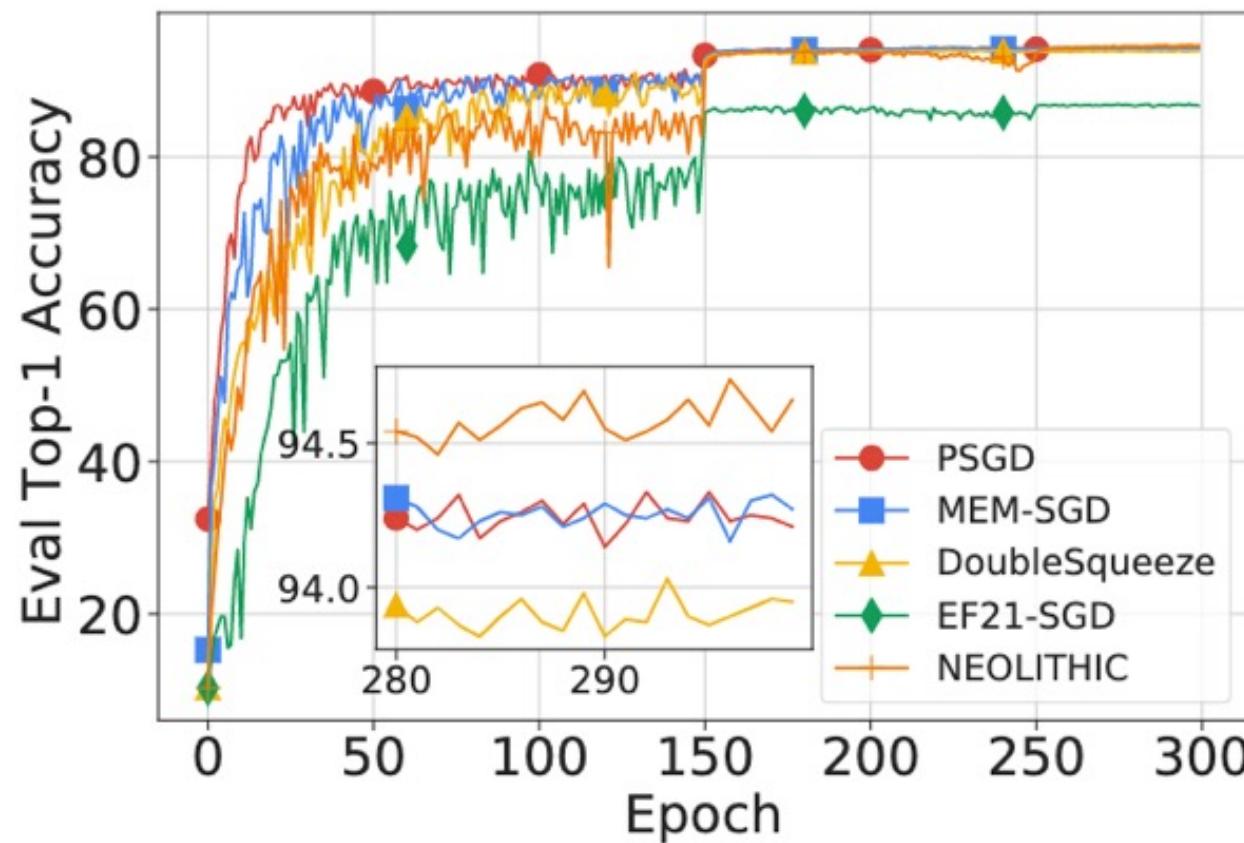


Contractive:

$$\|\mathcal{C}_k(v) - v\|^2 = \|v\|^2 - \|\mathcal{C}_k(v)\|^2 \leq (1 - k/d)\|v\|^2, \text{ i.e., } \delta = k/d \in (0, 1]$$

# Deep learning experiments

- 8 workers, **1% compression ratio** (top-k compressors), minibatch=128, R=2, ResNet18/ResNet20



# Deep learning experiments



Cifar10; 8 workers; **5% compression ratio**

METHODS	RESNET18	RESNET20
PSGD	$93.99 \pm 0.52$	$91.62 \pm 0.13$
MEM-SGD	$94.35 \pm 0.01$	$91.27 \pm 0.08$
DOUBLE-SQUEEZE	$94.11 \pm 0.14$	$90.73 \pm 0.02$
EF21-SGD	$87.37 \pm 0.49$	$65.82 \pm 4.86$
NEOLITHIC	<b><math>94.63 \pm 0.09</math></b>	<b><math>91.43 \pm 0.10</math></b>

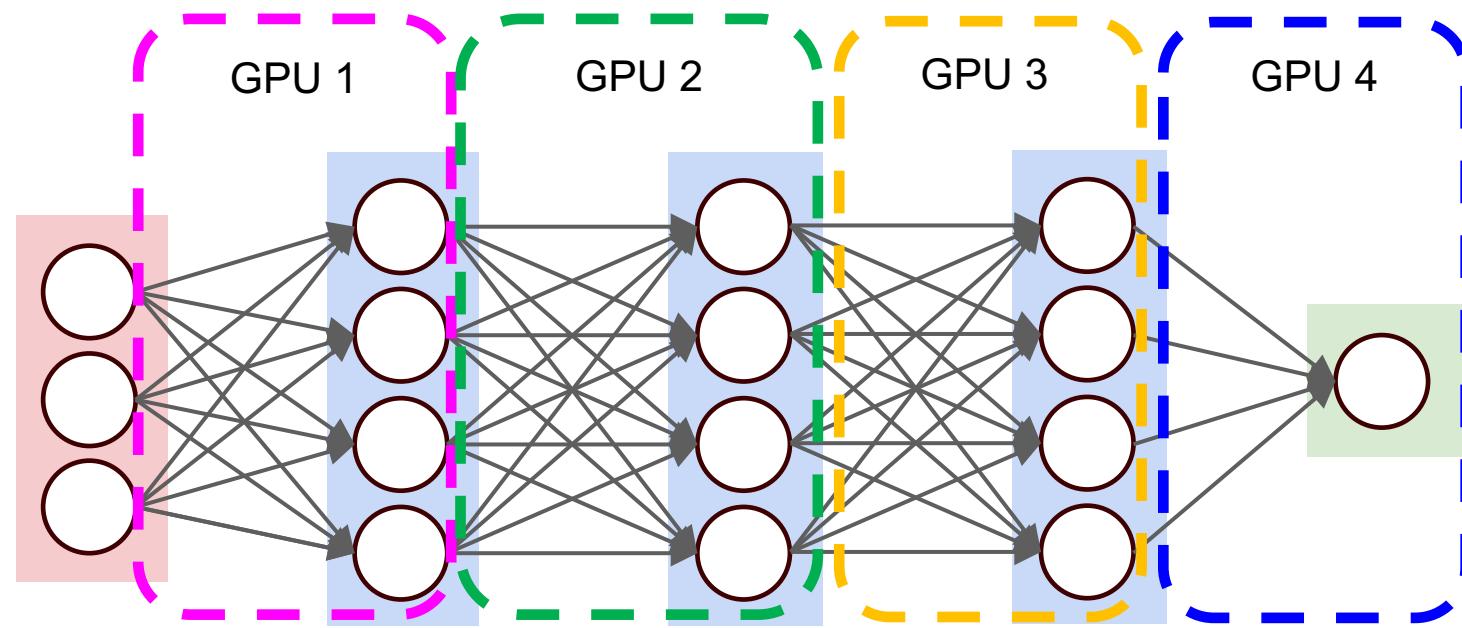
## PART 05

---

### Pipeline Parallelism

# Pipeline parallelism: architecture

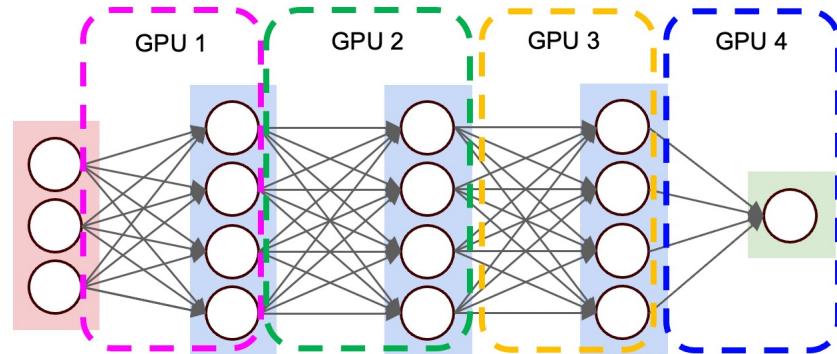
- Pipeline parallelism partitions layers and assigns them to different GPUs



- If models are partitioned uniformly, the parameters (memories) are cut down by  $N$  folds

# Model parallelism: no pipeline

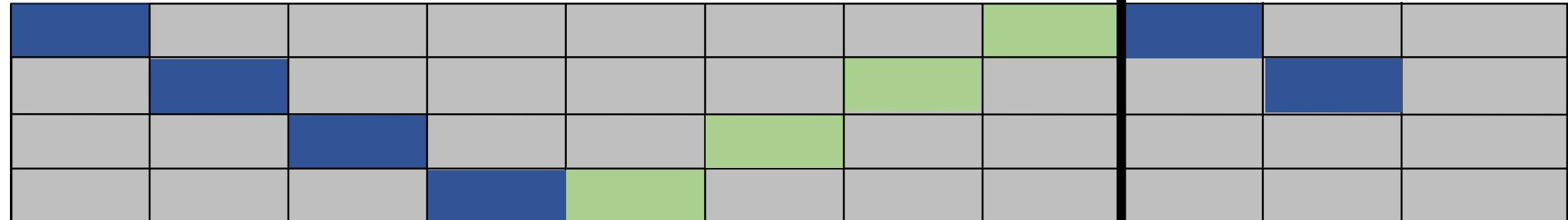
- Naïve model parallelism introduces substantial idle time. GPUs are not parallel.



Batch 1 starts

GPU 1  
GPU 2  
GPU 3  
GPU 4

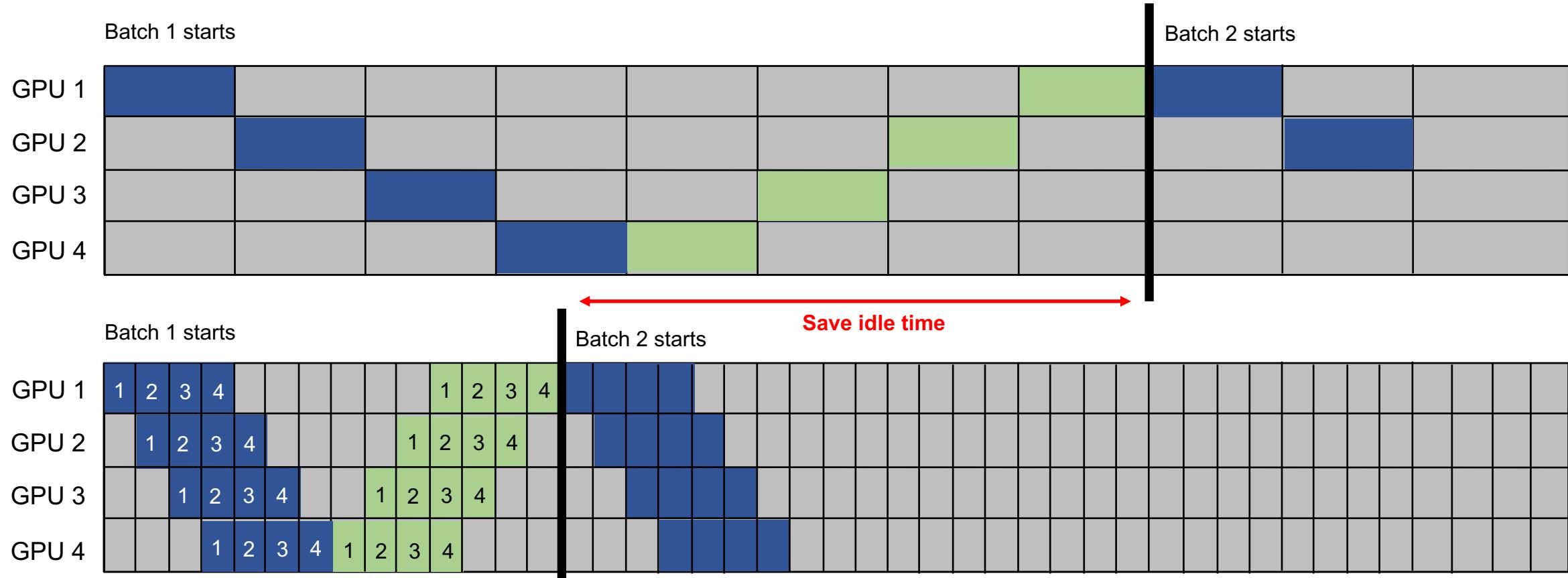
Batch 2 starts



(此处仅是示意图。实际中反向传播一般是前向时间的2倍)

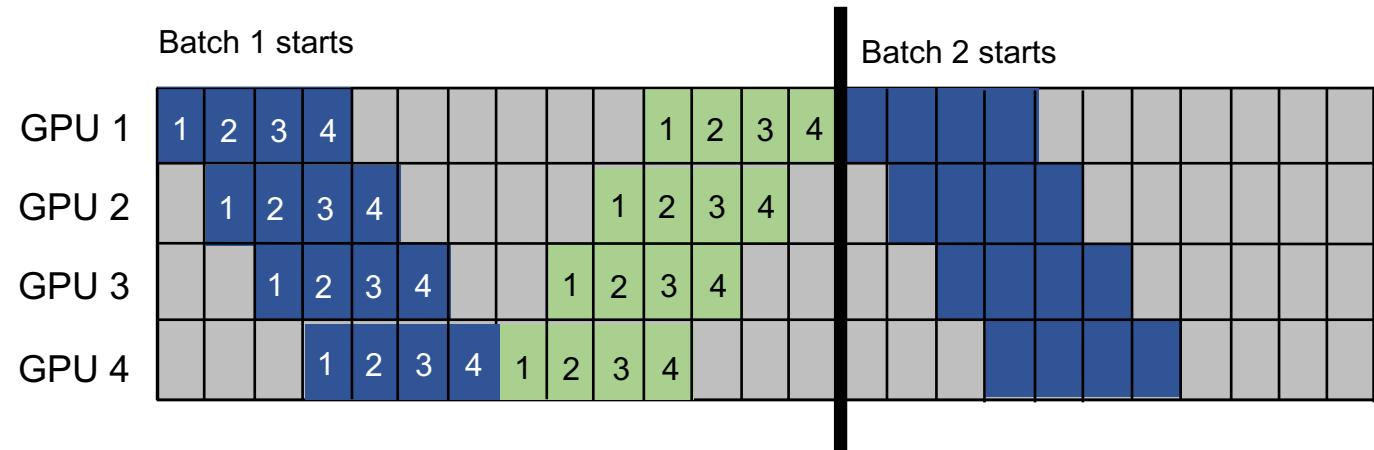
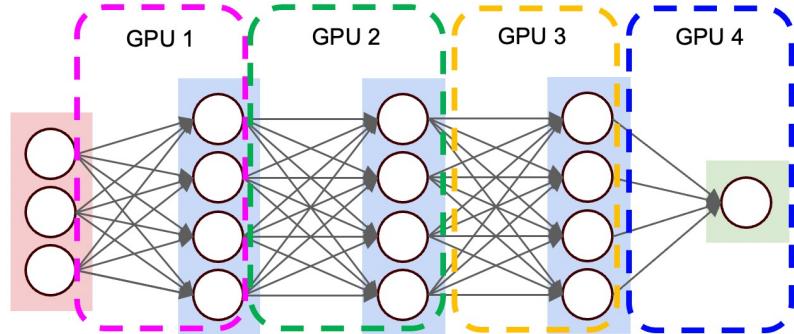
# Pipeline parallelism: pipeline reduces idle time

- Split one batch into a set of **micro-batches**



[GPipe: Efficient Training of Giant Neural Networks using Pipeline Parallelism]

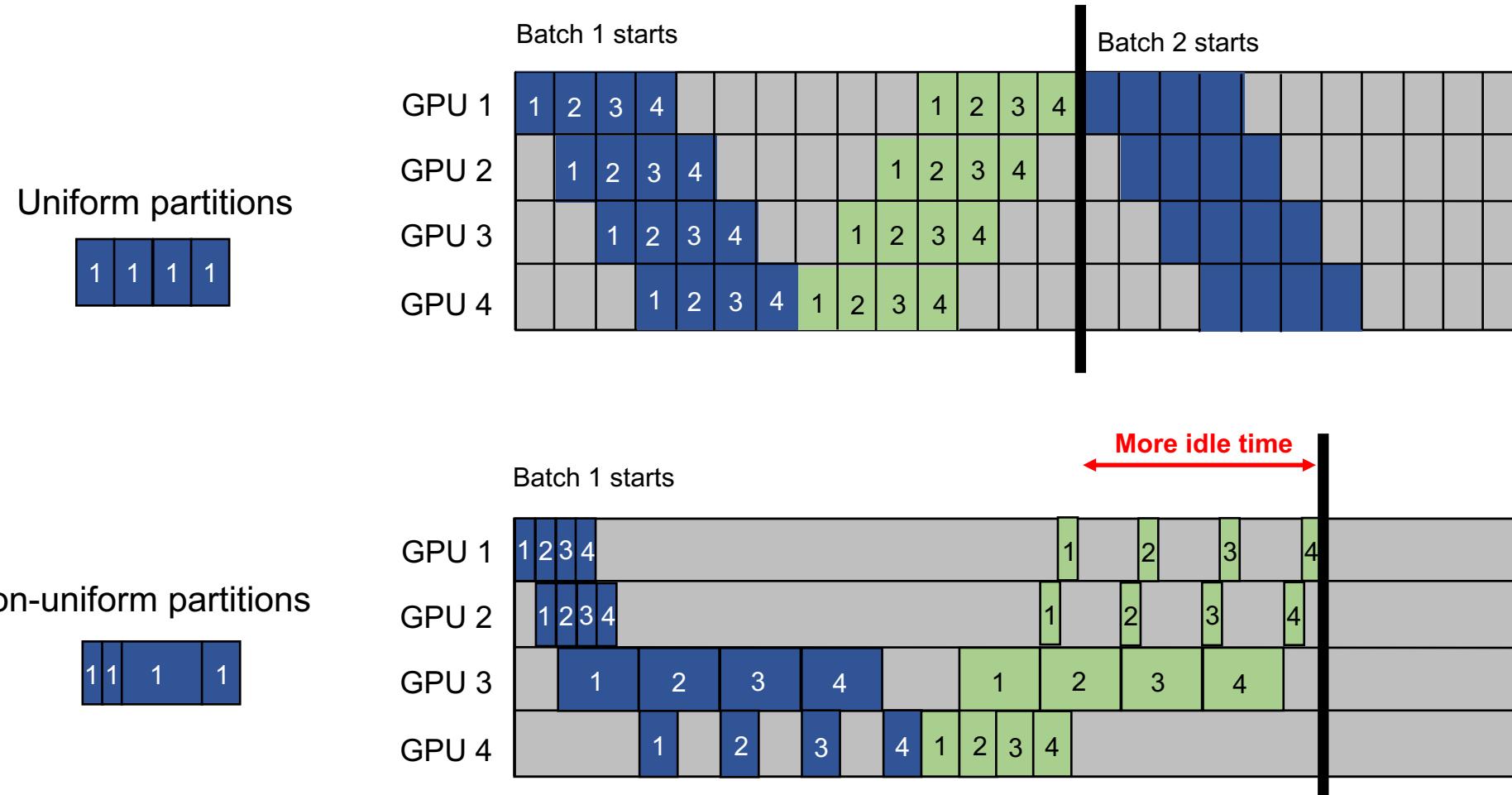
# Pipeline parallelism: idle time



- While idle time is reduced, it cannot be fully eliminated
- If each partition is of similar size, the idle time fraction is  $\frac{N - 1}{m}$     N: # GPUs    m: # micro-batches
- More micro-batches, less idle time (in line with the intuition)

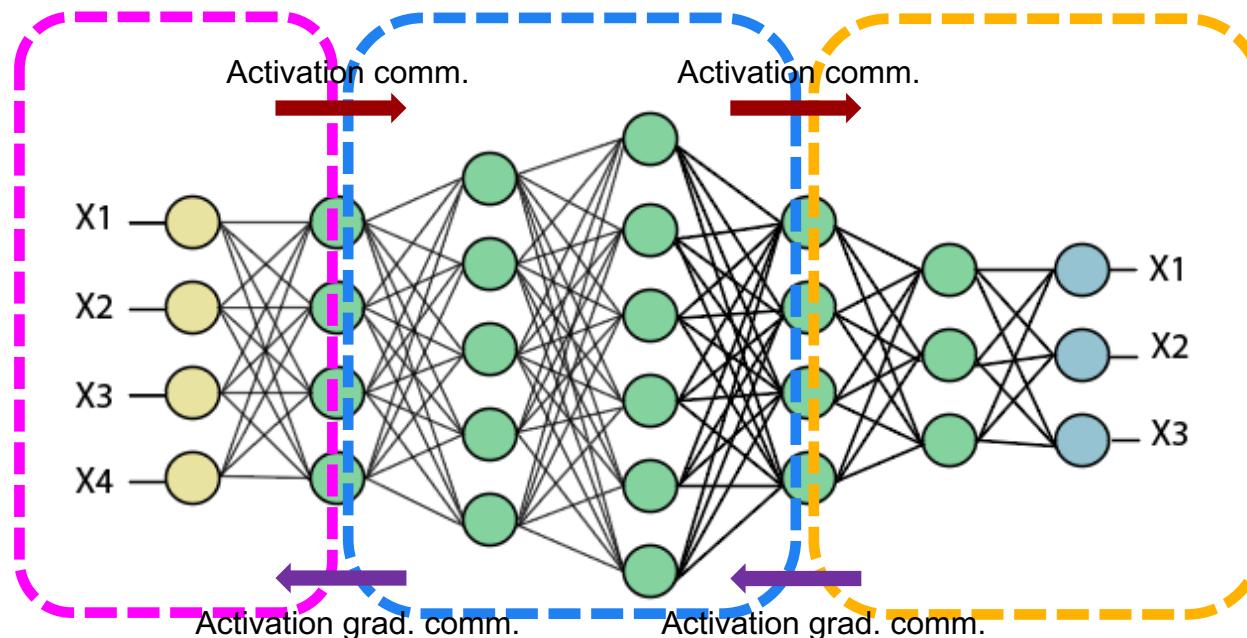
# Pipeline parallelism: idle time

- If each partition has very different parameter size, the idle time fraction will be enlarged



# Pipeline parallelism: communication time

- Pipeline parallelism incurs additional communication of activation and activation gradients
- Not all activations are to be communicated, but those at the partition points



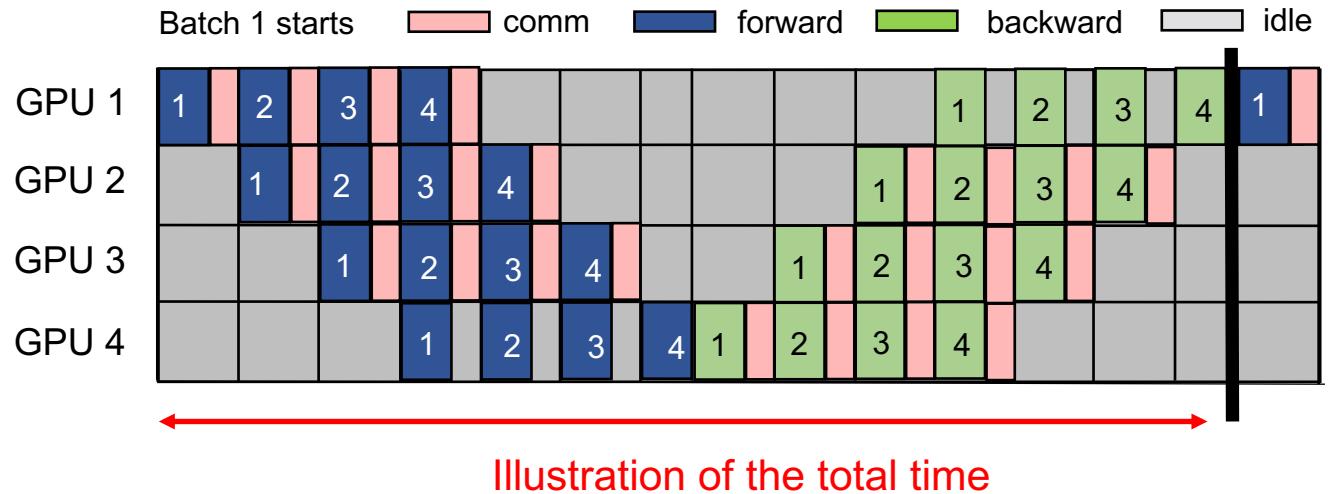
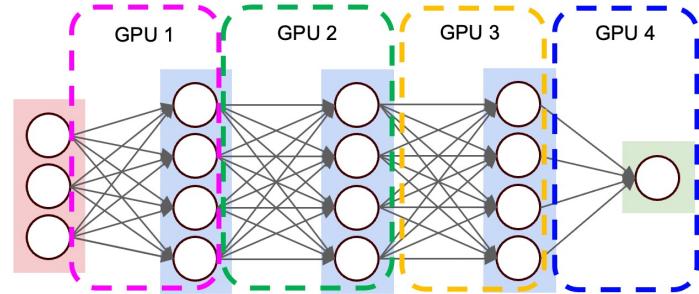
Micro-batch's communicated activation cost

$\leq$  Pipeline comm. cost  $\leq$

Full-batch's communicated activation cost

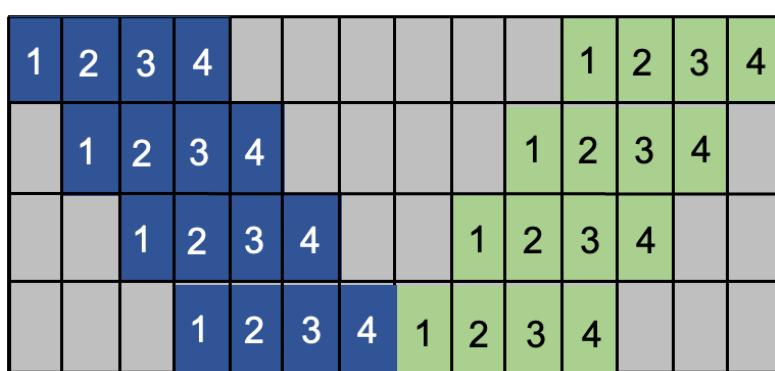
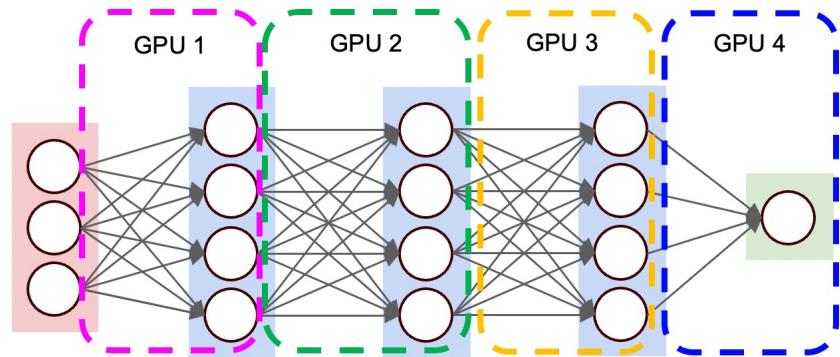
(can be modeled precisely)

# Pipeline parallelism: total time



- The above figure illustrates the very ideal scenario. While practical scenarios are more complicated (non-uniform partition), they share the same gist.
- Since each frame of comm., forward, backward, and idle time can be calculated, the total time can also be calculated

# Pipeline parallelism: summary



## Pros:

- Cuts the memory by N folds (if partition uniformly)
- Accelerates training due to micro-batch parallelism
- Less comm. overhead due to **micro-batch** (not full-batch) **activation** (not model) communication at **partitioned** layers (not all layers)

## Cons:

- Introduce idle time which cannot be fully removed
- Difficult to partition layers uniformly

Performance can be fully computed



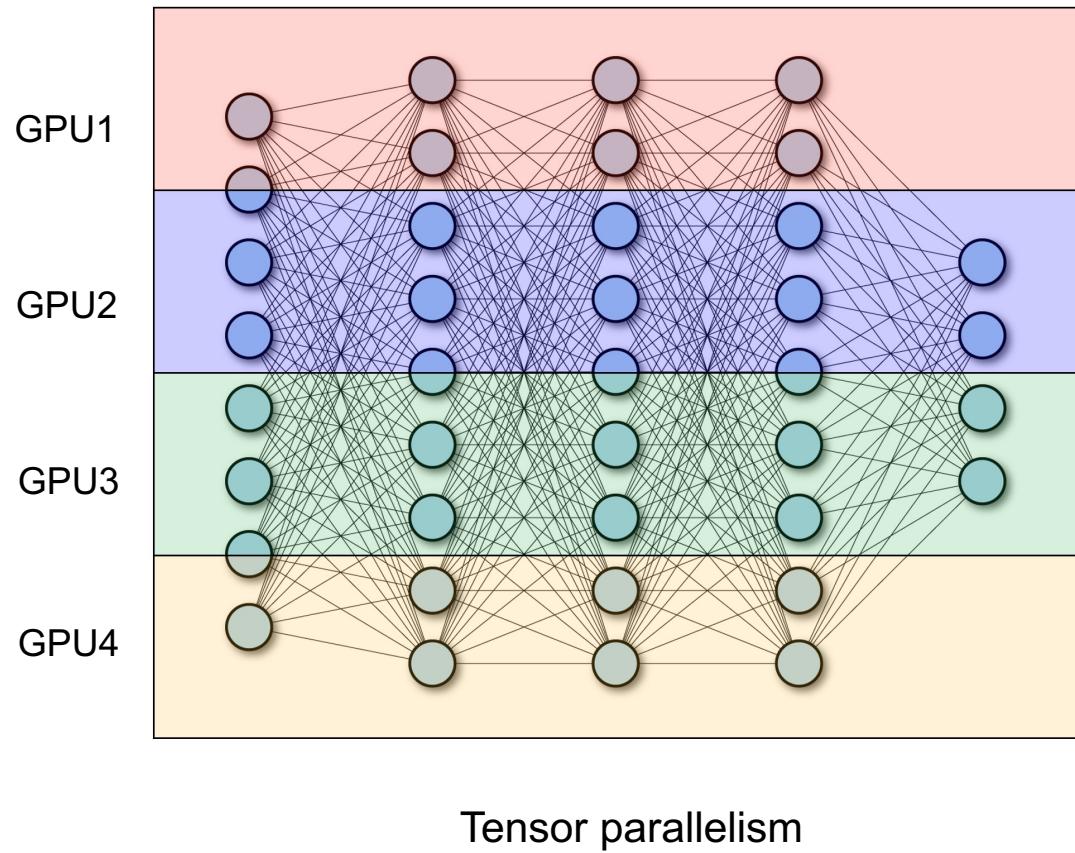
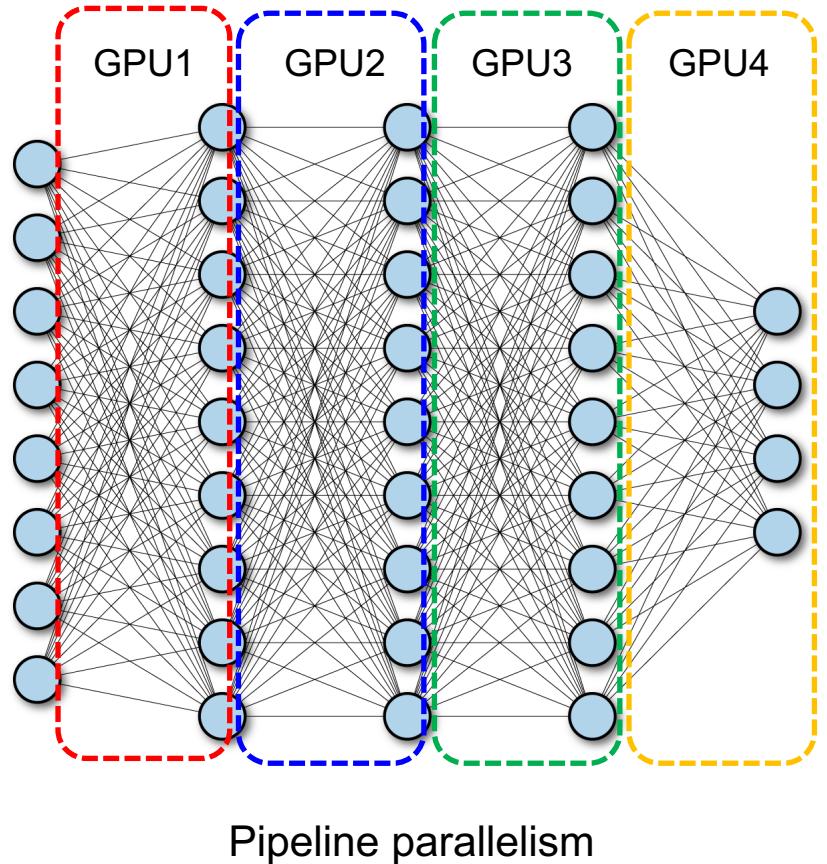
## PART 06

---

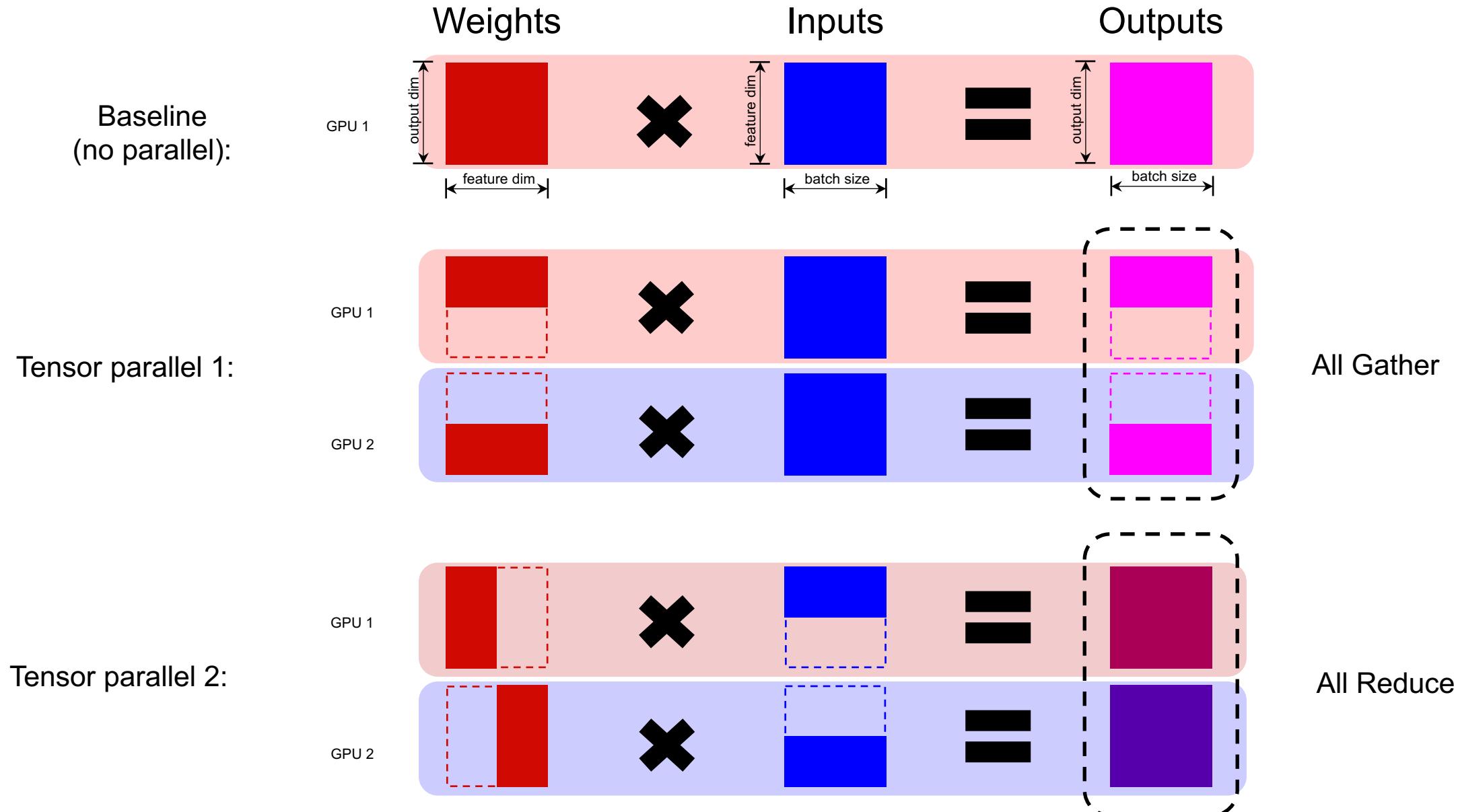
### Tensor Parallelism

# Tensor parallelism

- Tensor parallelism partitions tensors and assigns them to different GPUs (save memory)

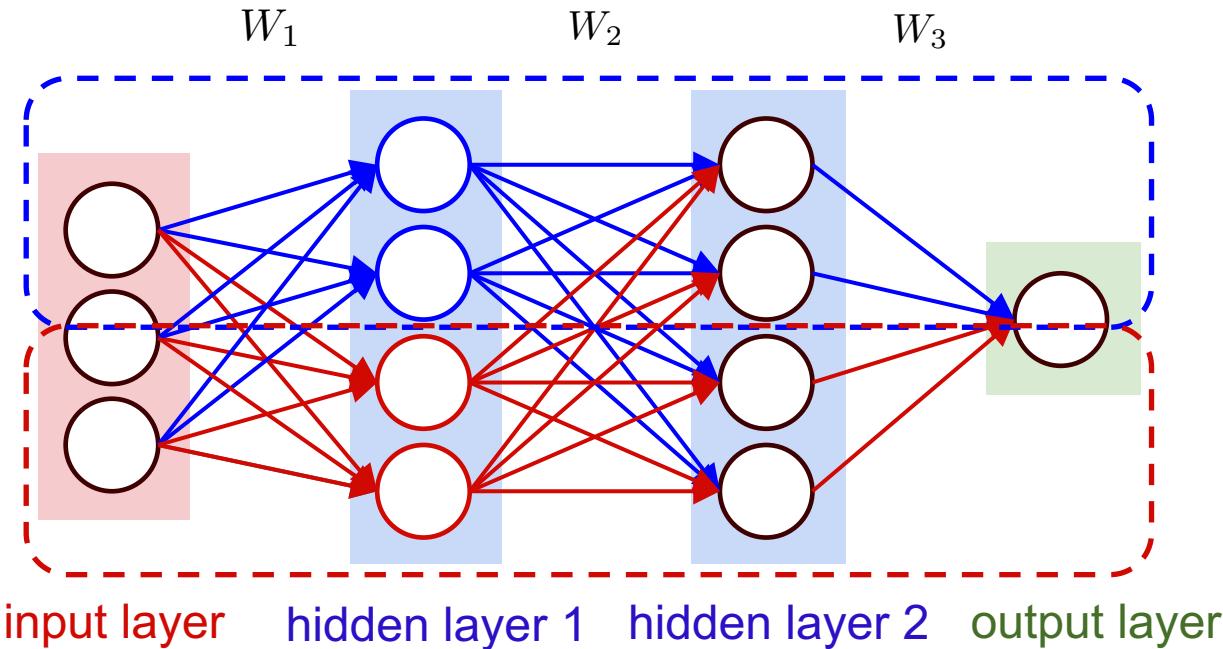


# Tensor parallelism is very flexible in partition manners



# Tensor parallelism: an example

$$y = W_3 W_2 W_1 x$$

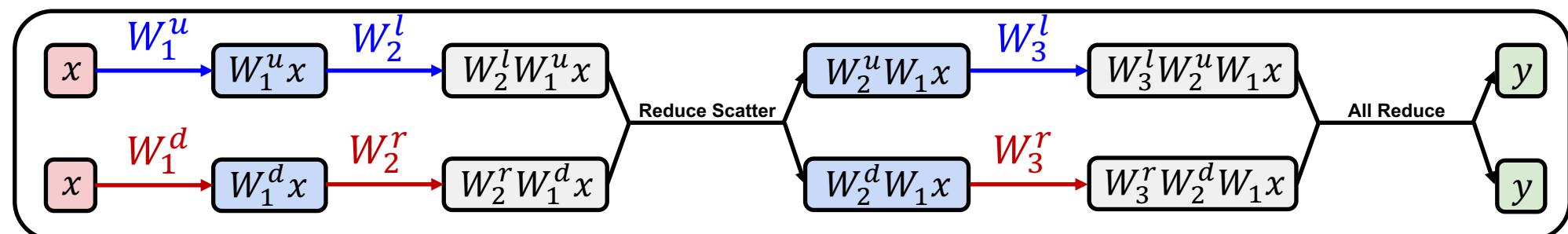


Partition 1:

$$\begin{bmatrix} W_1^u \\ W_1^d \end{bmatrix}$$

$$[W_2^\ell \quad W_2^r]$$

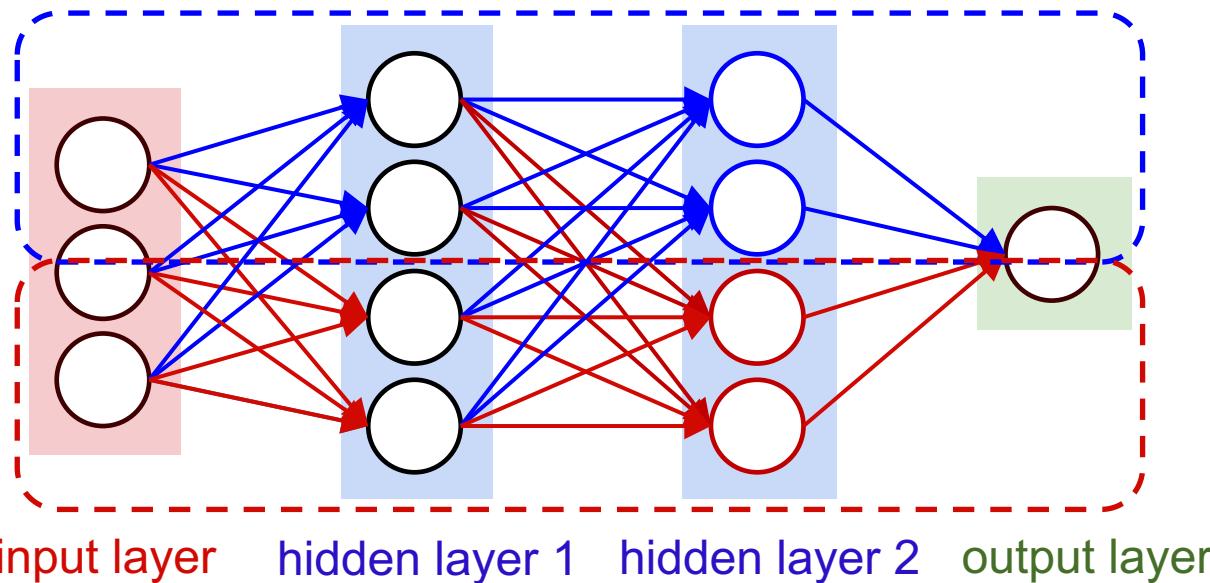
$$[W_3^\ell \quad W_3^r]$$



No idle time; only comm. and compt. time exist; almost each layer incurs comm.

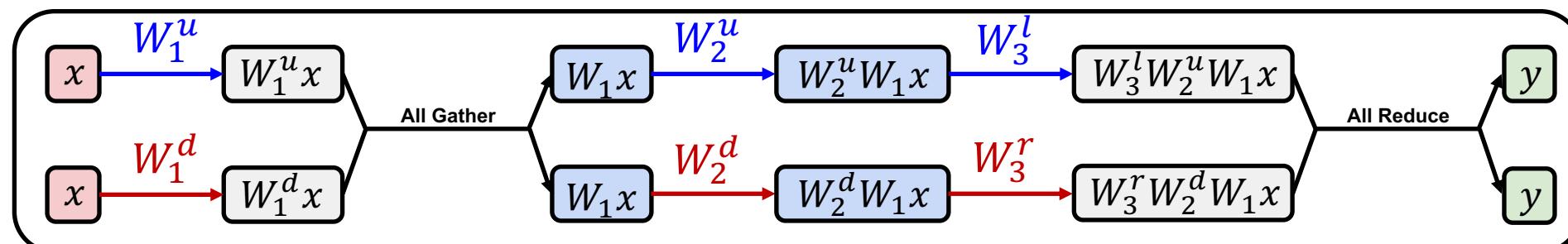
# Tensor parallelism: an example

$$y = W_3 W_2 W_1 x$$



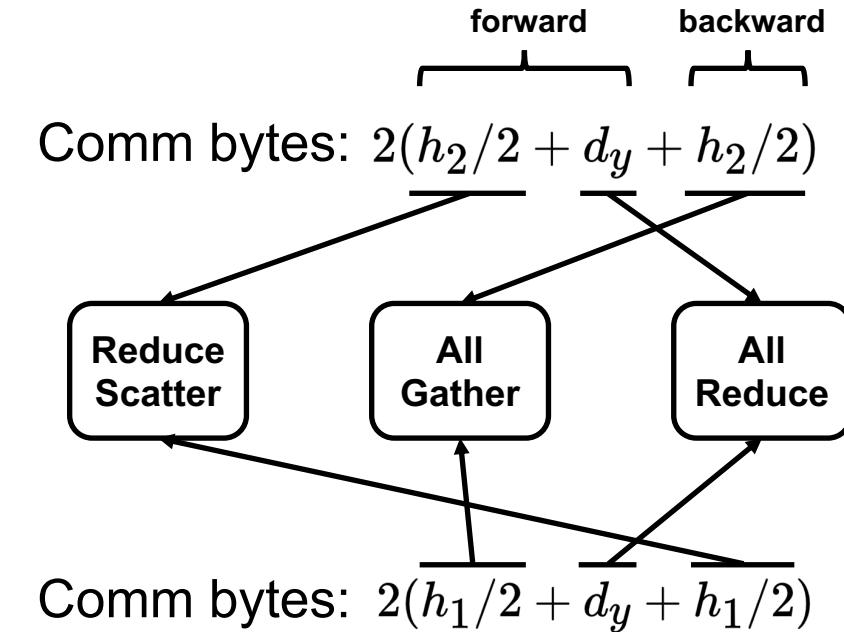
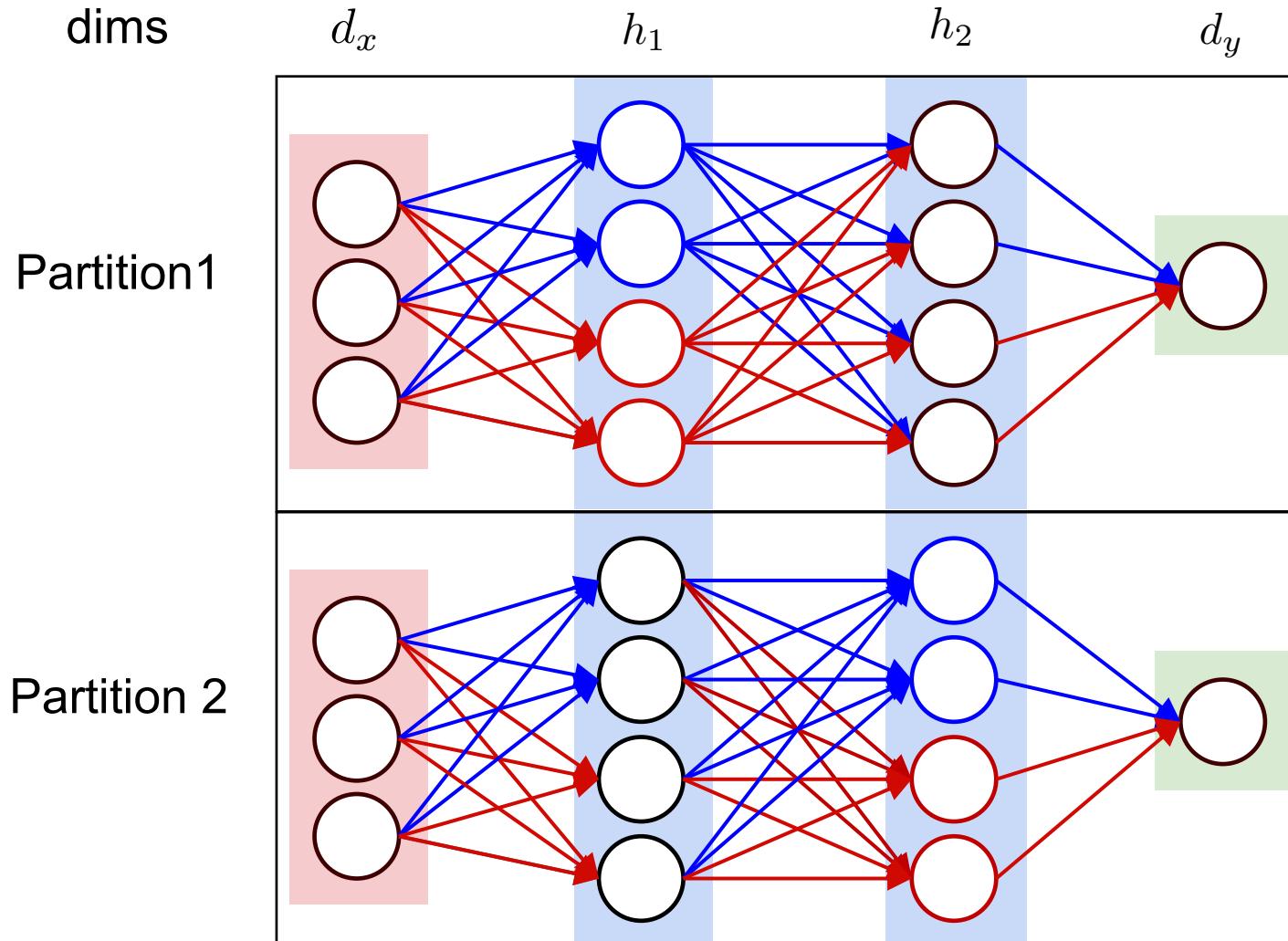
Partition 2:  $\begin{bmatrix} W_1^u \\ W_1^d \end{bmatrix}$

$$\begin{bmatrix} W_2^u \\ W_2^d \end{bmatrix} \quad [W_3^\ell \quad W_3^r]$$

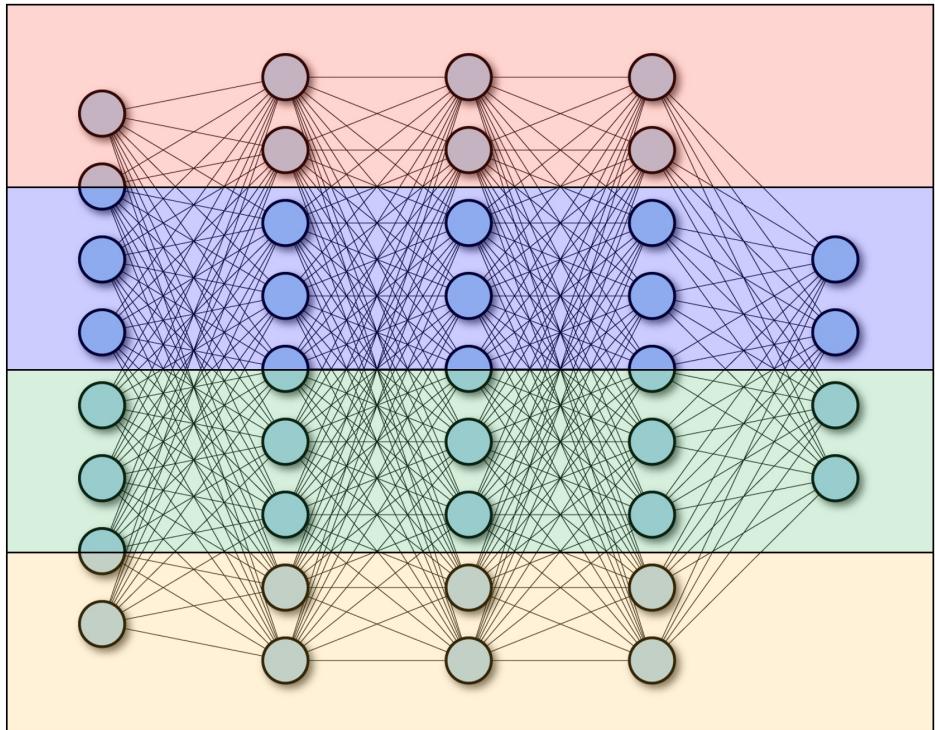


No idle time; only comm. and compt. time exist; almost each layer incurs comm.

# Different partitions incurs different communication cost



# Tensor parallelism: summary



## Pros:

- Cuts the memory by N folds (if partition uniformly)
- Accelerates training due to parallelism
- No idle time

## Cons:

- Introduce more comm. than pipeline parallelism
- Partition affects comm.; non-trivial to find good partition

The training time can be fully computed since both computation and communication cost are clear

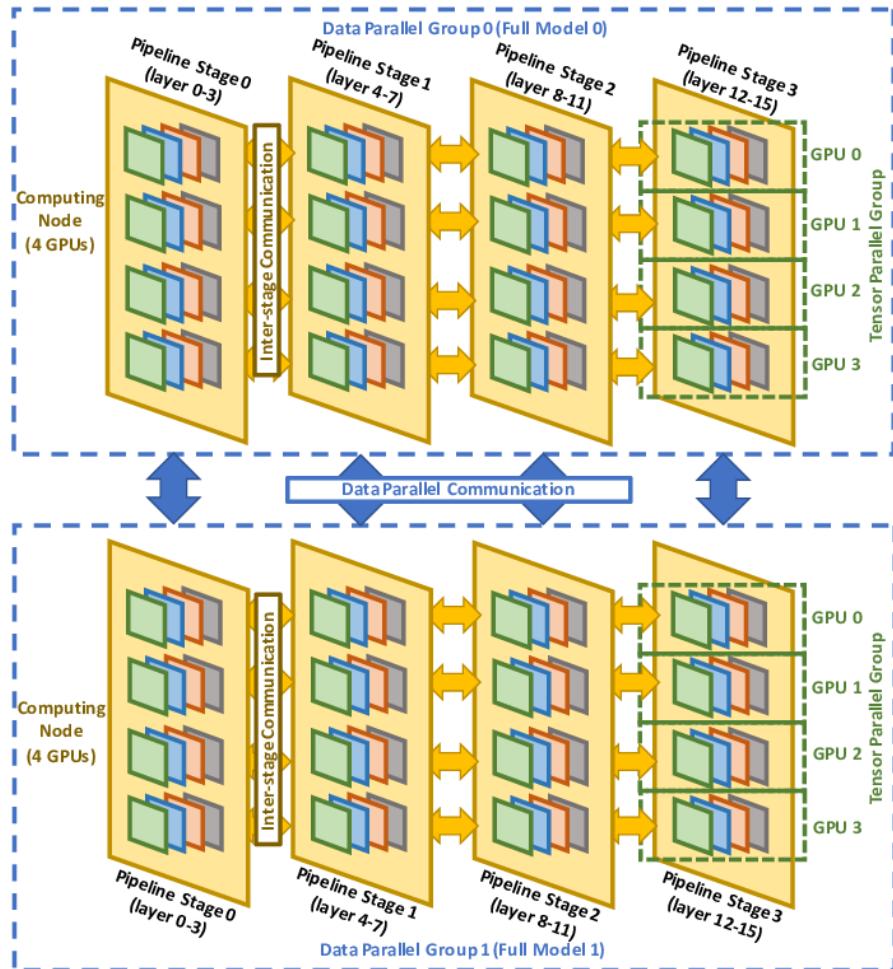


## PART 07

---

### 3D Parallelism

# 3D parallelism



- 3D parallelism is the most common approach to organize GPUs
- It is critical to determine how to combine data/pipeline/tensor parallelism
- Different combination will incur very different communication cost

# 3D parallelism: summary

---



- Data parallelism reduces token cost but maintain the big model
- Pipeline parallelism reduces memory cost; incurs less comm. cost; introduces idle time; sensitive to how layers are partitioned
- Tensor parallelism reduces memory cost; incurs big comm. cost; does not introduce time; sensitive to how tensors are partitioned
- Different 3D parallel combination will result in different LLM training performance
- Each parallelism's mechanism is very clear; training time can be fully computed

# Search for better 3D parallel policy

---



- Given a cluster of GPUs, a LLM task, and a specific 3D parallel policy, we can evaluate whether the policy is feasible, and what the performance (throughput, scalability, etc.) is
- In fact, we have a performance simulator!
- Given the simulator, we can easily search for a better 3D parallel policy using
  - Reinforcement learning
  - Dynamic programming
  - Integer programming
  - Zeroth-order optimization



# Thank you!

Kun Yuan homepage: <https://kunyuan827.github.io/>