



Kempner
INSTITUTE



HARVARD
UNIVERSITY

Large Language Model Distributed Training

October 18, 2024

Objectives

By the end of this workshop, you will be able to:

- Outline reasons to train models using more than one GPU.
- Understand different GPU collective communication primitives and their role in each parallel technique.
- Understand different parallelization techniques for distributed LLM training using GPUs.
- Train different transformers using OLMo (AI2 Open Large Language Model) in a distributed fashion on the HPC cluster.

Agenda

- 1 Why Going Distributed?
- 2 Intro to Distributed GPU Computing
- 3 Different Distributed LLM Training Techniques
- 4 Training a Large LLM on the Cluster

Why Distributed?

Data Size, Model Size or Both?

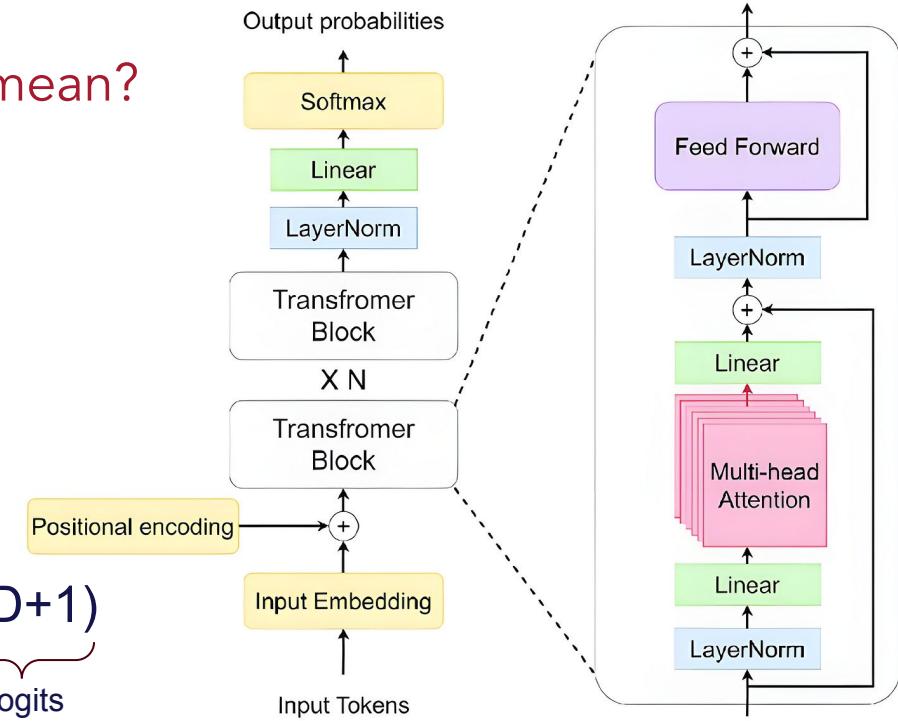
- Speedup Training
 - Model fits into a single GPU
 - Have a huge dataset to process
 - More GPUs - More computing resources
- Train Larger Models
 - We need to divide the model across multiple GPUs to be able to train it.

Decoder-only LLM Architecture

What do 1B, 7B, 70B, ... LLM sizes mean?

- V: vocabulary size
- B: batch size
- D: model dimension
- T: sequence length
- N: number of transformer blocks
- N_h : number of attention heads

$$P = \underbrace{N(12D^2 + 10D)}_{\text{Transformer Block}} + \underbrace{D(V+T)}_{\text{Embedding}} + \underbrace{2D}_{\text{Last Norm}} + \underbrace{V(D+1)}_{\text{Logits}}$$



Required Memory Estimation

Memory Contributors:

- Parameters
- Optimizer states
- Activations
- Others (Input/Target)

$$M_{\text{parameters}} = P$$

$$M_{\text{optimizer}} = 3P$$

$$M_{\text{activations}} = \underbrace{5NBTD}_{\text{FF}} + \underbrace{\text{BNN}_h T^2}_{\text{attention}} + \underbrace{2BVT}_{\text{embed}} + \underbrace{2BTD}_{\text{LNs, head}}$$

$$M_{\text{other}} = 2BT$$

Total Memory Requirement:

$$M = M_{\text{parameters}} + M_{\text{optimizer}} + M_{\text{activations}} + M_{\text{other}}$$

- V: vocabulary size
 - B: batch size
 - D: model dimension
 - T: sequence length
 - N: number of transformer blocks
 - N_h : number of attention heads
- $$P = N(12D^2 + 10D) + D(V+T) + 2D + V(D+1)$$

Floating-point Formats

Floating Point Formats

bfloat16: Brain Floating Point Format

Range: $\sim 1e^{-38}$ to $\sim 3e^{38}$



fp32: Single-precision IEEE Floating Point Format

Range: $\sim 1e^{-38}$ to $\sim 3e^{38}$



fp16: Half-precision IEEE Floating Point Format

Range: $\sim 5.96e^{-8}$ to 65504



<https://www.tensorops.ai/post/what-are-quantized-langs>

Required Memory Estimation

Just Loading Weights

- Parameters in bfloat16 (2 Bytes)
- Optimizer in float32 (4 Bytes)

$$M_{\text{model}} = P * 2 + 3P * 4 = 14 P \text{ (Bytes)}$$

GPUs:

H100 \Rightarrow 80 G
A100 \Rightarrow 40 G

Exercise: Assessing Memory Needs

One of the OLMo models we will be working with has the hyperparameters shown on the right.

- V (vocabulary size): 50280
- B (batch size): 256
- D (model dimension): 8192
- T (sequence length): 4096
- N (number of transformer blocks): 80
- N_h (number of attention heads): 64

- 1) Compute the number of parameters
- 2) Compute the memory that the model weights will take
- 3) Will this fit on one GPU? If not, how many A100 GPUs would you need? How many H100s?

GPUs:

H100 \Rightarrow 80 G
A100 \Rightarrow 40 G

Exercise Solution

$$P = N(12D^2 + 10D) + D(V+T) + 2D + V(D+1)$$

$$\begin{aligned} P &= 80 * (12 * 8192 ** 2 + 10 * 8192) + \\ &8192 * (50280 + 4096) + 2 * 8192 + \\ &50280 * (8192 + 1) \end{aligned}$$

$$P = 65288471656$$

P = 65 billion

$$M = 14P = 14 * 65288471656$$

M = 910 billion bytes = 910 GB

- V (vocabulary size): 50280
- B (batch size): 256
- D (model dimension): 8192
- T (sequence length): 4096
- N (number of transformer blocks): 80
- N_h (number of attention heads): 64

GPUs:

H100 \Rightarrow 80 G
A100 \Rightarrow 40 G

Required Memory Estimation

Just Loading Weights

- Parameters in bfloat16 (2 Bytes)
- Optimizer in float32 (4 Bytes)

$$M_{\text{model}} = P * 2 + 3P * 4 = 14 P \text{ (Bytes)}$$

Model Size (P)	Approx. memory used to train model (GB)
300M	4
1B	14
7B	98
13B	182
70B	980

GPUs:

H100 \Rightarrow 80 G
A100 \Rightarrow 40 G

Parameter Comparisons

	Vocabulary size (V)	Model dimension (D)	Sequence length (T)	Number of transformer blocks (N)	Number of attention heads (N_h)
OLMo 1B	32100	2046	2048	16	16
OLMo 7B	32100	4096	2048	32	32
OLMo 70B	50280	8192	4096	80	64

OLMo

Open Language Model

- A highly performant, truly open LLM and framework
 - 100% of ingredients are available to public including code, weights, checkpoints, training data and system logs.
 - To advance AI and study language models collectively
 - Decoder-only Architecture

Let's setup and install it on the cluster:

https://github.com/KempnerInstitute/OLMo/blob/main/README_KempnerInstitute.md

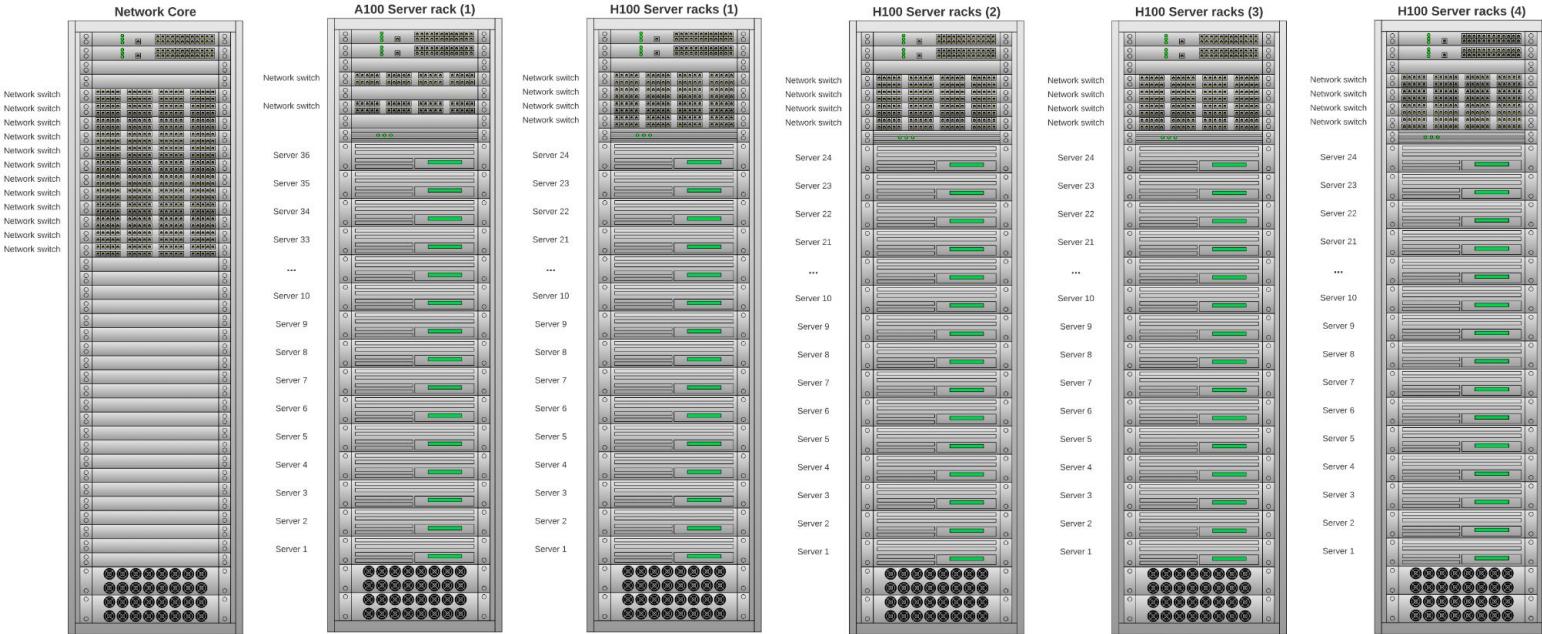
To request Kempner GPUs if you don't have access to Kempner cluster:

https://handbook.eng.kempnerinstitute.harvard.edu/s1_high_performance_computing/kempner_cluster/accessing_gpu_by_fasrc_users.html

Agenda

- 1 Why Going Distributed?
- 2 Intro to Distributed GPU Computing
- 3 Different Distributed LLM Training Techniques
- 4 Training a Large LLM on the Cluster

HPC Cluster



In Production (144 x A100 40 GB, 384 x H100 80 GB)

HPC Cluster - Computational Power

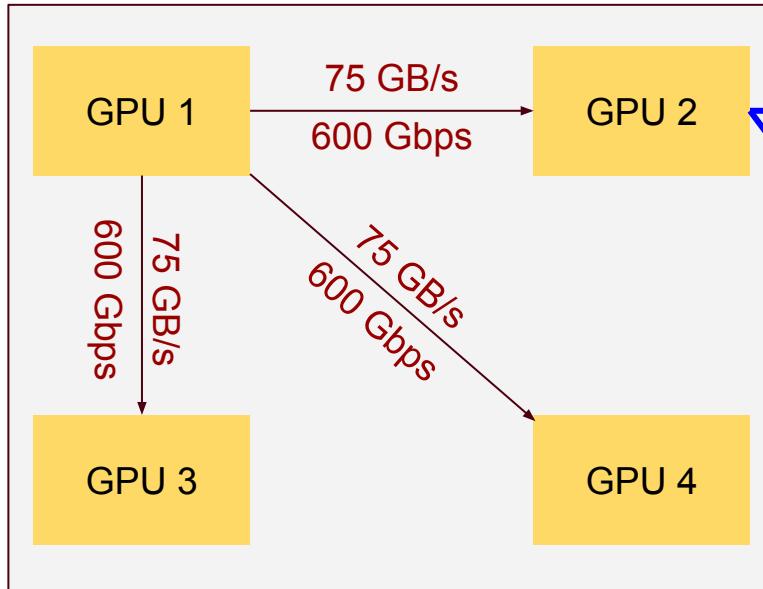
Compute Power (H100 GPUs)

- Number of Racks (N_{Rack}) = 4
- Number of Nodes per Rack (N_{Node}) = 24
- Number of GPUs per Node (N_{GPU}) = 4
- Total Number of GPUs = $N_{\text{Rack}} \times N_{\text{Node}} \times N_{\text{GPU}} = 384$ H100 GPUs
- Total Computational Power in FLOPs,
 - Total FLOPs (BFLOAT16 Tensor Core) = $384 \text{ GPU} \times 1979 \text{ TFLOPs / GPU} = 759,936$
TFLOPs = **759 PFLOPs**
 - Total FLOPs (FP32) = $384 \text{ GPU} \times 67 \text{ TFLOPs / GPU} = 25,728 \text{ TFLOPs} = \text{25 PFLOPs}$

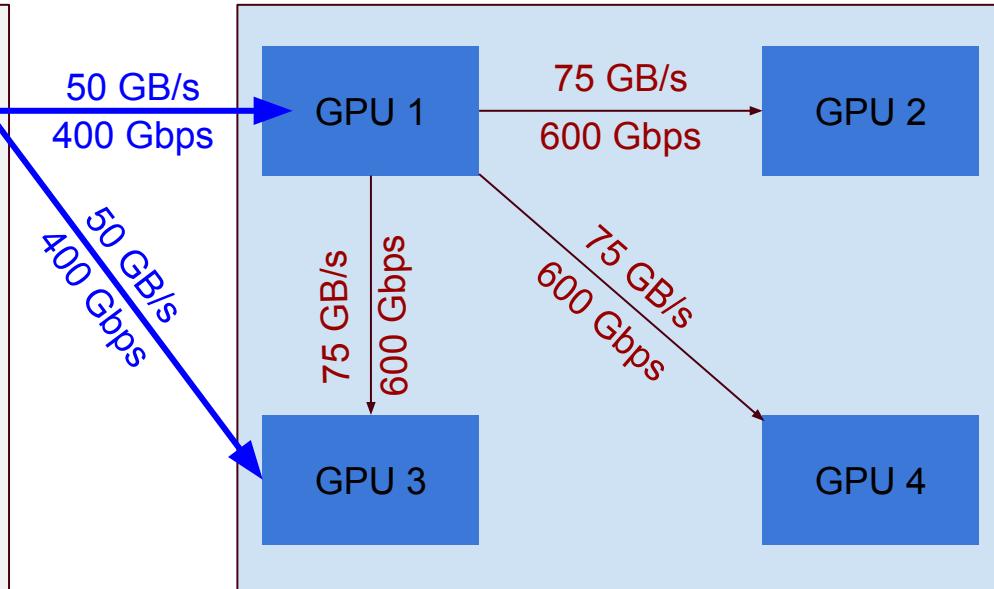
Technical Specifications	
	H100 SXM
FP64	34 teraFLOPS
FP64 Tensor Core	67 teraFLOPS
FP32	67 teraFLOPS
TF32 Tensor Core*	989 teraFLOPS
BFLOAT16 Tensor Core*	1,979 teraFLOPS
FP16 Tensor Core*	1,979 teraFLOPS
FP8 Tensor Core*	3,958 teraFLOPS
INT8 Tensor Core*	3,958 TOPS

GPU-to-GPU Communication

Node 1



Node 2



Inside Node (**NVLINK**): Each GPU talks to other three GPUs at 75 GB/s (single direction). This sums up to 900 GB/s all GPU-GPU bidirectional speed. $75 \text{ GB/s} * 6 * 2 = 900 \text{ GB/s}$

Outside Node (**InfiniBand Network NDR**): Each GPU communicates to other GPUs in another node at 400 Gbps (50 GB/s).

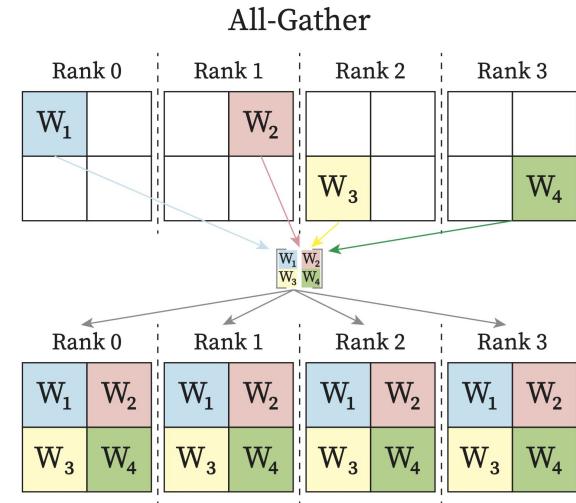
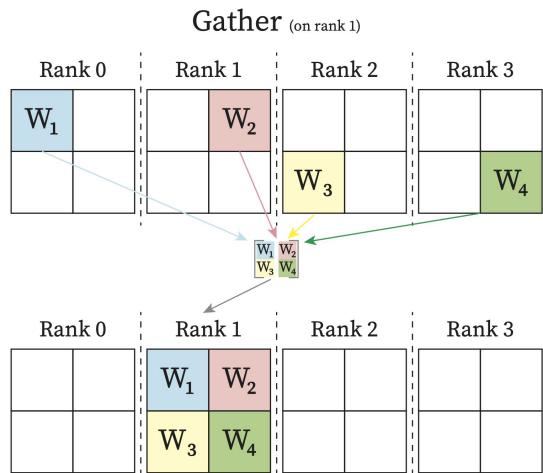
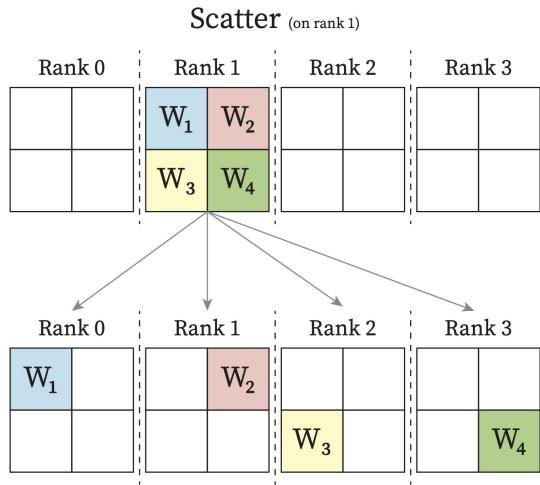
Inter-GPU Communication

NCCL

NVIDIA Collective Communication Library (**NCCL**, pronounced “NICKEL”) is used as backend in distributed strategies for NVIDIA GPUs

NCCL offers various collective **communication primitives**

Other NCCL Collective Primitives

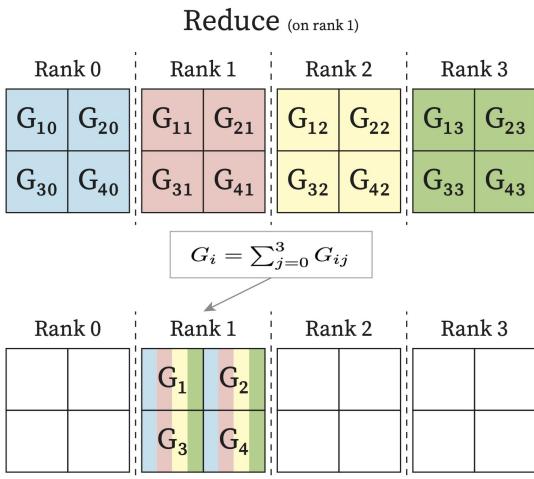


Scatter: From one rank data will be distributed across ranks.

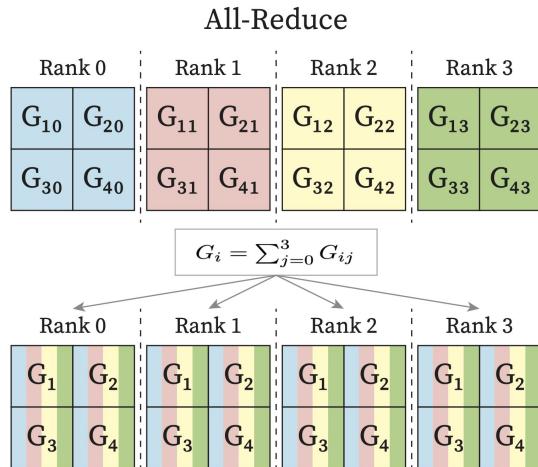
Gather: One rank will receive the aggregation of data from all ranks.

All-Gather: Each rank receives the aggregation of data from all ranks in the order of the ranks.

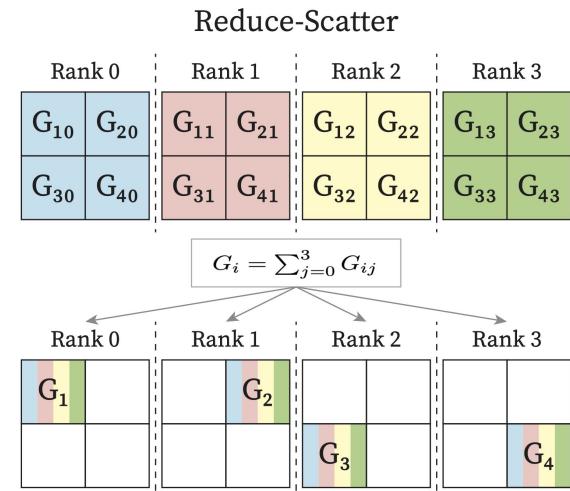
Other NCCL Collective Primitives



Reduce: One rank receives the reduction of input values across ranks.



All-Reduce: Each rank receives the reduction of input values across ranks.

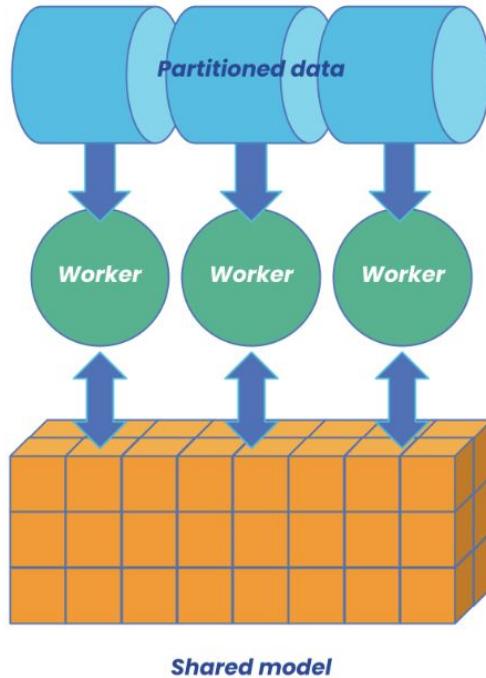


Reduce-Scatter: Input values are reduced across ranks, with each rank receiving a subpart of the result.

Agenda

- 1 Why Going Distributed?
- 2 Intro to Distributed GPU Computing
- 3 Different Distributed LLM Training Techniques
- 4 Training a Large LLM on the Cluster

Distributed Data Parallel Processing

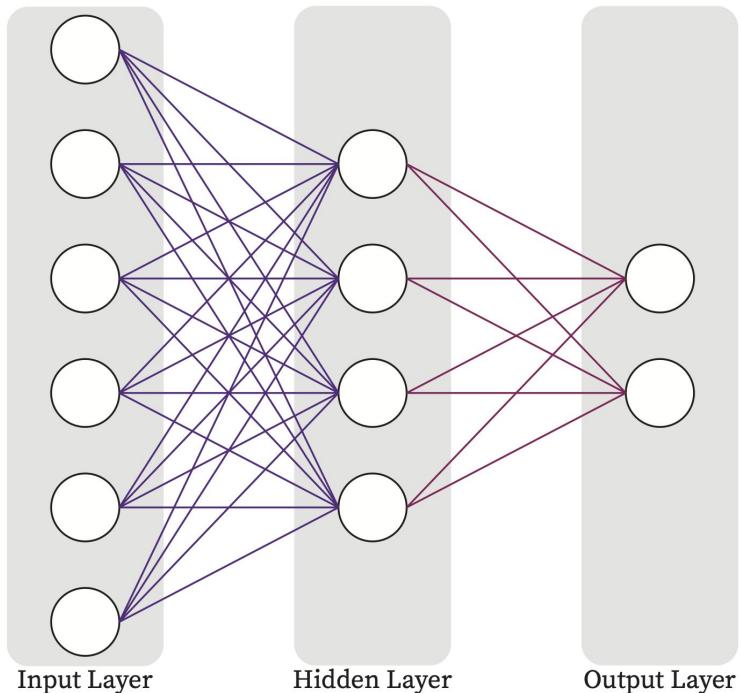


Most common approach to distributed training in machine learning

Each GPU trains a copy of the model.
Dataset is split into **different batches of data on each GPU**

<https://www.anyscale.com/blog/what-is-distributed-training>

Multi-layer Perceptron



```
class MLP(nn.Module):
    def __init__(self, in_feature, hidden_units, out_feature):
        super().__init__()

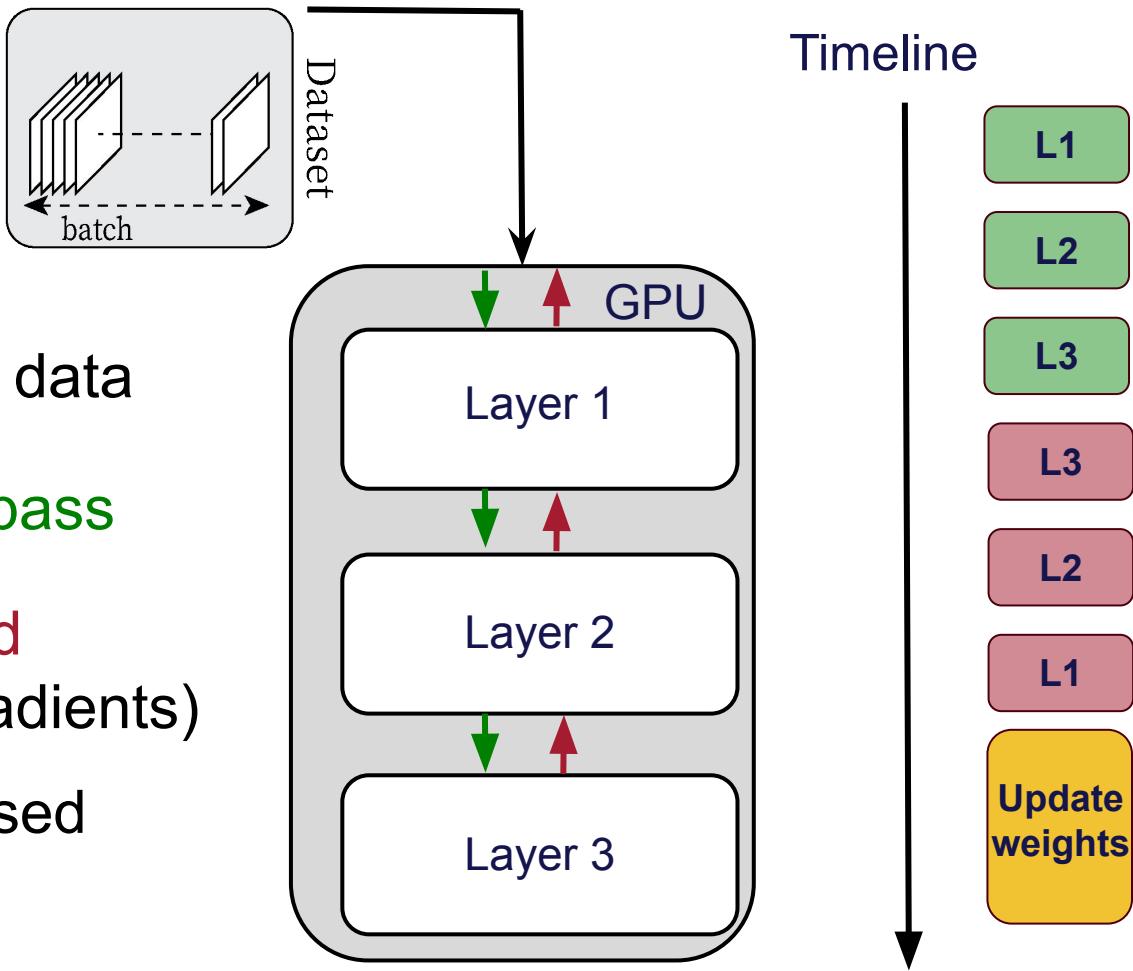
        self.hidden_layer = nn.Linear(in_feature, hidden_units)
        self.output_layer = nn.Linear(hidden_units, out_feature)

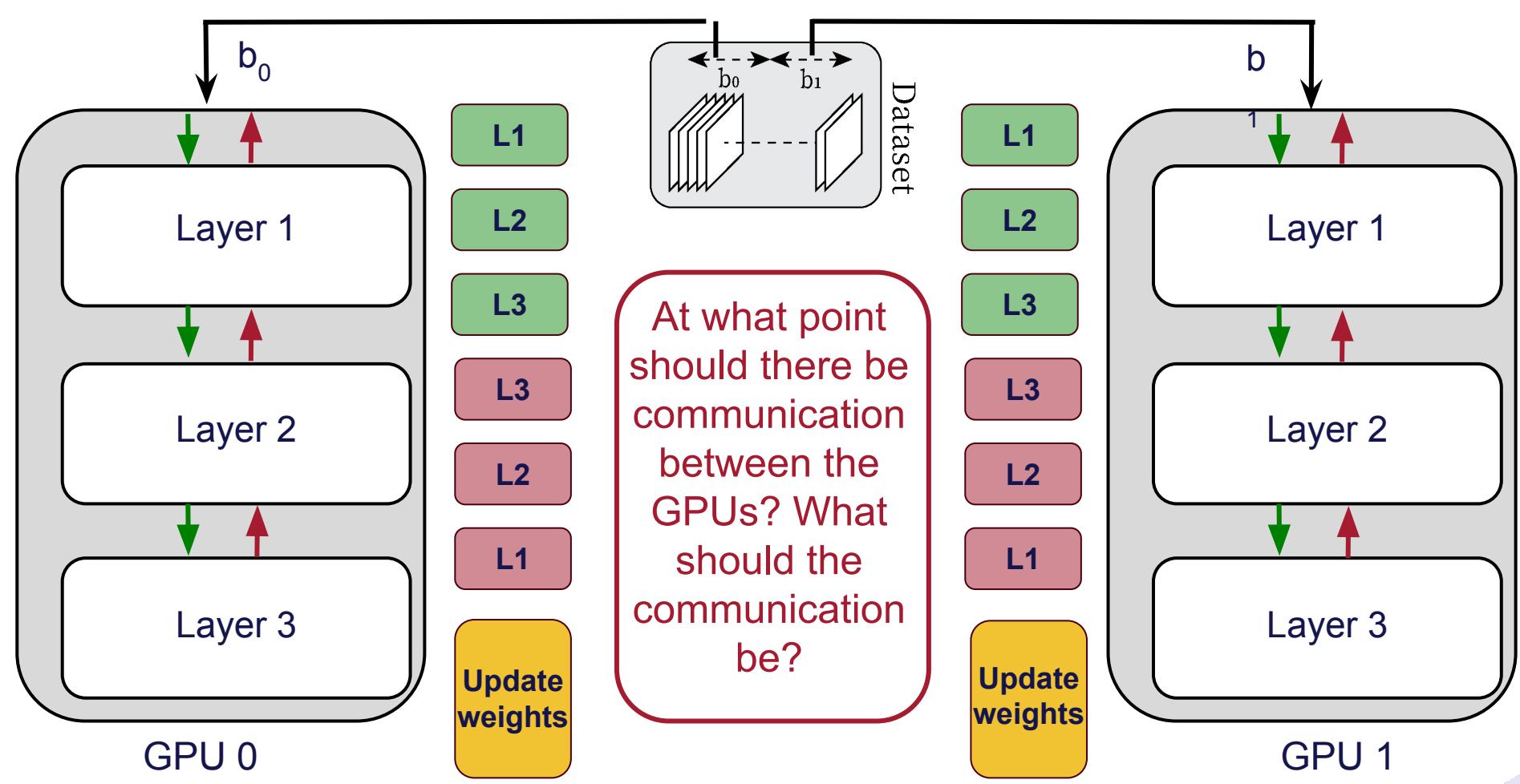
    def forward(self, x):
        x = self.hidden_layer(x)
        x = self.output_layer(x)
        return x
```

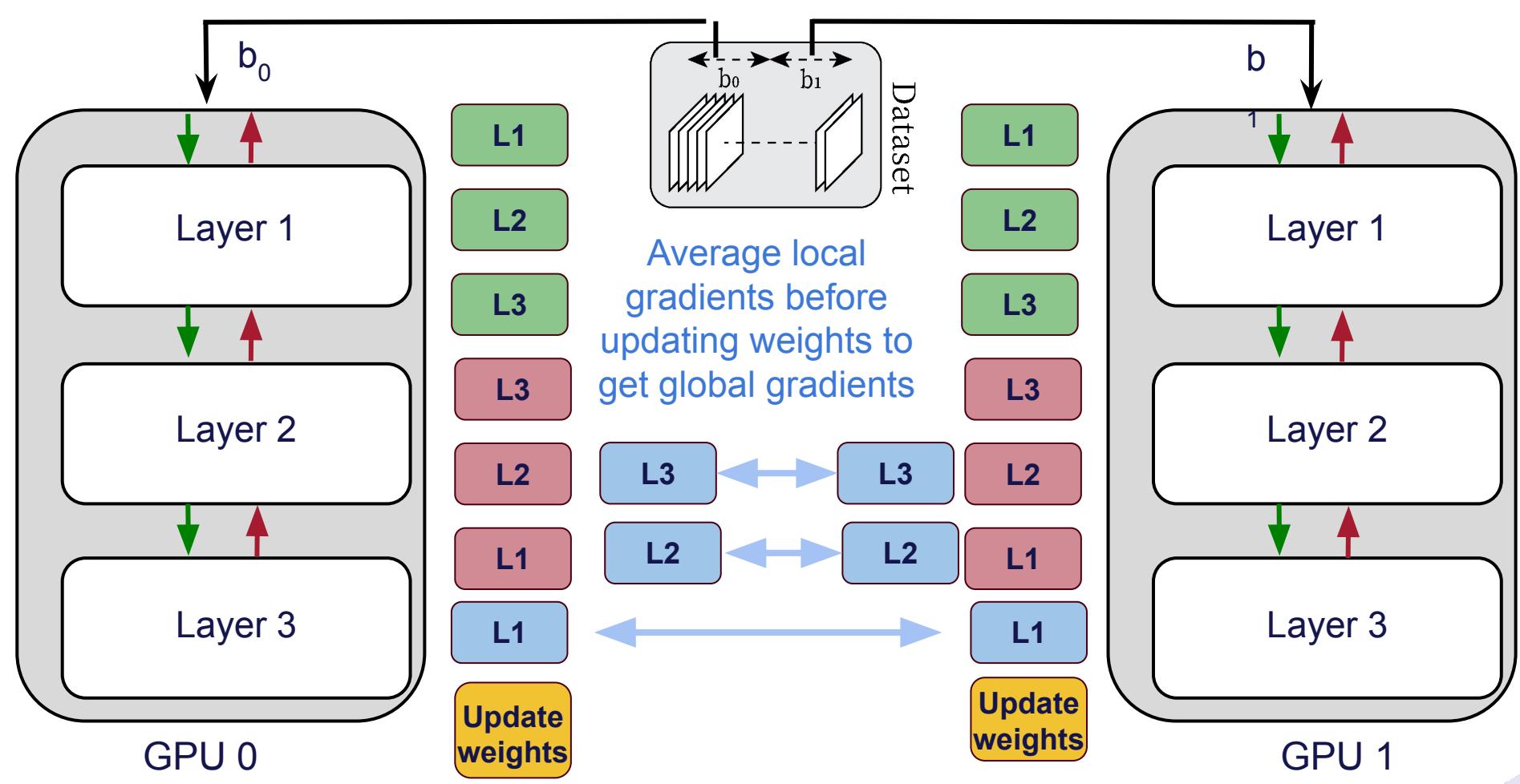
$$\begin{aligned} \mathbf{x}_i &= [x_1 \ x_2 \ x_3 \ x_4 \ x_5 \ x_6] & \mathbf{W}_1 &= \begin{bmatrix} w_1 & w_2 & w_3 & w_4 \\ w_5 & w_6 & w_7 & w_8 \\ w_9 & w_{10} & w_{11} & w_{12} \\ w_{13} & w_{14} & w_{15} & w_{16} \\ w_{17} & w_{18} & w_{19} & w_{20} \\ w_{21} & w_{22} & w_{23} & w_{24} \end{bmatrix} & \mathbf{W}_2 &= \begin{bmatrix} z_1 & z_2 \\ z_3 & z_4 \\ z_5 & z_6 \\ z_7 & z_8 \end{bmatrix} \\ \mathbf{y}_i &= [y_1 \ y_2] & \mathbf{b}_1 &= [b_1 \ b_2 \ b_3 \ b_4] & \mathbf{b}_2 &= [c_1 \ c_2] \end{aligned}$$

Single GPU MLP Training

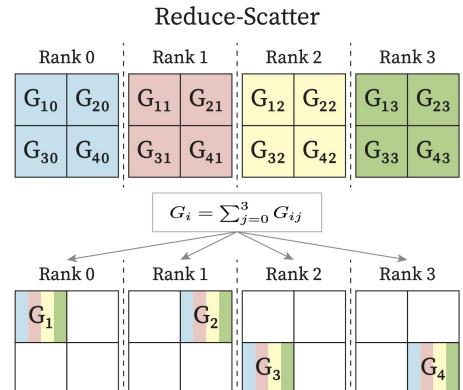
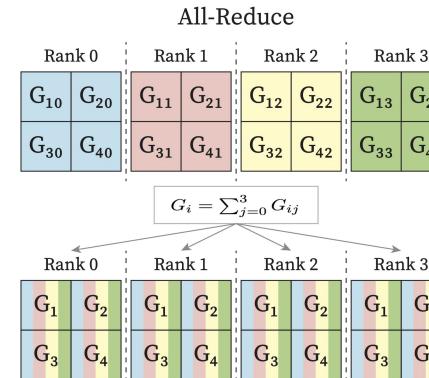
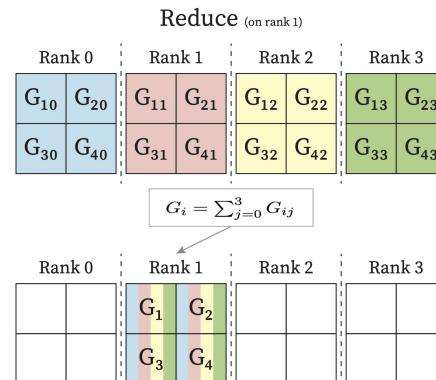
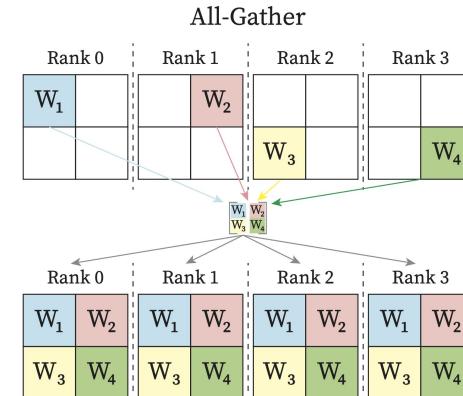
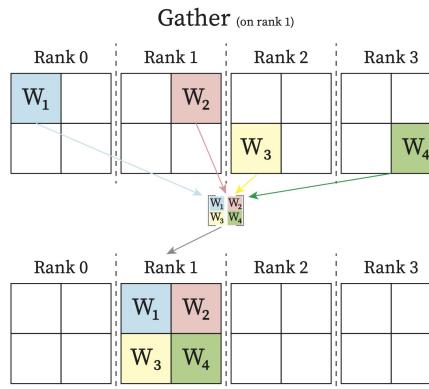
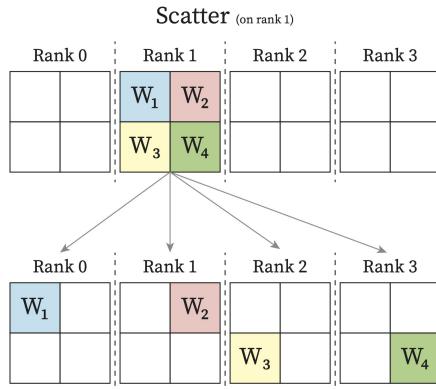
- 1) Model gets batch of data
- 2) Computes forward pass
- 3) Computes backward pass (computing gradients)
- 4) Updates weights based on gradients

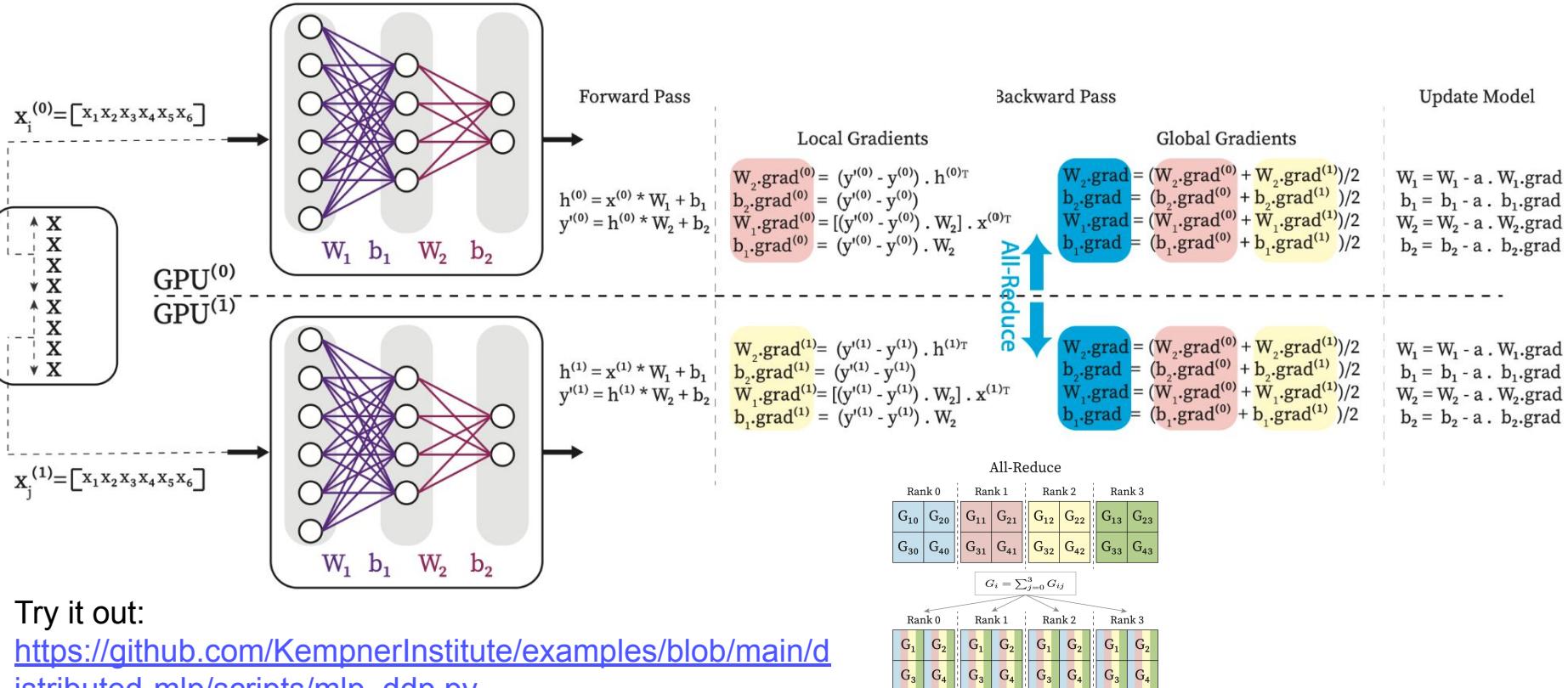






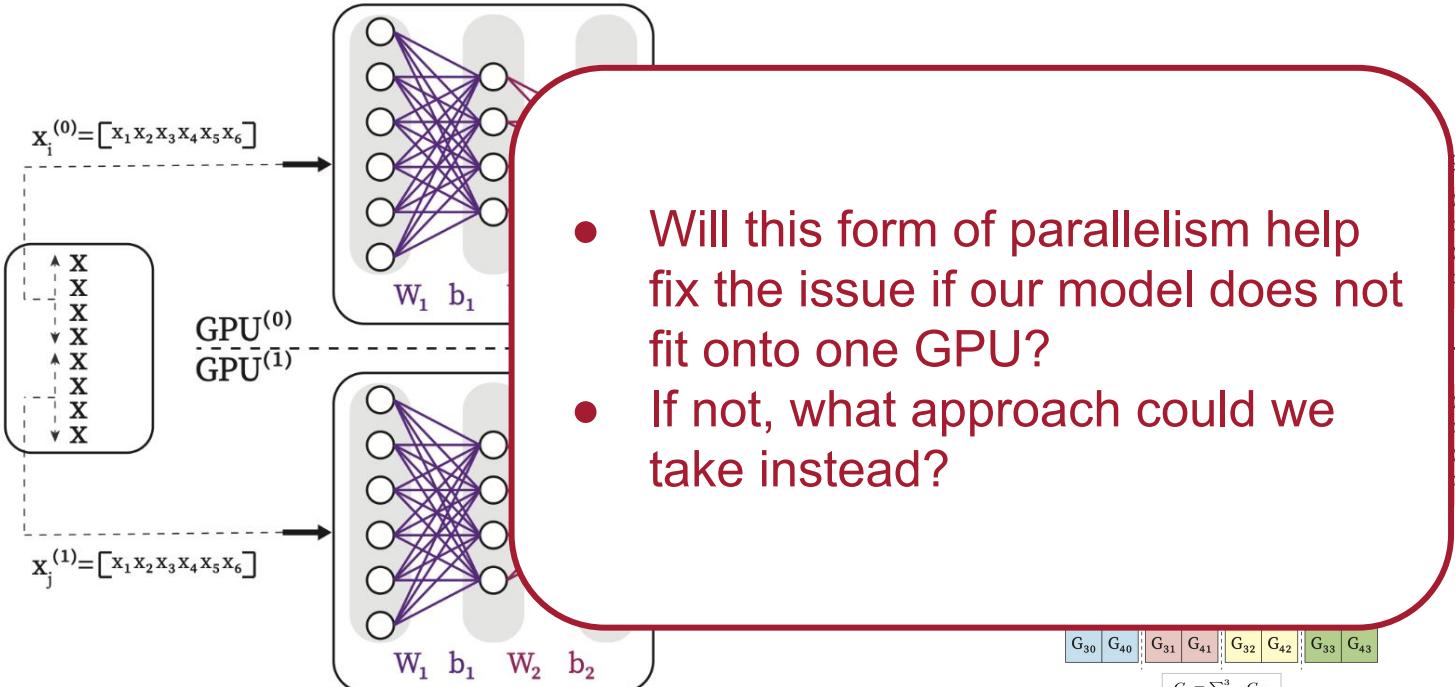
Which NCCL Collective Primitive?





Try it out:

https://github.com/KempnerInstitute/examples/blob/main/distributed-mlp/scripts/mlp_ddp.py



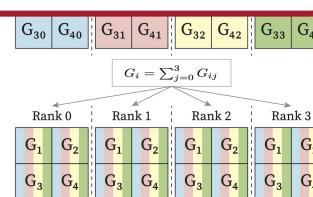
Try it out:

https://github.com/KempnerInstitute/examples/blob/main/distributed-mlp/scripts/mlp_ddp.py

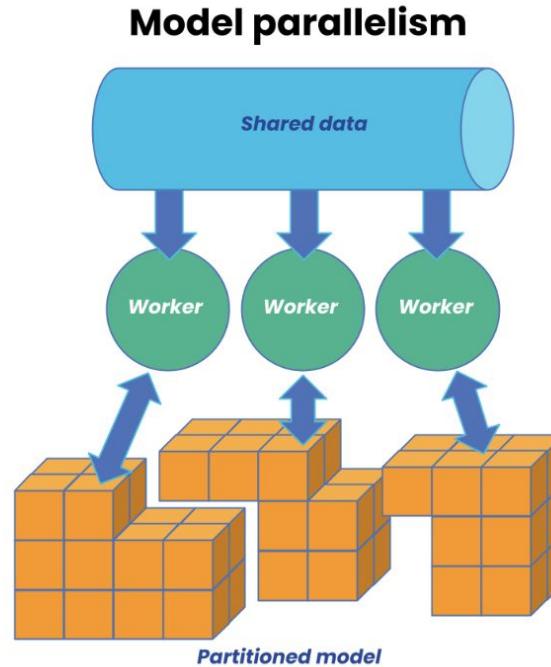
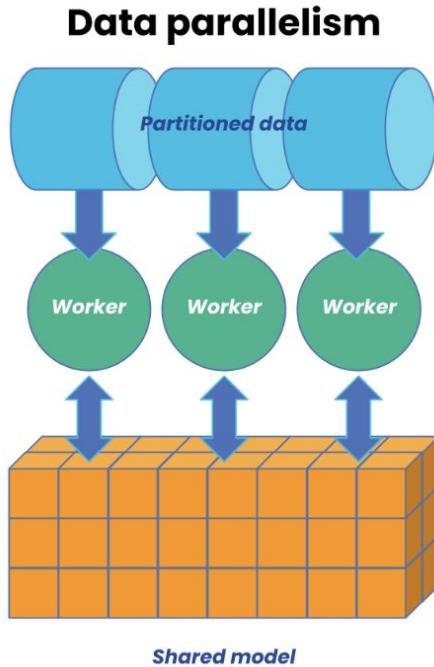
Update Model

Gradients $\overset{(0)}{g} + W_2 \cdot \text{grad}^{(1)}/2$ $\overset{(1)}{g} + b_2 \cdot \text{grad}^{(1)}/2$ $\overset{(0)}{g} + W_1 \cdot \text{grad}^{(1)}/2$ $\overset{(1)}{g} + b_1 \cdot \text{grad}^{(1)}/2$	$W_1 = W_1 - a \cdot W_1 \cdot \text{grad}$ $b_1 = b_1 - a \cdot b_1 \cdot \text{grad}$ $W_2 = W_2 - a \cdot W_2 \cdot \text{grad}$ $b_2 = b_2 - a \cdot b_2 \cdot \text{grad}$
---	--

$\overset{(0)}{g} + W_2 \cdot \text{grad}^{(1)}/2$ $\overset{(1)}{g} + b_2 \cdot \text{grad}^{(1)}/2$ $\overset{(0)}{g} + W_1 \cdot \text{grad}^{(1)}/2$ $\overset{(1)}{g} + b_1 \cdot \text{grad}^{(1)}/2$	$W_1 = W_1 - a \cdot W_1 \cdot \text{grad}$ $b_1 = b_1 - a \cdot b_1 \cdot \text{grad}$ $W_2 = W_2 - a \cdot W_2 \cdot \text{grad}$ $b_2 = b_2 - a \cdot b_2 \cdot \text{grad}$
--	--

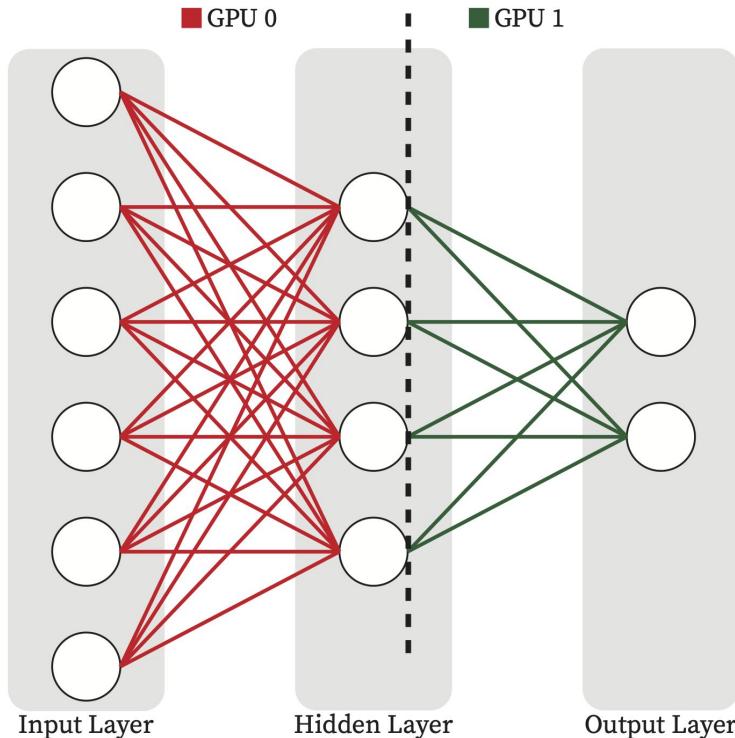


Model Parallelism



<https://www.anyscale.com/blog/what-is-distributed-training>

Model Parallelism

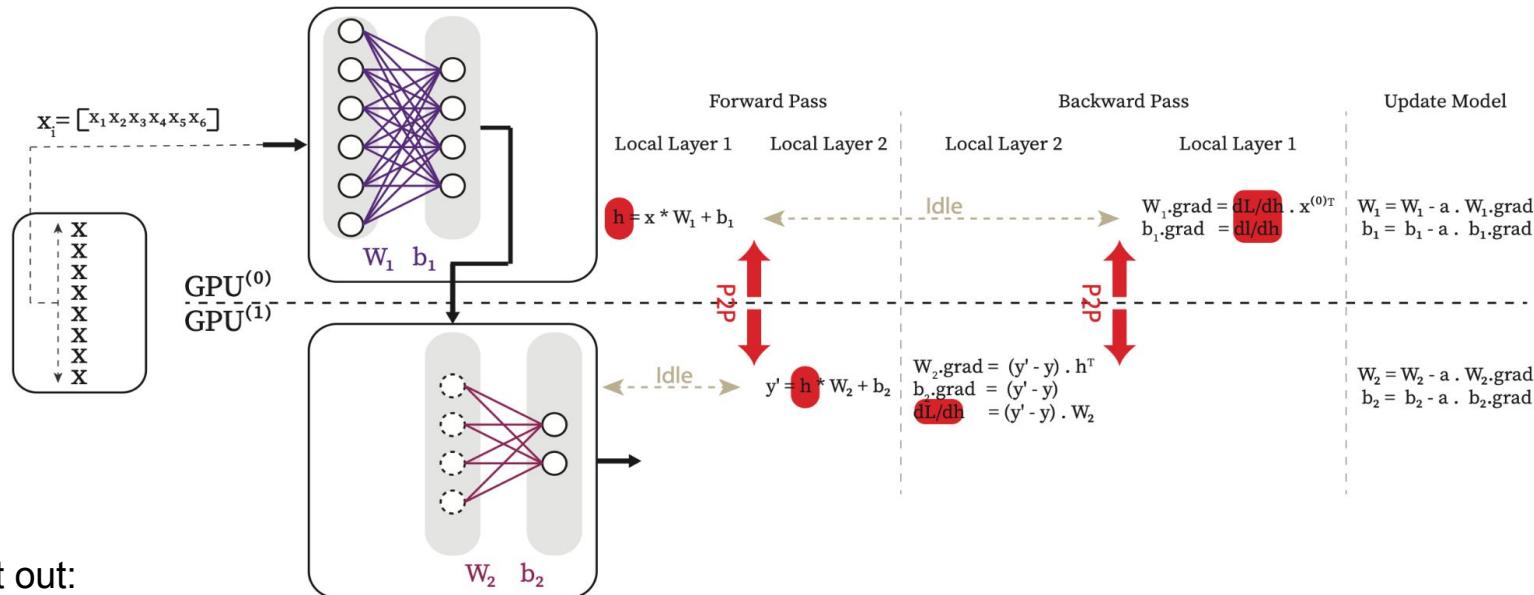


```
class MLP(nn.Module):
    def __init__(self, in_feature, hidden_units, out_feature):
        super().__init__()

        self.hidden_layer = nn.Linear(in_feature, hidden_units).to(0)
        self.output_layer = nn.Linear(hidden_units, out_feature).to(1)

    def forward(self, x):
        x = self.hidden_layer(x.to(0))
        x = self.output_layer(x.to(1))
        return x
```

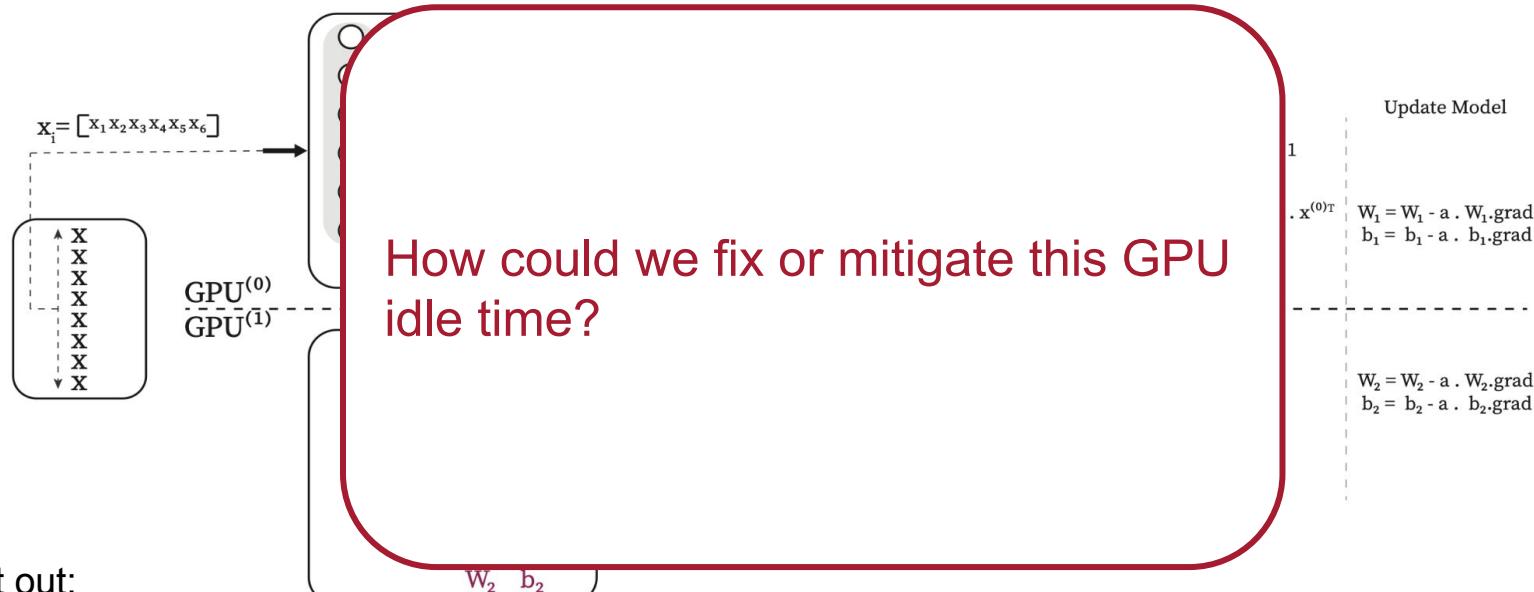
Model Parallelism and its Drawback



Try it out:

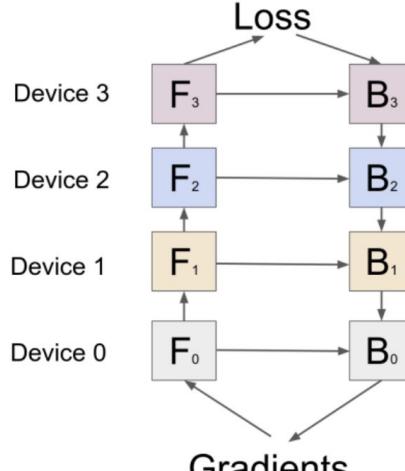
https://github.com/KempnerInstitute/examples/blob/main/distributed-mlp/scripts/mlp_model_parallel.py

Model Parallelism and its Drawback

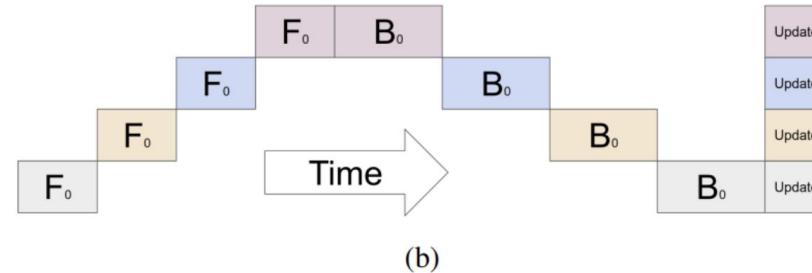


Pipeline Parallelism

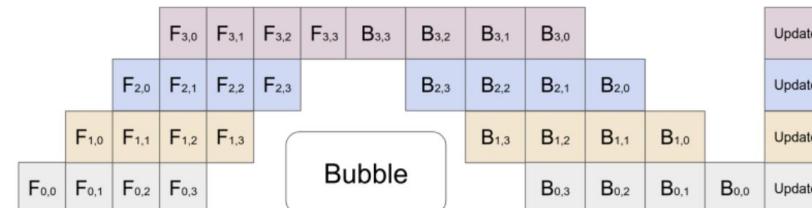
Pipeline Parallelism uses micro batches to reduce the idle time by adding overlaps



(a)



(b)

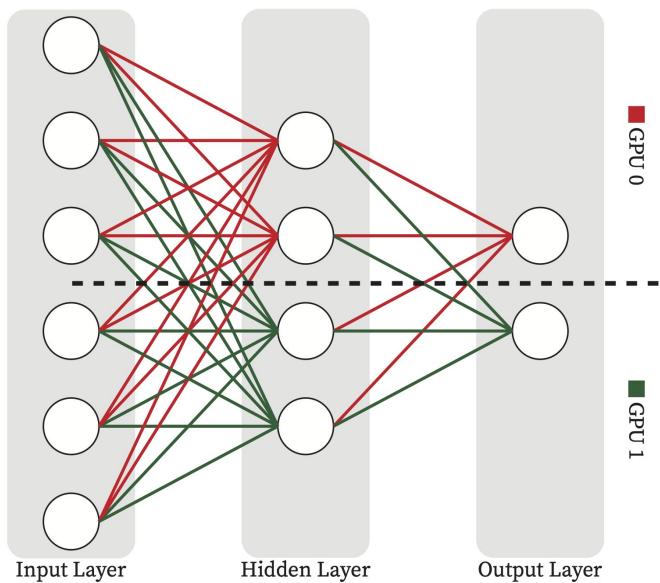
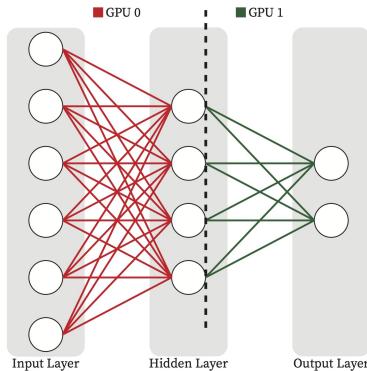


(c)

<https://medium.com/nerd-for-tech/an-overview-of-pipeline-parallelism-and-its-research-progress-7934e5e6d5b8>

Tensor Parallelism

Model/Pipeline Parallelism vs Tensor Parallelism

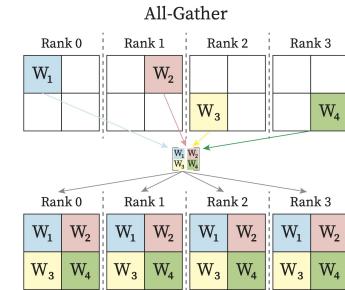
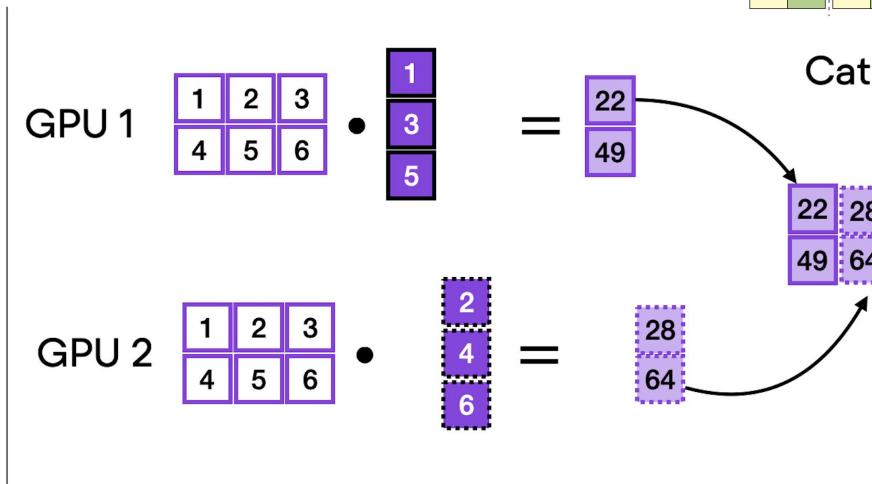


- All GPUs contribute in each layer computation
- Remove the GPU Idle time

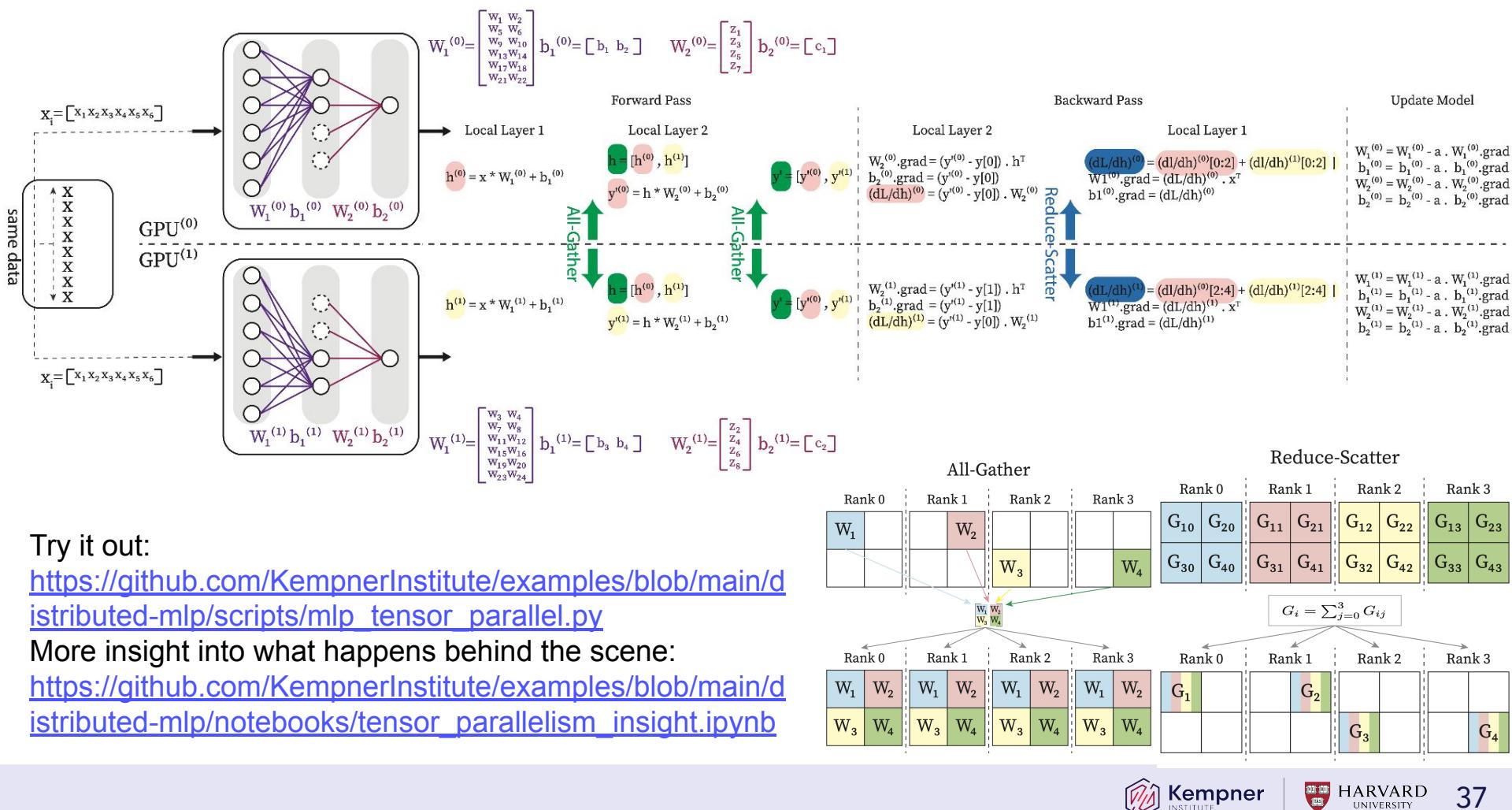
Tensor Parallelism

Splitting the weights column-wise between GPUs

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \cdot \begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \end{bmatrix} = \begin{bmatrix} 22 & 28 \\ 49 & 64 \end{bmatrix}$$



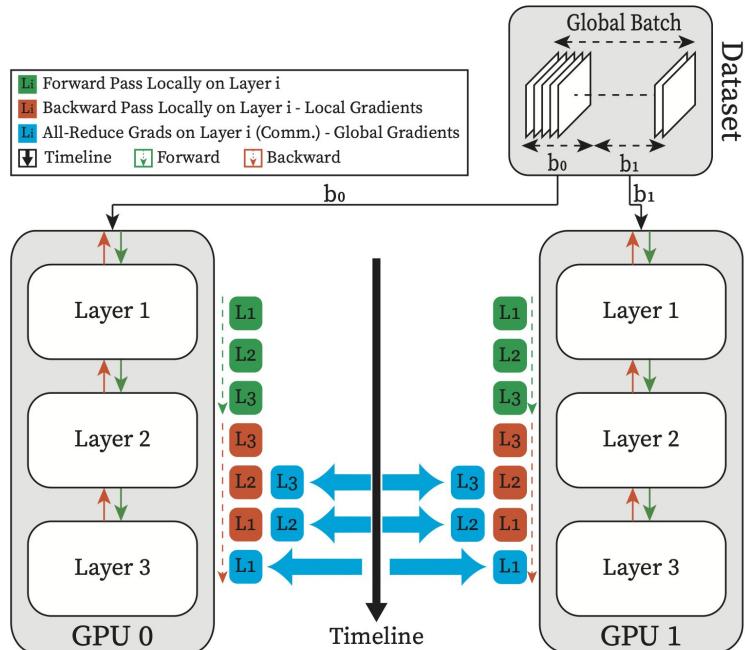
<https://magazine.sebastianraschka.com/p/accelerating-pytorch-model-training>



DDP vs FSDP: DDP

Each GPU has a copy of the model

1. Perform forward and backward passes locally
2. All-reduce gradients across GPUs (NCCL operation)
3. Update optimizer states and weights locally



FSDP vs DDP: Model Size

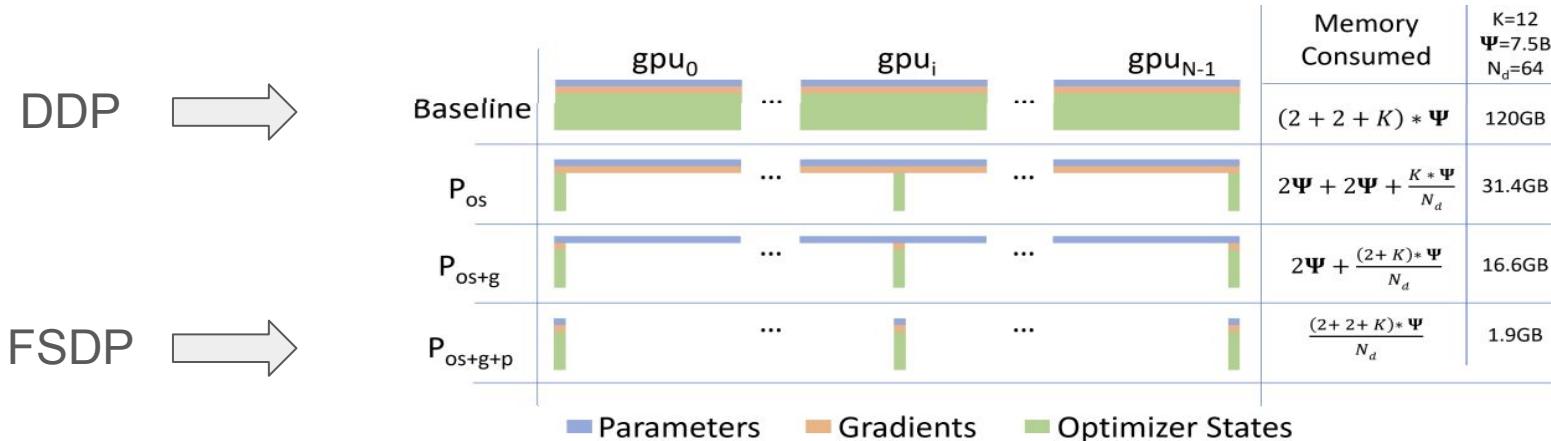


Figure 1: Comparing the per-device memory consumption of model states, with three stages of *ZeRO*-DP optimizations. Ψ denotes model size (number of parameters), K denotes the memory multiplier of optimizer states, and N_d denotes DP degree. In the example, we assume a model size of $\Psi = 7.5B$ and DP of $N_d = 64$ with $K = 12$ based on mixed-precision training with Adam optimizer.

Samyam R. et al, ZeRO: Memory Optimizations Toward Training Trillion Parameter Models. [arxiv](#)

DDP vs FSDP: FSDP

Each GPU has a shard of the model - 1D flatten parameters divided between GPUs

Forward:

1. All-gather all the weights across GPUs (NCCL operation)
2. Perform the forward pass locally
3. Release the collected weights to free memory

Backward:

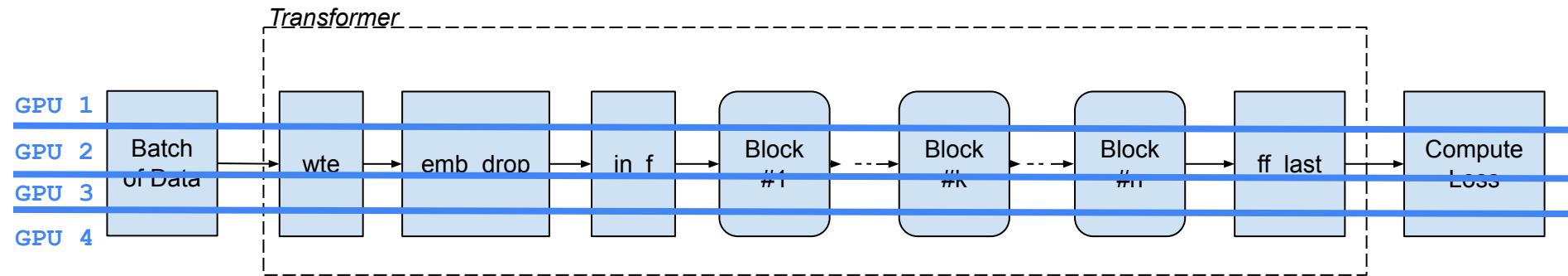
4. All-gather all the weights across GPUs (NCCL operation)
5. Perform the backward pass locally
6. Release the collected weights to free memory

Optimizer step and weight update:

7. Reduce-scatter gradients across GPUs (NCCL operation)

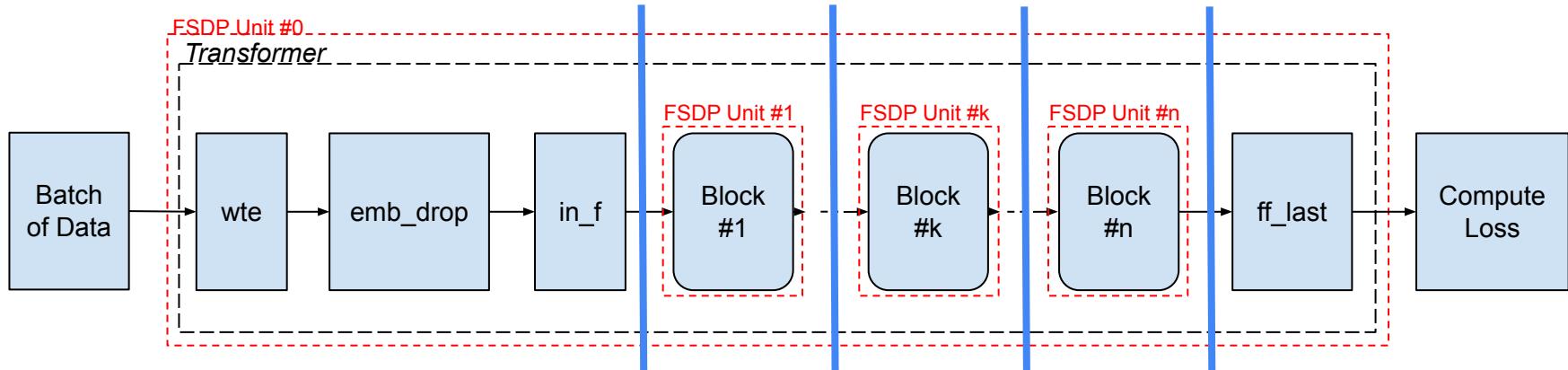
(Each GPUs will only compute their own partial of the gradient)
8. Update optimizer states and weights locally

FSDP: Overlap Computation and Communication



- FSDP fully shards all Parameters, Gradients and Optimizer states across the GPUs.
- Each All-gather, Forward pass, All-gather, Backward pass, Reduce-scatter, Optimizer and Weights update needs to be done sequentially.
 - No opportunity for overlapping computation and communication
 - Needs sort of dividing the model vertically into multiple subsets (aka units) to make this overlap possible

FSDP: Overlap Computation and Communication



- All-gathers and perform forward/backward pass is performed unit by unit
 - Helps with memory
 - Loads parameters only for the current unit
 - Needs to have enough memory to load the largest FSDP unit
 - Provide the computation and communication overlap
 - While unit #1 is performing forward pass, unit #2 all-gathers its parameters

FSDP: Overlap Computation and Communication

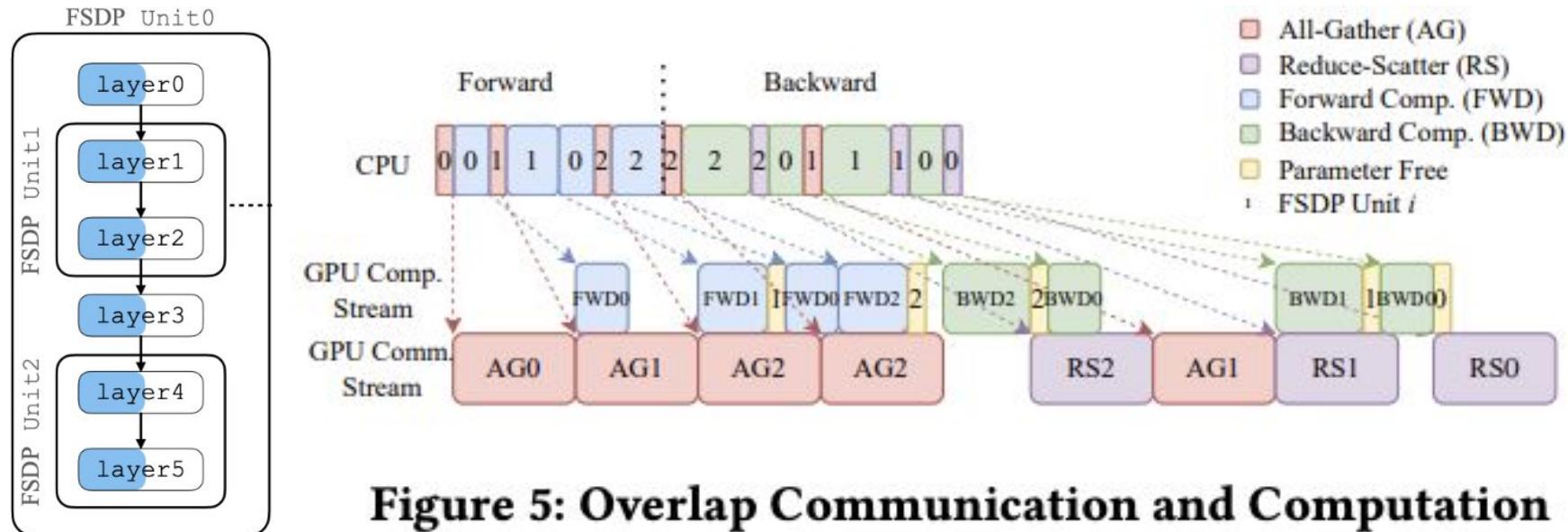
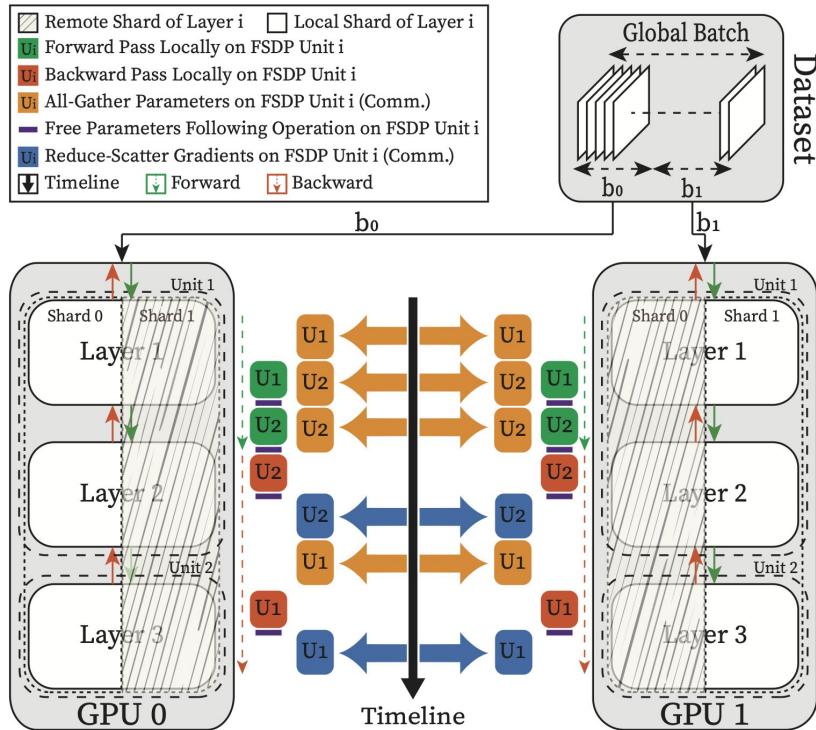
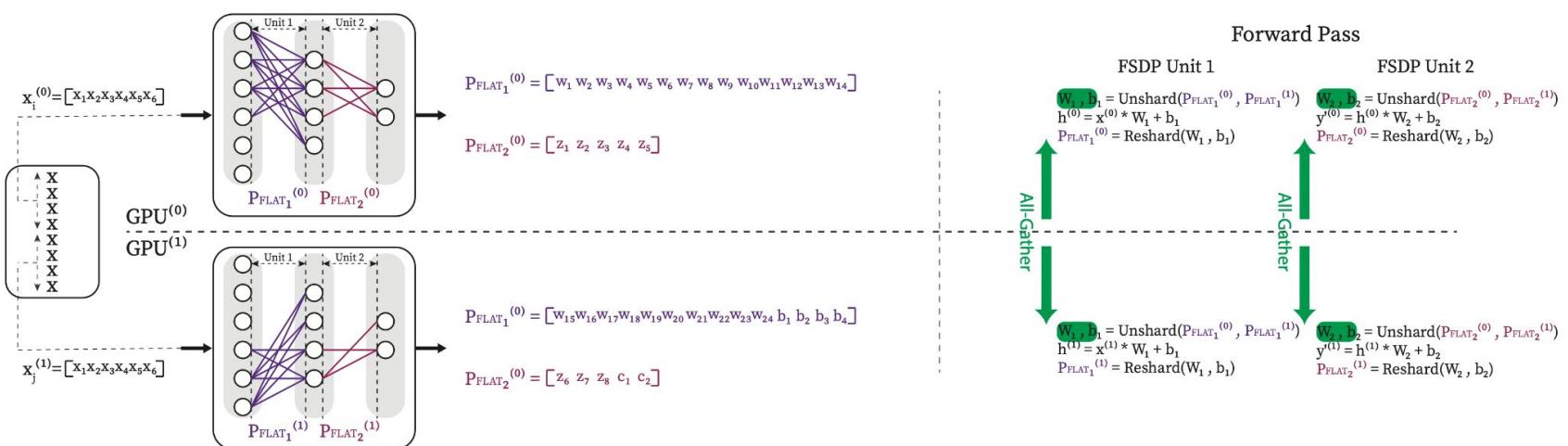


Figure 5: Overlap Computation and Communication

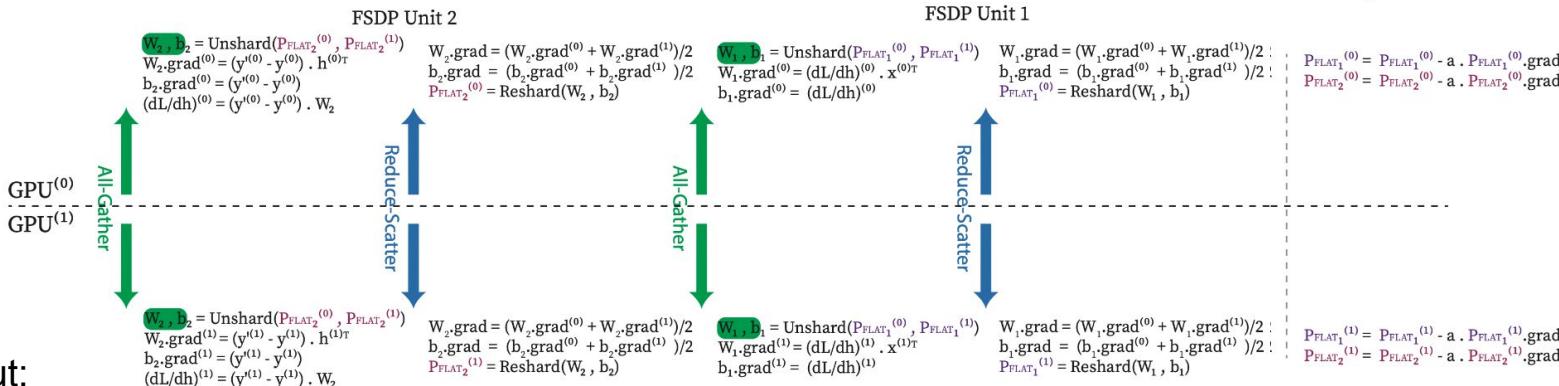
FSDP: Overlap Computation and Communication



https://handbook.eng.kempnerinstitute.harvard.edu/s5_ai_scaling_and_engineering/scalability/gpu_computing.html



Backward Pass



Try it out:

https://github.com/KempnerInstitute/examples/blob/main/distributed-mlp/scripts/mlp_fsdp.py

Agenda

- 1 Why Going Distributed?
- 2 Intro to Distributed GPU Computing
- 3 Different Distributed LLM Training Techniques
- 4 Training a Large LLM on the Cluster

OLMo: 1B - DDP

Distributed Data Parallelism

model:

d_model (D): 2048

n_layers (N): 16

n_heads (N_h): 16

max_sequence_length (T): 2048

vocab_size (V): 32100

ddp:

grad_sync_mode: batch

find_unused_params: false

distributed_strategy: ddp

$$P = N(12D^2 + 10D) + D(V+T) + 2D + V(D+1)$$

$$P = 0.944 \text{ B}$$

https://github.com/KempnerInstitute/OLMo/blob/main/configs/kempner_institute/1b_Olmo.yaml

Instruction:

https://github.com/KempnerInstitute/OLMo/blob/main/README_KempnerInstitute.md

OLMo: 7B - FSDP

Fully Sharded Data Parallelism

model:

d_model (D): 4096

n_layers (N): 32

n_heads (N_h): 32

max_sequence_length (T): 2048

vocab_size (V): 32100

fsdp:

wrapping_strategy: by_block

precision: mixed

sharding_strategy: FULL_SHARD

distributed_strategy: fsdp

$$P = N(12D^2 + 10D) + D(V+T) + 2D + V(D+1)$$

$$P = 6.715 \text{ B}$$

https://github.com/KempnerInstitute/OLMo/blob/main/configs/kempner_institute/7b_Olmo.yaml

Instruction:

https://github.com/KempnerInstitute/OLMo/blob/main/README_KempnerInstitute.md

Exercise: Try out running OLMo (Steps 1 and 2)

1) Install OLMo if you haven't already

(instructions: http://github.com/KempnerInstitute/OLMo/blob/main/README_KempnerInstitute.md)

2) Look through the config file for 7B FSDP

(Section 2.1: https://github.com/KempnerInstitute/OLMo/blob/main/README_KempnerInstitute.md)

Different Wrapping Policy

How to wrap modules into FSDP units

- **by_block**: put each OLMo block into its own FSDP unit.
- **by_block_and_size**: same as by_block but `wte` and `ff_out` will be wrapped separately.
- **by_block_group**: put each m OLMo blocks into its own FSDP unit.
- **size_based**: uses PyTorch's default size-based auto wrap policy. (Wraps any module above 100M size in its own FSDP unit)

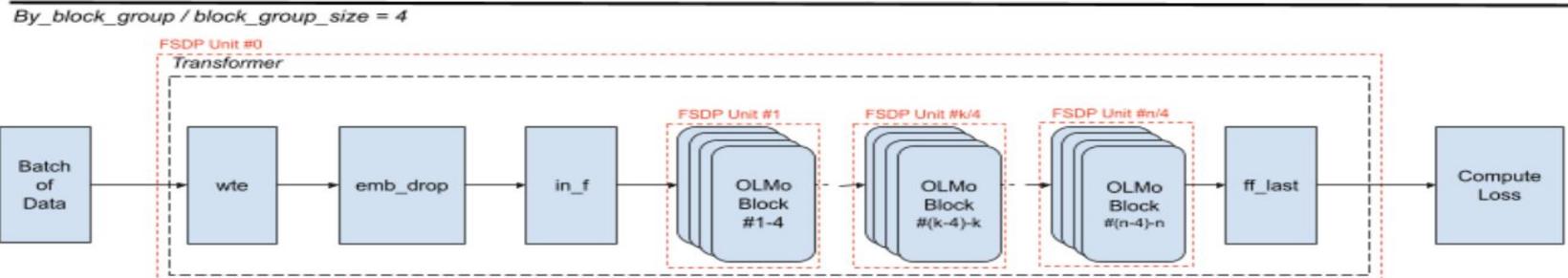
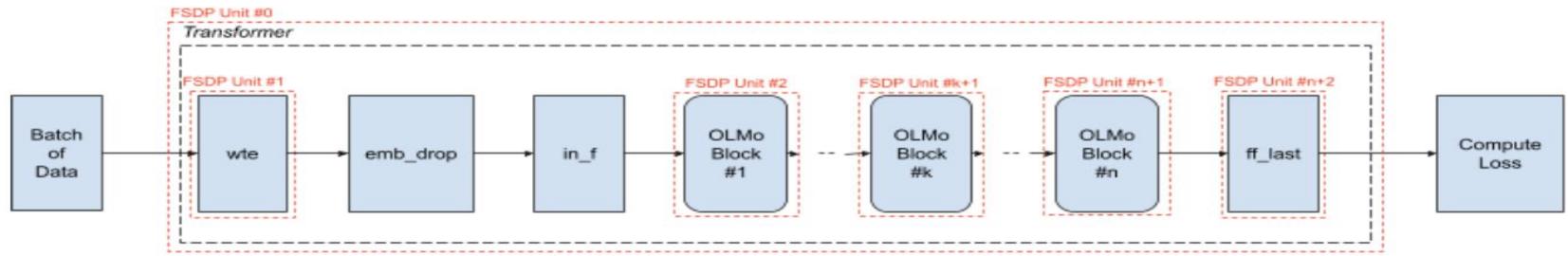
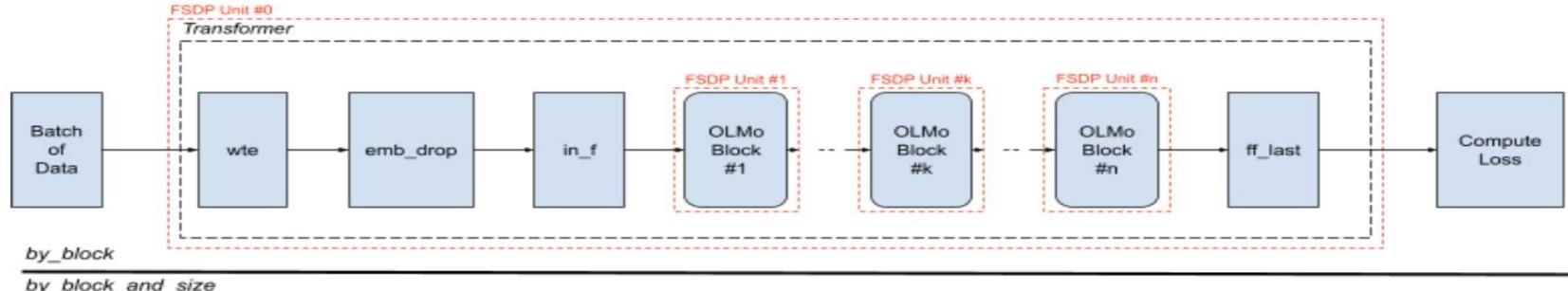


Figure 3: This shows how OLMo's FSDP wrapping policies form the FSDP units.

Different Sharding Strategies

May trade off memory saving and communication overhead

- **NO_SHARD**: as same as DDP
- **FULL_SHARD**: Parameters, gradients, and optimizer states are sharded.
- **SHARD_GRAD_OP**: Gradients and optimizer states are sharded during computation, and additionally, parameters are sharded outside computation meaning it keeps parameters unshared throughout the forward and backward computation.
- **HYBRID_SHARD**: Apply ``**FULL_SHARD**`` within a node, and replicate parameters across nodes. This results in reduced communication volume as expensive all-gathers and reduce-scatters are only done within a node, which can be more performant for medium-sized models.

Exercise: Try out running OLMo (Steps 3 and 4)

3) Update the [SLURM script](#). Pair up and run OLMo with FSDP! Check out the output logs

- You can use `tail -f <path-to-output-file>` to monitor the logs.

4) Try out different wrapping policies and sharding strategies.



Kempner
INSTITUTE



HARVARD
UNIVERSITY

Thank you

