

---

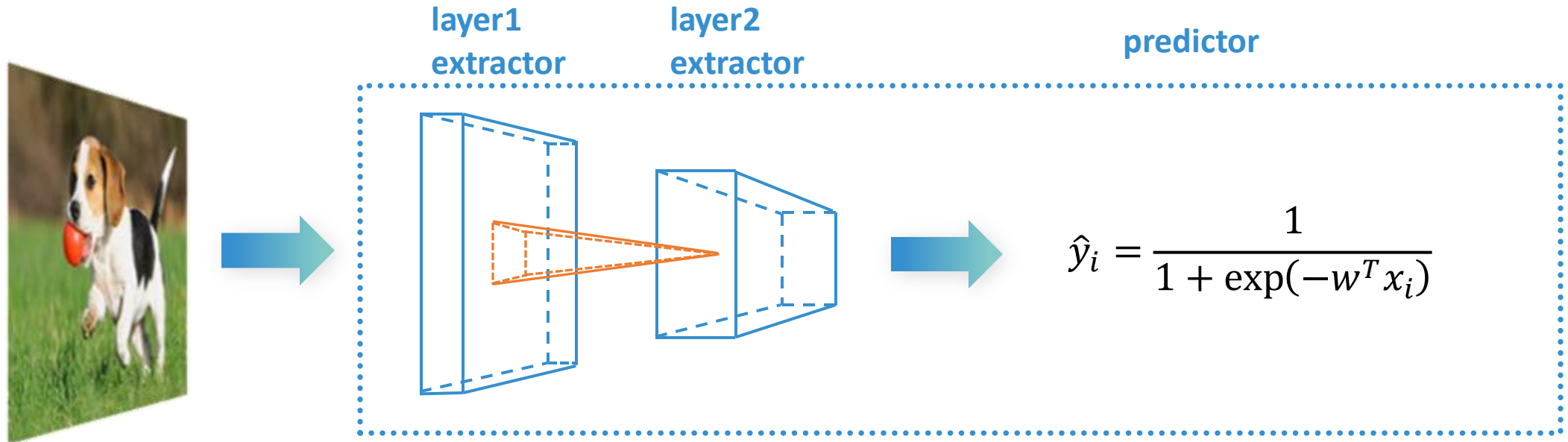
System for Artificial Intelligence

# Parallelization and Training I

Siyuan Feng  
Shanghai Innovation Institute

---

# Recap: DNN Training Overview



## Objective

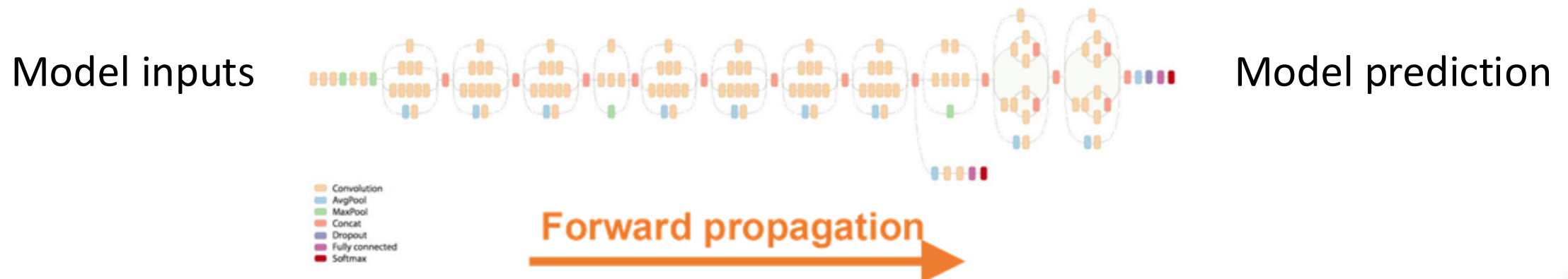
$$L(w) = \sum_{i=1}^n l(y_i, \hat{y}_i) + \lambda ||w||^2$$

## Training

$$w \leftarrow w - \eta \nabla_w L(w)$$

# DNN Training Process

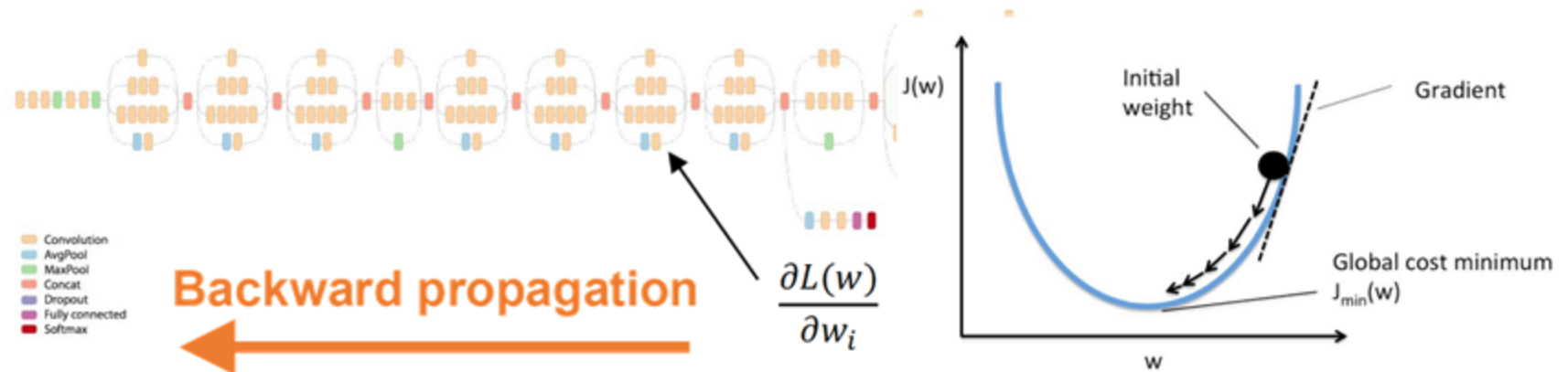
- Train ML models through many iterations of 3 stages
  - 1. Forward propagation:** apply model to a batch of input samples and run calculation through operators to produce a prediction
  - 2. Backward propagation:** run the model in reverse to produce error for each trainable weight
  - 3. Weight update:** use the loss value to update model weights



# DNN Training Process

- Train ML models through many iterations of 3 stages
  - 1. Forward propagation:** apply model to a batch of input samples and run calculation through operators to produce a prediction
  - 2. Backward propagation:** run the model in reverse to produce error for each trainable weight
  - 3. Weight update:** use the loss value to update model weights

Model inputs



# DNN Training Process

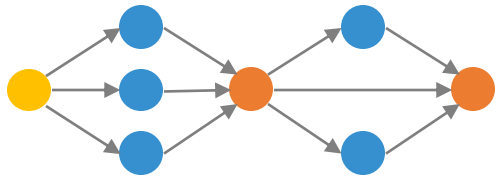
- Train ML models through many iterations of 3 stages
  - 1. Forward propagation:** apply model to a batch of input samples and run calculation through operators to produce a prediction
  - 2. Backward propagation:** run the model in reverse to produce error for each trainable weight
  - 3. Weight update:** use the loss value to update model weights

$$w_i := w_i - \gamma \frac{\partial L(w)}{\partial w_i} = w_i - \frac{\gamma}{n} \sum_{j=1}^n \boxed{\frac{\partial l_i(w)}{\partial w_i}} \quad \text{Gradients of individual samples}$$

# How can we parallelize DNN training?

$$w_i := w_i - \gamma \nabla L(w_i) = w_i - \frac{\gamma}{n} \sum_{j=1}^n \nabla L_j(w_i)$$

# Data Parallelism

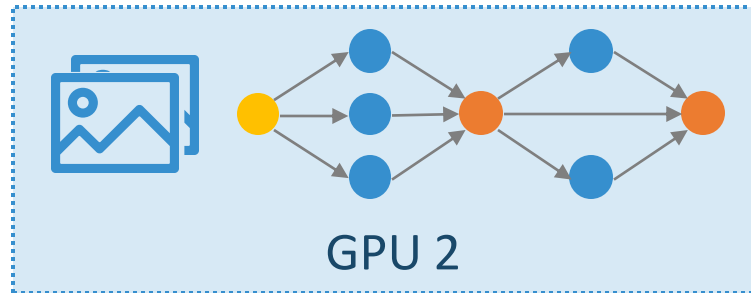
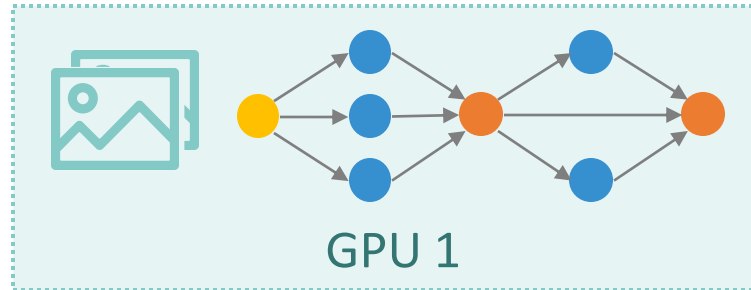
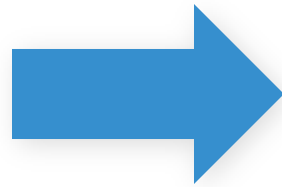


ML Model

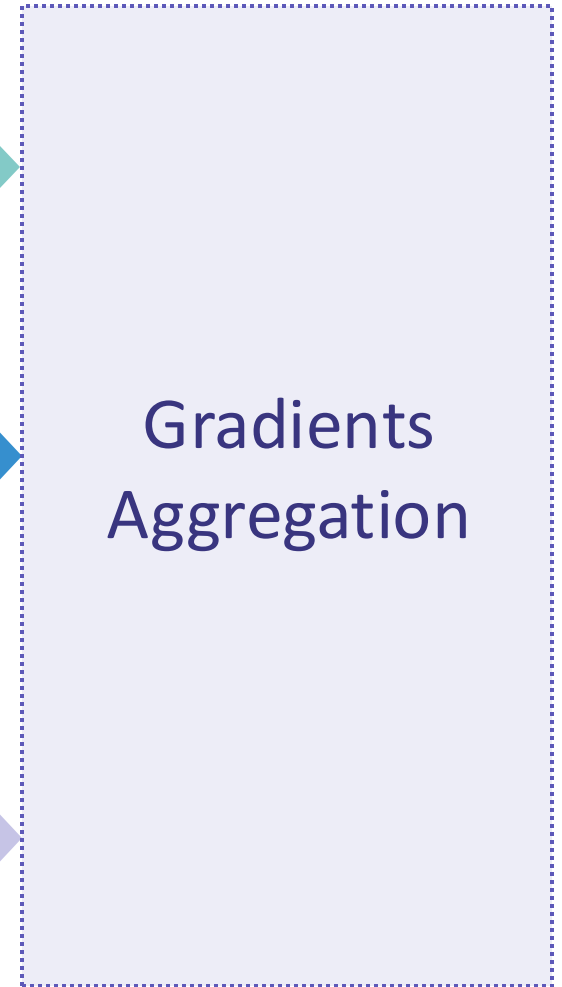
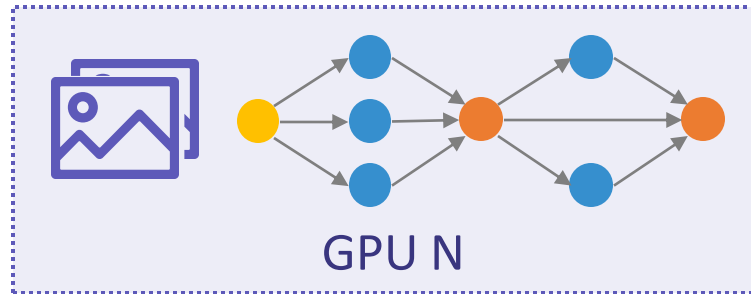


Training Dataset

$$w_i := w_i - \gamma \nabla L(w_i) = w_i - \frac{\gamma}{n} \sum_{j=1}^n \nabla L_j(w_i)$$



...

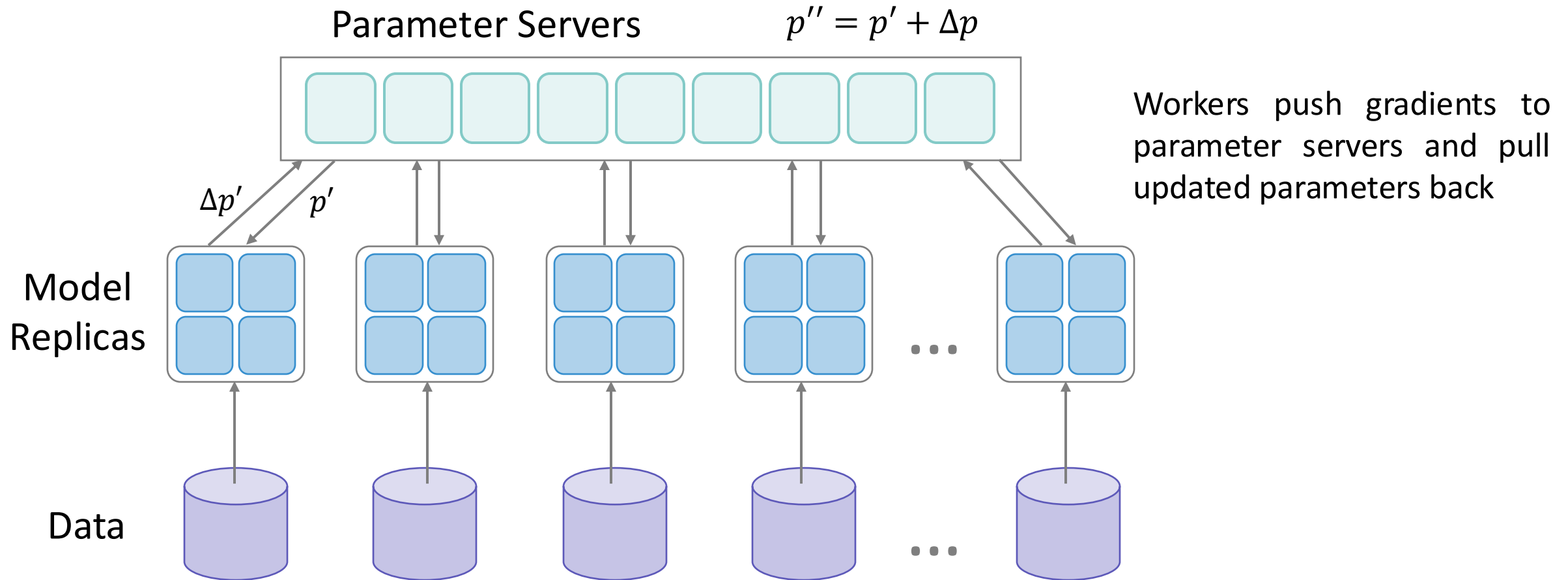


1. Partition training data into batches

2. Compute the gradients of each batch on a GPU

3. Aggregate gradients across GPUs

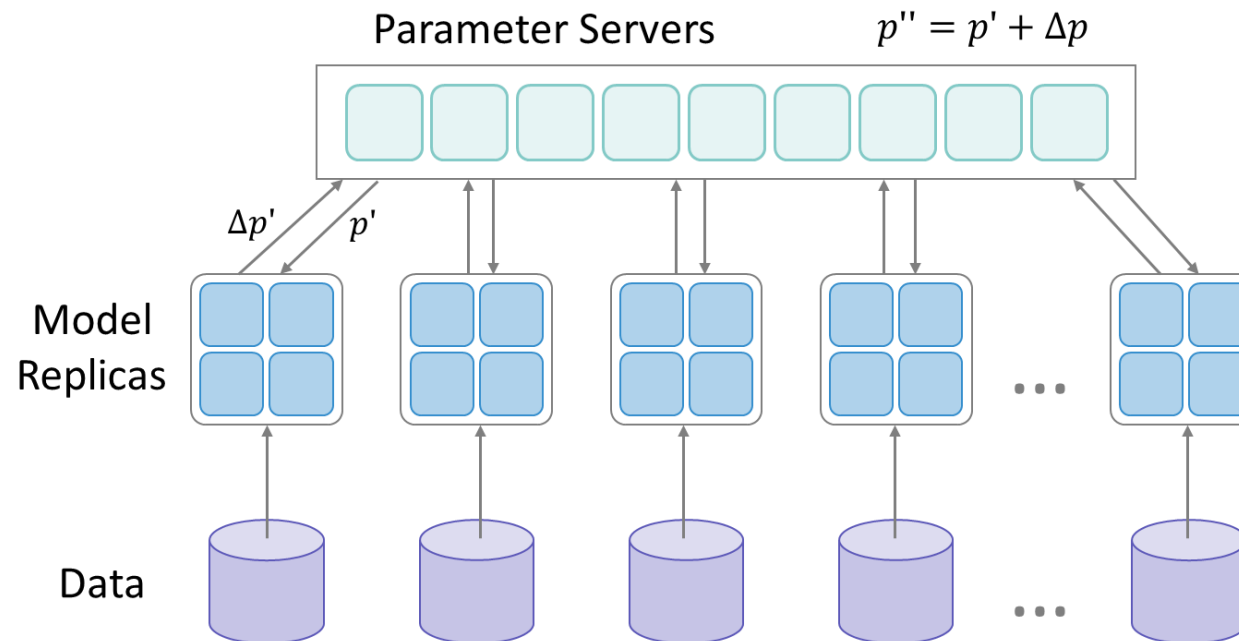
# Data Parallelism: Parameter Server





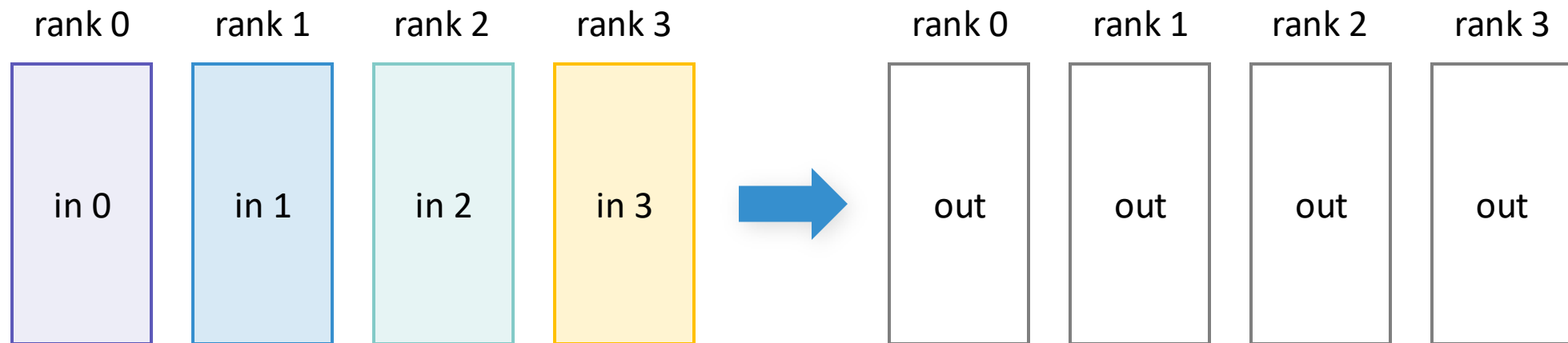
# Inefficiency of Parameter Server

- **Centralized communication:** all workers communicate with parameter servers for weights update; cannot scale to large numbers of workers
- How can we decentralize communication in DNN training?



# Inefficiency of Parameter Server

- **Centralized communication:** all workers communicate with parameter servers for weights update; cannot scale to large numbers of workers
- How can we decentralize communication in DNN training?
- **All-Reduce:** perform element-wise reduction across multiple devices



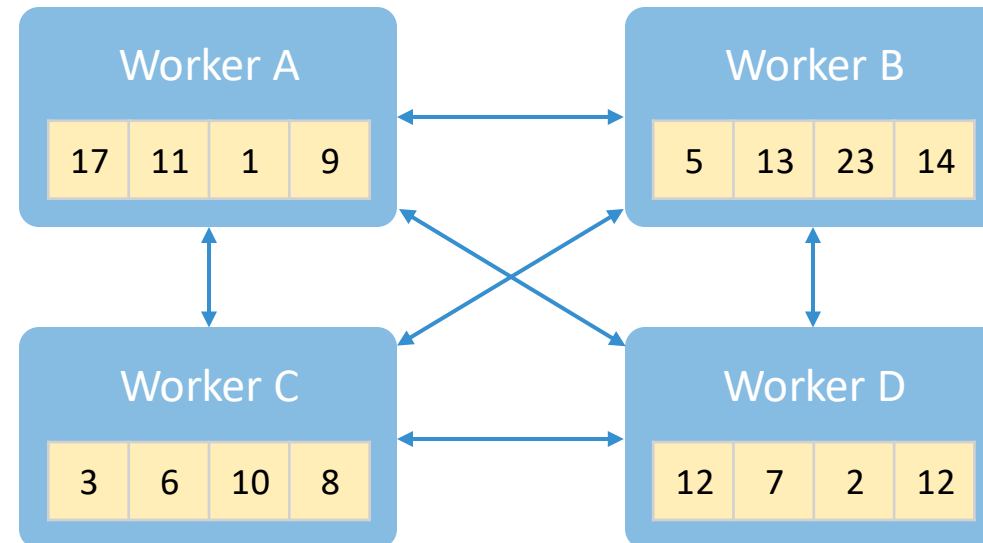
$$\text{out}[i] = \text{sum}(\text{in } X[i])$$

# Different Ways to Perform All-Reduce

- Naïve All-Reduce
- Ring All-Reduce
- Tree All-Reduce
- Butterfly All-Reduce

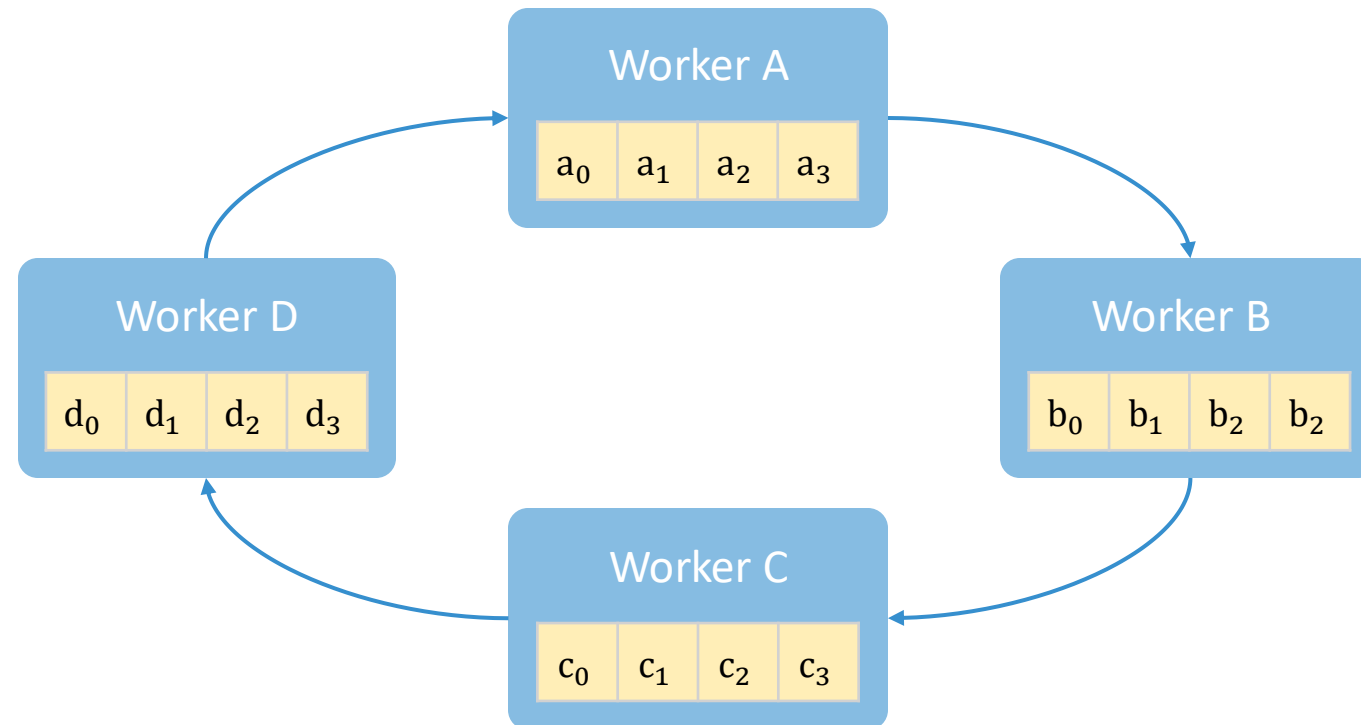
## Naïve All-Reduce

- Each worker can send its local gradients to all other workers
- If we have  $N$  workers and each worker contains  $M$  parameters
- Overall communication:  $N * (N-1) * M$  parameters
- **Issue:** each worker communicates with all other workers; same scalability issue as parameter server



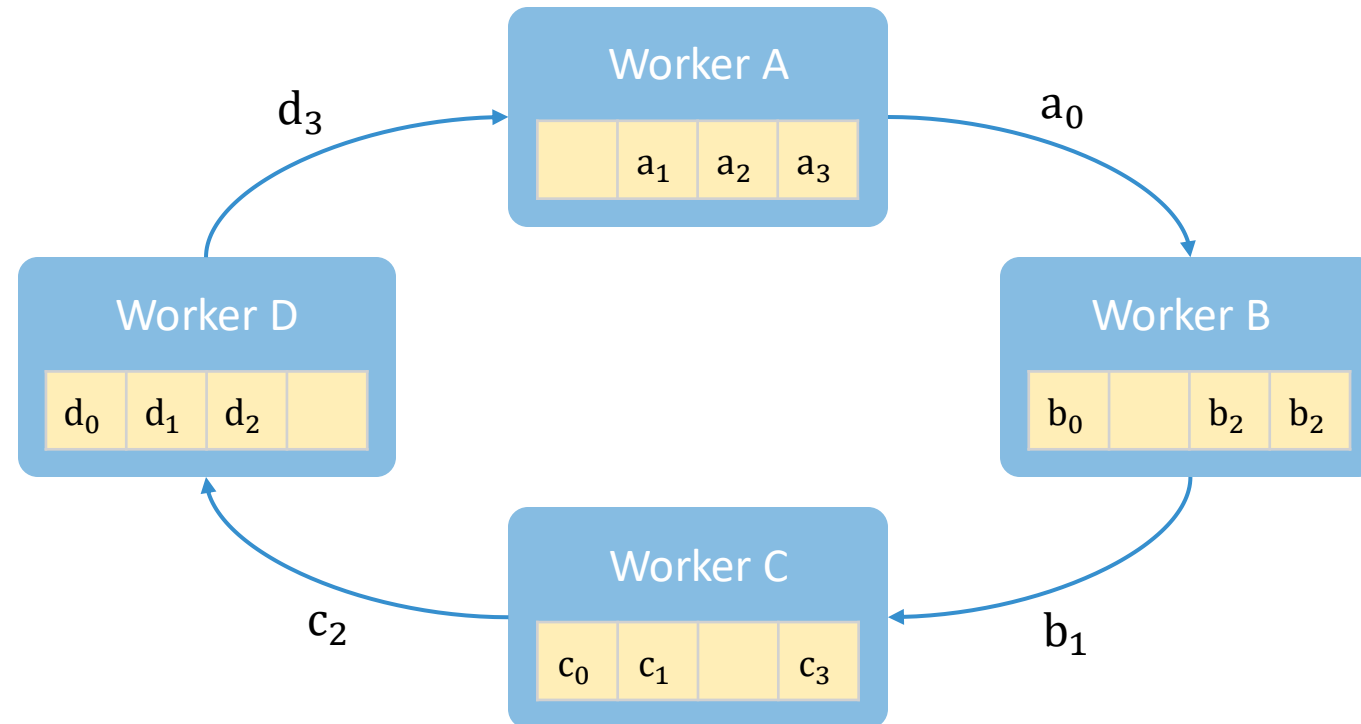
# Ring All-Reduce

- Construct a ring of  $N$  workers, divide  $M$  parameters into  $N$  slices
- Step 1 (Aggregation): each worker send one slice ( $M/N$  parameters) to the next worker on the ring; repeat  $N$  times



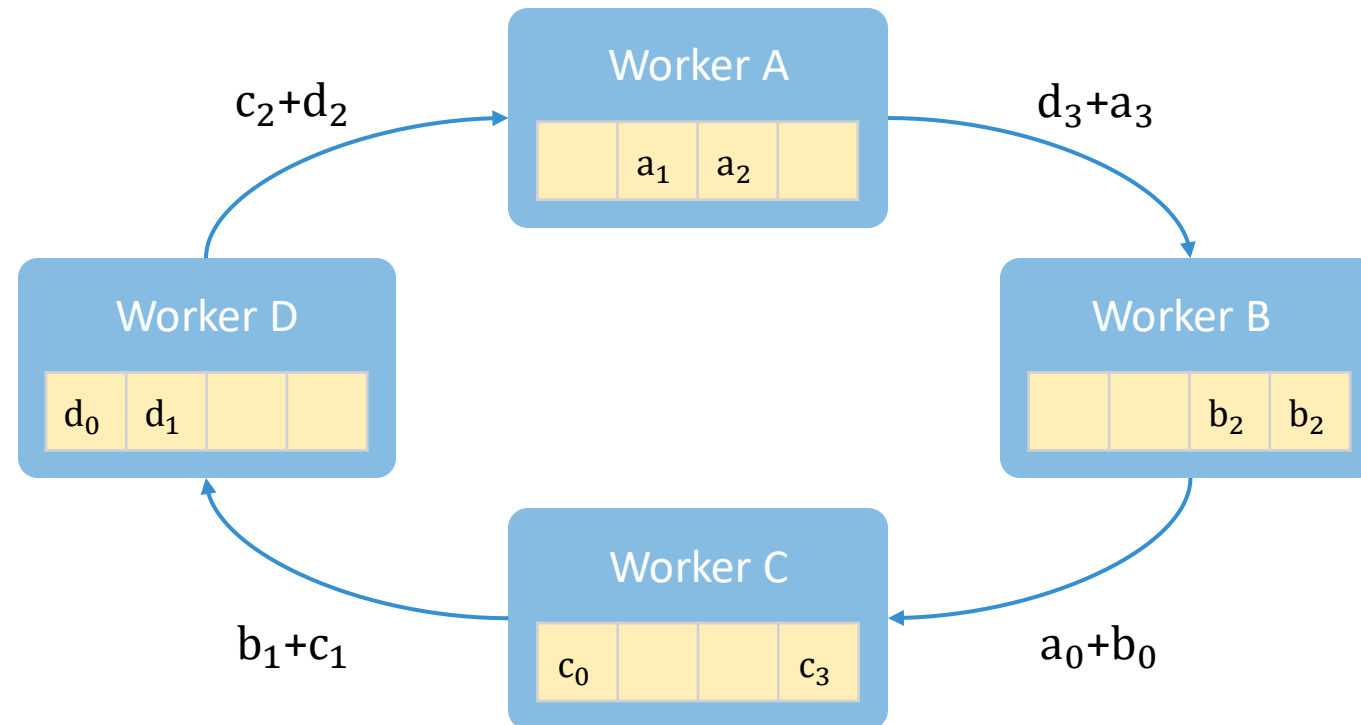
# Ring All-Reduce

- Construct a ring of  $N$  workers, divide  $M$  parameters into  $N$  slices
- Step 1 (Aggregation): each worker send one slice ( $M/N$  parameters) to the next worker on the ring; repeat  $N$  times



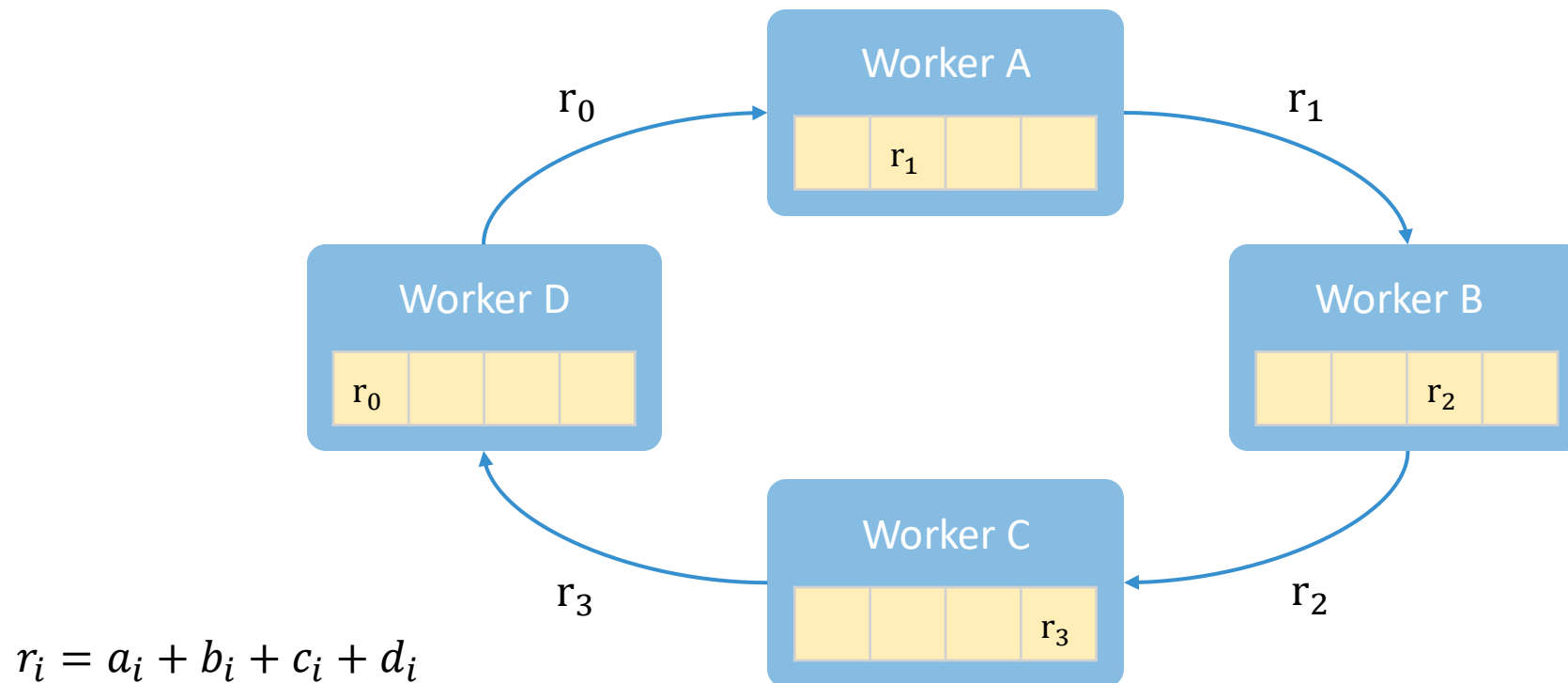
# Ring All-Reduce

- Construct a ring of  $N$  workers, divide  $M$  parameters into  $N$  slices
- Step 1 (Aggregation): each worker send one slice ( $M/N$  parameters) to the next worker on the ring; repeat  $N$  times



# Ring All-Reduce

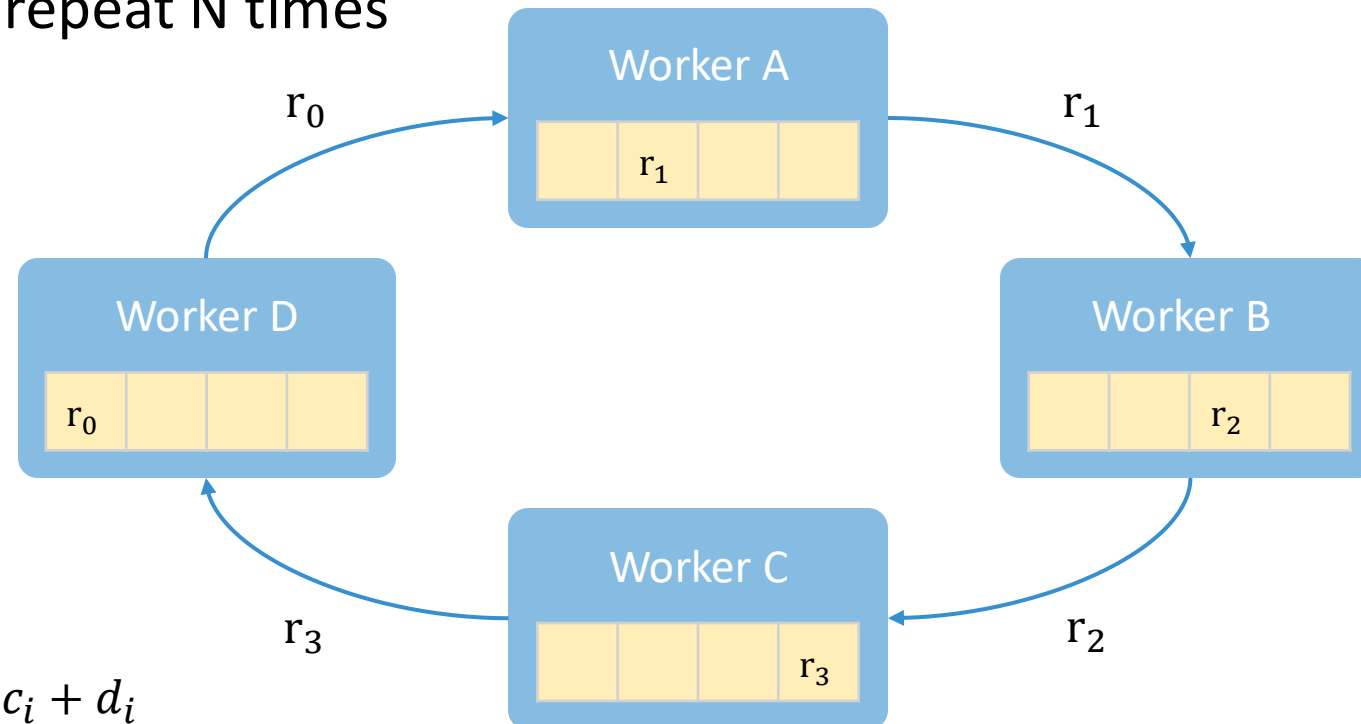
- Construct a ring of N workers, divide M parameters into N slices
- Step 1 (Aggregation): each worker send one slice (M/N parameters) to the next worker on the ring; repeat N times
- After step 1, each worker has the aggregated version of M/N parameters





# Ring All-Reduce

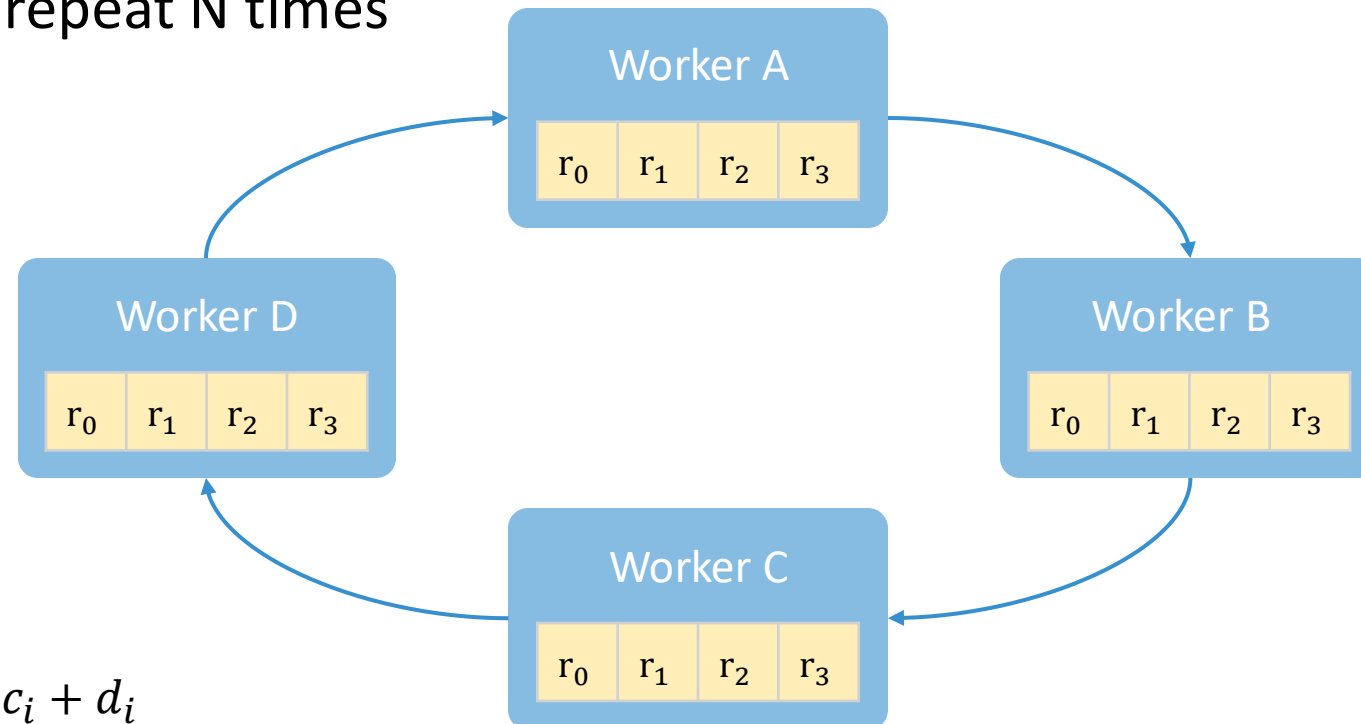
- Construct a ring of N workers, divide M parameters into N slices
- Step 1 (Aggregation): each worker send one slice (M/N parameters) to the next worker on the ring; repeat N times
- Step 2 (Broadcast): each worker send one slice of aggregated parameters to the next worker; repeat N times



$$r_i = a_i + b_i + c_i + d_i$$

# Ring All-Reduce

- Construct a ring of N workers, divide M parameters into N slices
- Step 1 (Aggregation): each worker send one slice (M/N parameters) to the next worker on the ring; repeat N times
- Step 2 (Broadcast): each worker send one slice of aggregated parameters to the next worker; repeat N times



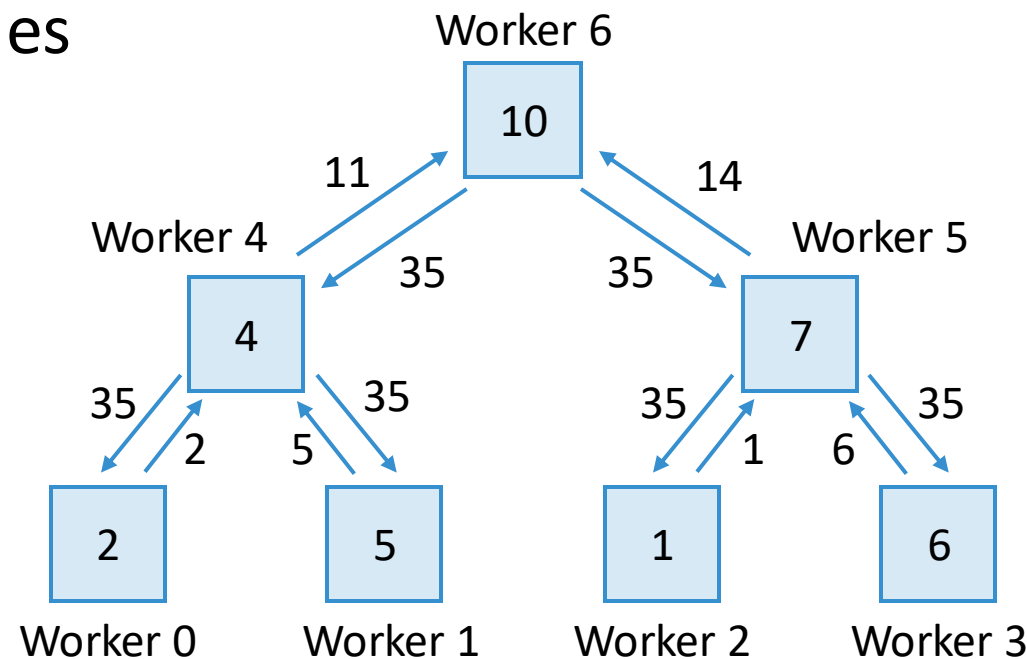
$$r_i = a_i + b_i + c_i + d_i$$

## Ring All-Reduce

- Construct a ring of  $N$  workers, divide  $M$  parameters into  $N$  slices
- Step 1 (Aggregation): each worker send one slice ( $M/N$  parameters) to the next worker on the ring; repeat  $N$  times
- Step 2 (Broadcast): each worker send one slice of aggregated parameters to the next worker; repeat  $N$  times
- Overall communication:  $2 * M * N$  parameters
  - Aggregation:  $M * N$  parameters
  - Broadcast:  $M * N$  parameters

## Tree All-Reduce

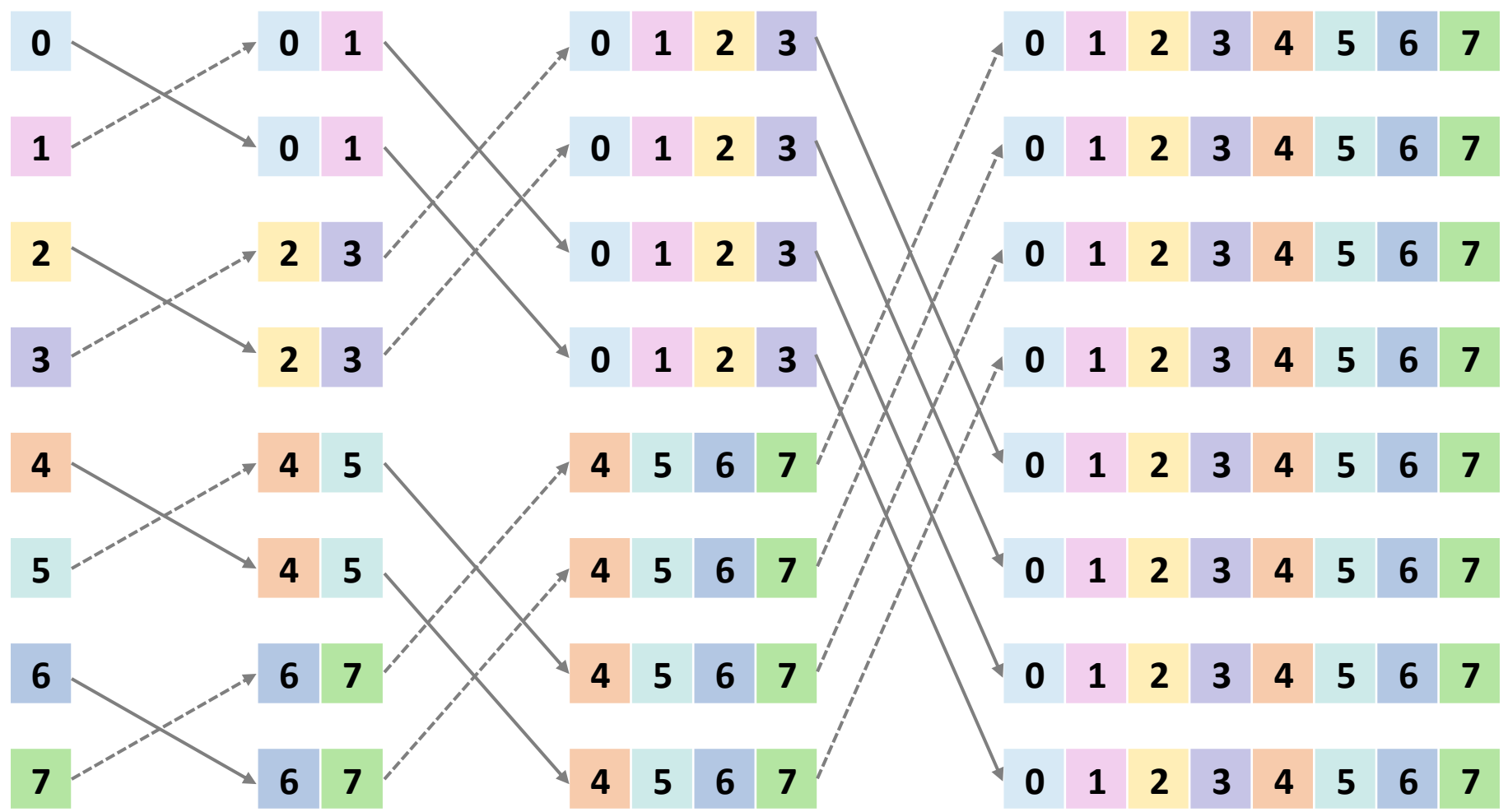
- Construct a tree of  $N$  workers;
- Step 1 (Aggregation): each worker sends  $M$  parameters to its parent; repeat  $\log(N)$  times
- Step 2 (Broadcast): each worker sends  $M$  parameters to its children; repeat  $\log(N)$  times



## Tree All-Reduce

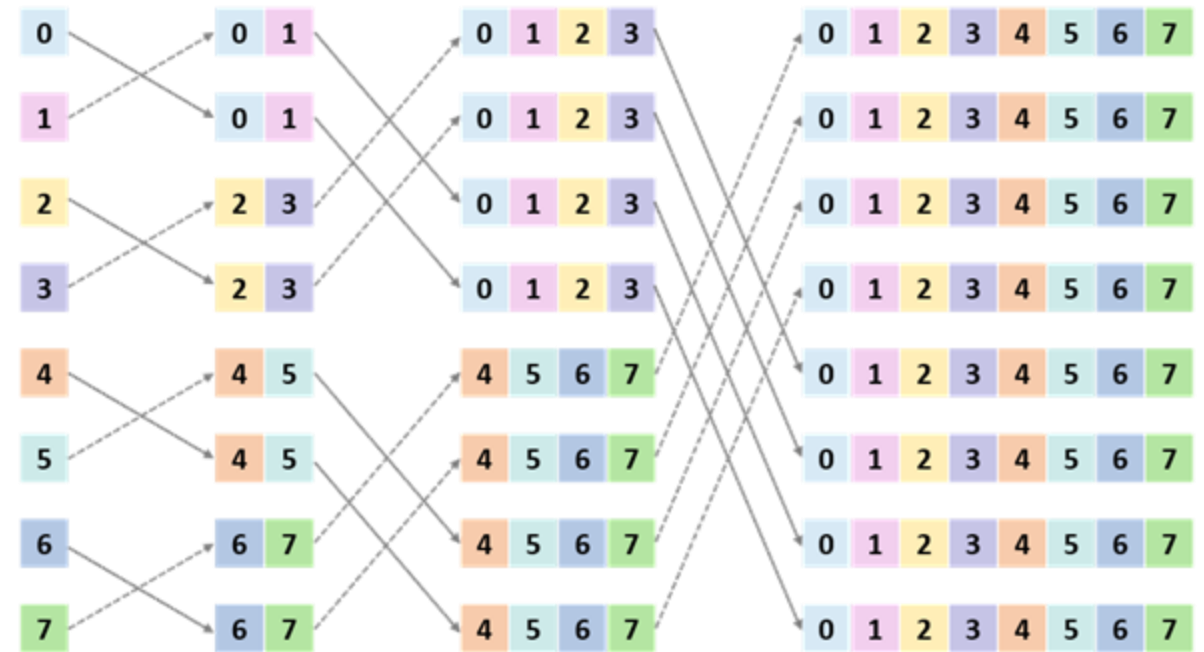
- Construct a tree of  $N$  workers;
- Step 1 (Aggregation): each worker sends  $M$  parameters to its parent; repeat  $\log(N)$  times
- Step 2 (Broadcast): each worker sends  $M$  parameters to its children; repeat  $\log(N)$  times
- Overall communication:  $2 * N * M$  parameters
  - Aggregation:  $M * N$  parameters
  - Broadcast:  $M * N$  parameters

# Butterfly Network



# Butterfly All-Reduce

- Repeat  $\log(N)$  times:
  1. Each worker sends  $M$  parameters to its target node in the butterfly network
  2. Each worker aggregates gradients locally
- Overall communication:  $N * M * \log(N)$  parameters



# Comparing different All-Reduce Methods

	Parameter Server	Naïve All-Reduce	Ring All-Reduce	Tree All-Reduce	Butterfly All-Reduce
Overall communication	$2 \times N \times M$	$N^2 \times M$	$2 \times N \times M$	$2 \times N \times M$	$N \times M \times \log N$

- Question: Ring All-Reduce is more efficient and scalable than Tree All-Reduce and Parameter Server, why?

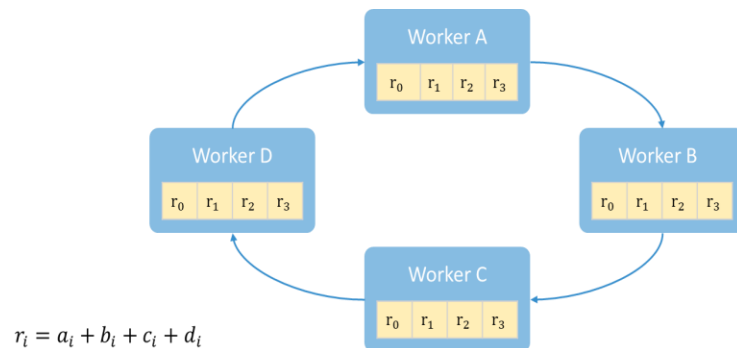


# Ring v.s. Tree v.s. Parameter Server

Ring All-Reduce:

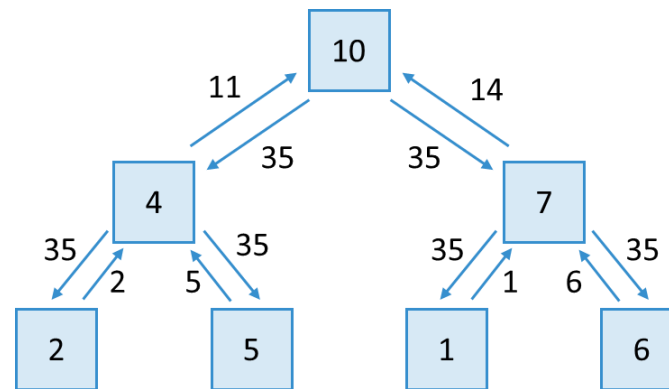
- Best latency
- Balanced workload across workers
- More scalable since each worker

sends  $2 \cdot M$  parameters (independent to the number of workers)



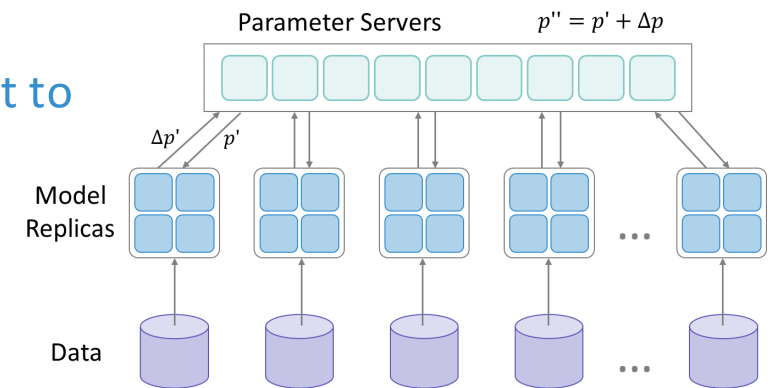
Each worker sends  $M/N$  parameters per iteration; repeat for  $2 \cdot N$  iterations

**Latency:**  $M/N \cdot (2 \cdot N) / \text{bandwidth}$



Each worker sends  $M$  parameters per iteration; repeat for  $2 \cdot \log(N)$  iterations

**Latency:**  $M \cdot 2 \cdot \log(N) / \text{bandwidth}$

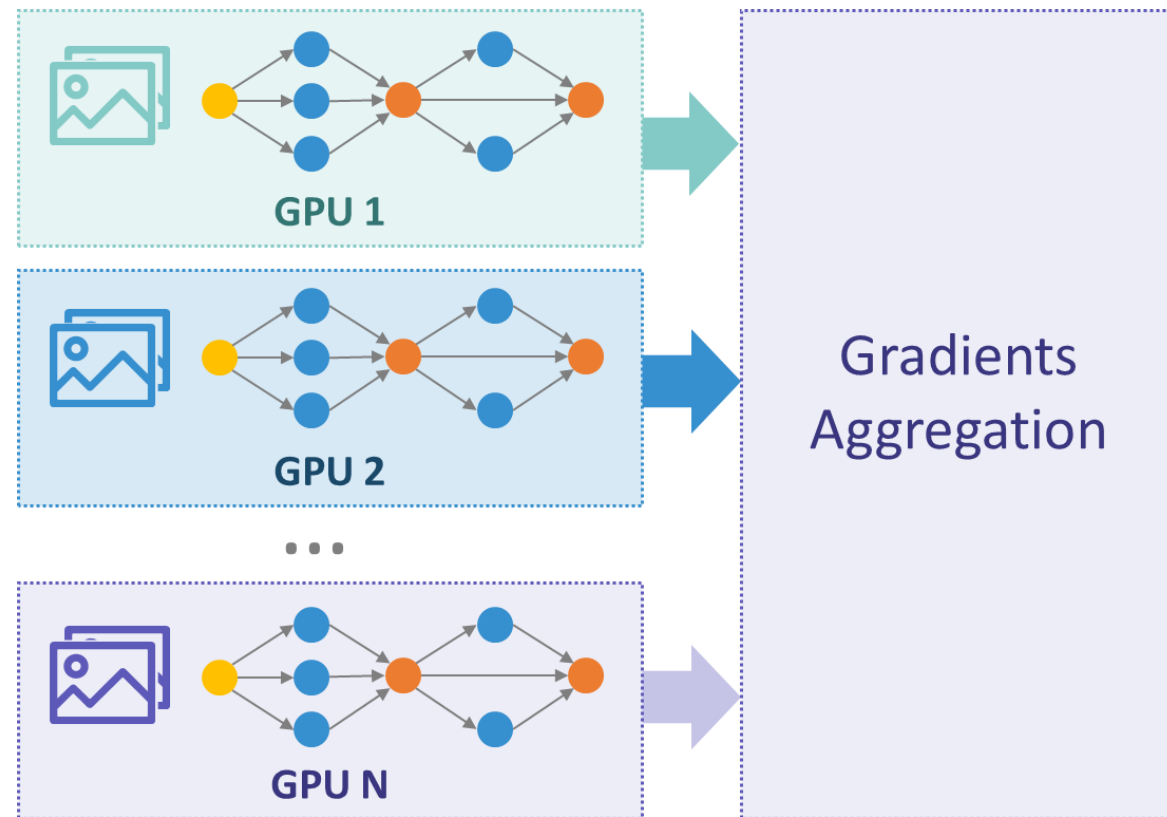


All workers send  $M$  parameters to parameter servers and receive  $M$  parameters from servers

**Latency:**  $M \cdot N / \text{bandwidth}$

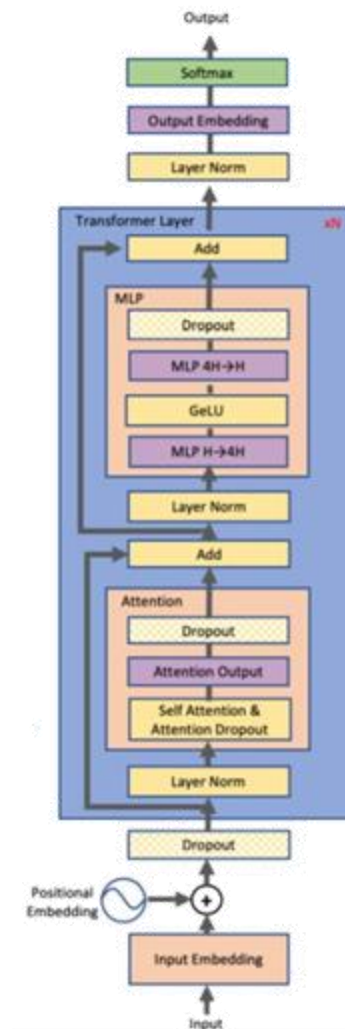
## An Issue with Data Parallelism

- Each GPU saves a replica of the entire model
- Cannot train large models that exceed GPU device memory



# Large Model Training Challenges

	BertLarge	GPT-2	Turing 17.2 NLG	GPT-3
Parameters	0.32B	1.5B	17.2B	175B
Layers	24	48	78	96
Hidden Dimension	1024	1600	4256	12288
Relative Computation	1x	4.7x	<b>54x</b>	<b>547x</b>
Memory Footprint	5.12GB	24GB	<b>275GB</b>	<b>2800GB</b>



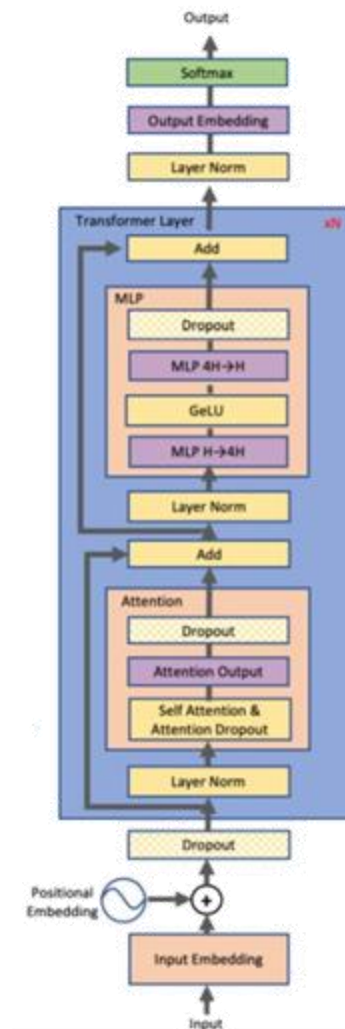
# Large Model Training Challenges

	BertLarge	GPT-2	Turing 17.2 NLG	GPT-3
Parameters	0.32B	1.5B	17.2B	175B
Layers	24	48	78	96
Hidden Dimension	1024	1600	4256	12288
Relative Computation	1x	4.7x	<b>54x</b>	<b>547x</b>
Memory Footprint	5.12GB	24GB	<b>275GB</b>	<b>2800GB</b>

NVIDIA V100 GPU memory capacity: 16G/32G

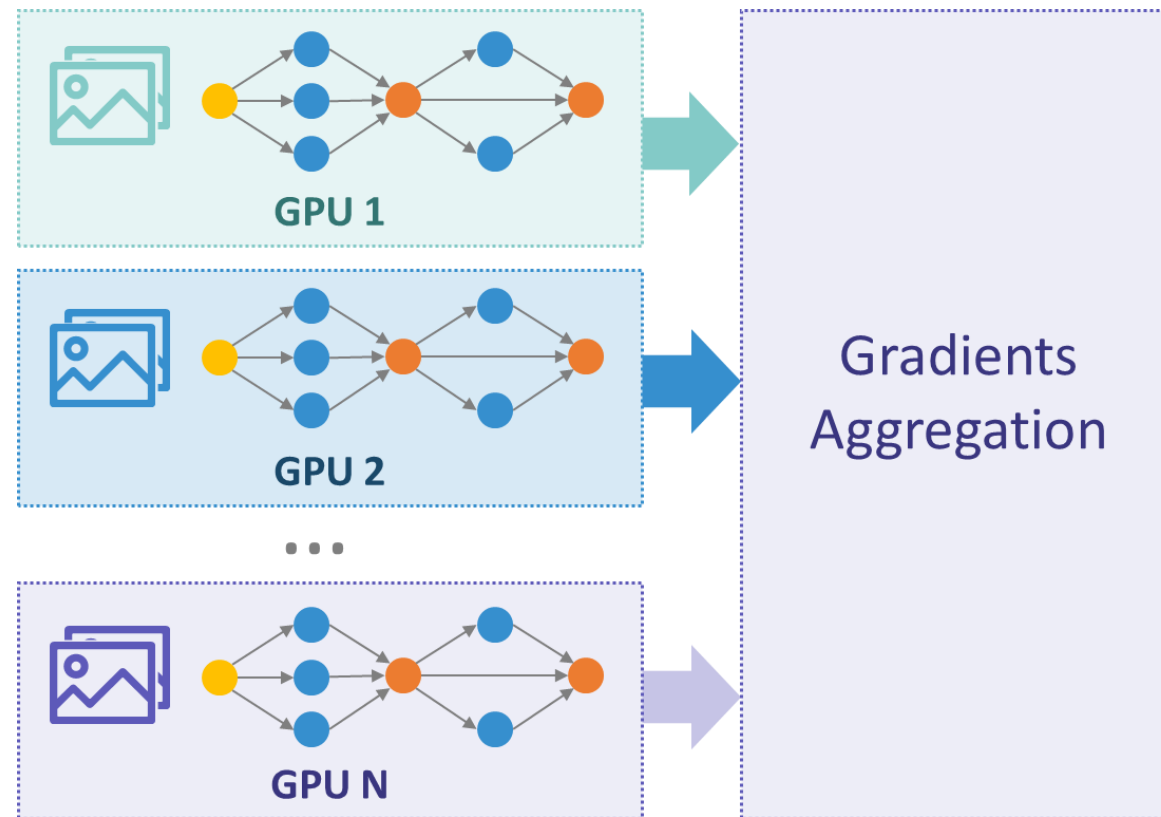
NVIDIA A100 GPU memory capacity: 40G/80G

**Out of Memory**



# ZeRO: Zero Redundancy Optimizer

- Eliminating data redundancy in data parallel training
- A widely used technique for data parallel training of large models



# Revisit: Stochastic Gradient Descent

For  $t = 1$  to  $T$

$$\Delta w = \eta x \frac{1}{b} \sum_{i=1}^b \nabla \left( \text{loss}(f_w(x_i, y_i)) \right) \quad // \text{ compute derivative and update}$$

Backward pass      Forward pass

$w \leftarrow w + \Delta w$  // apply update

End

# Adaptive Learning Rates (Adam)

For  $t = 1$  to  $T$

$$g = \frac{1}{b} \sum_{i=1}^b \nabla \left( \text{loss}(f_w(x_i, y_i)) \right)$$

$$\Delta w = \text{adam}(g)$$

$w \leftarrow w + \Delta w$  // apply update

End

$$\nu_t = \beta_1 * \nu_{t-1} + (1 - \beta_1) * g_t$$

$$s_t = \beta_2 * s_{t-1} + (1 - \beta_2) * g_t^2$$

$$\Delta \omega_t = -\eta \frac{\nu_t}{\sqrt{s_t} + \epsilon} * g_t$$

$g_t$  : Gradient at time  $t$  along  $\omega^j$

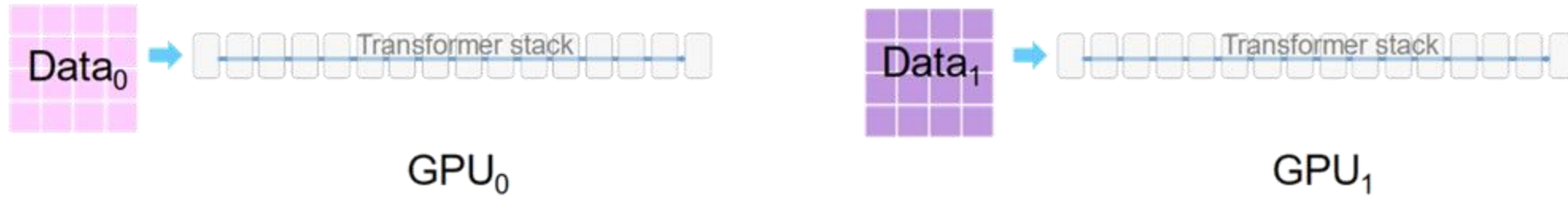
$\nu_t$  : Exponential Average of gradients along  $\omega_j$

$s_t$  : Exponential Average of squares of gradients along  $\omega_j$

$\beta_1, \beta_2$  : Hyperparameters

[1] Kingma and Ba, "Adam: A Method for Stochastic Optimization", 2014,  
<https://arxiv.org/abs/1412.6980>

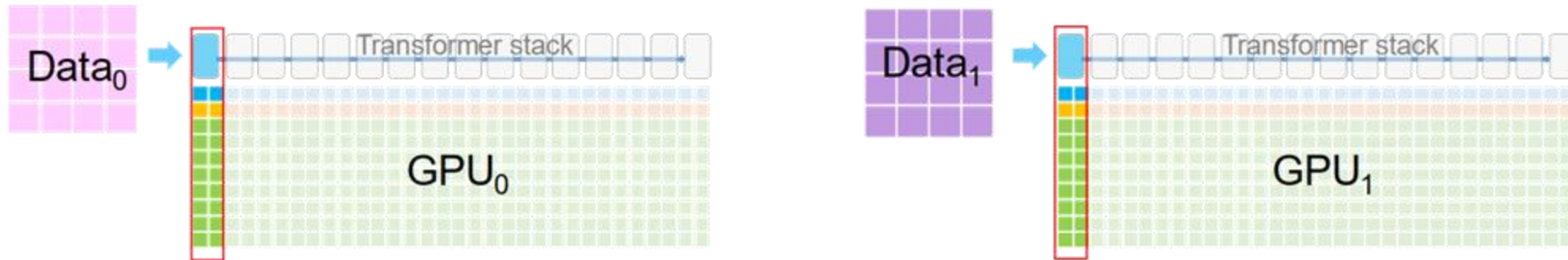
# Understanding Memory Consumption



A 16-layer transformer model  = 1 layer



# Understanding Memory Consumption



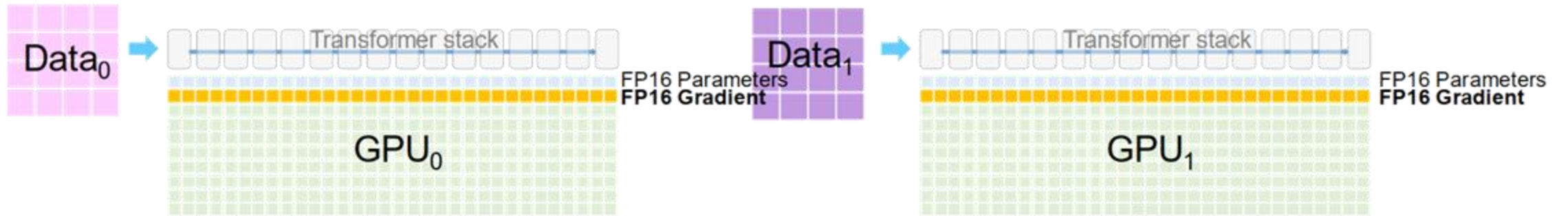
Each cell  represents GPU memory used by its corresponding transformer layer 

# Understanding Memory Consumption



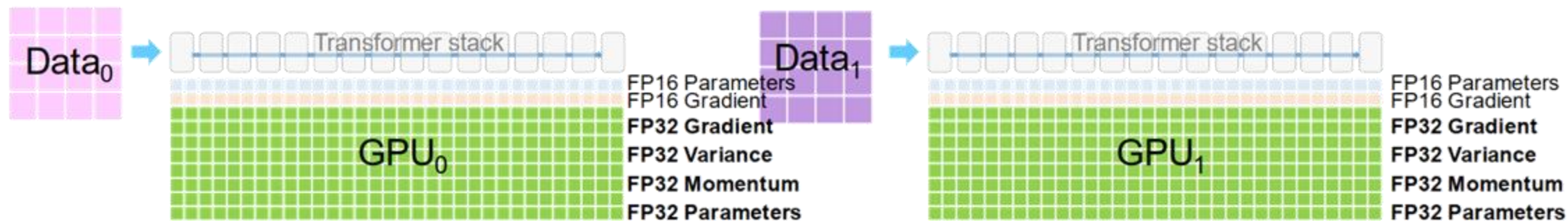
- FP16 parameter

# Understanding Memory Consumption



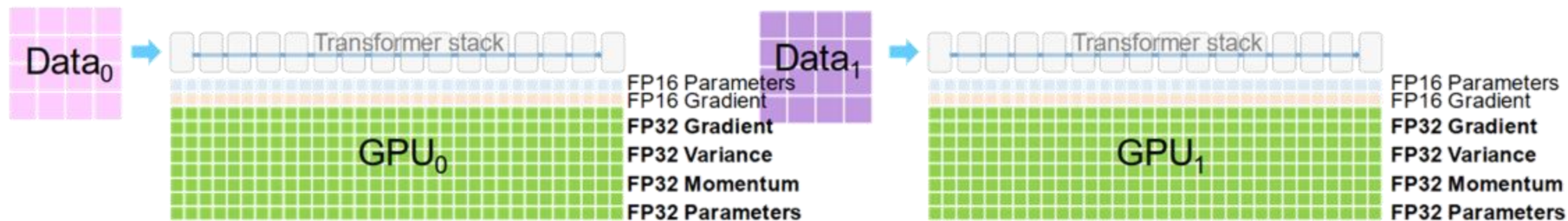
- FP16 parameter
- FP16 Gradients

# Understanding Memory Consumption



- FP16 parameter
- FP16 Gradients
- FP32 Optimizer States
  - Gradients, Variance, Momentum, Parameters

# Understanding Memory Consumption



- FP16 parameter: **2M bytes**
- FP16 Gradients: **2M bytes**
- FP32 Optimizer States : **16M bytes**
  - Gradients, Variance, Momentum, Parameters

Example 1B parameter model -> 20GB/GPU

Memory consumption doesn't include:

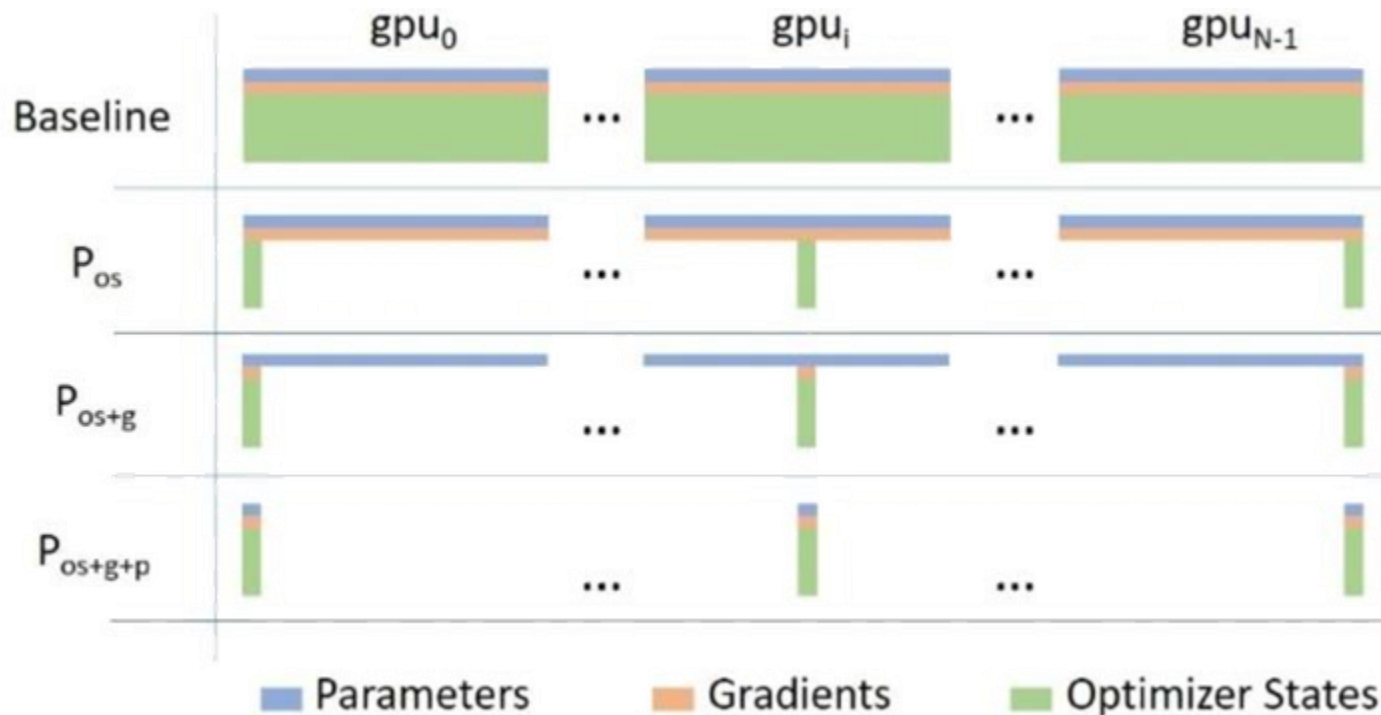
- Input batch + activations

M = number of parameters in the model

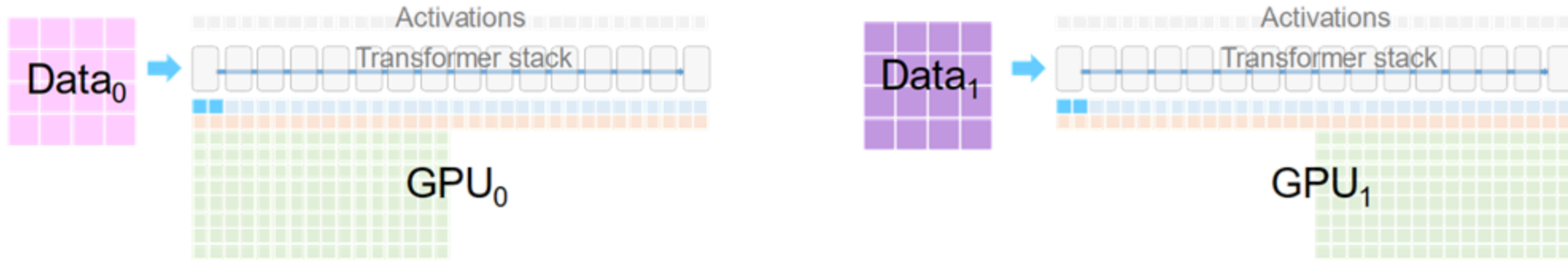


# ZeRO-DP: ZeRO powered Data Parallelism

- ZeRO removes the redundancy across data parallel process
- Stage 1: partitioning optimizer states
- Stage 2: partitioning gradients
- Stage 3: partitioning parameters

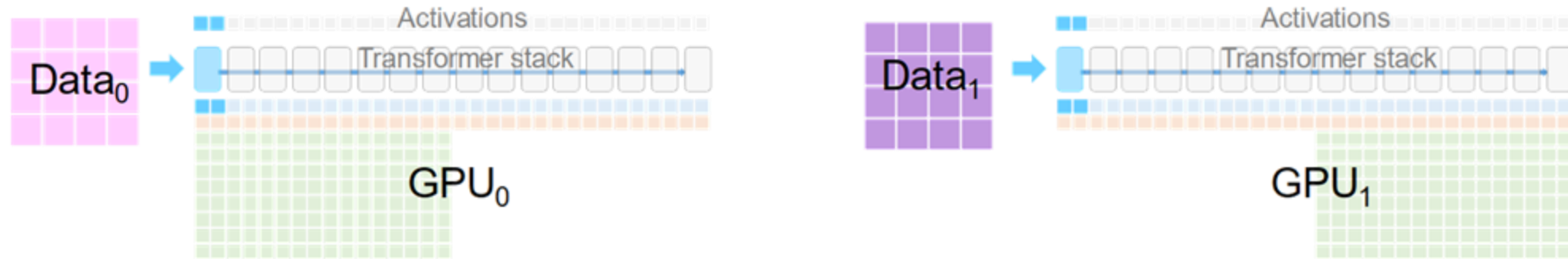


# ZeRO Stage 1: Partitioning Optimizer States



- ZeRO Stage 1
- Partitions optimizer states across GPUs

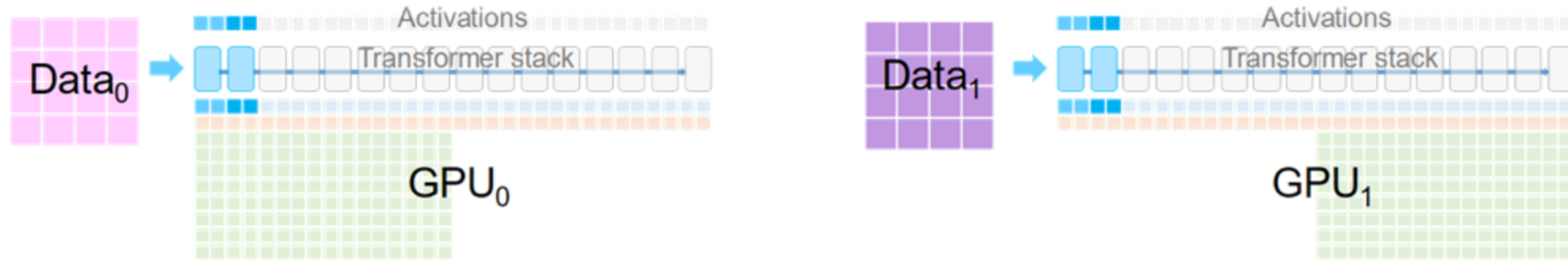
# ZeRO Stage 1: Partitioning Optimizer States



- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks

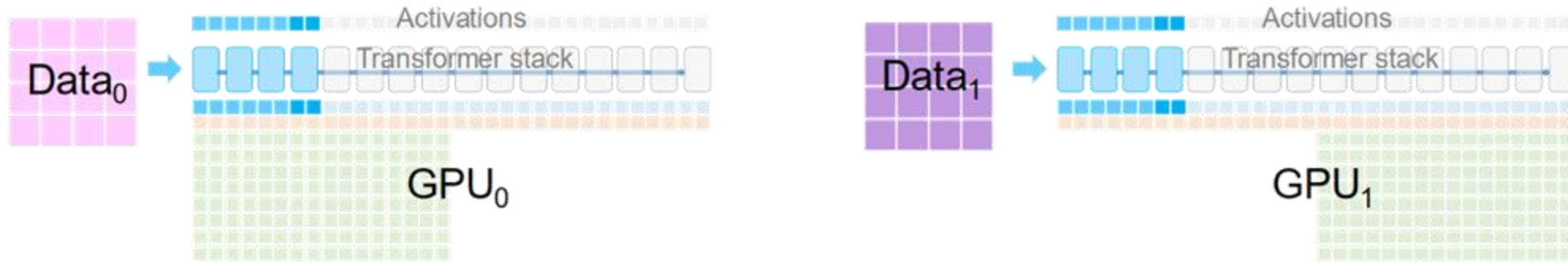


# ZeRO Stage 1: Partitioning Optimizer States



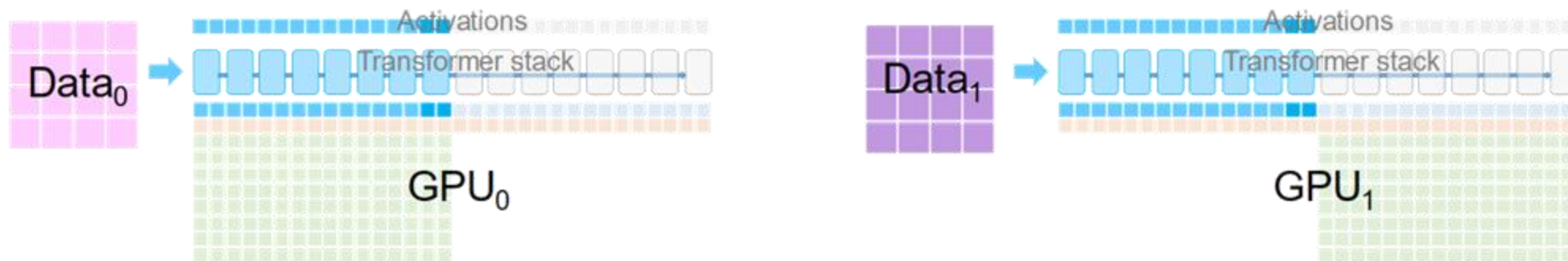
- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks

# ZeRO Stage 1: Partitioning Optimizer States



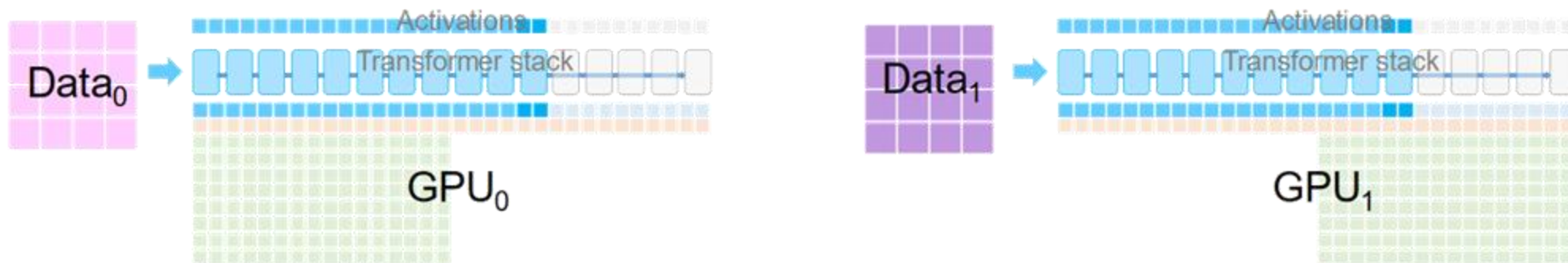
- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks

# ZeRO Stage 1: Partitioning Optimizer States



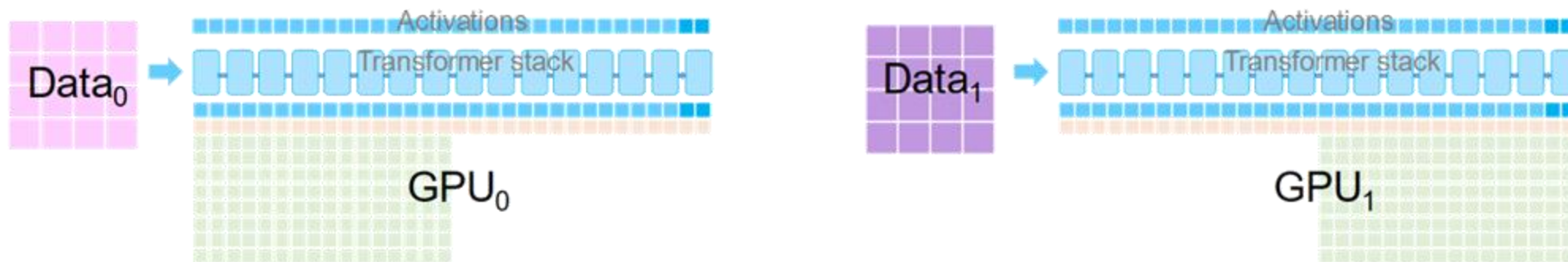
- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks

# ZeRO Stage 1: Partitioning Optimizer States



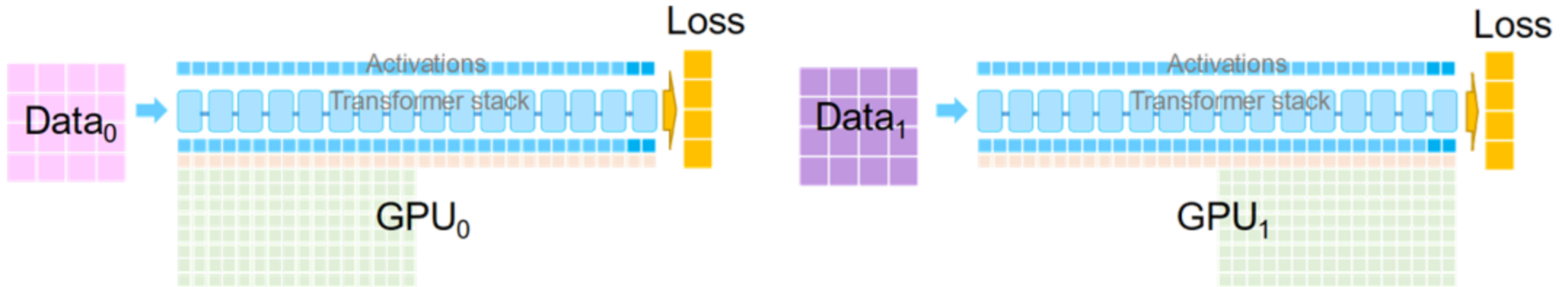
- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks

# ZeRO Stage 1: Partitioning Optimizer States



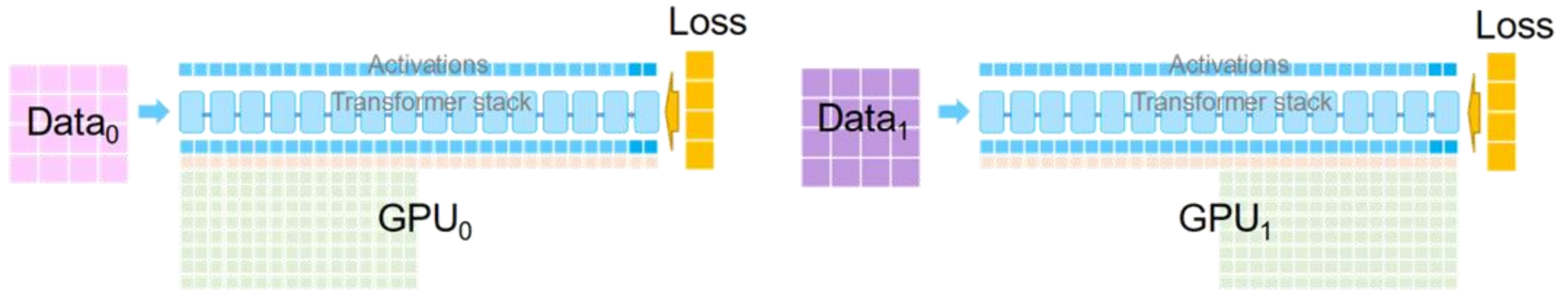
- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks

# ZeRO Stage 1: Partitioning Optimizer States



- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks

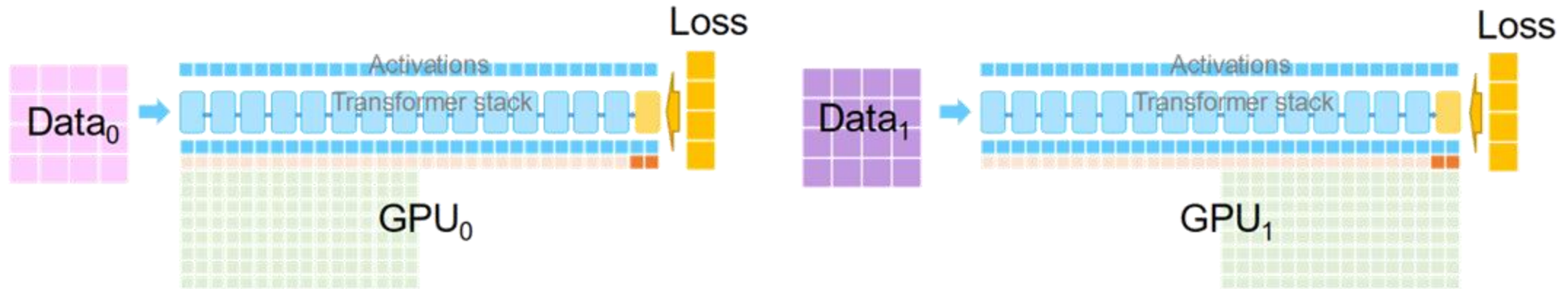
# ZeRO Stage 1: Partitioning Optimizer States



- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks
- Backward propagation to generate FP16 gradients



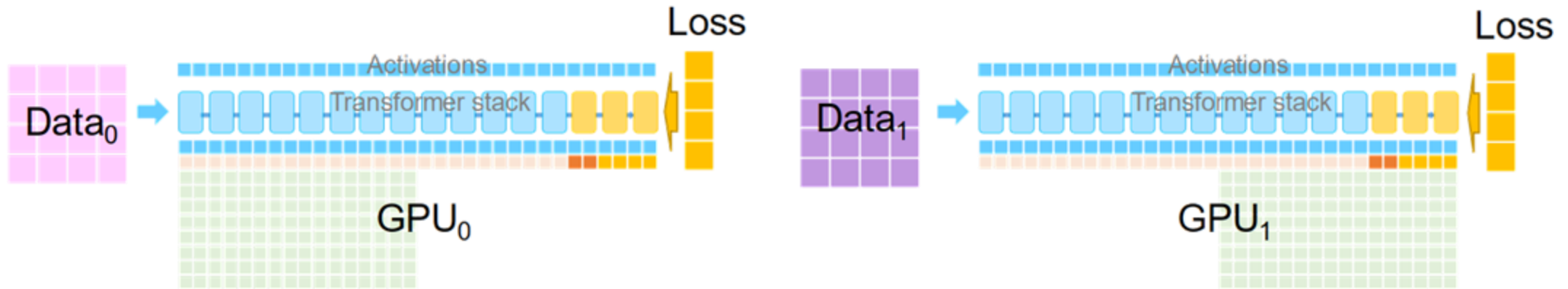
# ZeRO Stage 1: Partitioning Optimizer States



- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks
- Backward propagation to generate FP16 gradients

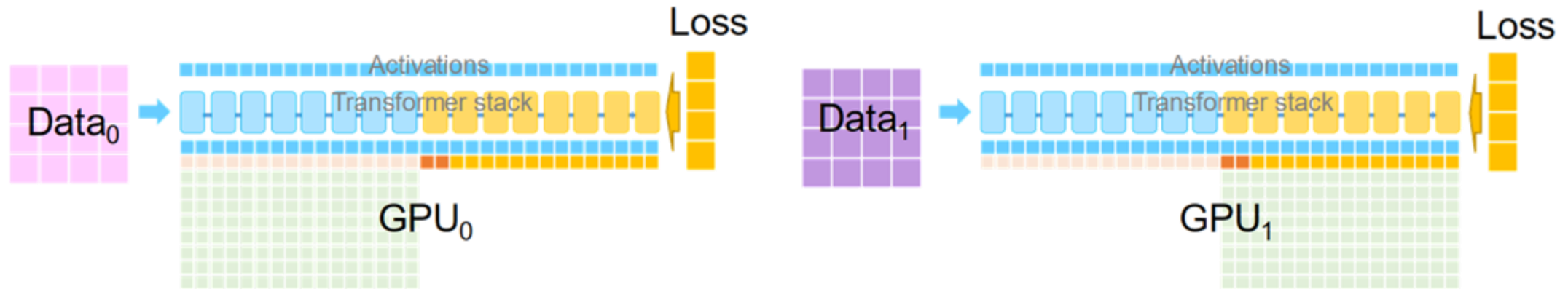


# ZeRO Stage 1: Partitioning Optimizer States



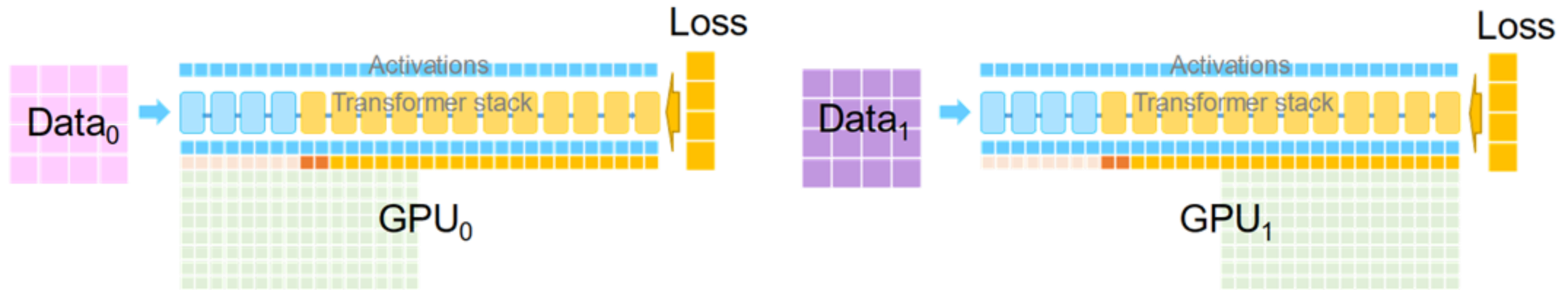
- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks
- Backward propagation to generate FP16 gradients

# ZeRO Stage 1: Partitioning Optimizer States



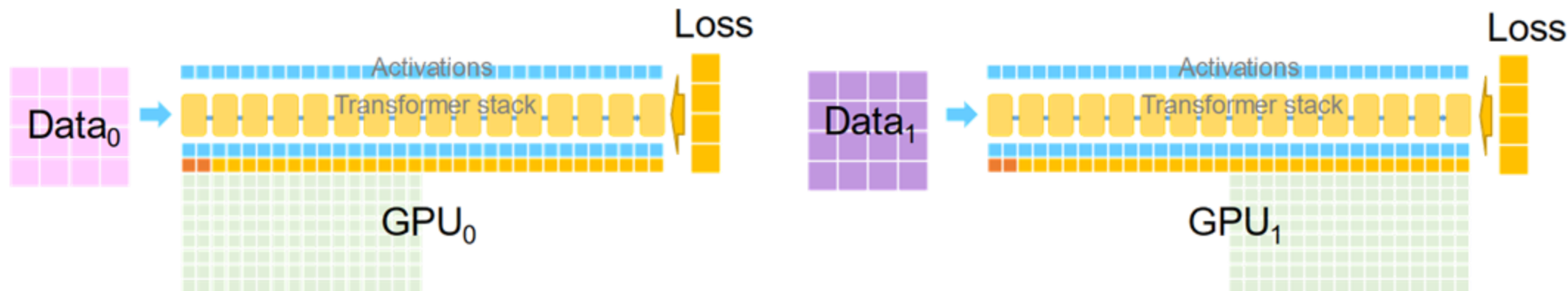
- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks
- Backward propagation to generate FP16 gradients

# ZeRO Stage 1: Partitioning Optimizer States



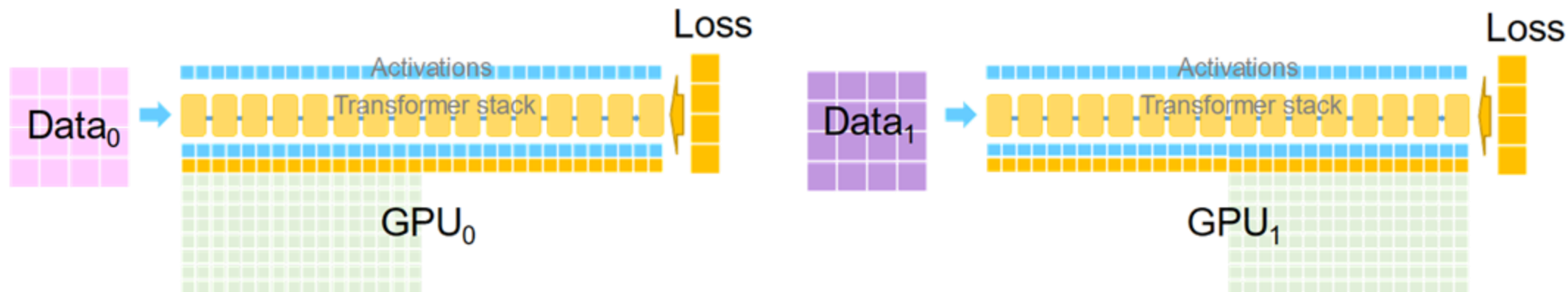
- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks
- Backward propagation to generate FP16 gradients

# ZeRO Stage 1: Partitioning Optimizer States



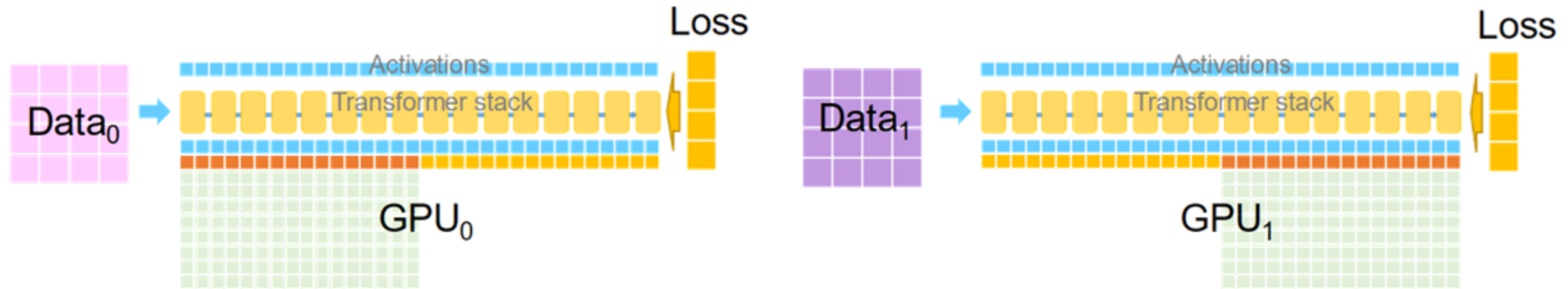
- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks
- Backward propagation to generate FP16 gradients

# ZeRO Stage 1: Partitioning Optimizer States



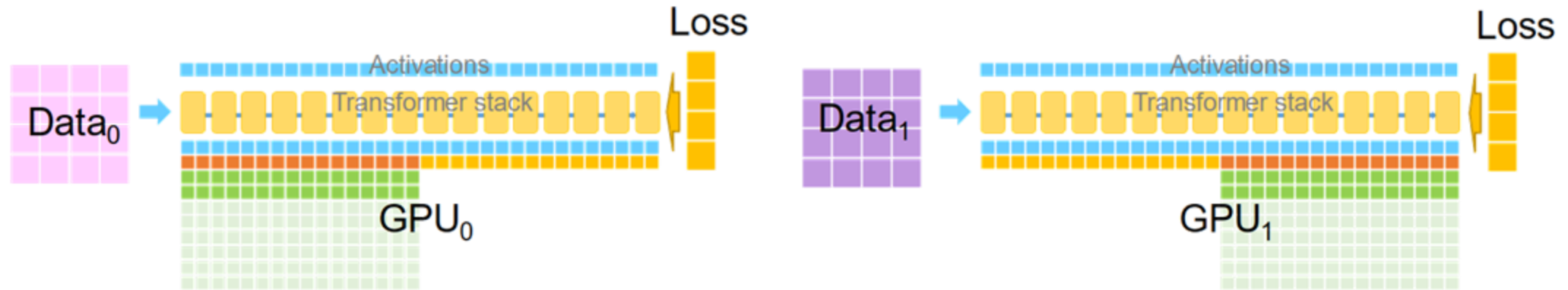
- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks
- Backward propagation to generate FP16 gradients and All-Reduce to average

# ZeRO Stage 1: Partitioning Optimizer States



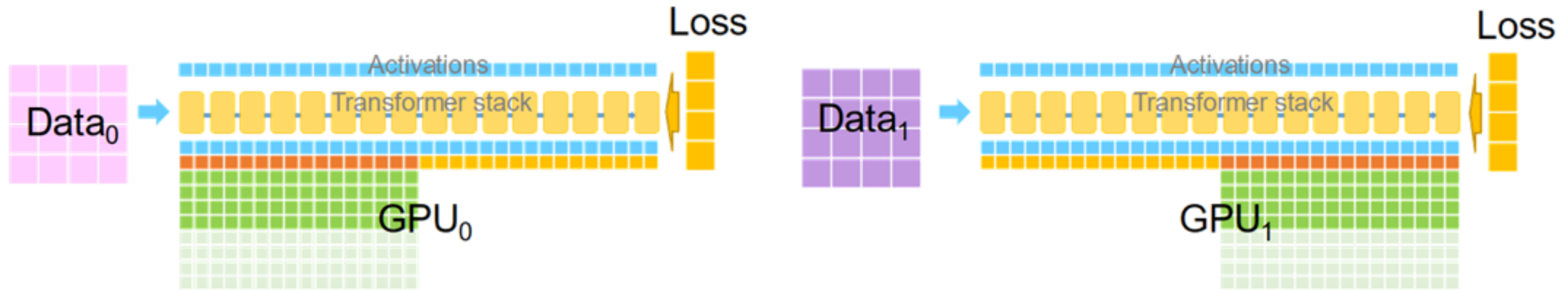
- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks
- Backward propagation to generate FP16 gradients and All-Reduce to average

# ZeRO Stage 1: Partitioning Optimizer States



- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks
- Backward propagation to generate FP16 gradients and All-Reduce to average
- Update the FP32 weights with ADAM optimizer

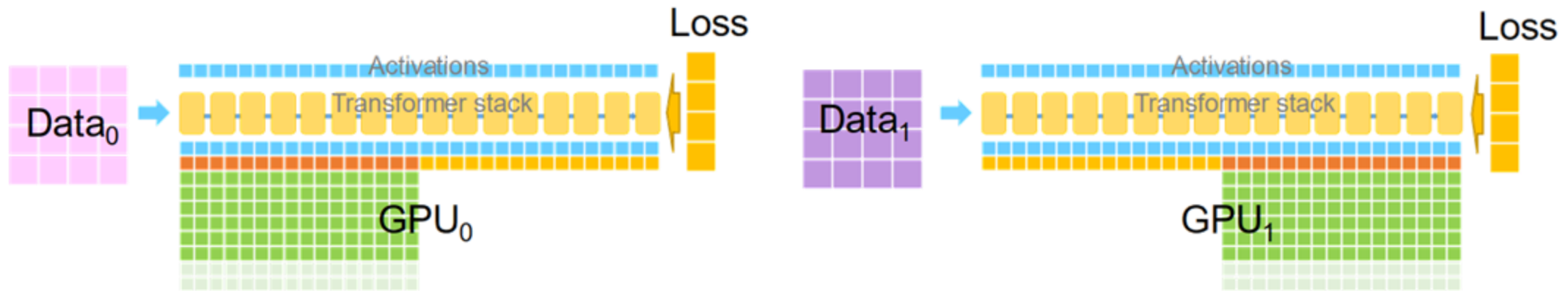
# ZeRO Stage 1: Partitioning Optimizer States



- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks
- Backward propagation to generate FP16 gradients and All-Reduce to average
- Update the FP32 weights with ADAM optimizer

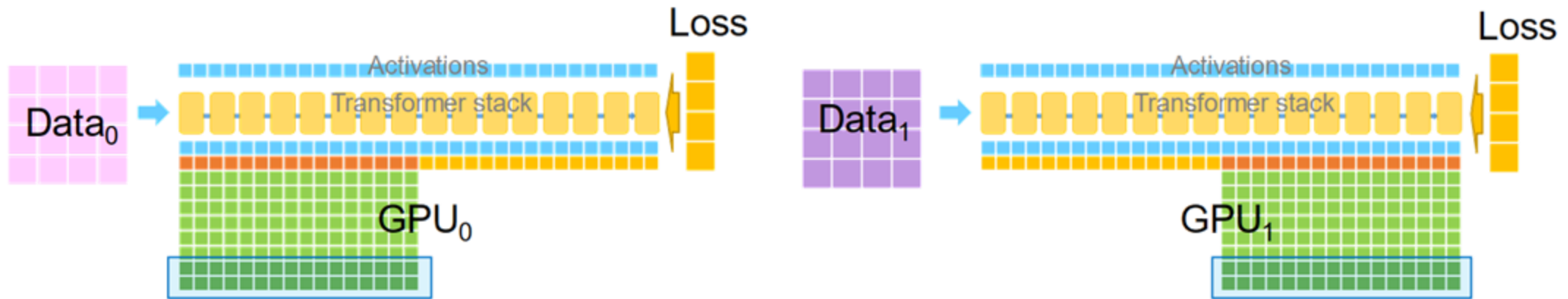


# ZeRO Stage 1: Partitioning Optimizer States



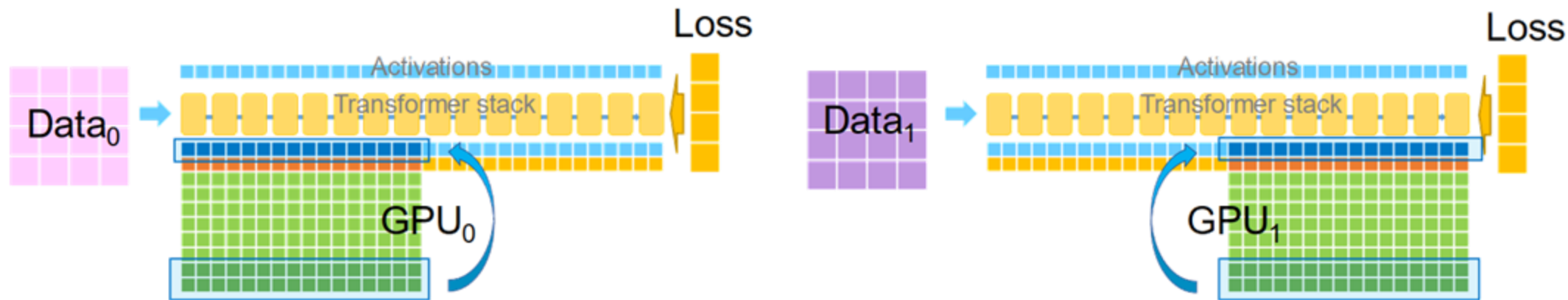
- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks
- Backward propagation to generate FP16 gradients and All-Reduce to average
- Update the FP32 weights with ADAM optimizer

# ZeRO Stage 1: Partitioning Optimizer States



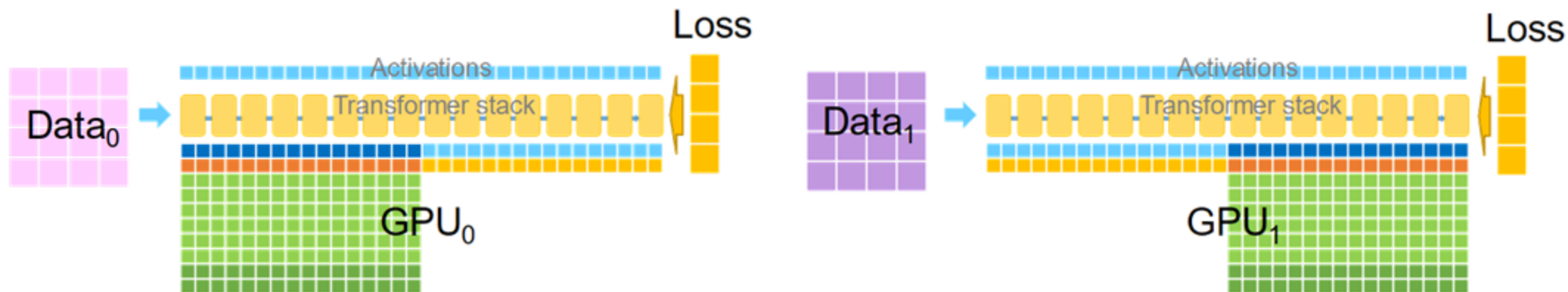
- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks
- Backward propagation to generate FP16 gradients and All-Reduce to average
- Update the FP32 weights with ADAM optimizer

# ZeRO Stage 1: Partitioning Optimizer States



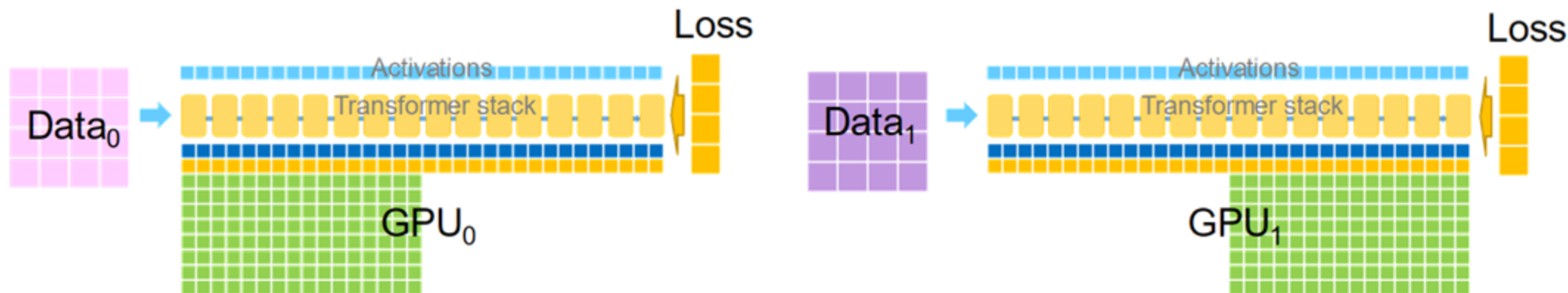
- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks
- Backward propagation to generate FP16 gradients and All-Reduce to average
- Update the FP32 weights with ADAM optimizer
- Update the FP16 weights

# ZeRO Stage 1: Partitioning Optimizer States



- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks
- Backward propagation to generate FP16 gradients and All-Reduce to average
- Update the FP32 weights with ADAM optimizer
- Update the FP16 weights
-

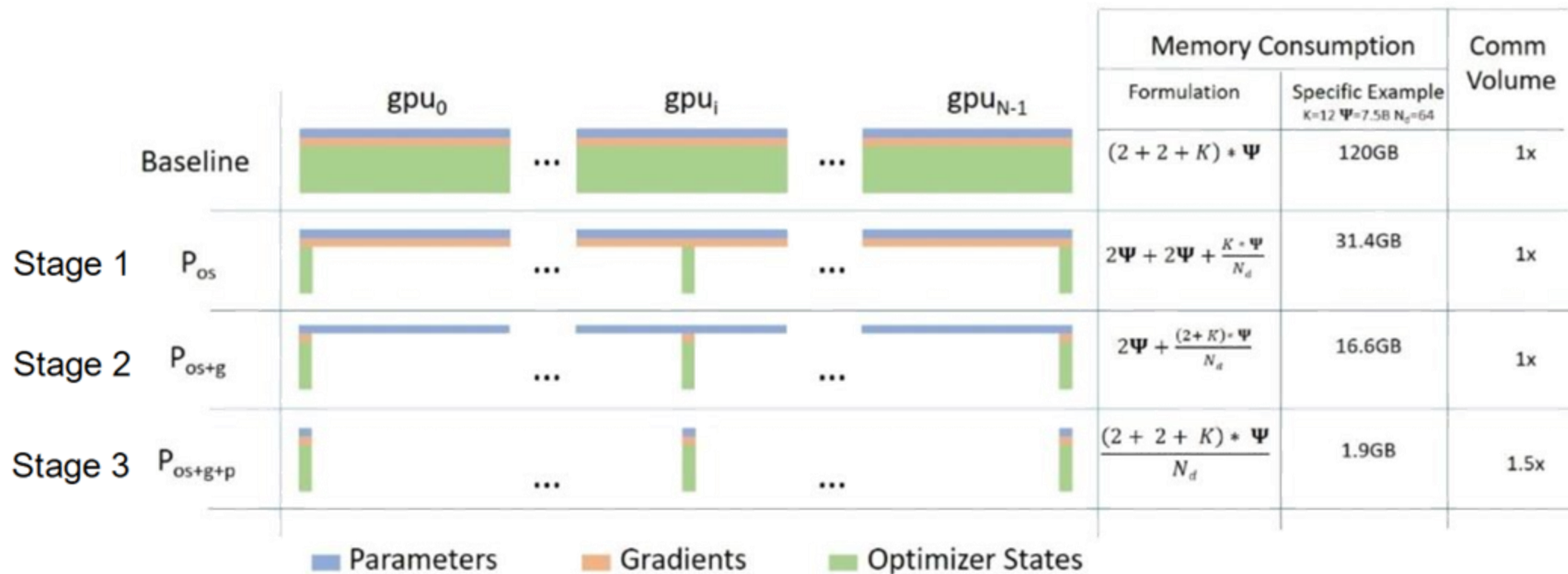
# ZeRO Stage 1: Partitioning Optimizer States



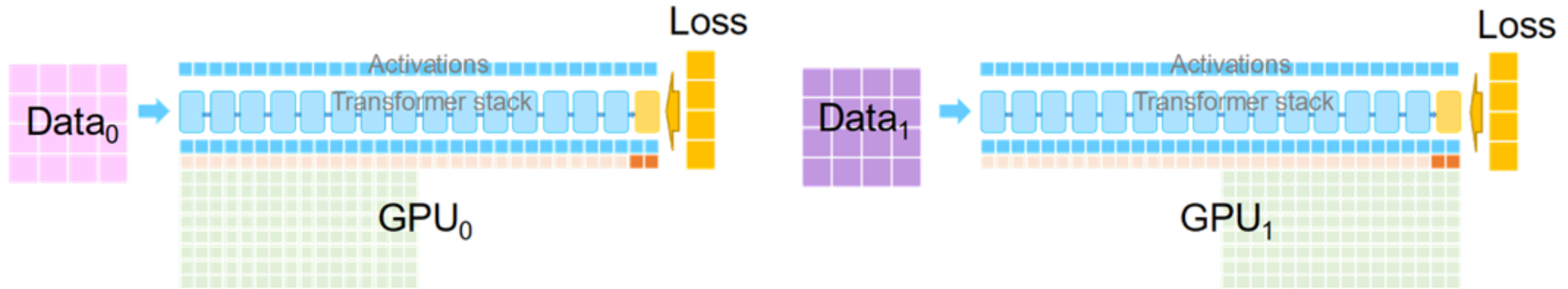
- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks
- Backward propagation to generate FP16 gradients and All-Reduce to average
- Update the FP32 weights with ADAM optimizer
- Update the FP16 weights
- All Gather the FP16 weights to complete the iteration

# ZeRO: Zero Redundancy Optimizer

- Progressive memory savings and communication volume

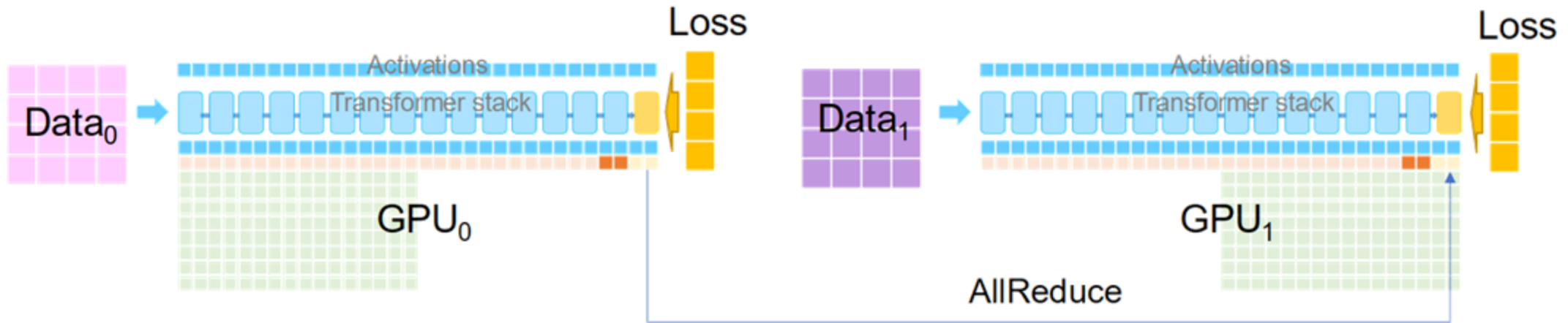


## ZeRO Stage 2: Partitioning Gradients



- Partitioning gradients across GPUs
- The forward process remains the same as stage 1

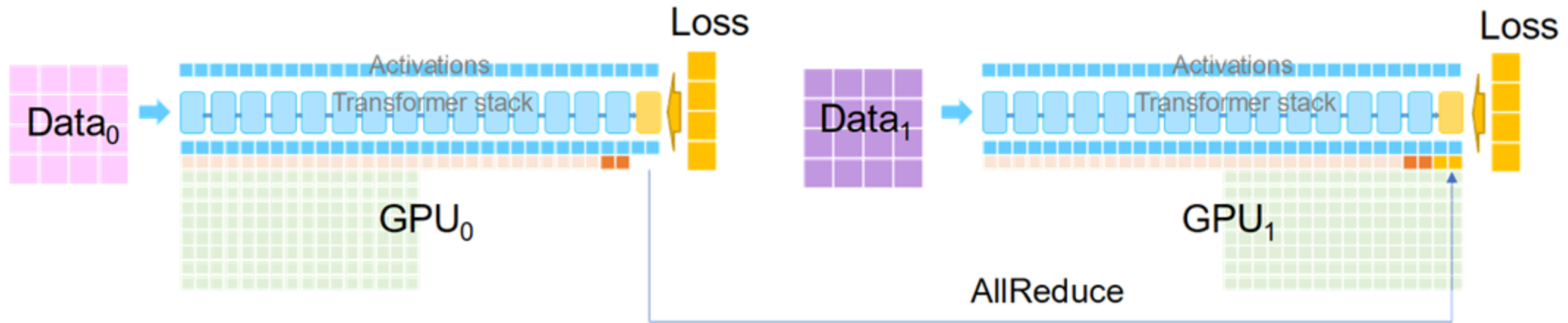
## ZeRO Stage 2: Partitioning Gradients



- Partitioning gradients across GPUs
- Perform All-Reduce right after back propagation of each layer

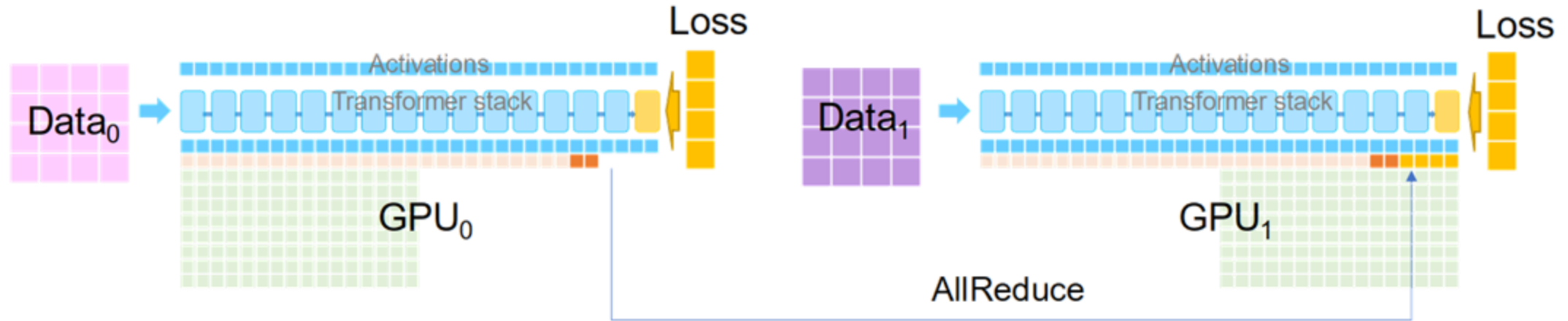


## ZeRO Stage 2: Partitioning Gradients



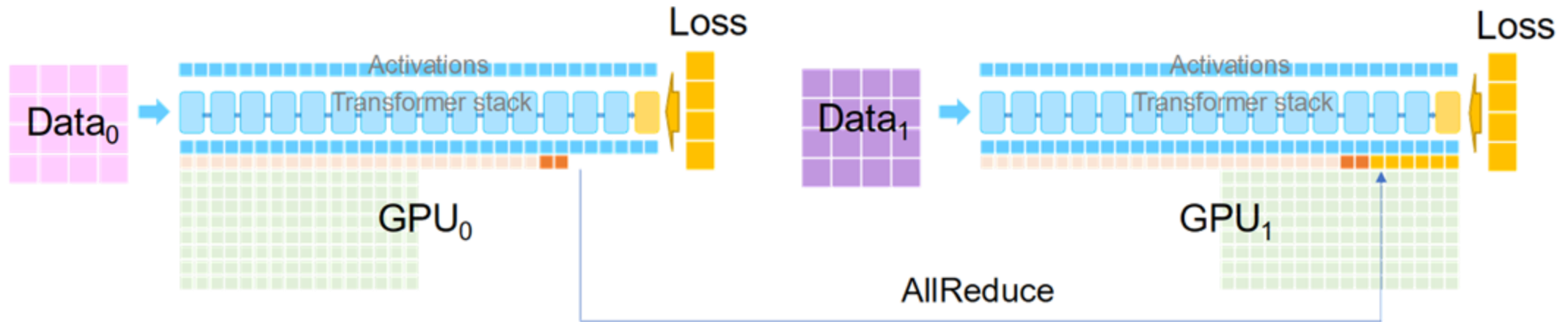
- Partitioning gradients across GPUs
- Only one GPU keeps the gradients after All-Reduce

## ZeRO Stage 2: Partitioning Gradients



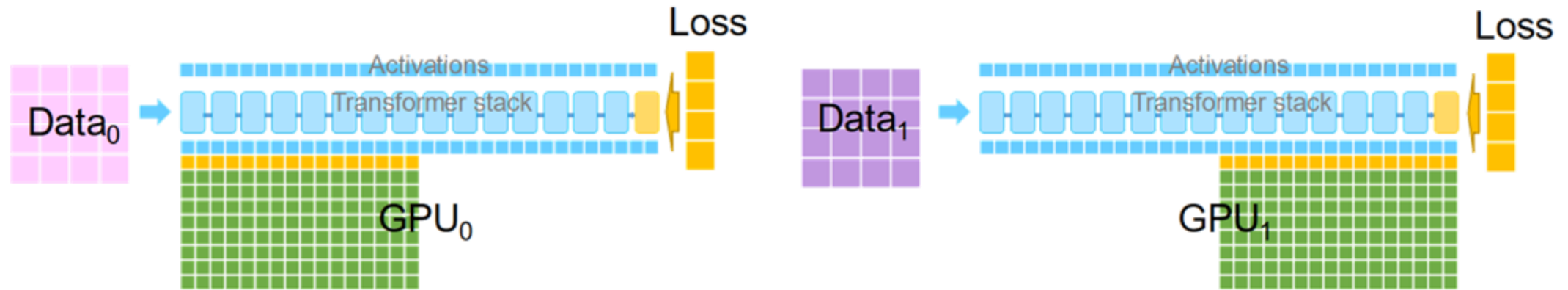
- Partitioning gradients across GPUs
- Reduce gradients on GPUs responsible for updating parameters

## ZeRO Stage 2: Partitioning Gradients



- Partitioning gradients across GPUs
- Reduce gradients on GPUs responsible for updating parameters

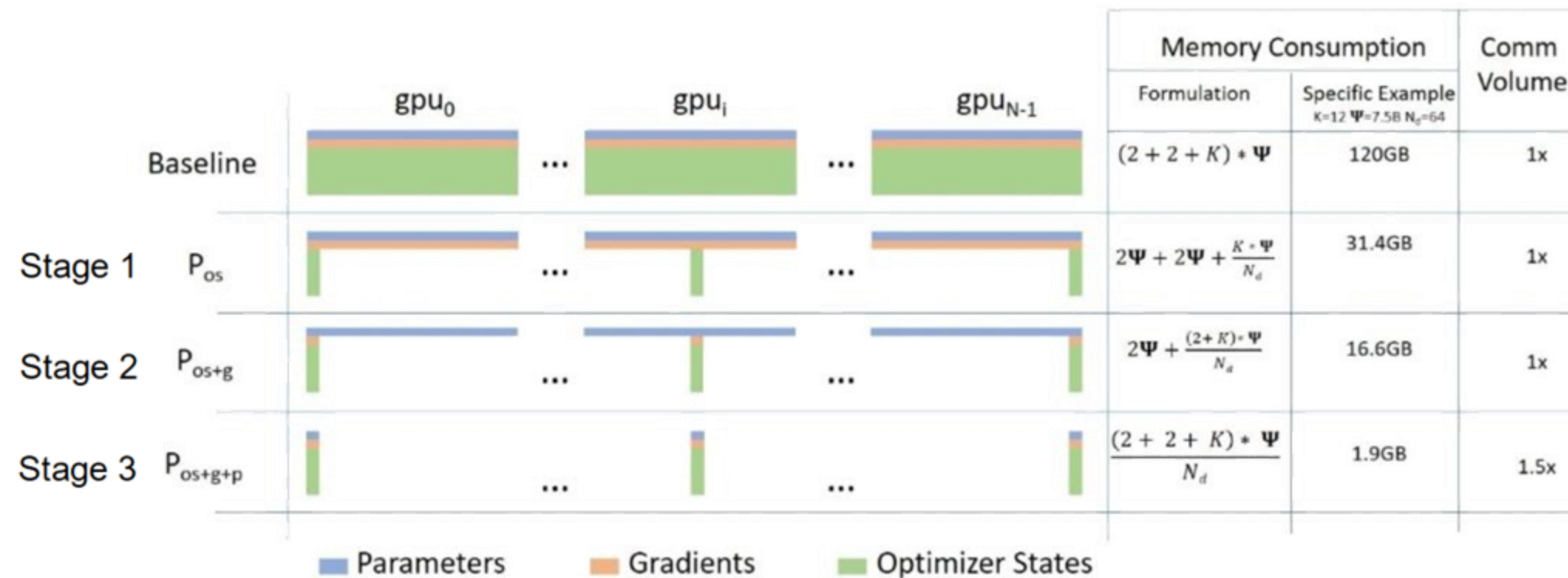
## ZeRO Stage 2: Partitioning Gradients



- Partitioning gradients across GPUs
- Reduce gradients on GPUs responsible for updating parameters

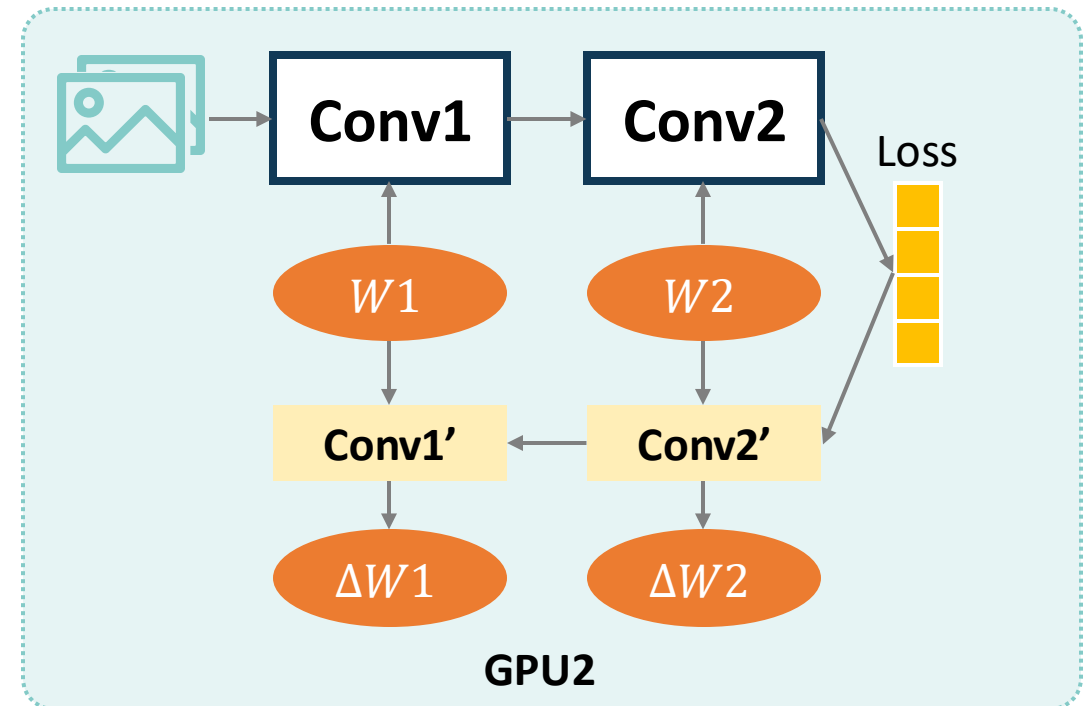
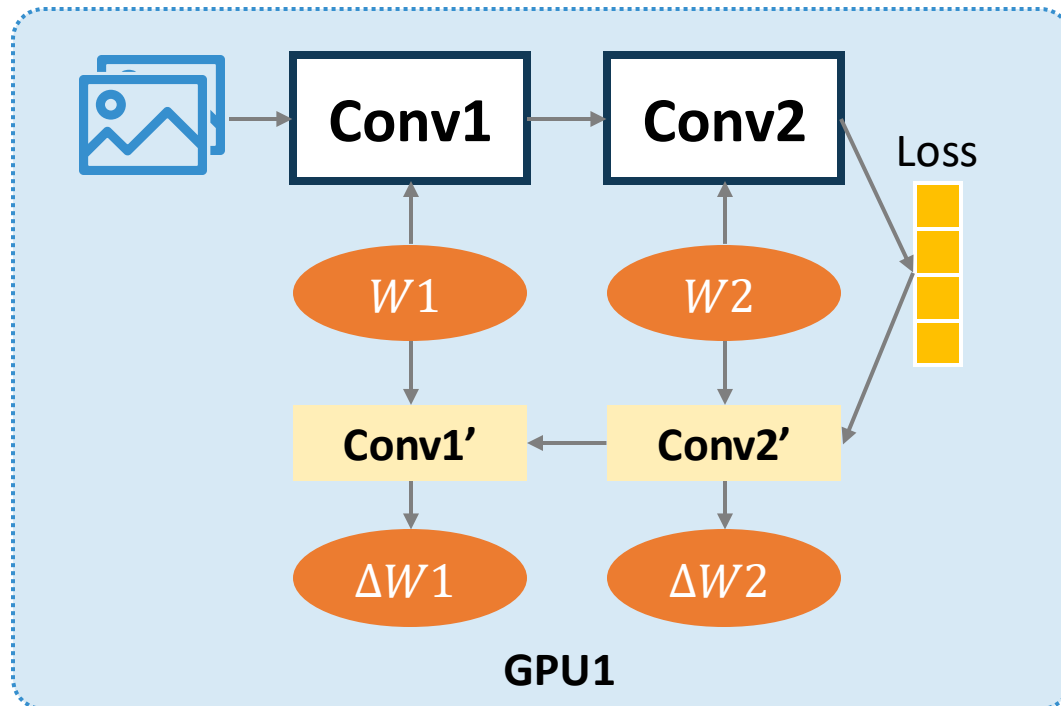
# ZeRO: Zero Redundancy Optimizer

- Progressive memory savings and communication volume
- Turning NLR 17.2B is powered by Stage 1 and Megatron



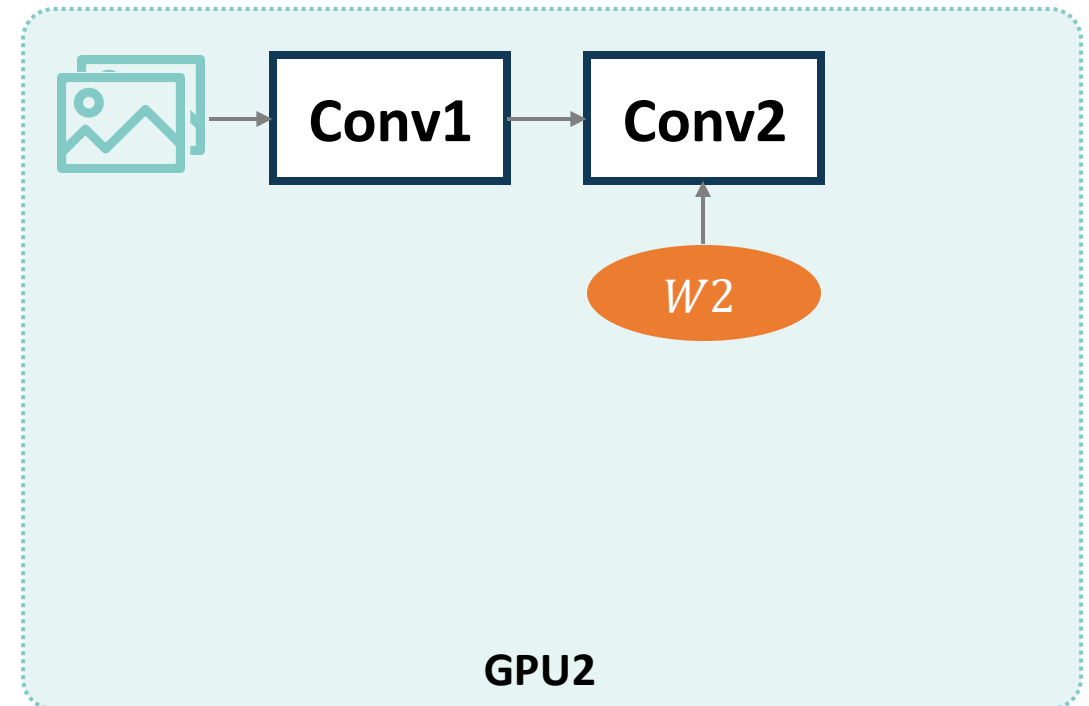
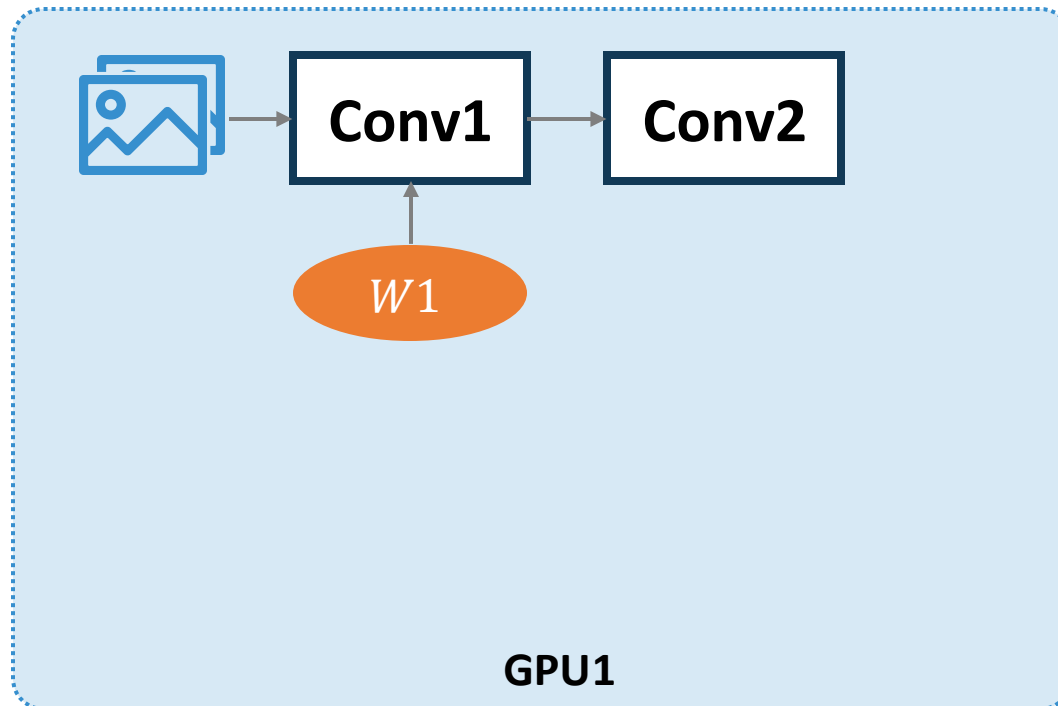
## ZeRO Stage 3: Partitioning Parameters

- In data parallel training, all GPUs keep **all** parameters during training



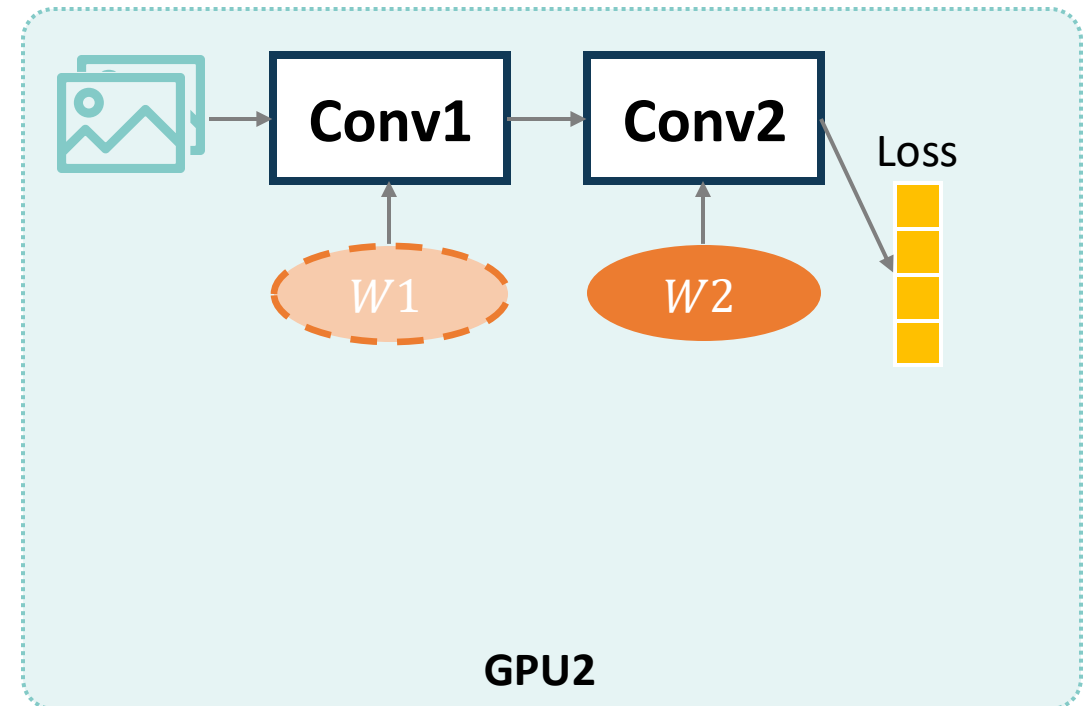
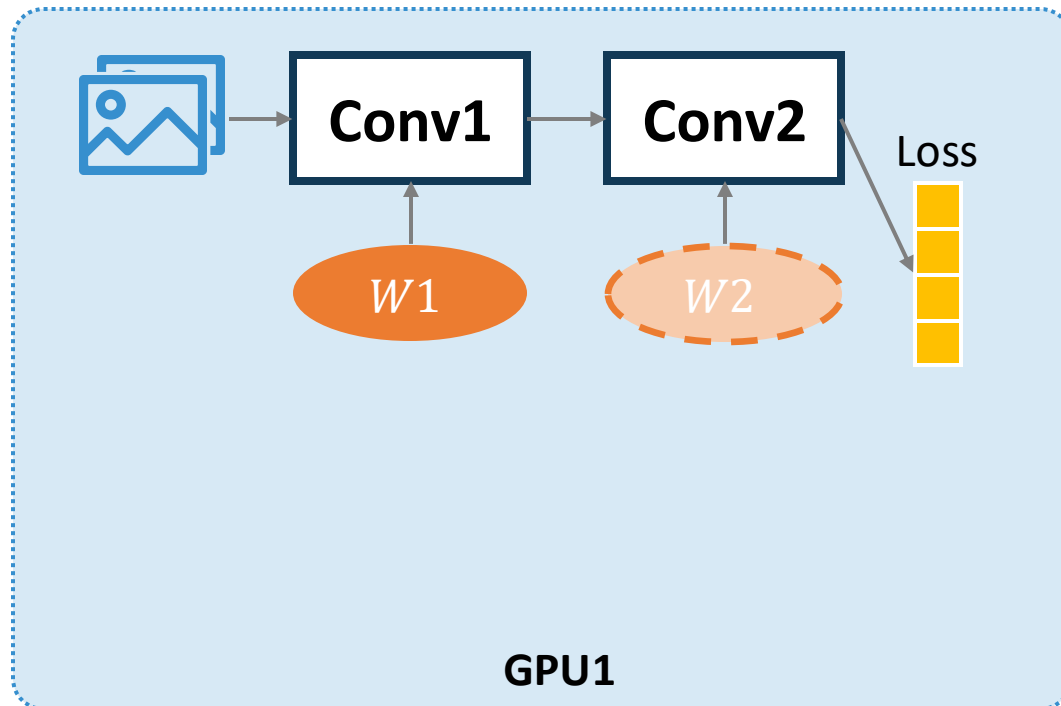
## ZeRO Stage 3: Partitioning Parameters

- In ZeRO, model parameters are partitioned across GPUs



## ZeRO Stage 3: Partitioning Parameters

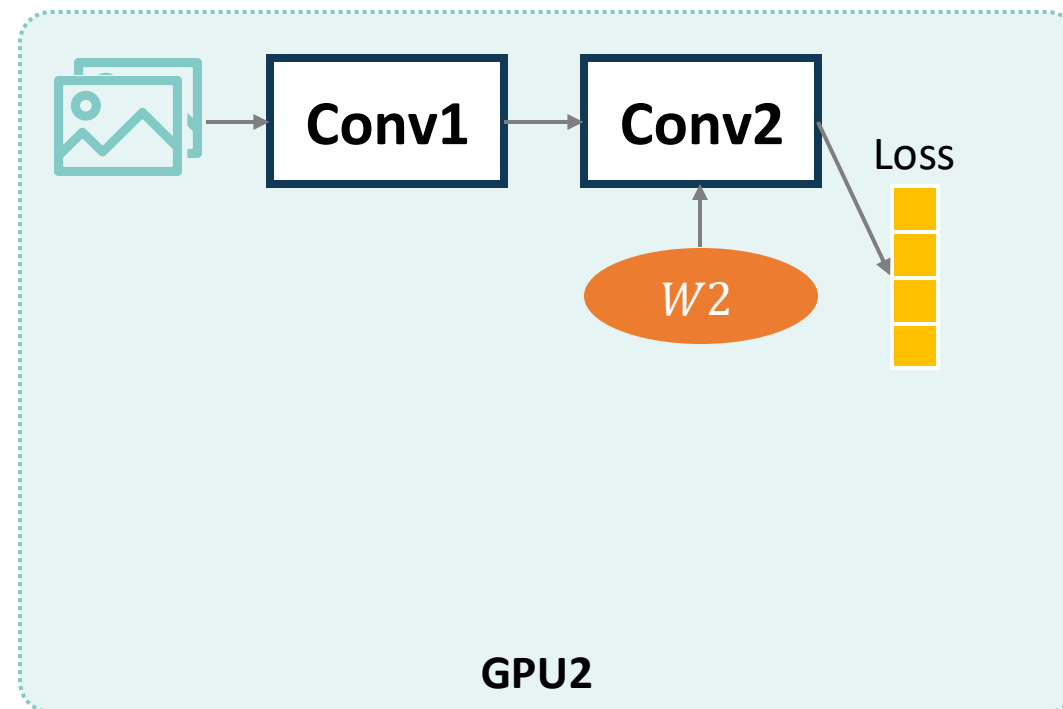
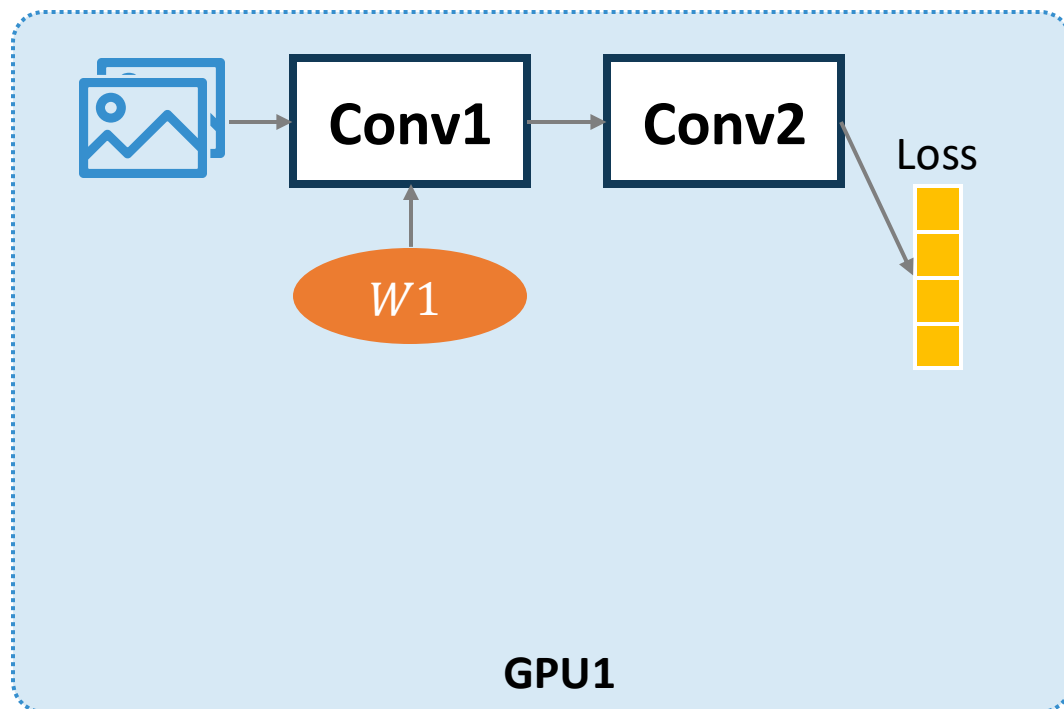
- In ZeRO, model parameters are partitioned across GPUs
- GPUs broadcast their parameters during forward





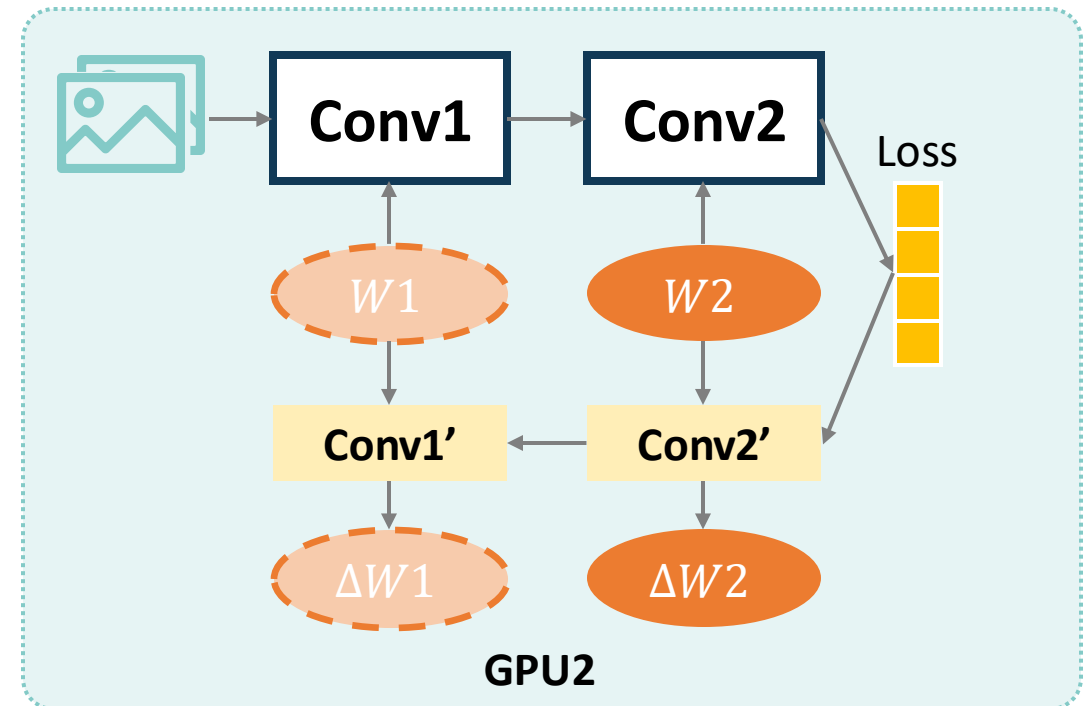
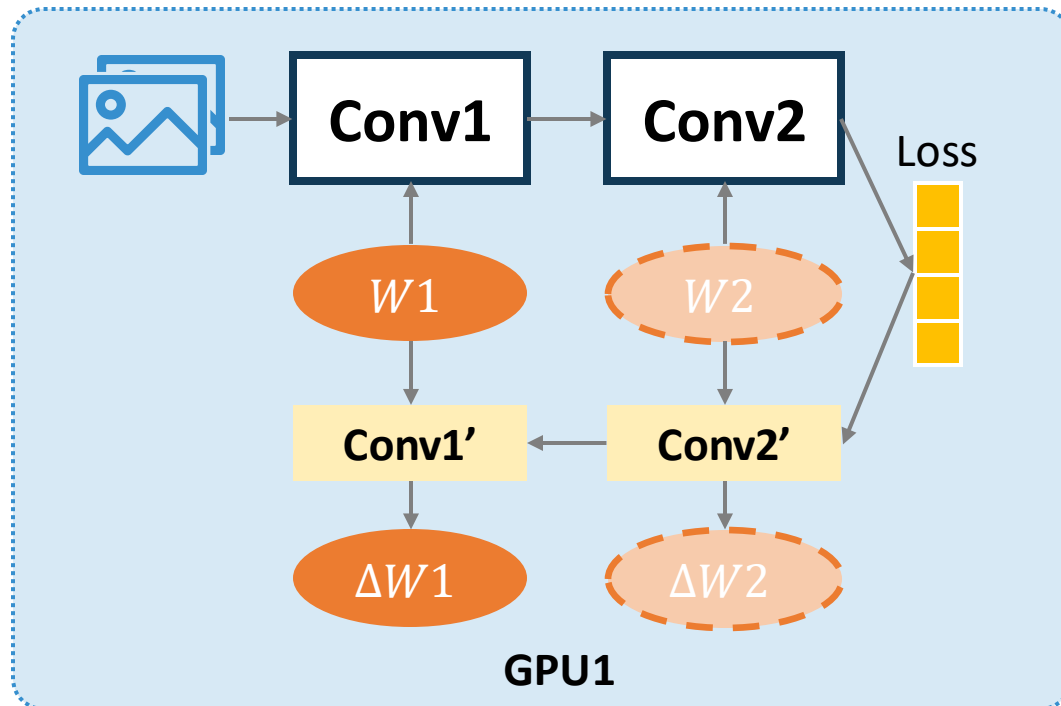
## ZeRO Stage 3: Partitioning Parameters

- In ZeRO, model parameters are partitioned across GPUs
- Parameters are discarded right after use



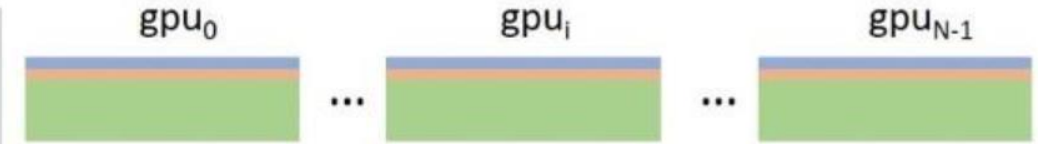



## ZeRO Stage 3: Partitioning Parameters

- In ZeRO, model parameters are partitioned across GPUs
- GPUs broadcast their parameters again during backward



# ZeRO: Zero Redundancy Optimizer

- ZeRO has three different stages
- Progressive memory savings and communication volume

			Memory Consumption		Comm Volume
			Formulation	Specific Example $K=12 \ \Psi=7.5B \ N_d=64$	
Baseline			$(2 + 2 + K) * \Psi$	120GB	1x
Stage 1	$P_{os}$		$2\Psi + 2\Psi + \frac{K * \Psi}{N_d}$	31.4GB	1x
Stage 2	$P_{os+g}$		$2\Psi + \frac{(2+K)*\Psi}{N_d}$	16.6GB	1x
Stage 3	$P_{os+g+p}$		$\frac{(2 + 2 + K) * \Psi}{N_d}$	1.9GB	1.5x

■ Parameters   
 ■ Gradients   
 ■ Optimizer States

# Acknowledgement

The development of this course, including its structure, content, and accompanying presentation slides, has been significantly influenced and inspired by the excellent work of instructors and institutions who have shared their materials openly. We wish to extend our sincere acknowledgement and gratitude to the following courses, which served as invaluable references and a source of pedagogical inspiration:

- Machine Learning Systems[15-442/15-642], by **Tianqi Chen** and **Zhihao Jia** at **CMU**.
- Advanced Topics in Machine Learning (Systems)[CS6216], by **Yao Lu** at **NUS**

While these materials provided a foundational blueprint and a wealth of insightful examples, all content herein has been adapted, modified, and curated to meet the specific learning objectives of our curriculum. Any errors, omissions, or shortcomings found in these course materials are entirely our own responsibility. We are profoundly grateful for the contributions of the educators listed above, whose dedication to teaching and knowledge-sharing has made the creation of this course possible.

---



System for Artificial Intelligence

Thanks

Siyuan Feng  
Shanghai Innovation Institute



---