



# CS 498: Machine Learning System Spring 2026

Minjia Zhang

The Grainger College of Engineering

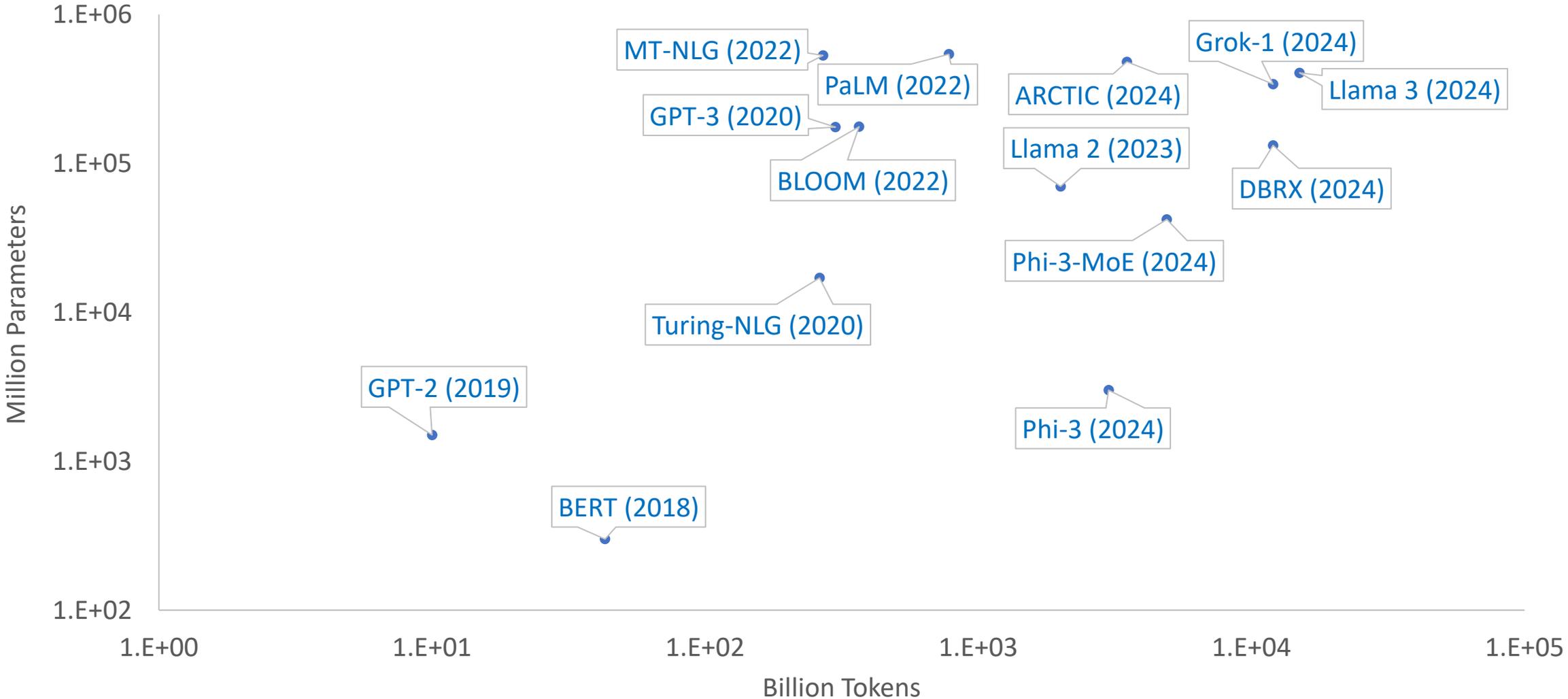
## **ZeRO-Style Data Parallelism (Fully-Sharded Data Parallelism)**

- Motivation
- ZeRO capability overview
- Understanding Memory Consumption
- ZeRO-DP: ZeRO powered data parallelism

### **Learning Objectives:**

- Understand more deeply the memory bottlenecks of DDP training
- Explain how ZeRO progressively removes memory redundancy

# Evolution of AI Models



# State-of-art and its limitations



	<b>Max Parameter (in billions)</b>	<b>Max Parallelism</b>	<b>Compute Efficiency</b>	<b>Usability (Model Rewrite)</b>
<b>Data Parallel (DP)</b>	Approx. 1.2	>1000	Very Good	Great

\*Mixed precision Adam on Cluster of DGX-2 with NVIDIA 32 GB V100 GPUs

# State-of-art and its limitations



	<b>Max Parameter (in billions)</b>	<b>Max Parallelism</b>	<b>Compute Efficiency</b>	<b>Usability (Model Rewrite)</b>
<b>Data Parallel (DP)</b>	Approx. 1.2	>1000	Very Good	Great

Democratization of large model training requires easy to use solutions that can support large model sizes and parallelism degree while remaining highly efficient

\*Mixed precision Adam on Cluster of DGX-2 with NVIDIA 32 GB V100 GPUs

# State-of-art and its limitations



	<b>Max Parameter (in billions)</b>	<b>Max Parallelism</b>	<b>Compute Efficiency</b>	<b>Usability (Model Rewrite)</b>
<b>Data Parallel (DP)</b>	Approx. 1.2	>1000	Very Good	Great
<b>Model Parallel (MP)</b>	Approx. 20	Approx. 16	Good	Needs Model Rewrite

Democratization of large model training requires easy to use solutions that can support large model sizes and parallelism degree while remaining highly efficient

\*Mixed precision Adam on Cluster of DGX-2 with NVIDIA 32 GB V100 GPUs

# State-of-art and its limitations



	<b>Max Parameter (in billions)</b>	<b>Max Parallelism</b>	<b>Compute Efficiency</b>	<b>Usability (Model Rewrite)</b>
<b>Data Parallel (DP)</b>	Approx. 1.2	>1000	Very Good	Great
<b>Model Parallel (MP)</b>	Approx. 20	Approx. 16	Good	Needs Model Rewrite
<b>MP + DP</b>	Approx. 20	> 1000	Good	Needs Model Rewrite

Democratization of large model training requires easy to use solutions that can support large model sizes and parallelism degree while remaining highly efficient

\*Mixed precision Adam on Cluster of DGX-2 with NVIDIA 32 GB V100 GPUs

# State-of-art and its limitations



	<b>Max Parameter (in billions)</b>	<b>Max Parallelism</b>	<b>Compute Efficiency</b>	<b>Usability (Model Rewrite)</b>
<b>Data Parallel (DP)</b>	Approx. 1.2	>1000	Very Good	Great
<b>Model Parallel (MP)</b>	Approx. 20	Approx. 16	Good	Needs Model Rewrite
<b>MP + DP</b>	Approx. 20	> 1000	Good	Needs Model Rewrite
<b>Pipeline Parallel (PP)</b>	Approx. 100	Approx. 128	Very Good	Needs Model Rewrite

Democratization of large model training requires easy to use solutions that can support large model sizes and parallelism degree while remaining highly efficient

\*Mixed precision Adam on Cluster of DGX-2 with NVIDIA 32 GB V100 GPUs

# State-of-art and its limitations



Democratization of large model training requires easy to use solutions that can support large model sizes and parallelism degree while remaining highly efficient

	<b>Max Parameter (in billions)</b>	<b>Max Parallelism</b>	<b>Compute Efficiency</b>	<b>Usability (Model Rewrite)</b>
<b>Data Parallel (DP)</b>	Approx. 1.2	>1000	Very Good	Great
<b>Model Parallel (MP)</b>	Approx. 20	Approx. 16	Good	Needs Model Rewrite
<b>MP + DP</b>	Approx. 20	> 1000	Good	Needs Model Rewrite
<b>Pipeline Parallel (PP)</b>	Approx. 100	Approx. 128	Very Good	Needs Model Rewrite
<b>PP + DP</b>	Approx. 100	> 1000	Very Good	Needs Model Rewrite

\*Mixed precision Adam on Cluster of DGX-2 with NVIDIA 32 GB V100 GPUs

# State-of-art and its limitations



Democratization of large model training requires easy to use solutions that can support large model sizes and parallelism degree while remaining highly efficient

	<b>Max Parameter (in billions)</b>	<b>Max Parallelism</b>	<b>Compute Efficiency</b>	<b>Usability (Model Rewrite)</b>
<b>Data Parallel (DP)</b>	Approx. 1.2	>1000	Very Good	Great
<b>Model Parallel (MP)</b>	Approx. 20	Approx. 16	Good	Needs Model Rewrite
<b>MP + DP</b>	Approx. 20	> 1000	Good	Needs Model Rewrite
<b>Pipeline Parallel (PP)</b>	Approx. 100	Approx. 128	Very Good	Needs Model Rewrite
<b>PP + DP</b>	Approx. 100	> 1000	Very Good	Needs Model Rewrite
<b>MP + PP + DP</b>	> 1000	> 1000	Very Good	Needs Significant Model Rewrite

\*Mixed precision Adam on Cluster of DGX-2 with NVIDIA 32 GB V100 GPUs

# Motivation

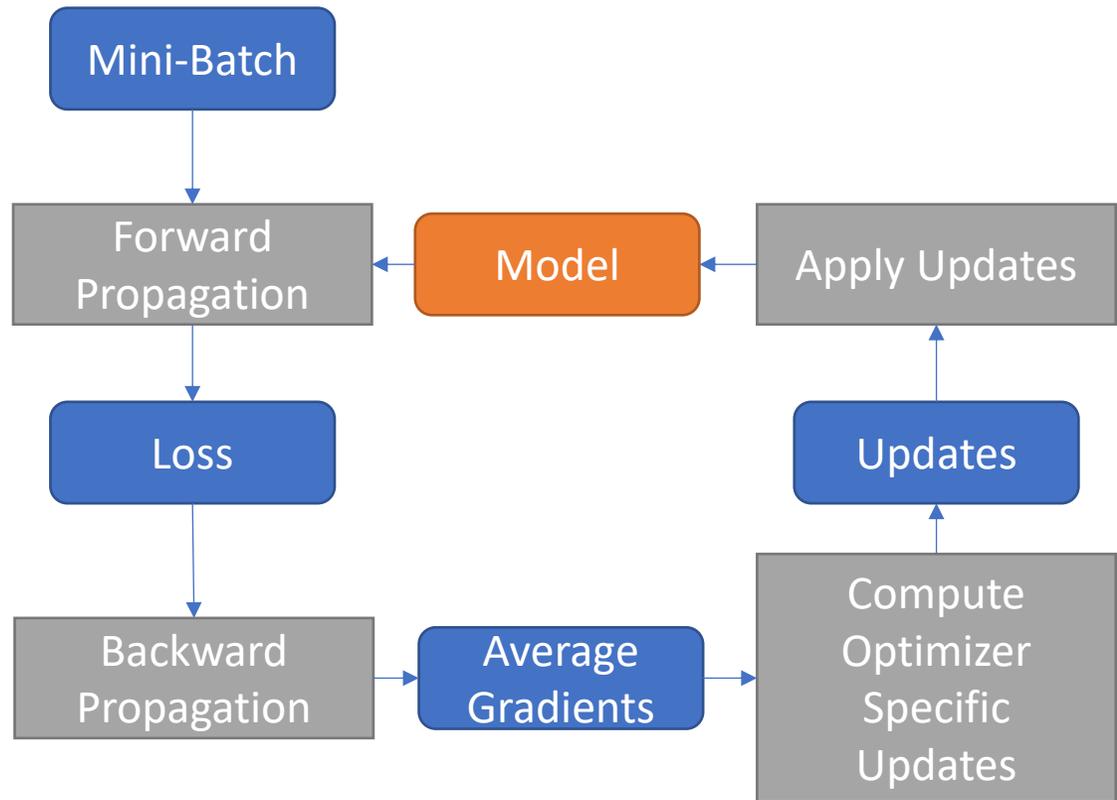


Democratization of large model training requires easy to use solutions that can support large model sizes and parallelism degree while remaining highly efficient

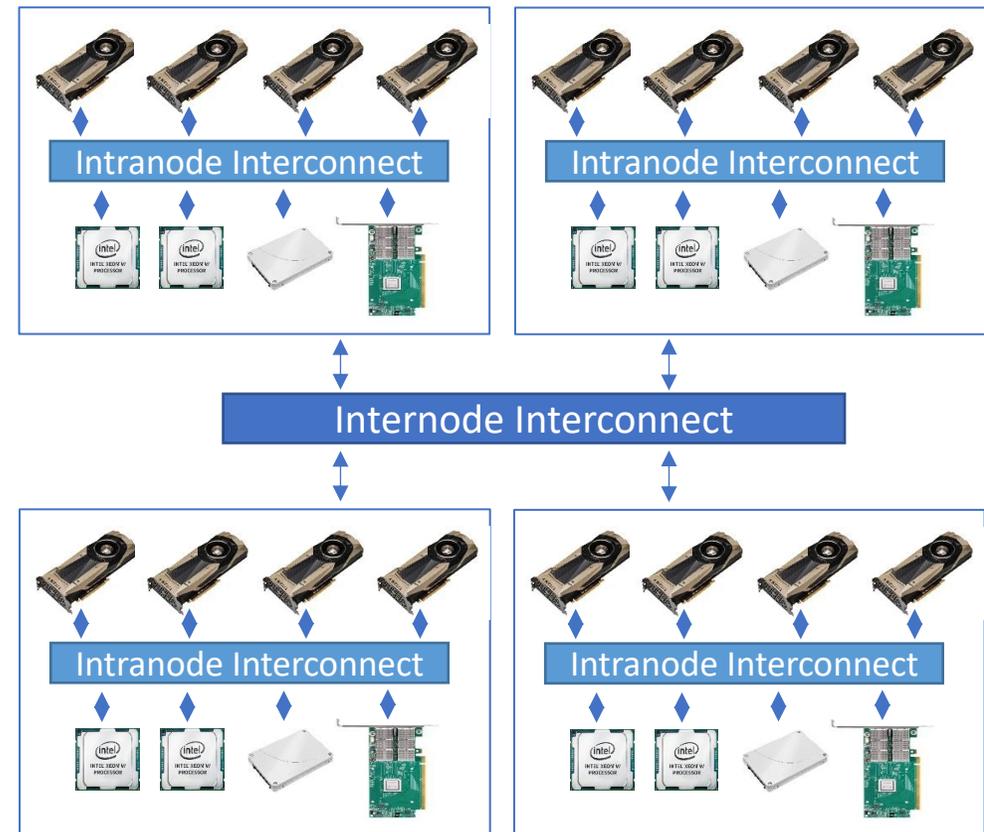
	<b>Max Parameter (in billions)</b>	<b>Max Parallelism</b>	<b>Compute Efficiency</b>	<b>Usability (Model Rewrite)</b>
<b>Data Parallel (DP)</b>	Approx. 1.2	>1000	Very Good	Great
<b>Model Parallel (MP)</b>	Approx. 20	Approx. 16	Good	Needs Model Rewrite
<b>MP + DP</b>	Approx. 20	> 1000	Good	Needs Model Rewrite
<b>Pipeline Parallel (PP)</b>	Approx. 100	Approx. 128	Very Good	Needs Model Rewrite
<b>PP + DP</b>	Approx. 100	> 1000	Very Good	Needs Model Rewrite
<b>MP + PP + DP</b>	> 1000	> 1000	Very Good	Needs Significant Model Rewrite
<b>ZeRO</b>	> 1000	> 1000	Very Good	Great

- Motivation
- ZeRO capability overview
- **Understanding Memory Consumption**
- ZeRO-DP: ZeRO powered data parallelism
- Evaluation

# Distributed Data Parallel Training Overview

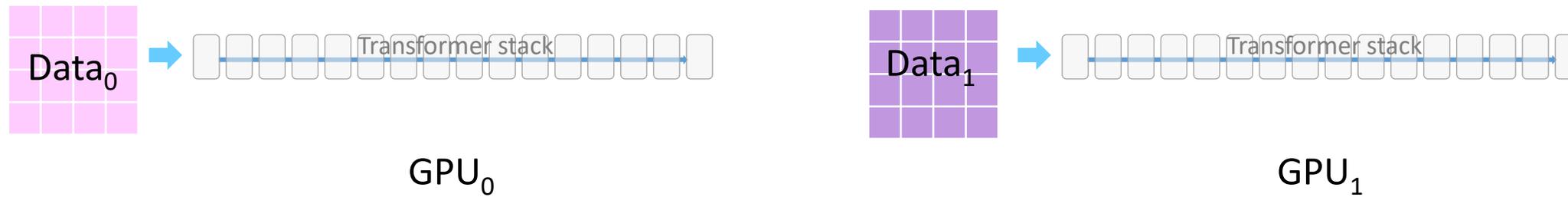


Data Parallel Training Loop



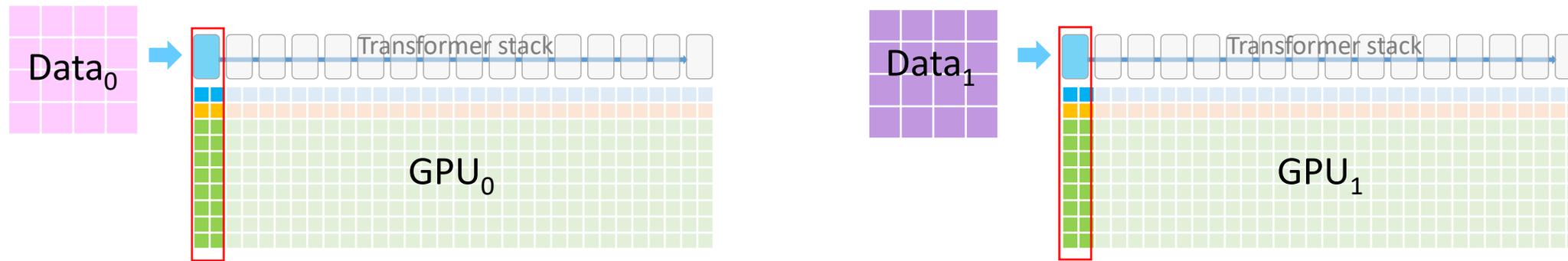
Distributed GPU Cluster

# Understanding Memory Consumption



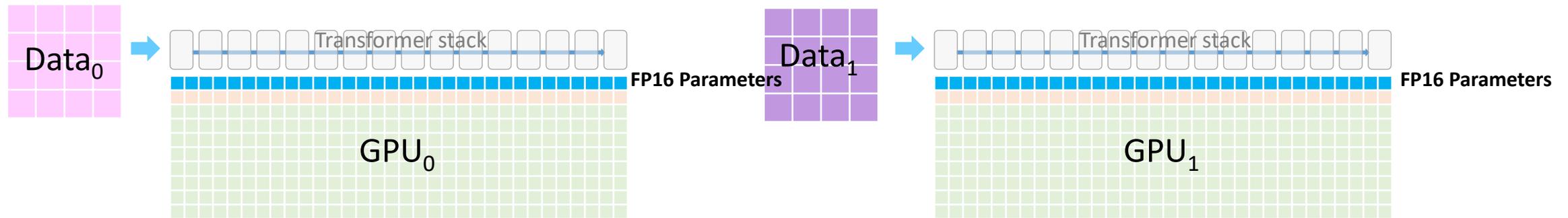
A 16-layer transformer model  = 1 layer

# Understanding Memory Consumption



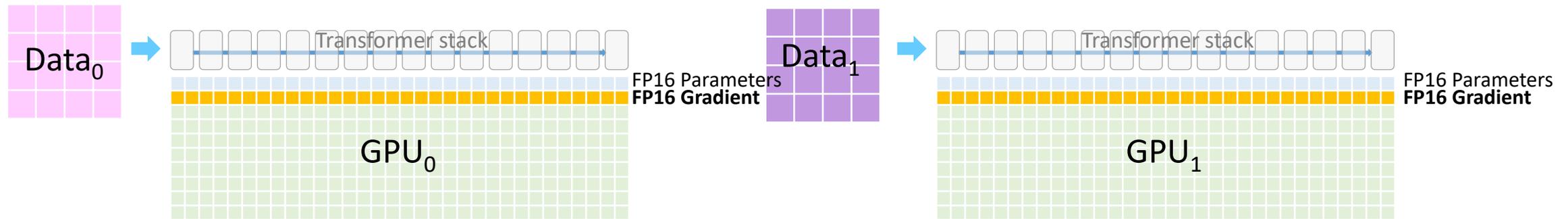
Each cell   represents GPU memory used by its corresponding transformer layer 

# Understanding Memory Consumption



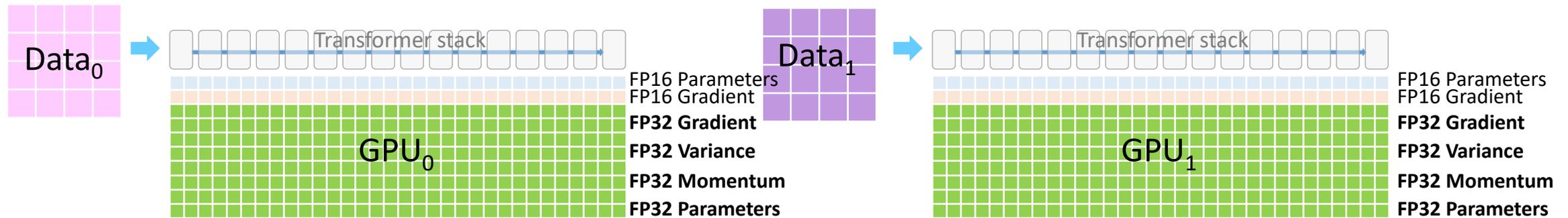
- FP16 parameter

# Understanding Memory Consumption



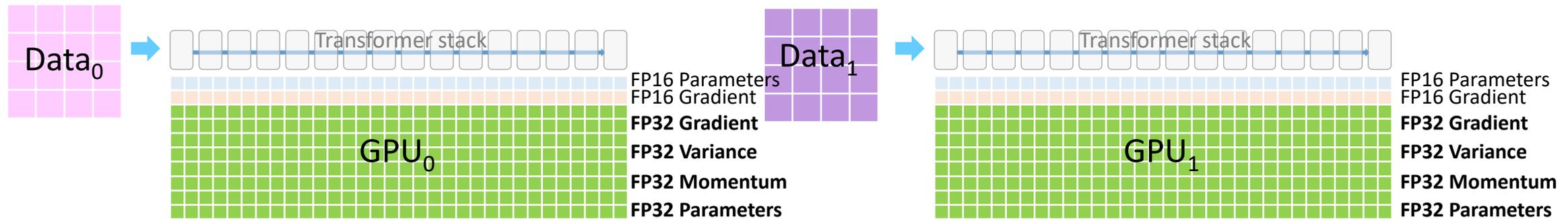
- FP16 parameter
- FP16 Gradients

# Understanding Memory Consumption



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

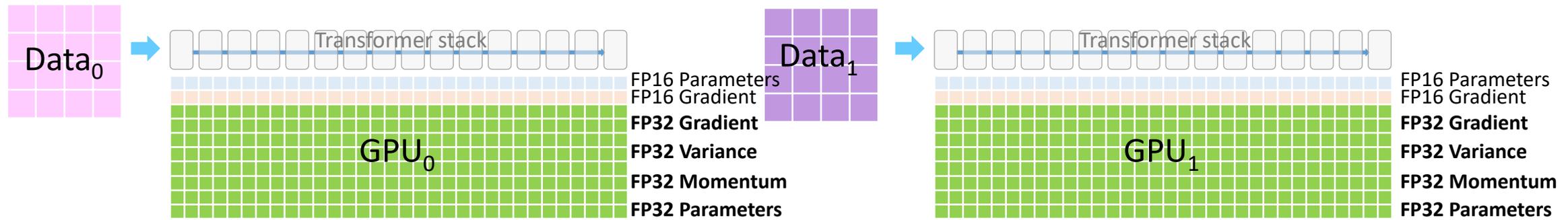
# Understanding Memory Consumption



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

M = number of parameters in the model

# Understanding Memory Consumption



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

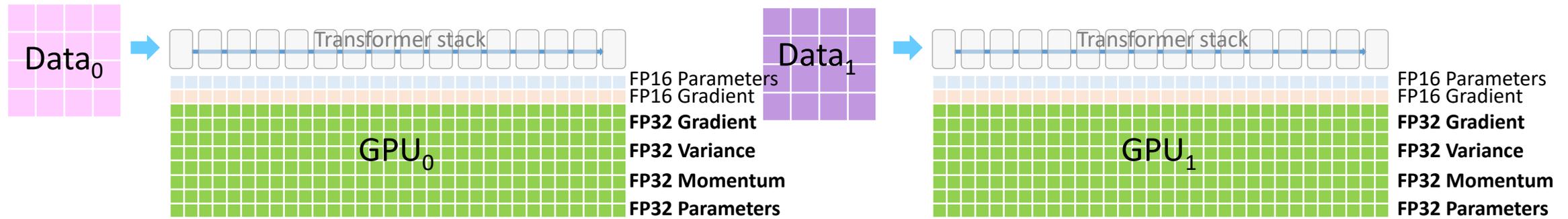
Example 1B parameter model -> 20GB/GPU

Memory consumption doesn't include:

- Input batch + activations

M = number of parameters in the model

# Understanding Memory Consumption



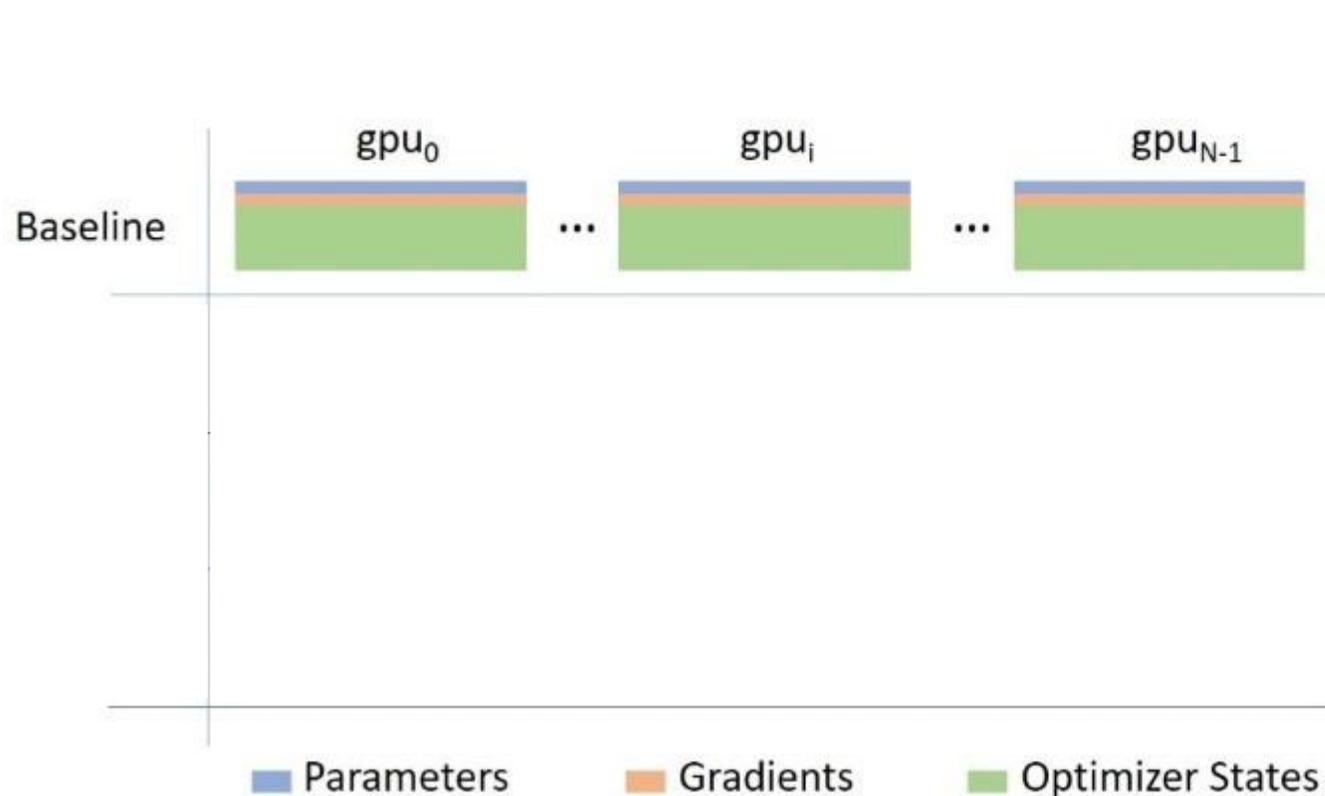
- FP16 parameter : **2M bytes**
- FP16 Gradients : **2M bytes**
- FP32 Optimizer States : **16M bytes**
  - Gradients, Variance, Momentum, Parameters

Example 1B parameter model -> 20GB/GPU

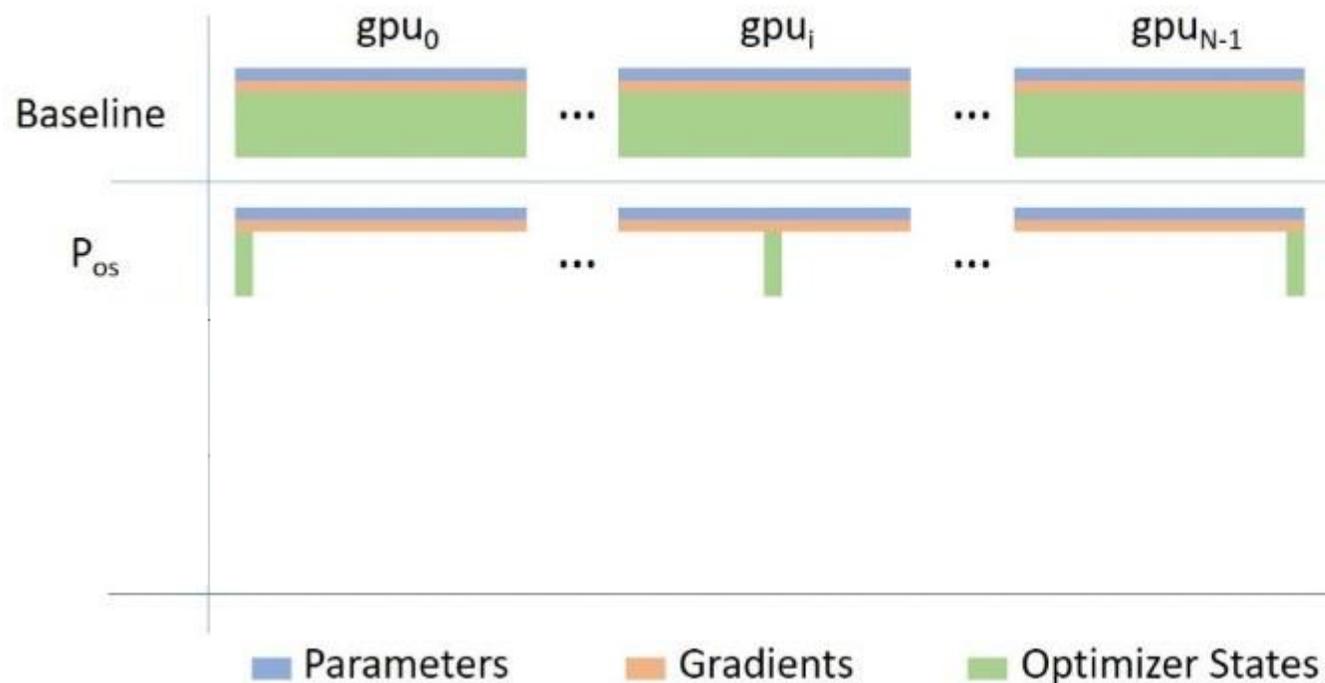
M = number of parameters in the model

- Motivation
- ZeRO capability overview
- Understanding Memory Consumption
- **ZeRO-DP: ZeRO powered data parallelism**
- Evaluation

- ZeRO removes the redundancy across data parallel process
- Partitioning optimizer states, gradients and parameters (3 stages)

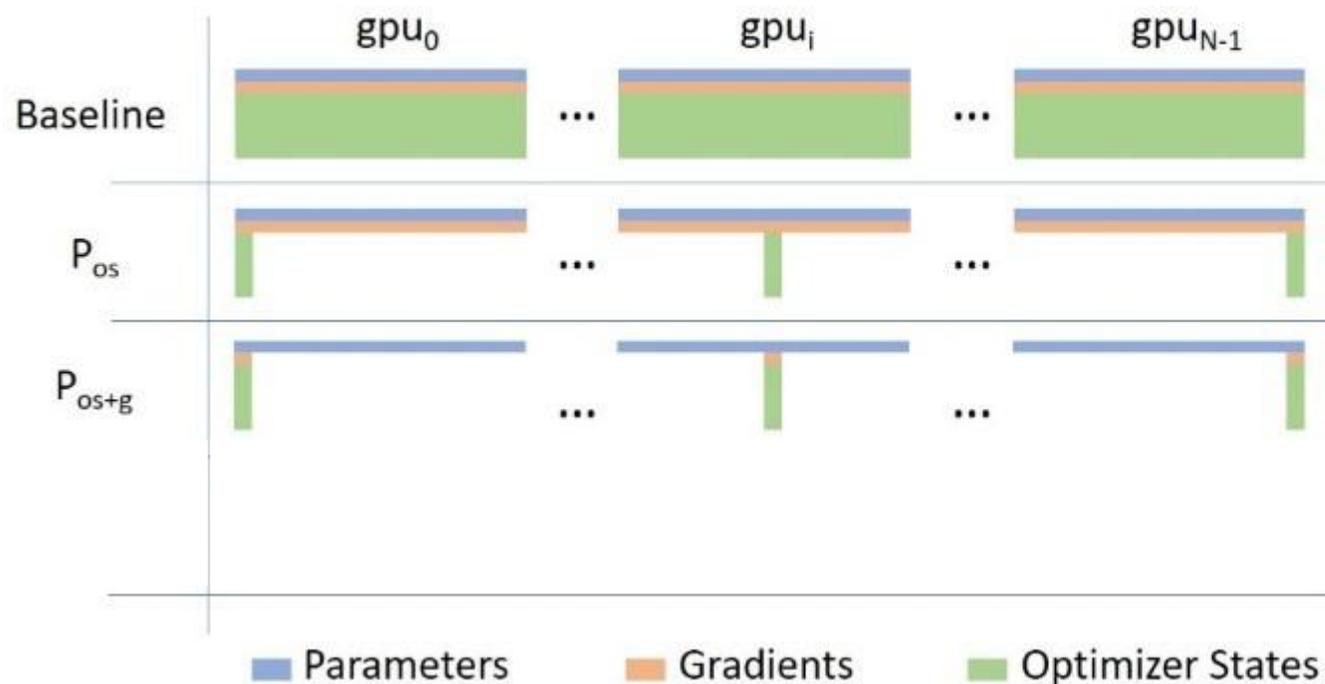


- ZeRO removes the redundancy across data parallel process
- Partitioning optimizer states, gradients and parameters (3 stages)



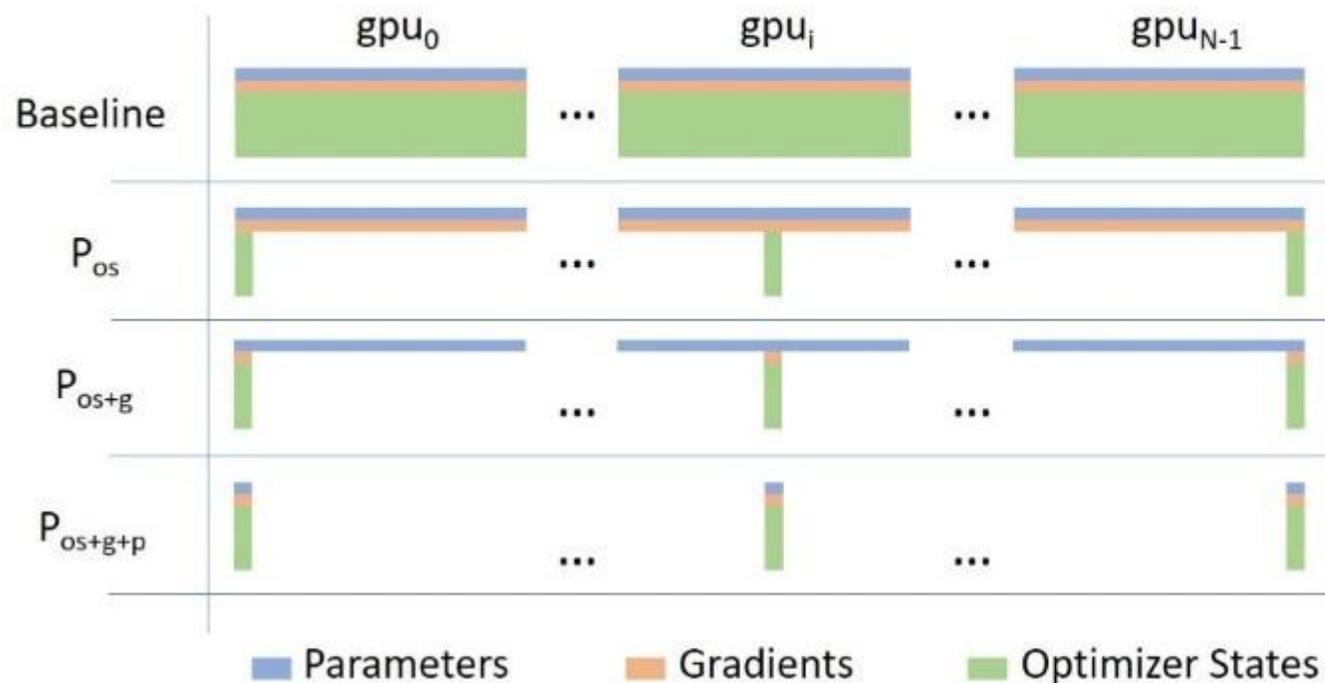
**Stage 1 ( $P_{os}$ )**

- ZeRO removes the redundancy across data parallel process
- Partitioning optimizer states, gradients and parameters (3 stages)



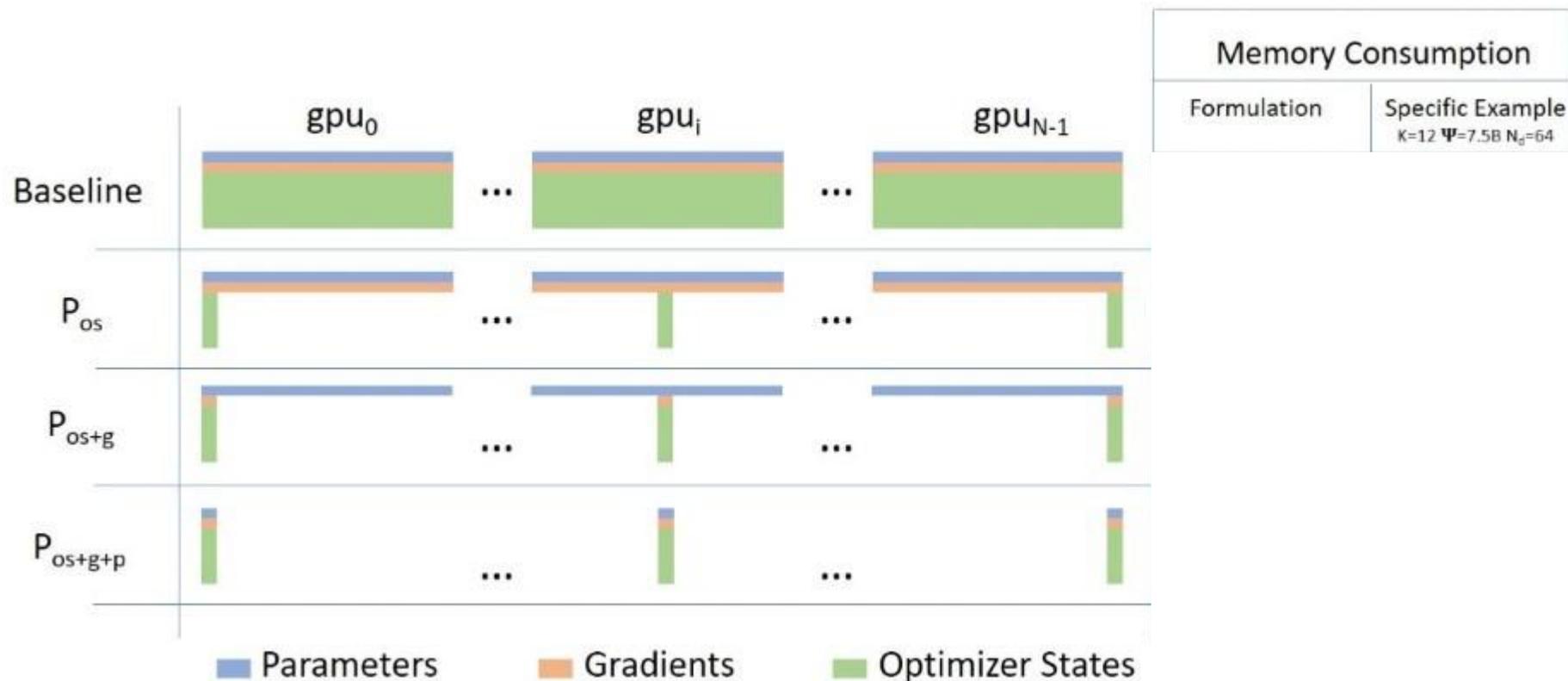
**Stage 2 (P<sub>os+g</sub>)**

- ZeRO removes the redundancy across data parallel process
- Partitioning optimizer states, gradients and parameters (3 stages)

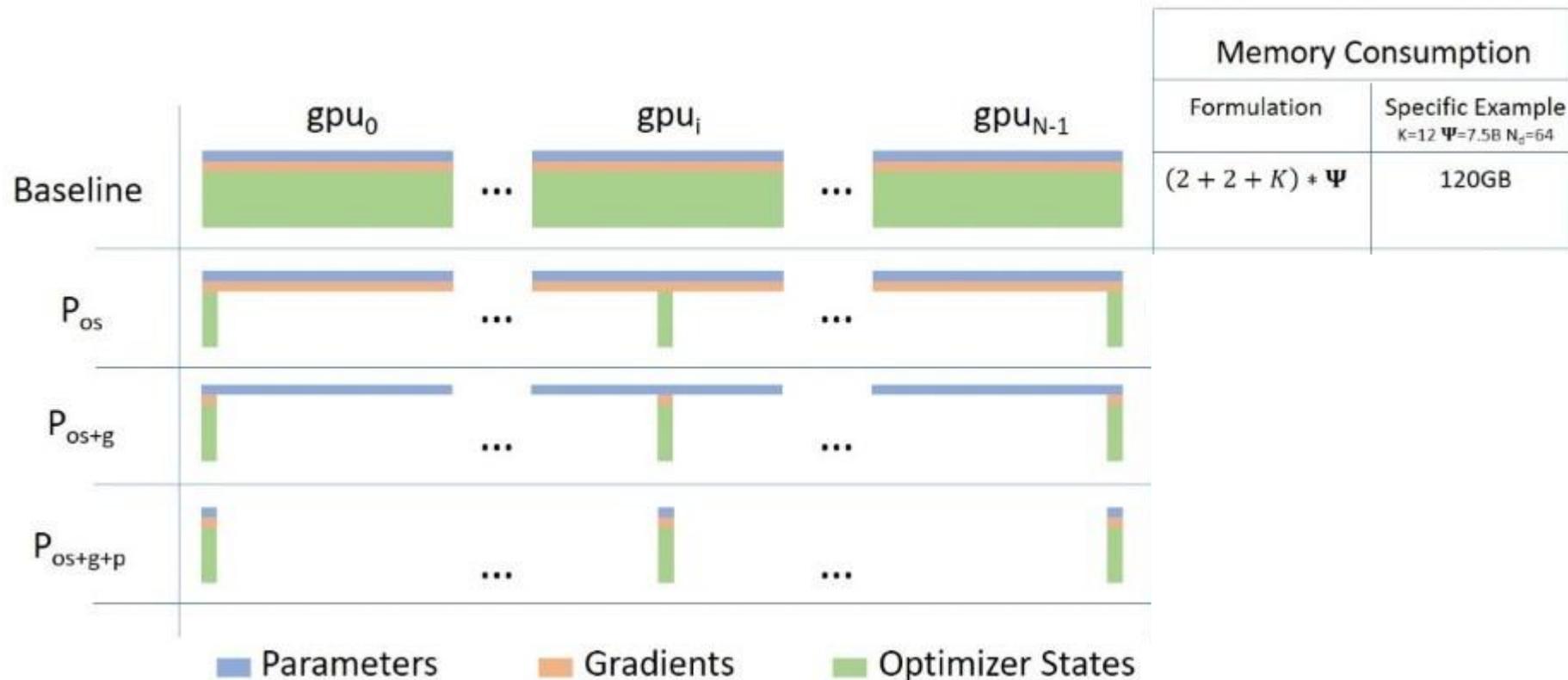


**Stage 3 (P<sub>os+g+p</sub>)**

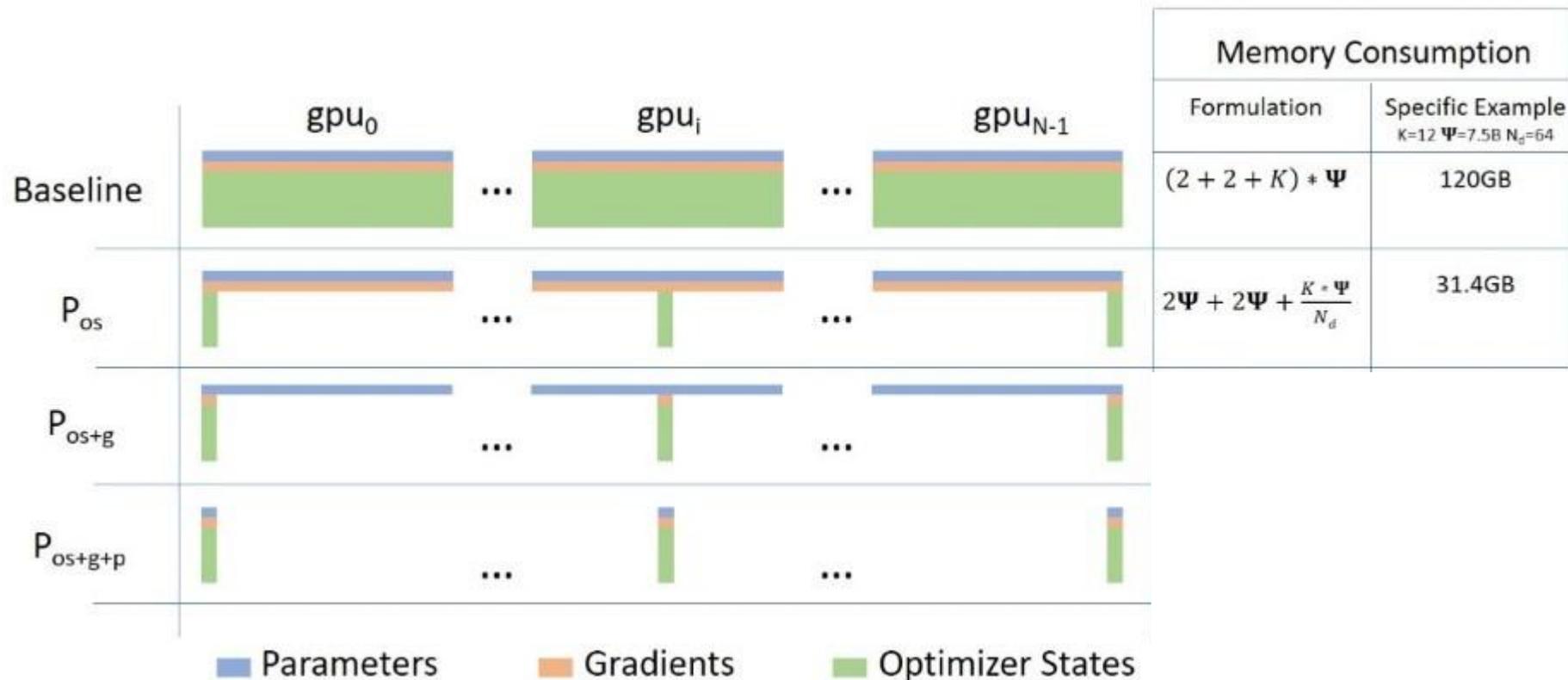
- ZeRO removes the redundancy across data parallel process
- Partitioning optimizer states, gradients and parameters (3 stages)



- ZeRO removes the redundancy across data parallel process
- Partitioning optimizer states, gradients and parameters (3 stages)



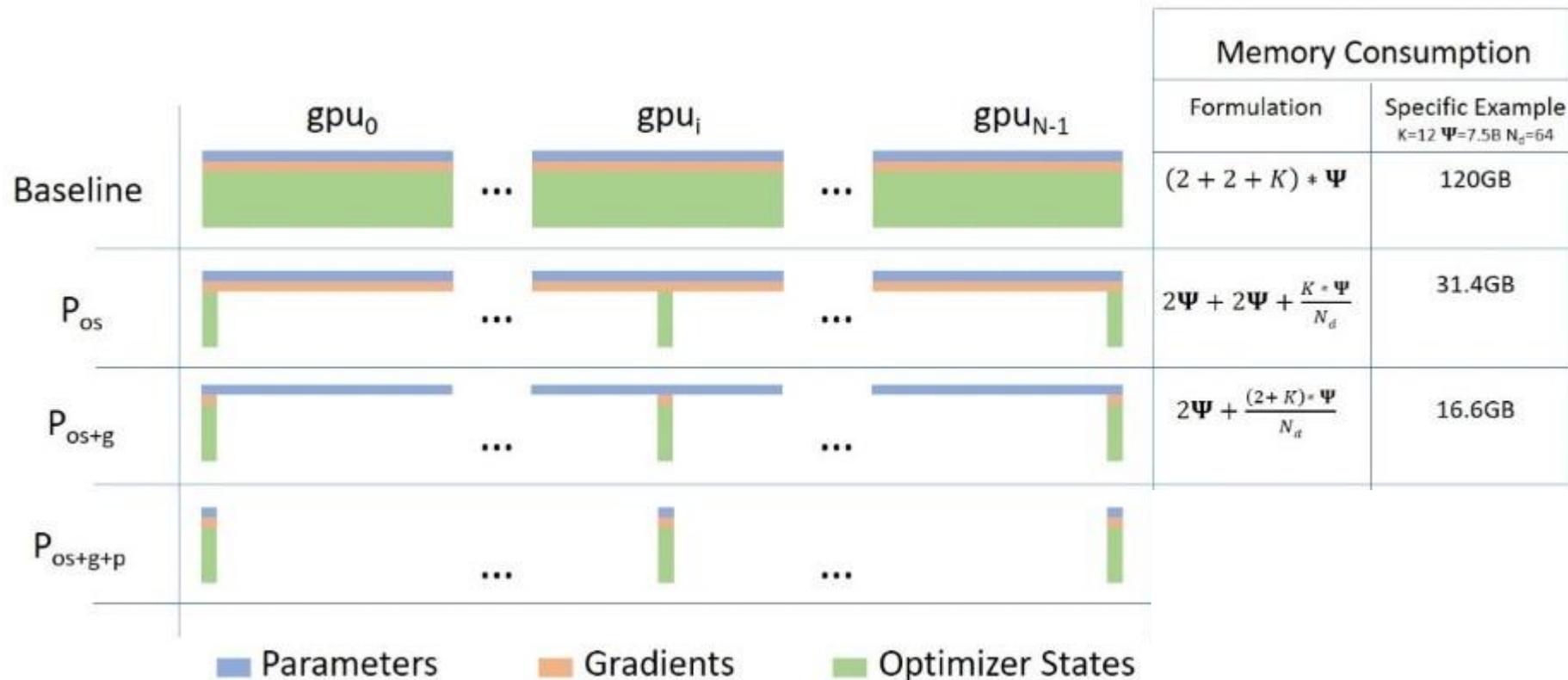
- ZeRO removes the redundancy across data parallel process
- Partitioning optimizer states, gradients and parameters (3 stages)



# ZeRO-DP: ZeRO powered Data Parallelism



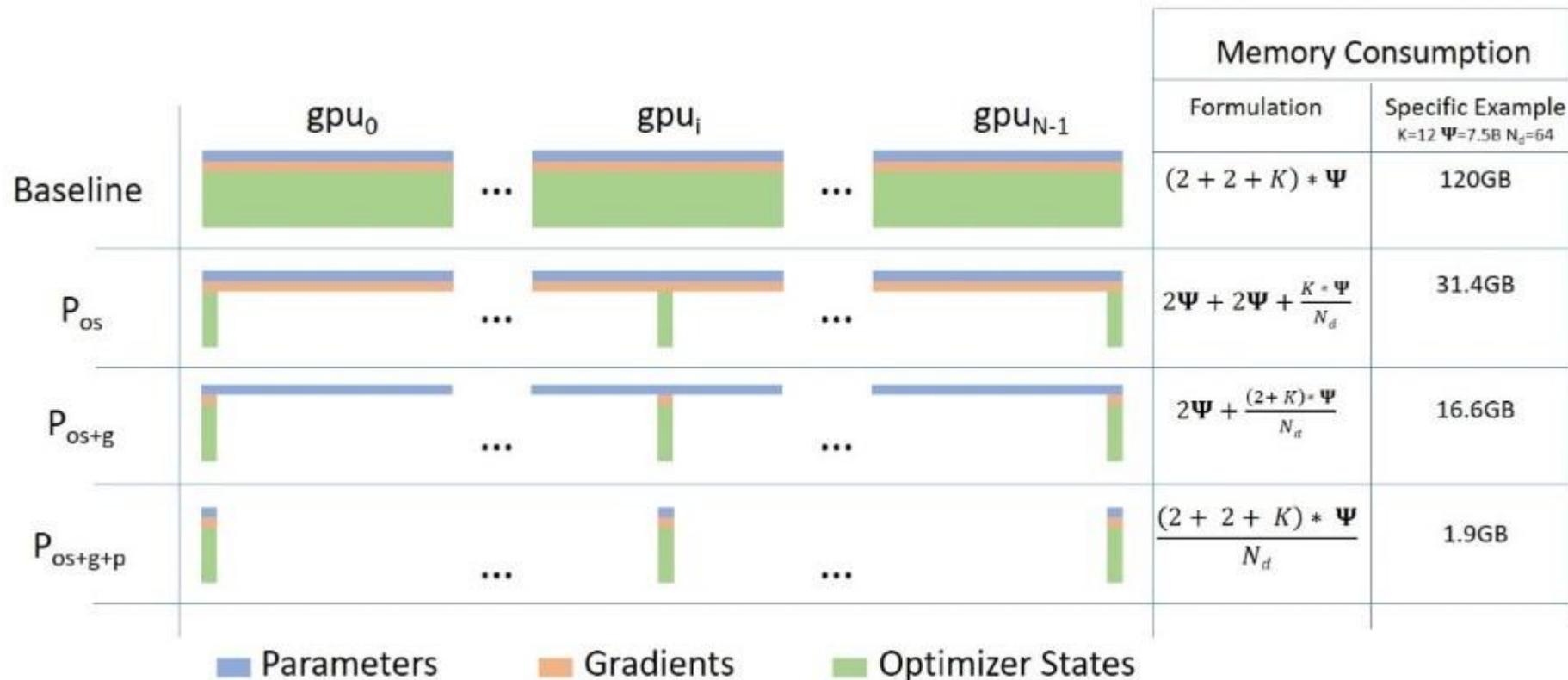
- ZeRO removes the redundancy across data parallel process
- Partitioning optimizer states, gradients and parameters (3 stages)



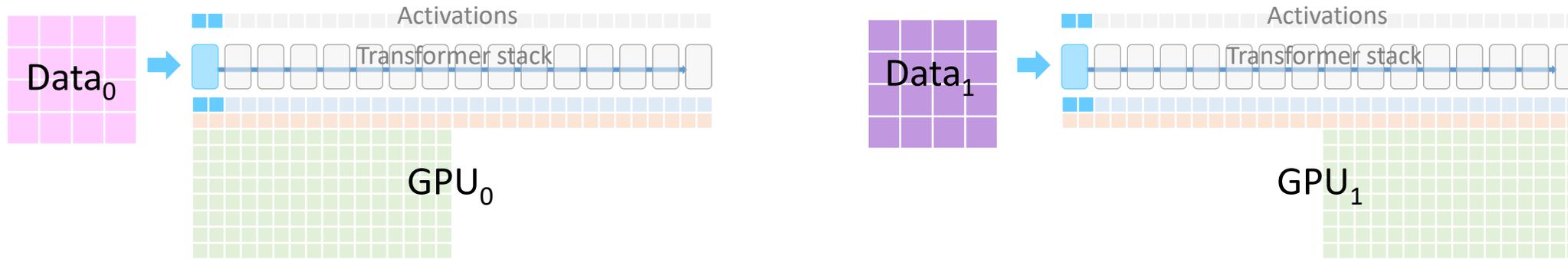
# ZeRO-DP: ZeRO powered Data Parallelism



- ZeRO removes the redundancy across data parallel process
- Partitioning optimizer states, gradients and parameters (3 stages)

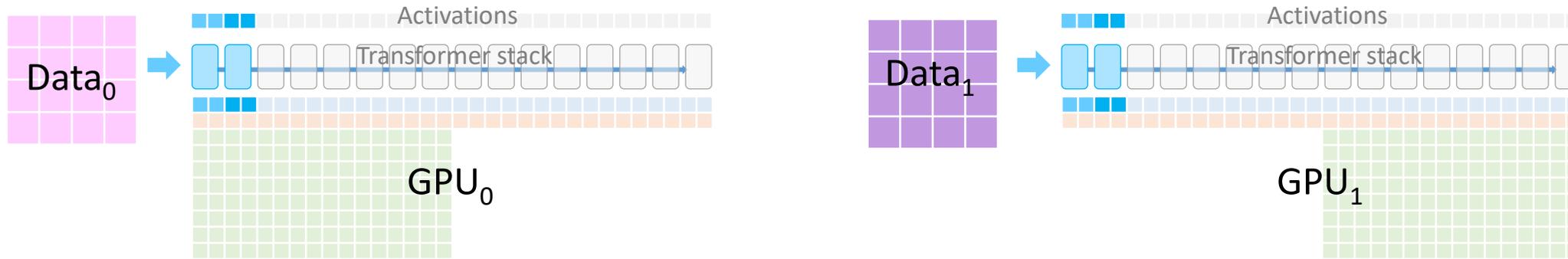


# ZeRO-DP: ZeRO powered Data Parallelism



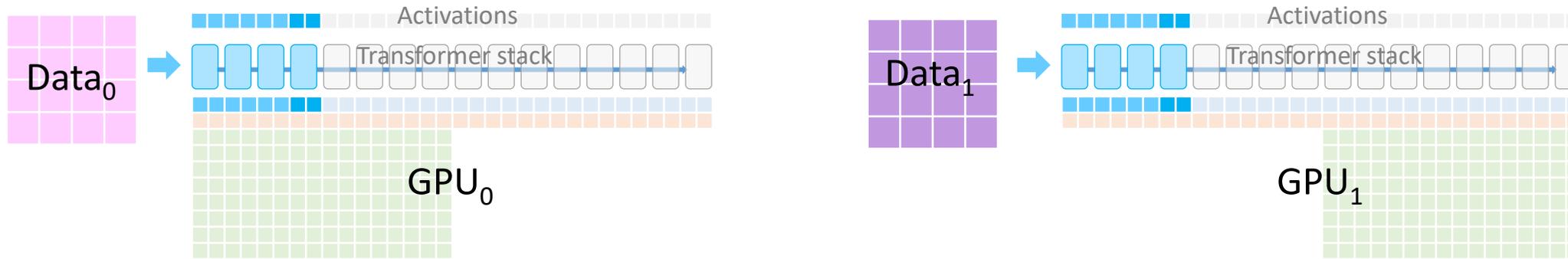
- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks

# ZeRO-DP: ZeRO powered Data Parallelism



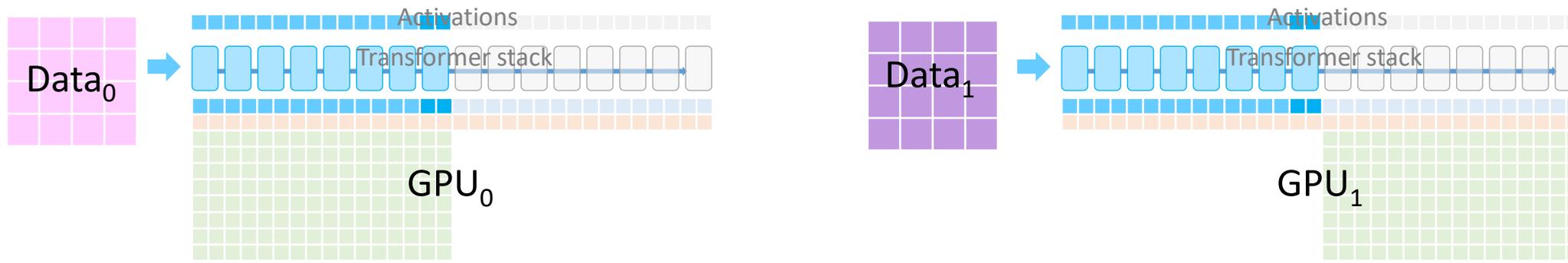
- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks

# ZeRO-DP: ZeRO powered Data Parallelism



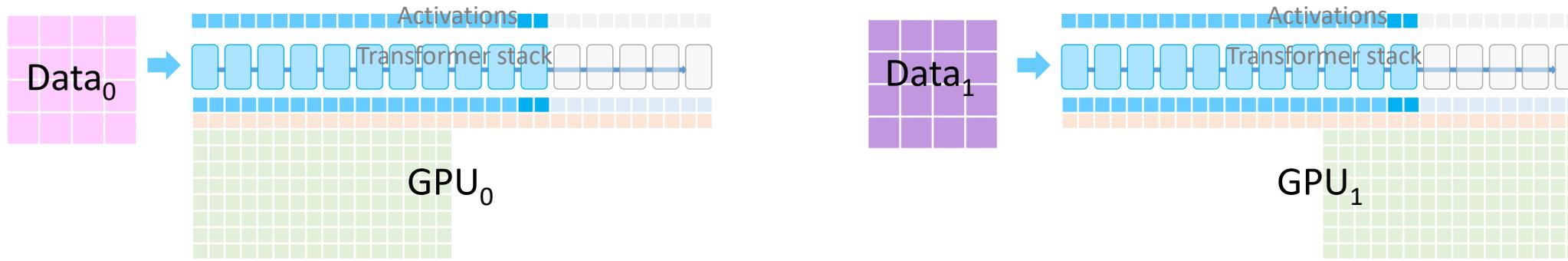
- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks

# ZeRO-DP: ZeRO powered Data Parallelism



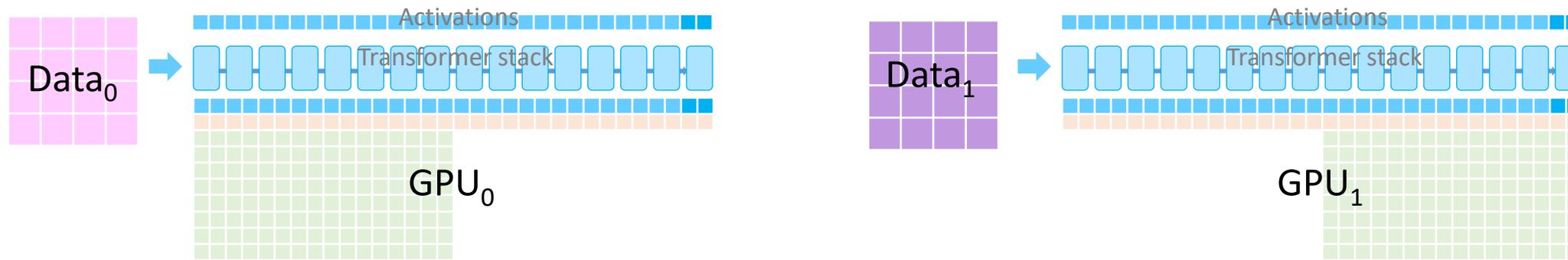
- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks

# ZeRO-DP: ZeRO powered Data Parallelism



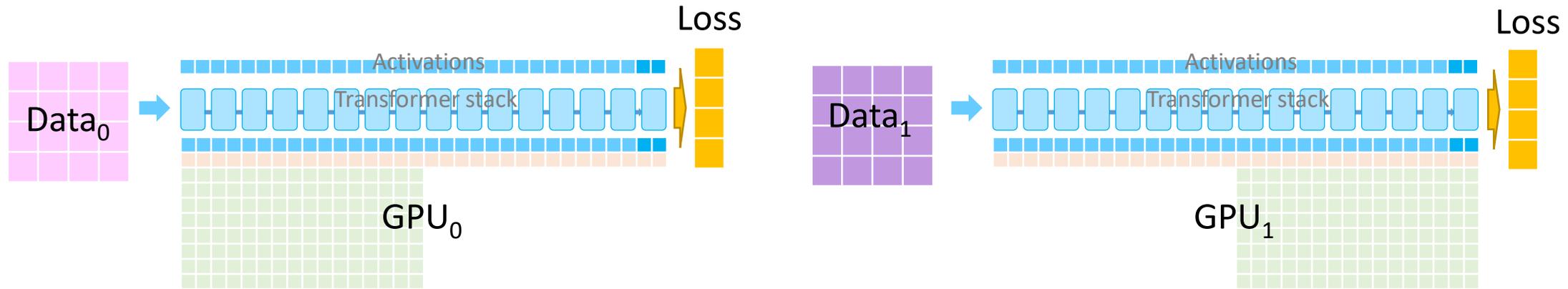
- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks

# ZeRO-DP: ZeRO powered Data Parallelism



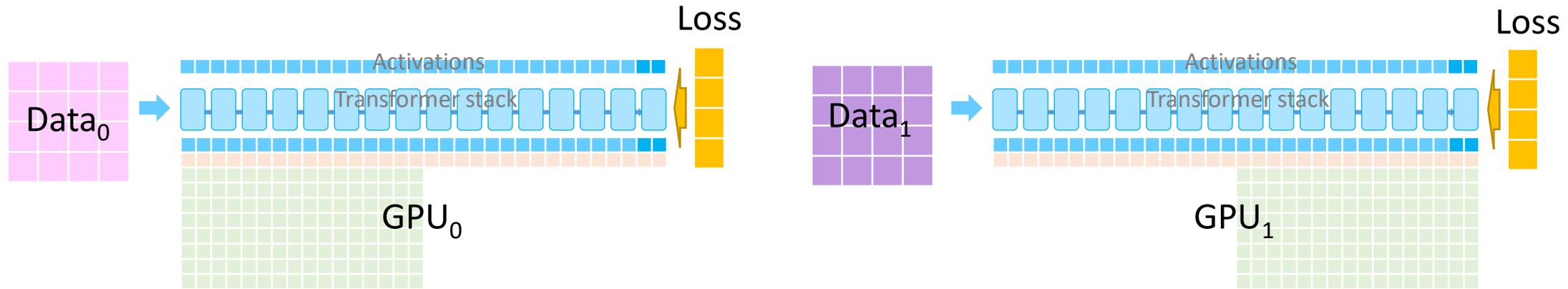
- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks

# ZeRO-DP: ZeRO powered Data Parallelism



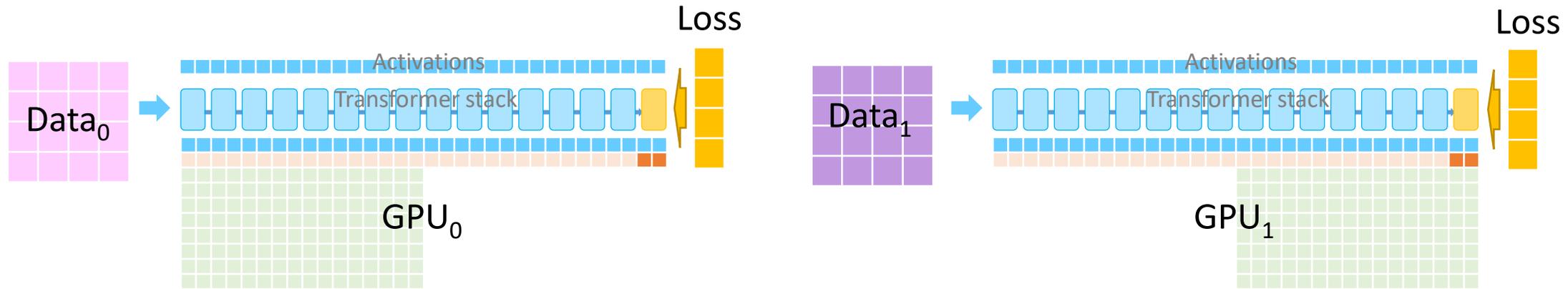
- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks

# ZeRO-DP: ZeRO powered Data Parallelism



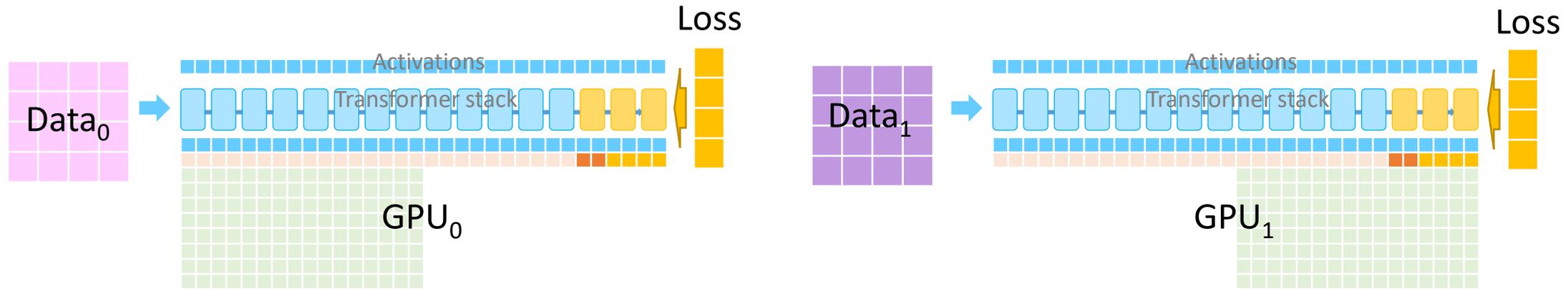
- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks
- Backward propagation to generate FP16 gradients

# ZeRO-DP: ZeRO powered Data Parallelism



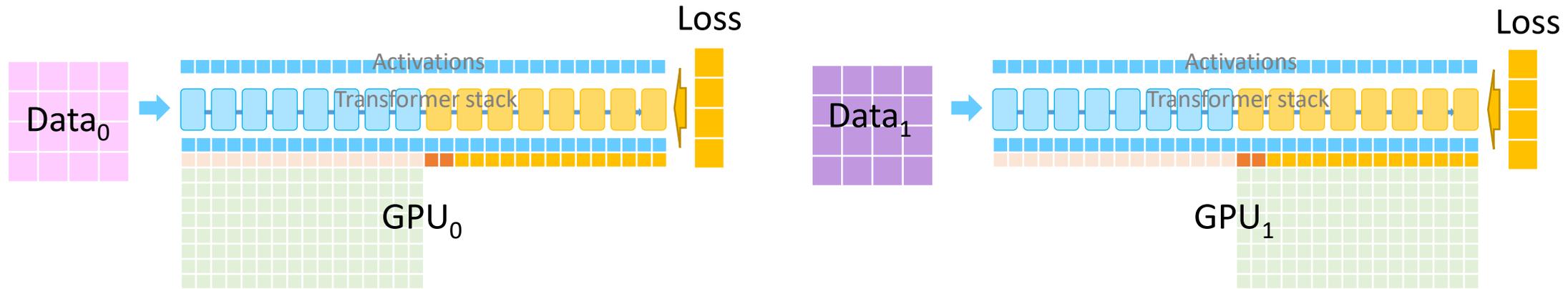
- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks
- Backward propagation to generate FP16 gradients

# ZeRO-DP: ZeRO powered Data Parallelism



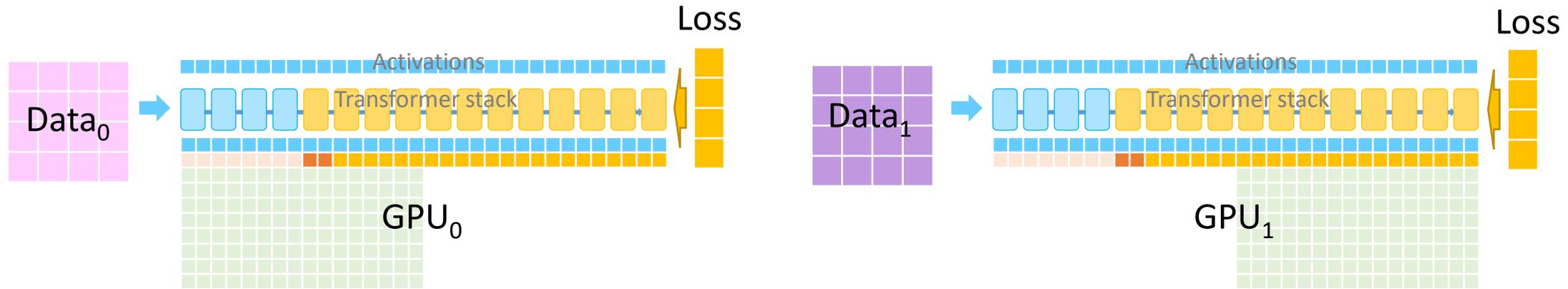
- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks
- Backward propagation to generate FP16 gradients

# ZeRO-DP: ZeRO powered Data Parallelism



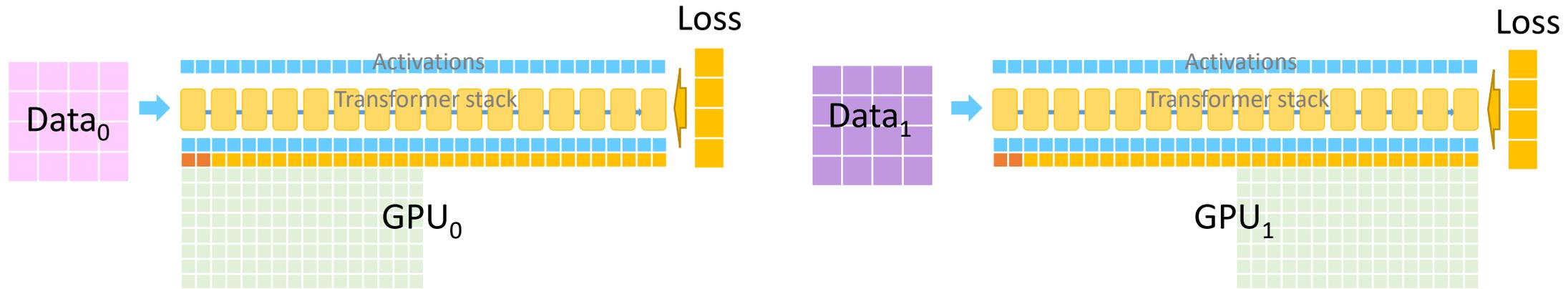
- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks
- Backward propagation to generate FP16 gradients

# ZeRO-DP: ZeRO powered Data Parallelism



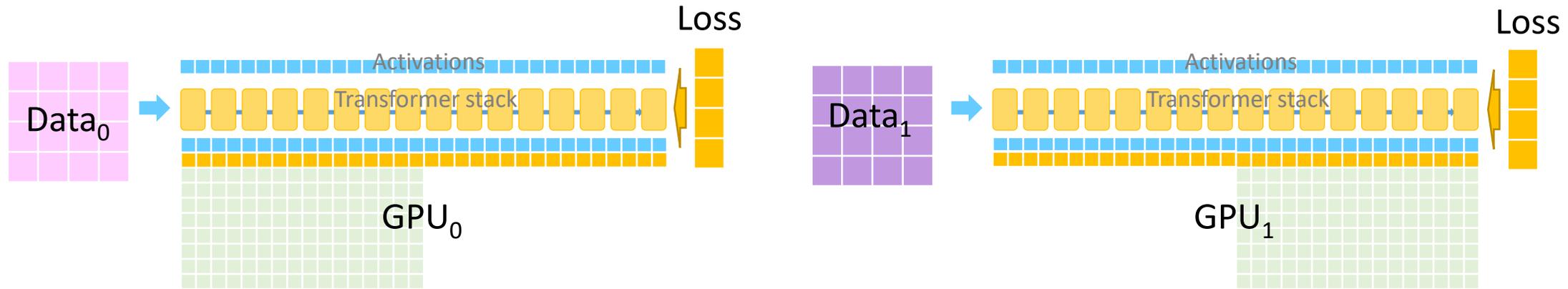
- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks
- Backward propagation to generate FP16 gradients

# ZeRO-DP: ZeRO powered Data Parallelism



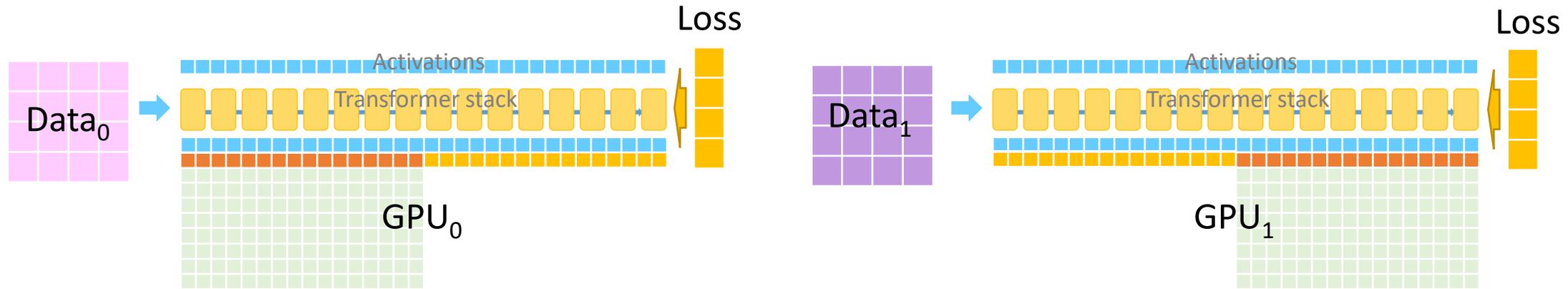
- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks
- Backward propagation to generate FP16 gradients

# ZeRO-DP: ZeRO powered Data Parallelism

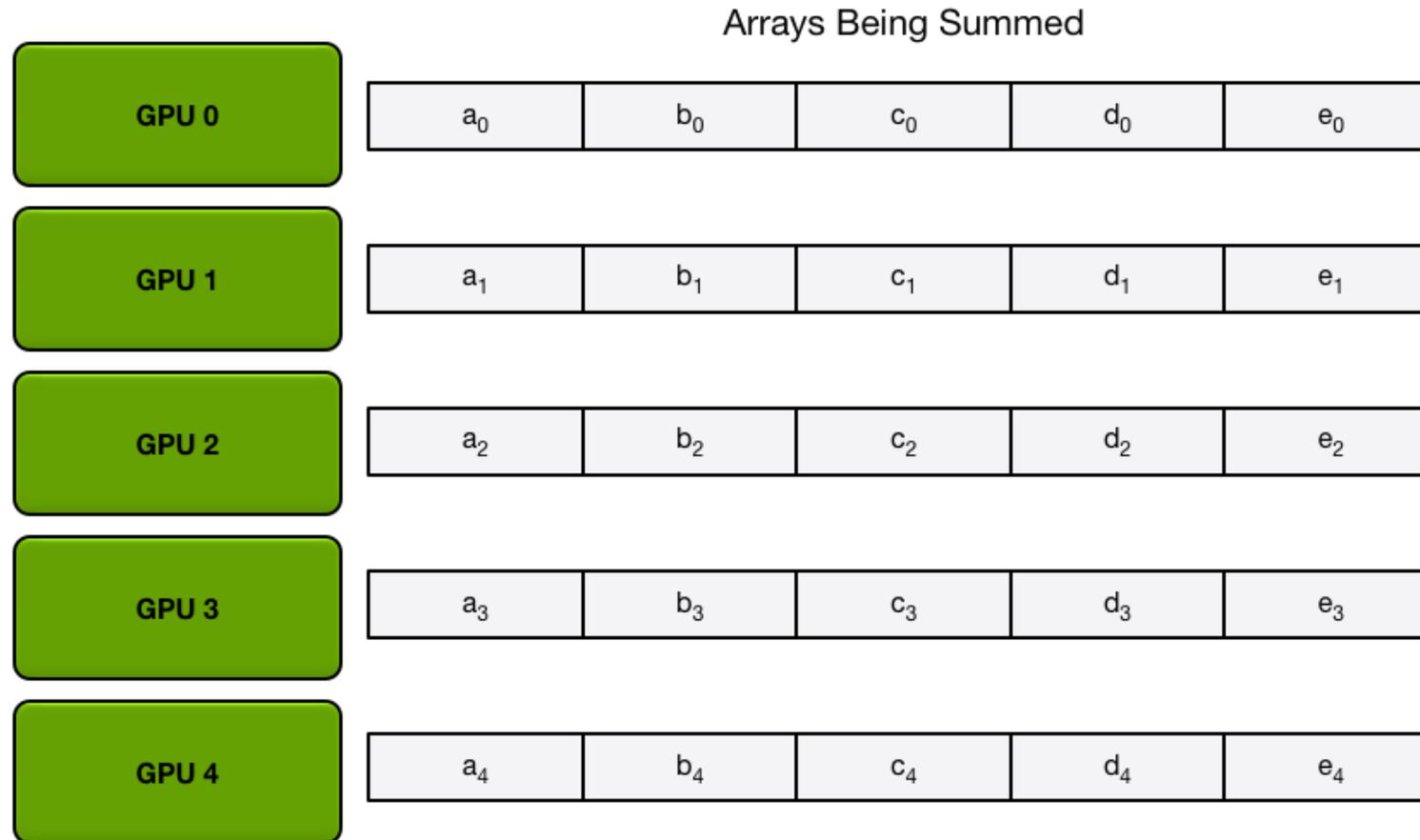


- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks
- Backward propagation to generate FP16 gradients and reduce scatter to average

# ZeRO-DP: ZeRO powered Data Parallelism

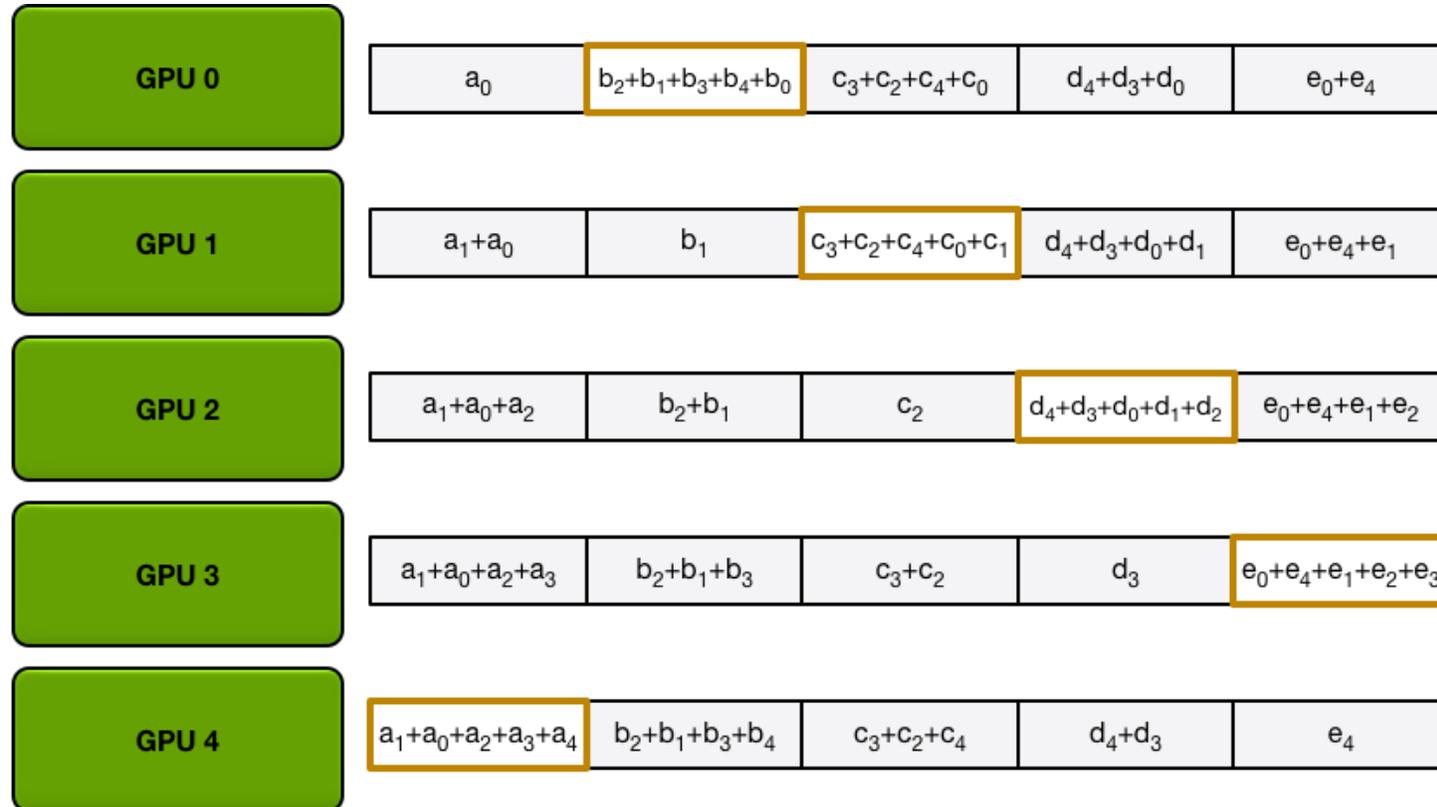


- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks
- Backward propagation to generate FP16 gradients and reduce scatter to average



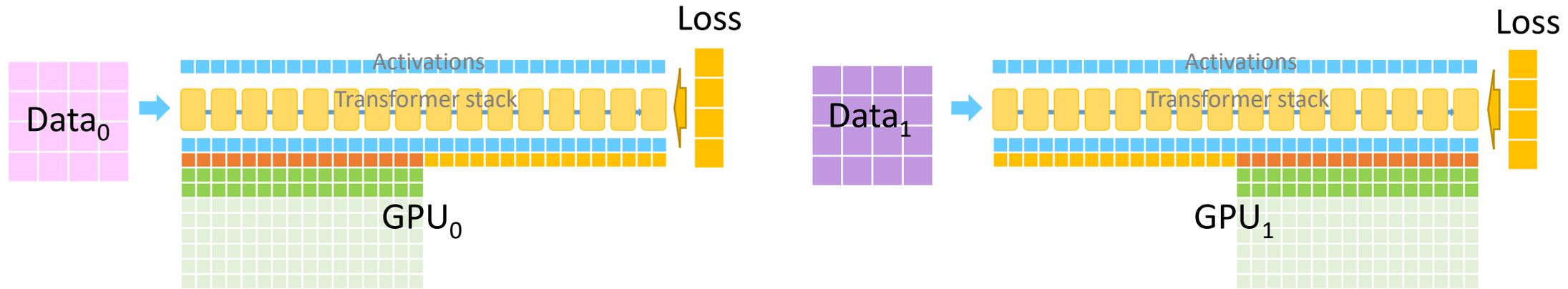
Partitioning of an array into N chunks

# Recall: Reduce-Scatter



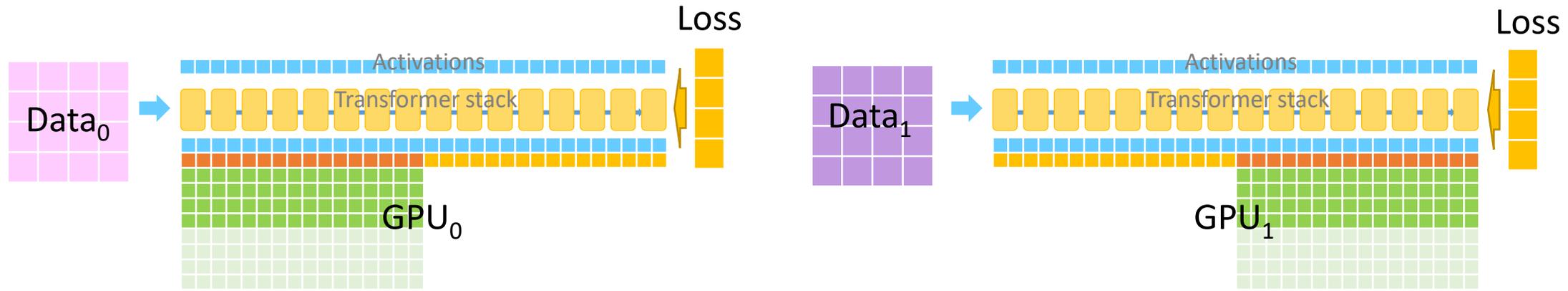
Reduce-scatter data transfer (after iteration 4)

# ZeRO-DP: ZeRO powered Data Parallelism



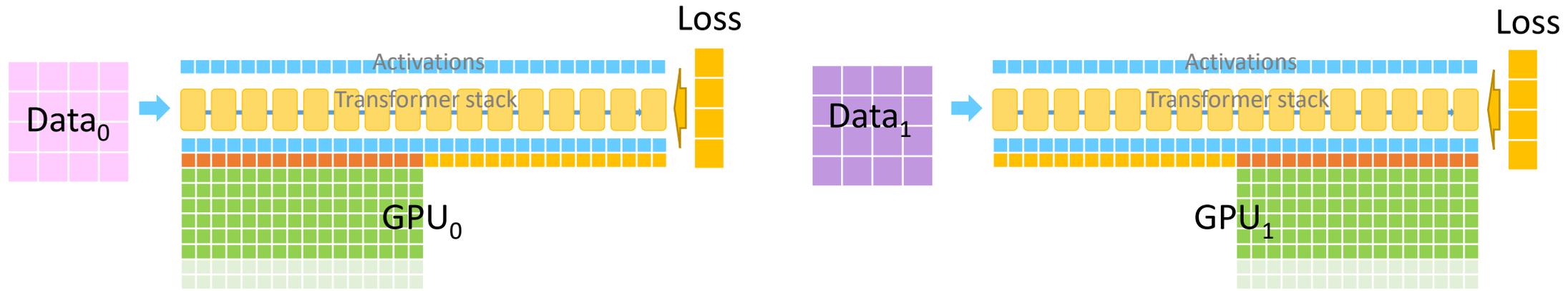
- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks
- Backward propagation to generate FP16 gradients and reduce scatter to average
- Update the FP32 weights with ADAM optimizer

# ZeRO-DP: ZeRO powered Data Parallelism



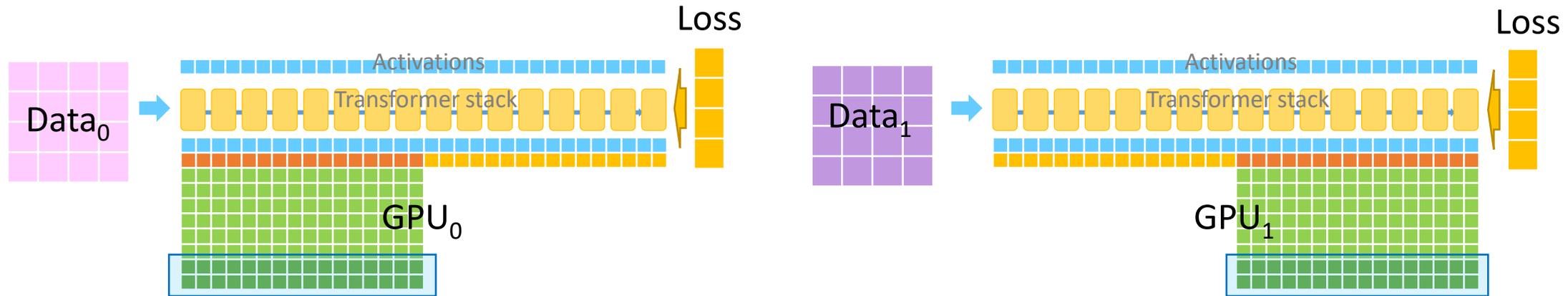
- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks
- Backward propagation to generate FP16 gradients and reduce scatter to average
- Update the FP32 weights with ADAM optimizer

# ZeRO-DP: ZeRO powered Data Parallelism



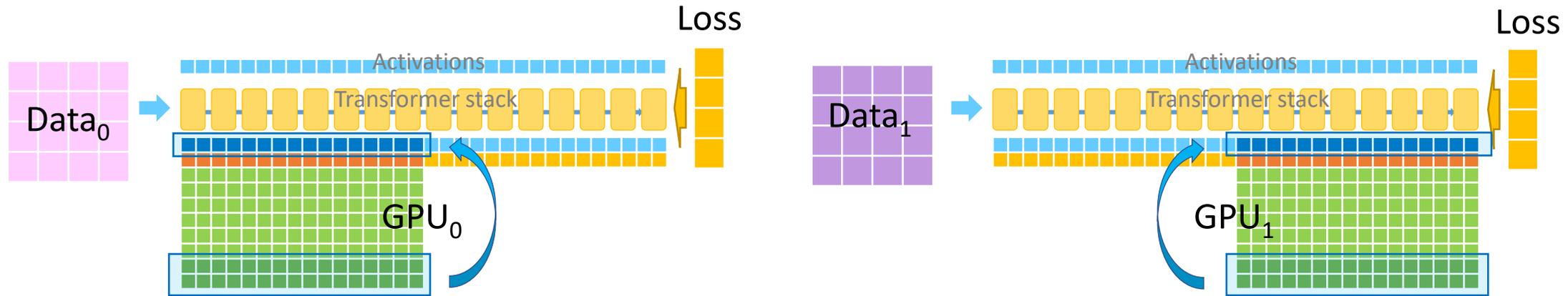
- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks
- Backward propagation to generate FP16 gradients and reduce scatter to average
- Update the FP32 weights with ADAM optimizer

# ZeRO-DP: ZeRO powered Data Parallelism



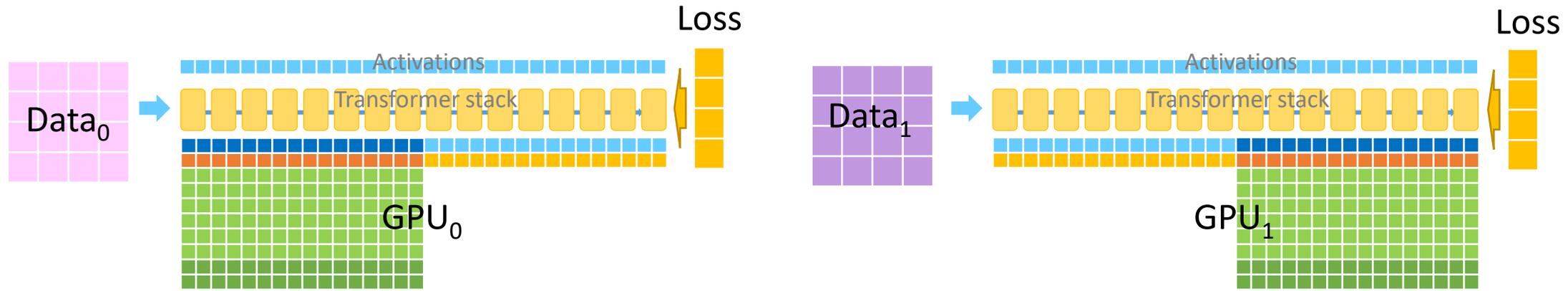
- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks
- Backward propagation to generate FP16 gradients and reduce scatter to average
- Update the FP32 weights with ADAM optimizer

# ZeRO-DP: ZeRO powered Data Parallelism



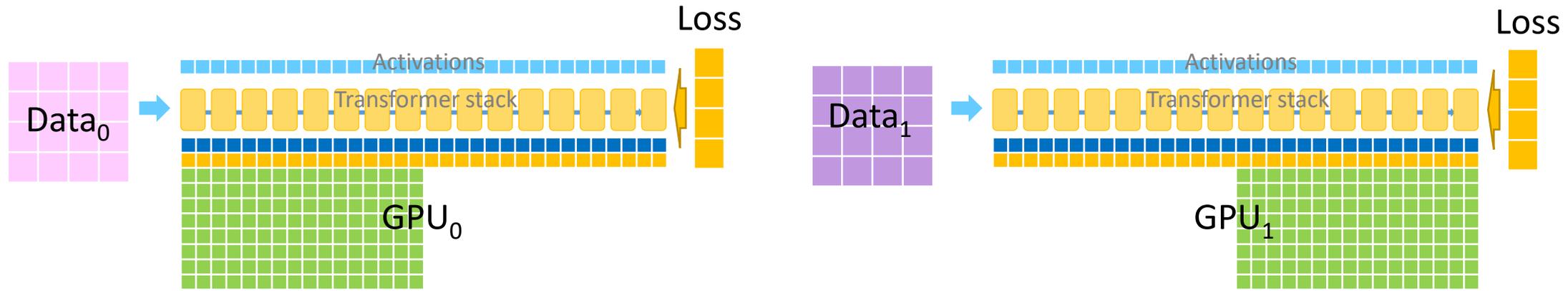
- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks
- Backward propagation to generate FP16 gradients and reduce scatter to average
- Update the FP32 weights with ADAM optimizer
- Update the FP16 weights

# ZeRO-DP: ZeRO powered Data Parallelism



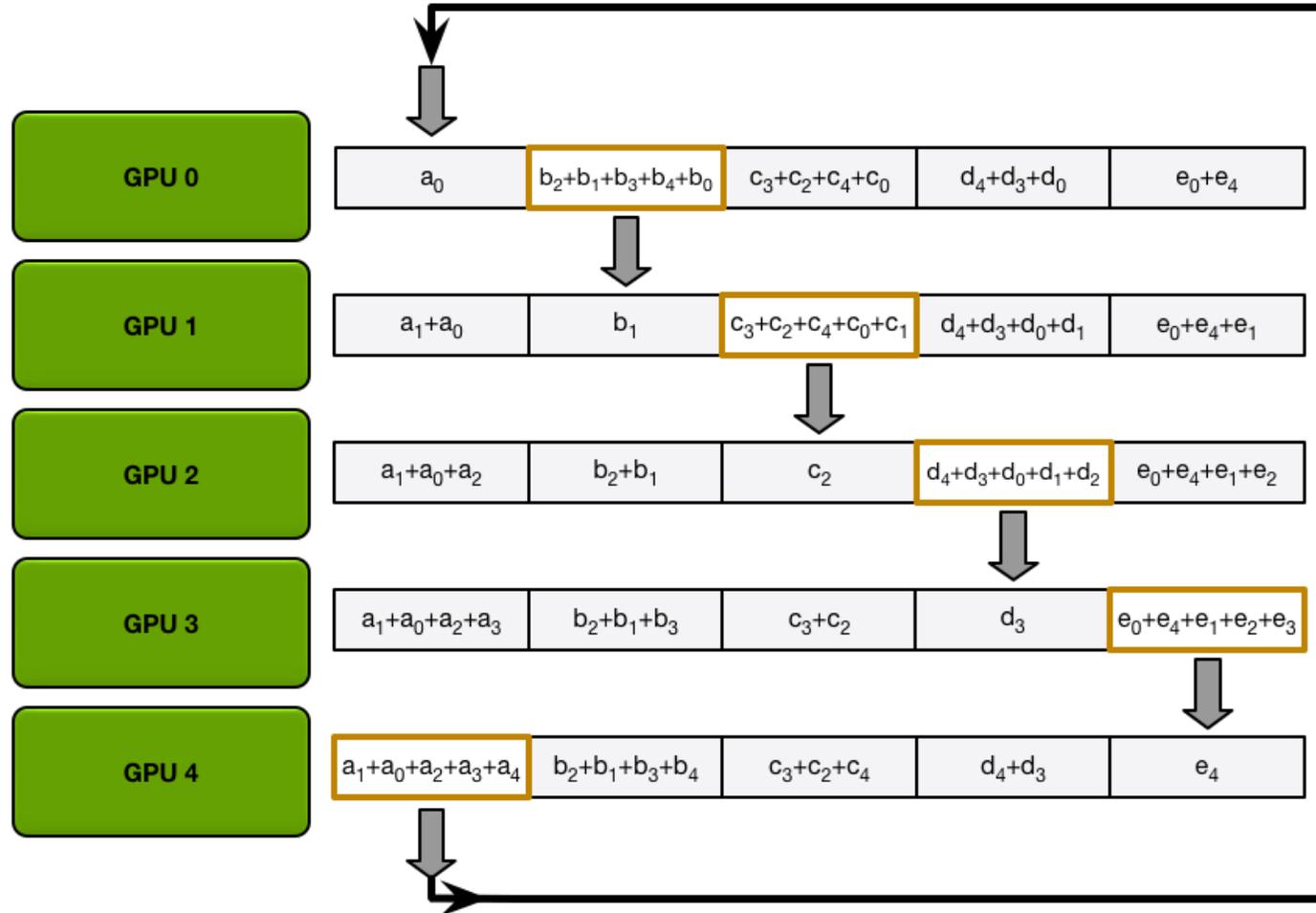
- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks
- Backward propagation to generate FP16 gradients and reduce scatter to average
- Update the FP32 weights with ADAM optimizer
- Update the FP16 weights
- All Gather the FP16 weights to complete the iteration

# ZeRO-DP: ZeRO powered Data Parallelism

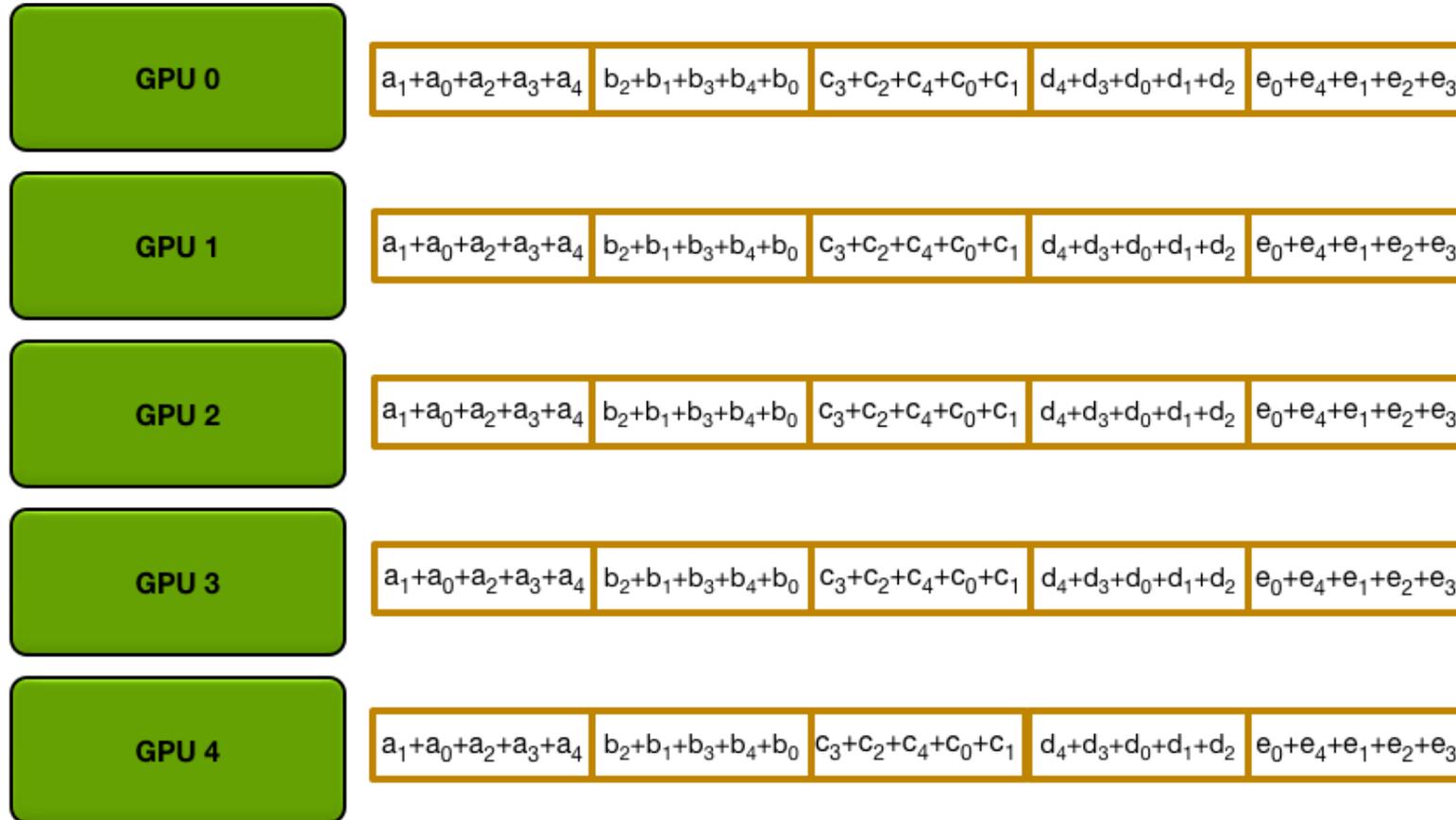


- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks
- Backward propagation to generate FP16 gradients and reduce scatter to average
- Update the FP32 weights with ADAM optimizer
- Update the FP16 weights
- All Gather the FP16 weights to complete the iteration

# Recall: Allgather

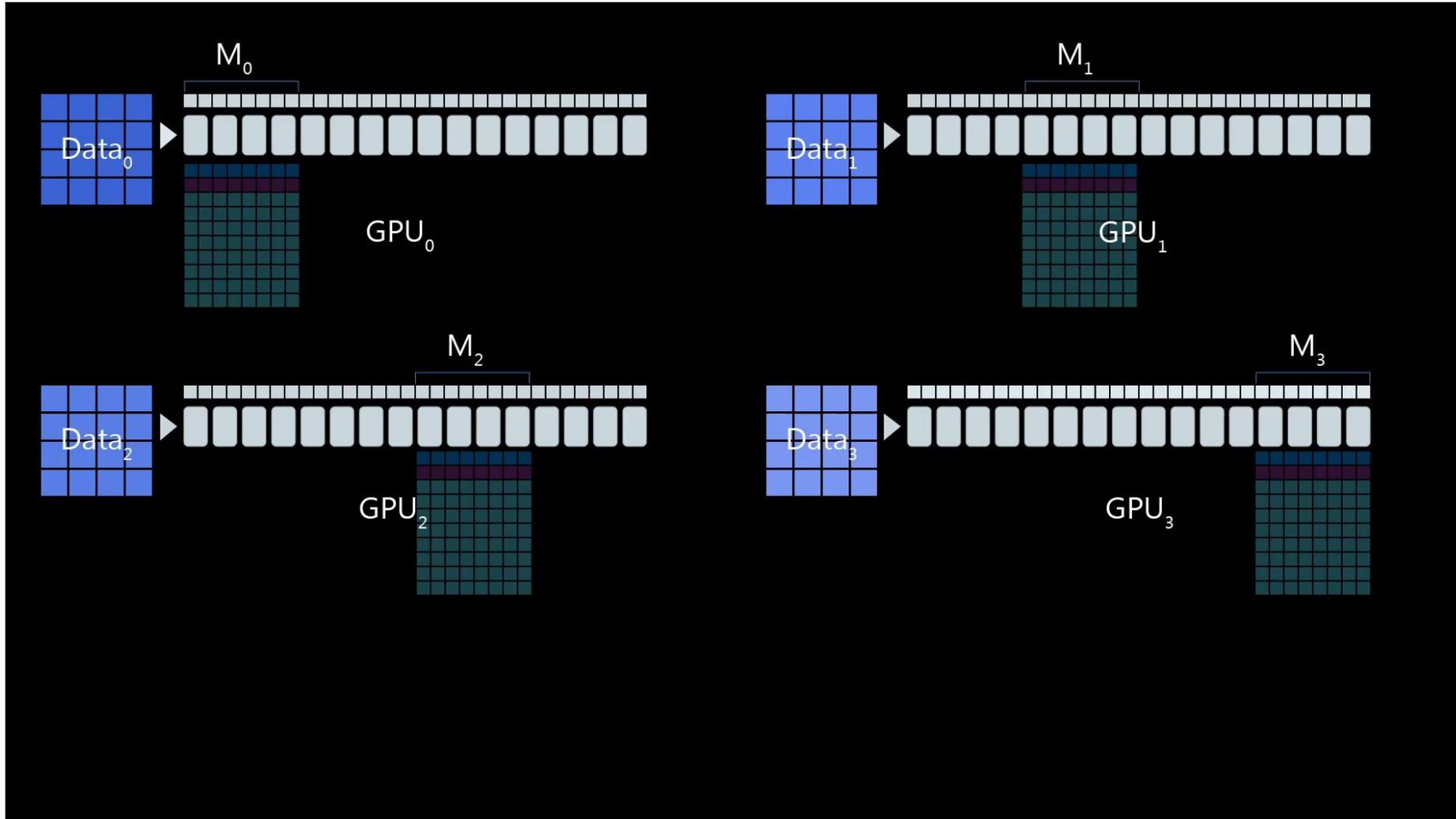


All processes can obtain the complete array by circulating again without reduction operations.



Final state after allgather transfer

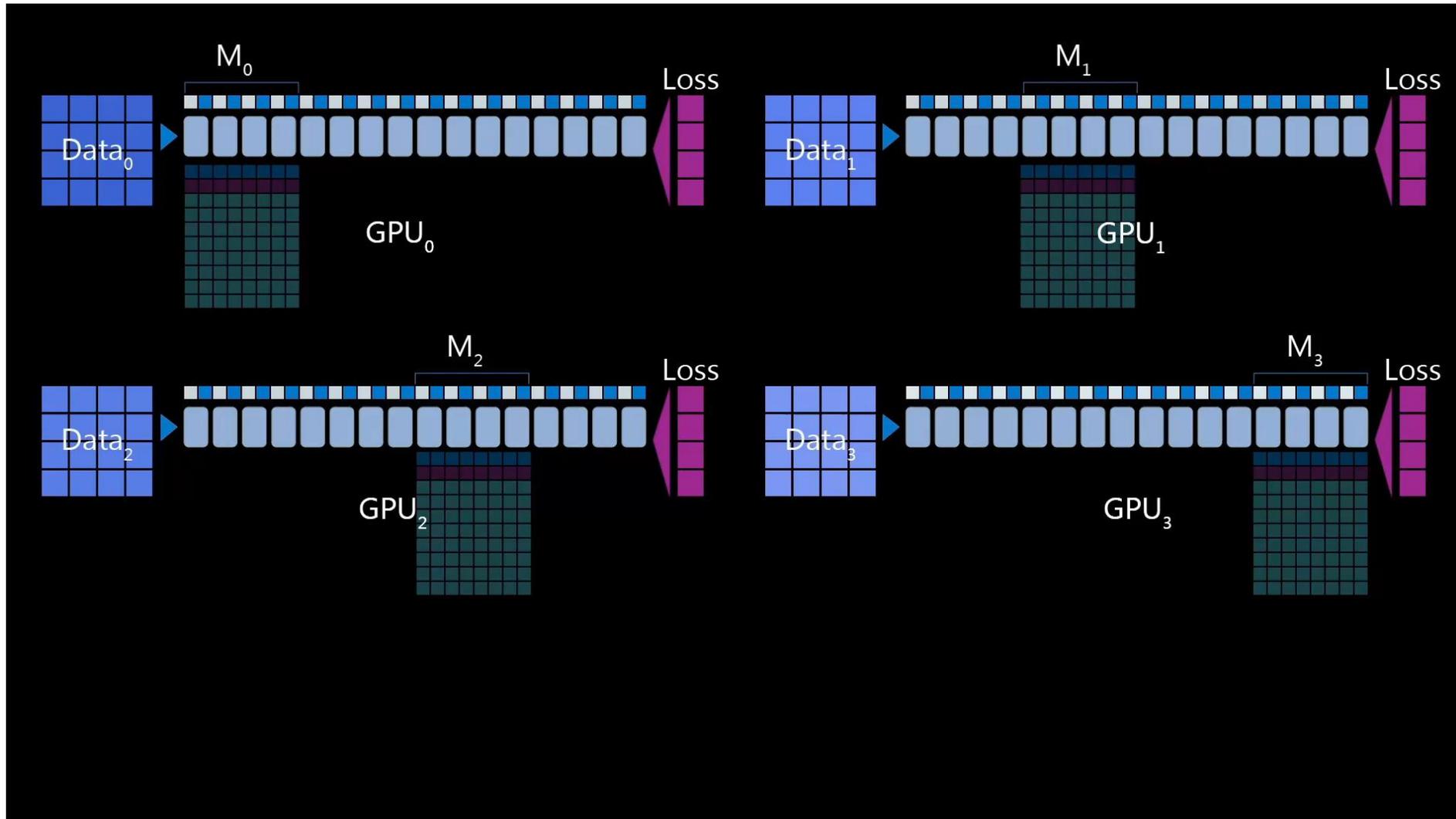
# ZeRO-DP : Stage 3 Forward Propagation



Legend for memory buffer components:

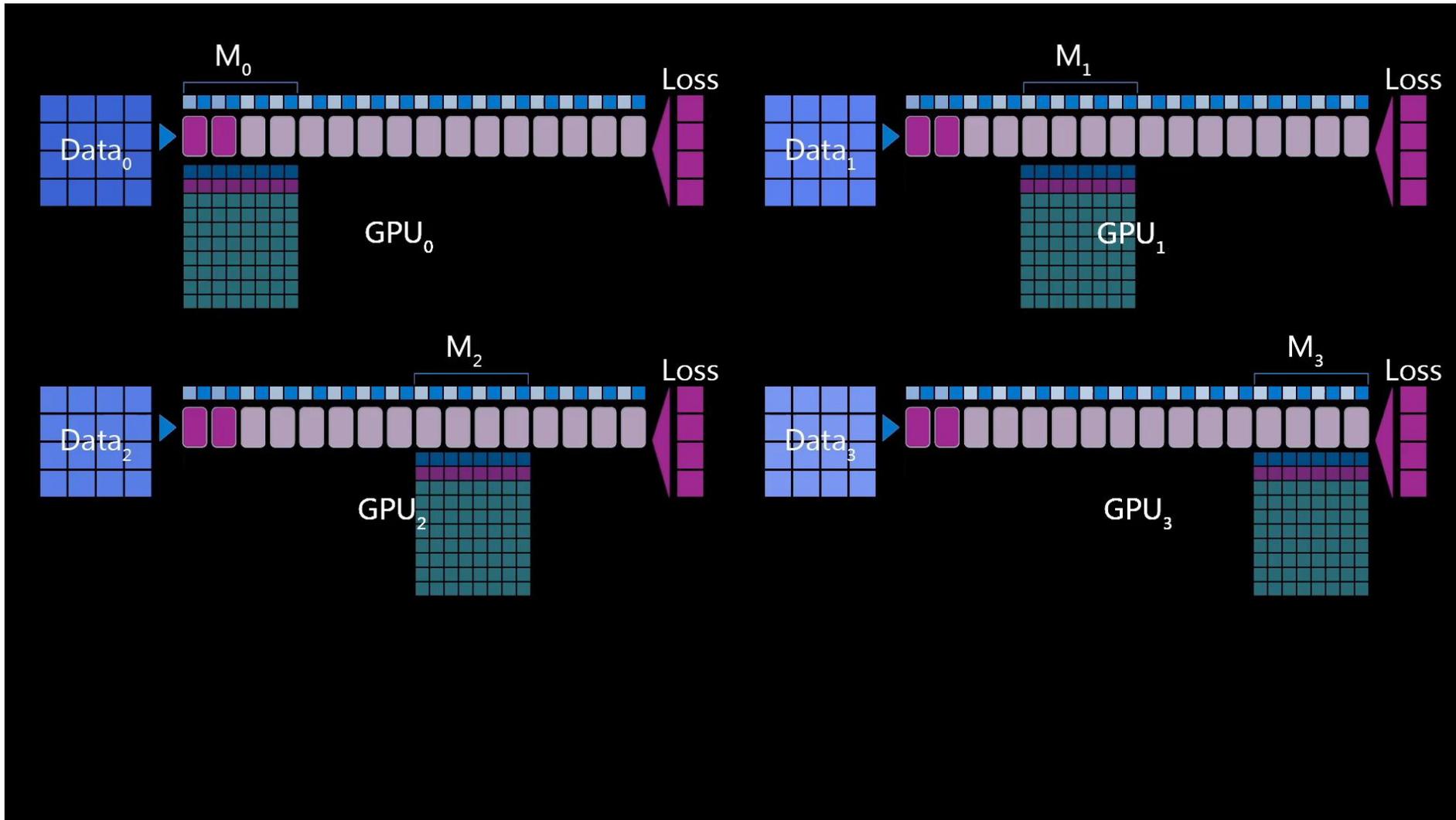
- fp16 params
- fp16 grads
- fp32 grads
- fp32 momentum
- fp32 variance
- fp32 params

# ZeRO-DP : Stage 3 Backward Propagation



- fp16 params
- fp16 grads
- fp32 grads
- fp32 momentum
- fp32 variance
- fp32 params

# ZeRO-DP : Stage 3 Optimizer Step



- fp16 params
- fp16 grads
- fp32 grads
- fp32 momentum
- fp32 variance
- fp32 params

- Progressive memory savings and Communication Volume

	Memory Reduction with N GPUs
Data Parallel	1x
ZeRO Stage 1 ( $P_{os}$ )	4x
ZeRO Stage 2 ( $P_{os+g}$ )	8x
ZeRO Stage 3 ( $P_{os+g+p}$ )	Nx

\*Mixed precision Adam on Cluster of DGX-2 with NVIDIA 32 GB V100 GPUs

- Progressive memory savings and Communication Volume

	Memory Reduction with N GPUs	Max params with ZeRO only (in billions)
Data Parallel	1x	1.2
ZeRO Stage 1 ( $P_{os}$ )	4x	7
ZeRO Stage 2 ( $P_{os+g}$ )	8x	14
ZeRO Stage 3 ( $P_{os+g+p}$ )	Nx	>1000

\*Mixed precision Adam on Cluster of DGX-2 with NVIDIA 32 GB V100 GPUs

- Progressive memory savings and Communication Volume

	Memory Reduction with N GPUs	Max params with ZeRO only (in billions)	Max params with ZeRO and model parallelism (in billions)
Data Parallel	1x	1.2	20
ZeRO Stage 1 ( $P_{os}$ )	4x	7	100
ZeRO Stage 2 ( $P_{os+g}$ )	8x	14	200
ZeRO Stage 3 ( $P_{os+g+p}$ )	Nx	>1000	>1000

\*Mixed precision Adam on Cluster of DGX-2 with NVIDIA 32 GB V100 GPUs

- Progressive memory savings and Communication Volume

	Memory Reduction with N GPUs	Max params with ZeRO only (in billions)	Max params with ZeRO and model parallelism (in billions)	Comm Volume
Data Parallel	1x	1.2	20	1x
ZeRO Stage 1 ( $P_{os}$ )	4x	7	100	1x
ZeRO Stage 2 ( $P_{os+g}$ )	8x	14	200	1x
ZeRO Stage 3 ( $P_{os+g+p}$ )	Nx	>1000	>1000	1.5x

\*Mixed precision Adam on Cluster of DGX-2 with NVIDIA 32 GB V100 GPUs



```
# construct torch.nn.Module
model = MyModel()

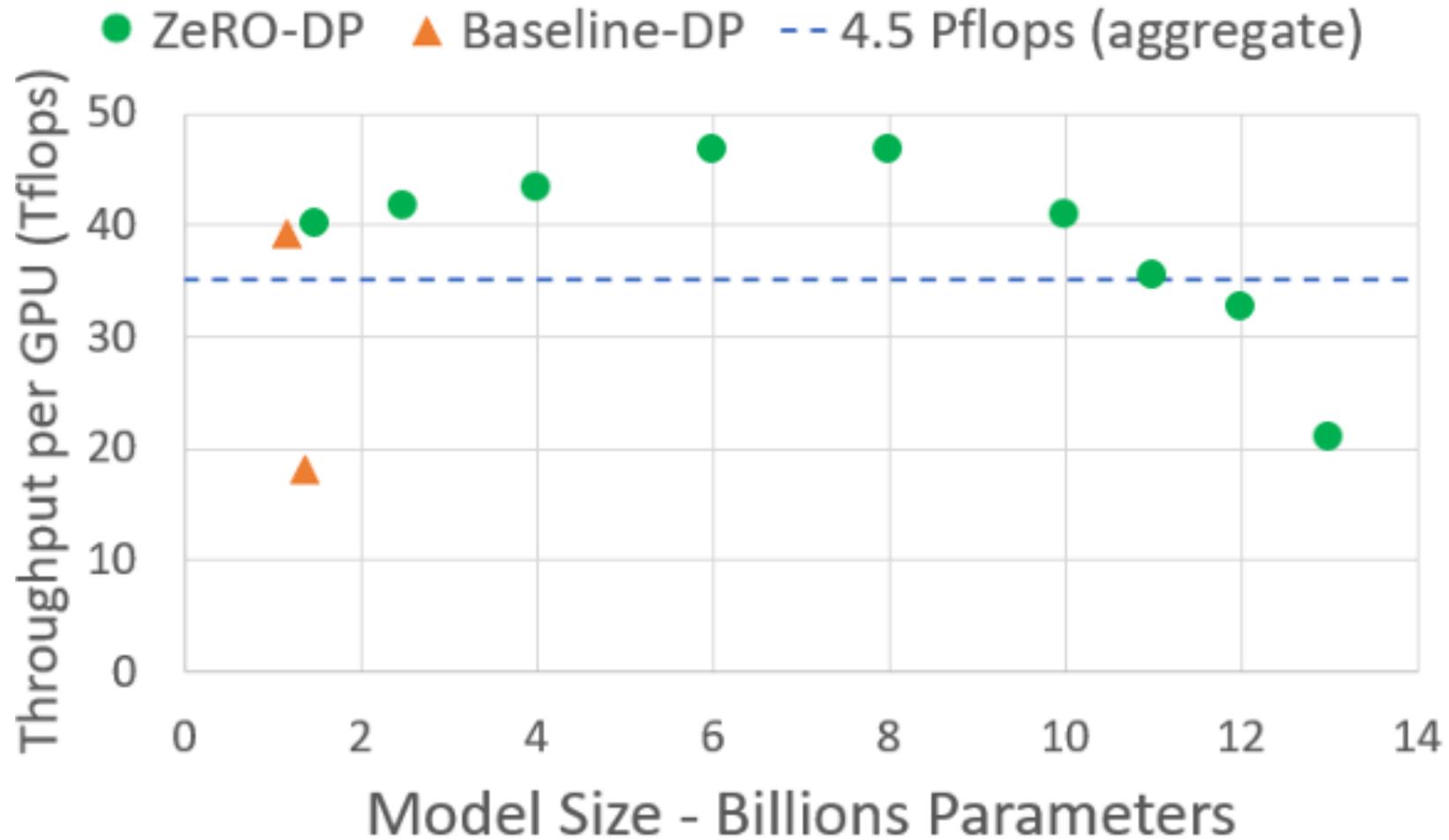
# wrap w. DeepSpeed engine
engine, *_ = deepspeed.initialize(
    model=model,
    config=ds_config

# training-loop w.r.t. engine
for batch in data_loader:
    loss = engine(batch)
    engine.backward(loss)
    engine.step()
```

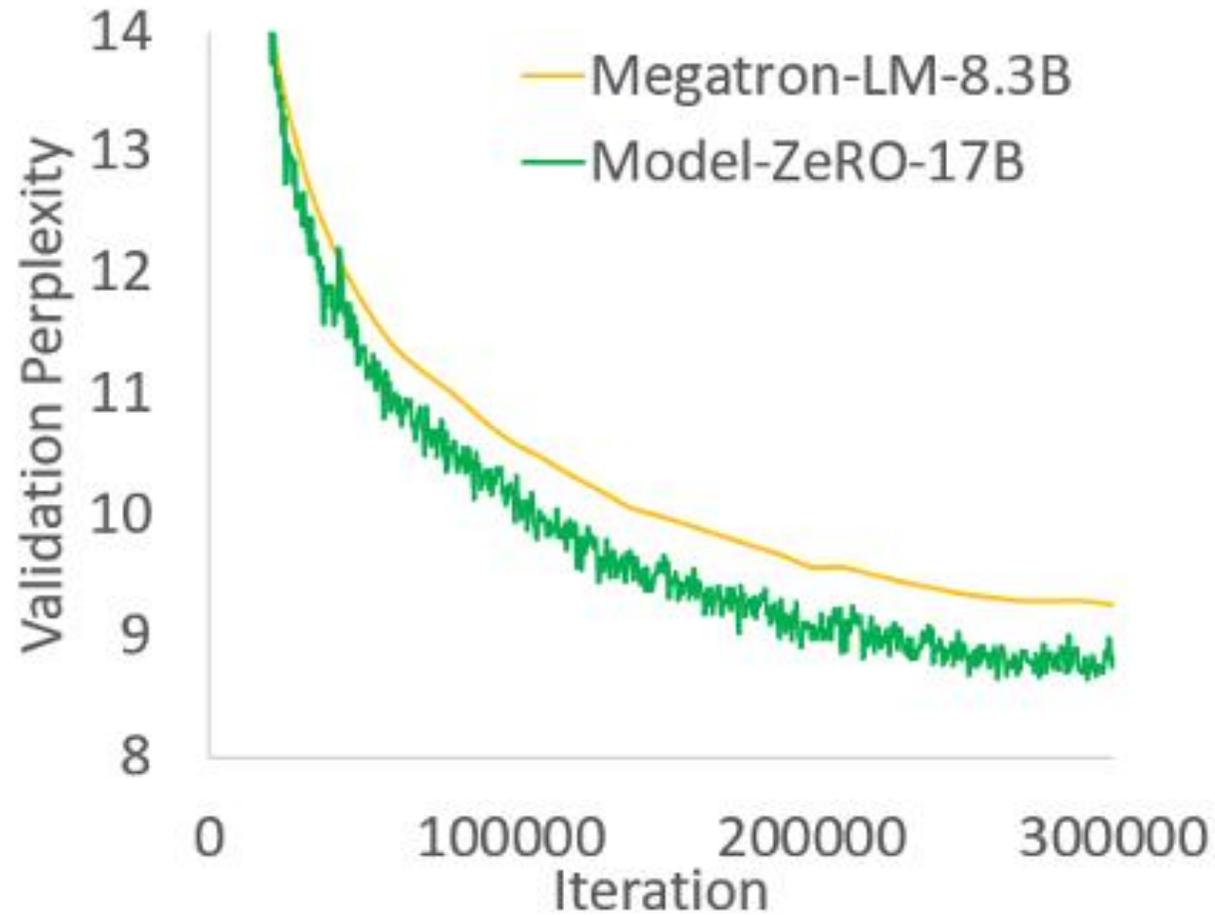


```
ds_config = {
    "optimizer": {
        "type": "Adam",
        "params": {"lr": 0.001}
    },
    "zero": {
        "stage": 3,
        "offload_optimizer": {
            "device": "[cpu|nvme]"
        },
        "offload_param": {
            "device": "[cpu|nvme]"
        }
    }
}
```

# ZeRO Performance: Max Model Throughput



# ZeRO Performance: Turning-NLG 17B



Two guest lectures from industry

- March 10, Meta, Zechun Liu (online)
- March 12, Amazon, Yida Wang (online)

# Questions?