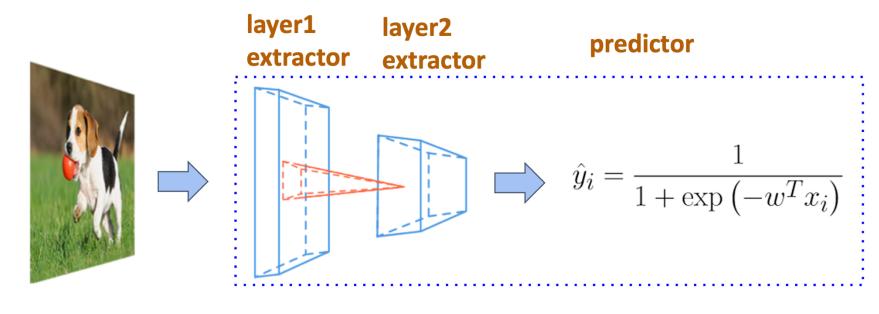
15-779 Lecture 12: Parallelization Part 1: Data Parallelism and Zero Redundancy

Zhihao Jia

Carnegie Mellon University

Recap: DNN Training Overview



Objective

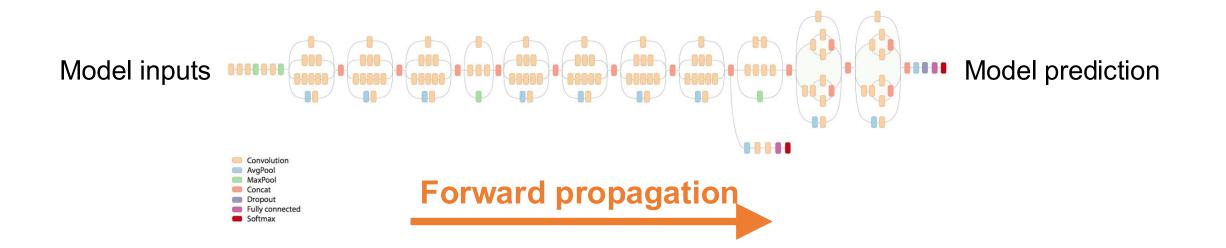
Training

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

DNN Training Process

Train ML models through many iterations of 3 stages

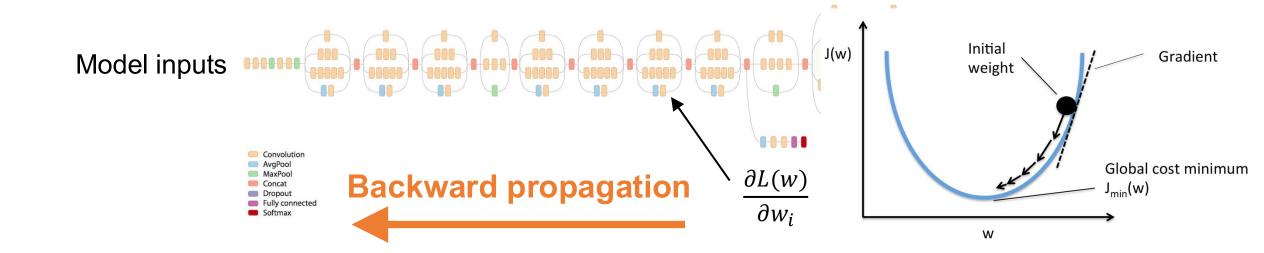
- Forward propagation: apply model to a batch of input samples and run calculation through operators to produce a prediction
- Backward propagation: run the model in reverse to produce error for each trainable weight
- 3. Weight update: use the loss value to update model weights



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 a gradient 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

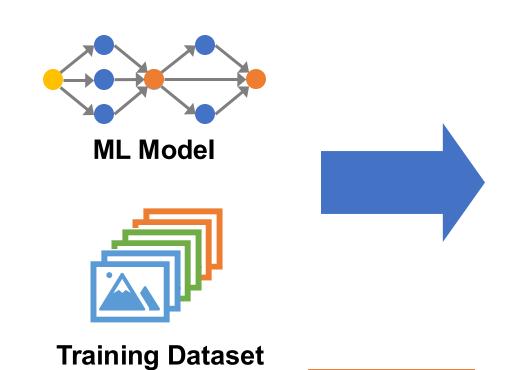
- 1. Forward propagation: apply model to a batch of input samples and run calculation through operators to produce a prediction
- 2. Backward propagation: run the model in reverse to produce a gradient for each trainable weight
- 3. Weight update: use the gradients to update model weights

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

How can we parallelize DNN training?

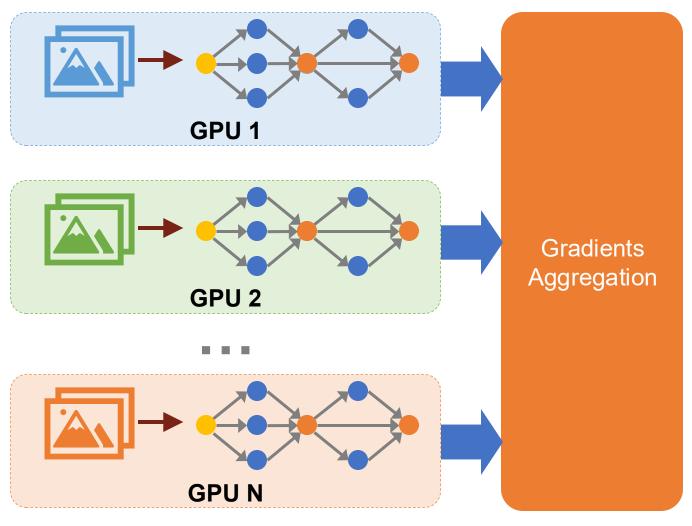
$$w_i \coloneqq w_i - \gamma \nabla L(w_i) = w_i - \frac{\gamma}{n} \sum_{j=1}^n \nabla L_j(w_i)$$

Data Parallelism



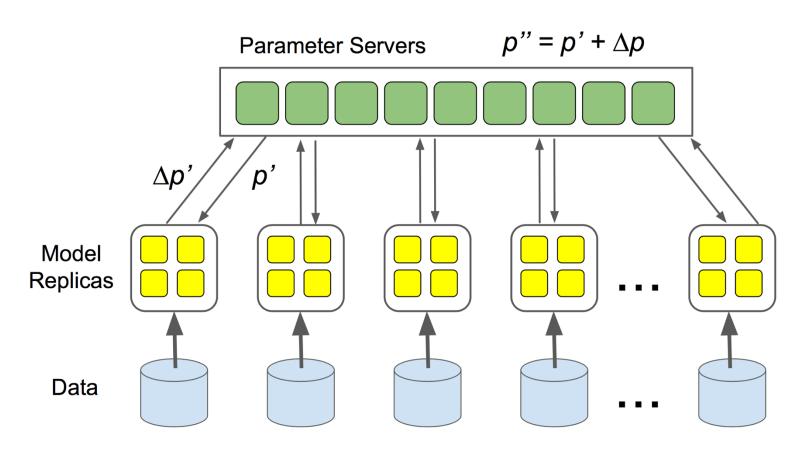
1. Partition training data into batches

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



- 2. Compute the gradients of each batch on a GPU
- 3. Aggregate gradients across GPUs

Data Parallelism: Parameter Server

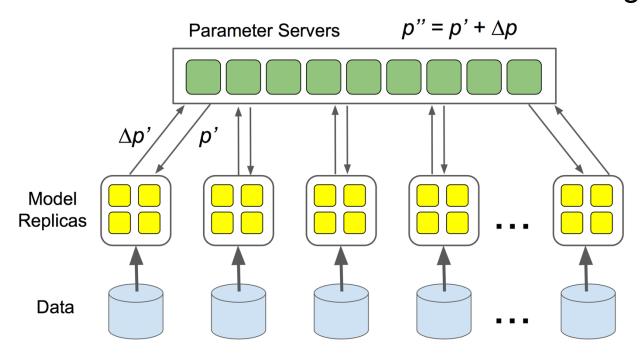


Workers push gradients to parameter servers and pull updated parameters back

Inefficiency of Parameter Server

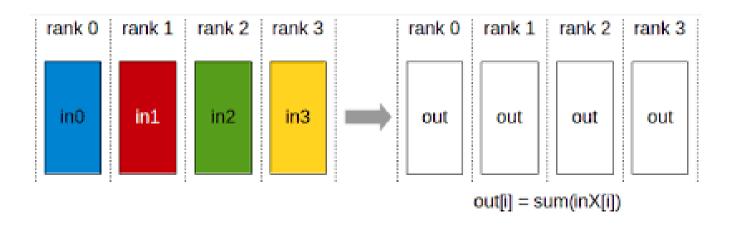
 Centralized communication: all workers communicate with parameter servers for weights update; cannot scale to large numbers of workers

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?
- AllReduce: perform element-wise reduction across multiple devices

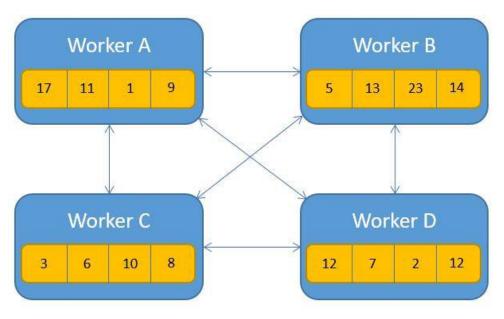


Different Ways to Perform AllReduce

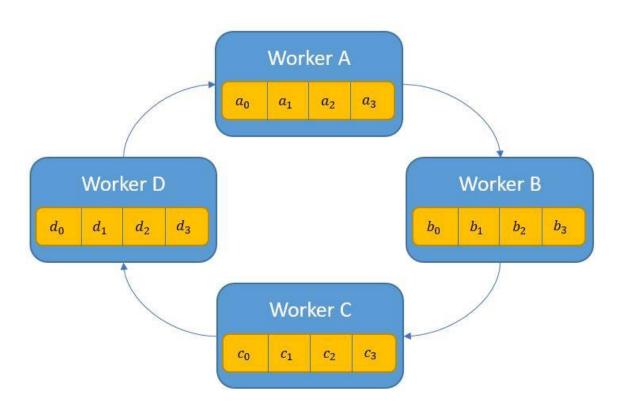
- Naïve AllReduce
- Ring AllReduce
- Tree AllReduce
- Butterfly AllReduce

Naïve AllReduce

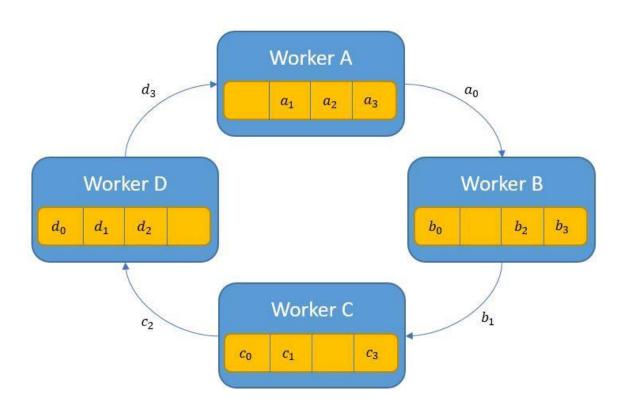
- Each worker can send its local gradients to all other workers
- If we have N workers and each worker contains M parameters
- Overall communication: N * (N-1) * M parameters
- Issue: each worker communicates with all other workers; same scalability issue as parameter server



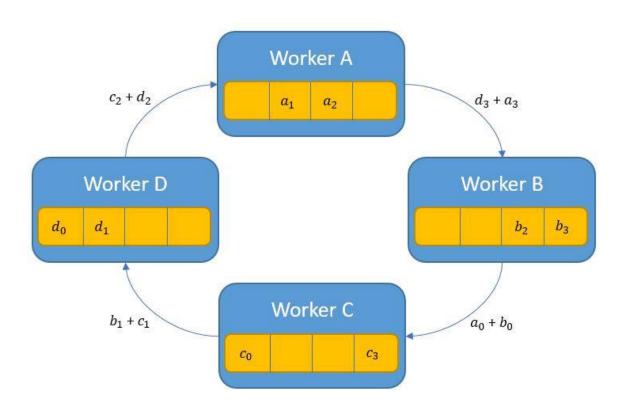
- Construct a ring of N workers, divide M parameters into N slices
- Step 1 (Aggregation): each worker send one slice (M/N parameters) to the next worker on the ring; repeat N times



- 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

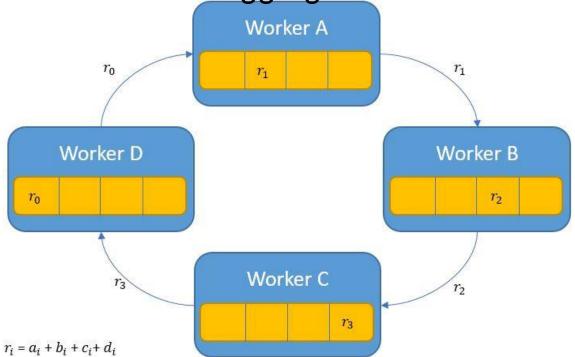


- Construct a ring of N workers, divide M parameters into N slices
- Step 1 (Aggregation): each worker send one slice (M/N parameters) to the next worker on the ring; repeat N times



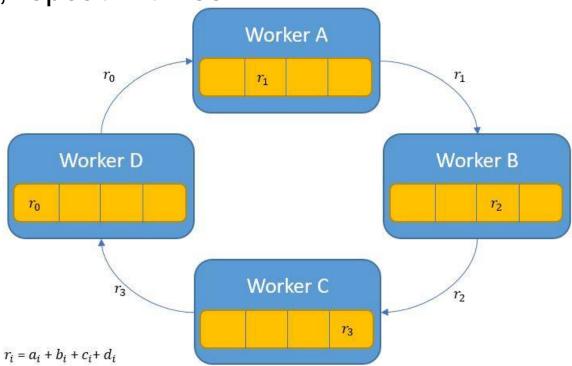
- Construct a ring of N workers, divide M parameters into N slices
- Step 1 (Aggregation): each worker send one slice (M/N parameters) to the next worker on the ring; repeat N times

After step 1, each worker has the aggregated version of M/N parameters



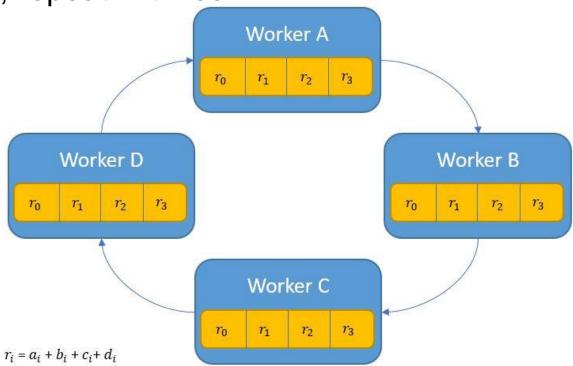
- Construct a ring of N workers, divide M parameters into N slices
- Step 1 (Aggregation): each worker send one slice (M/N parameters) to the next worker on the ring; repeat N times

 Step 2 (Broadcast): each worker send one slice of aggregated parameters to the next worker; repeat N times



- 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



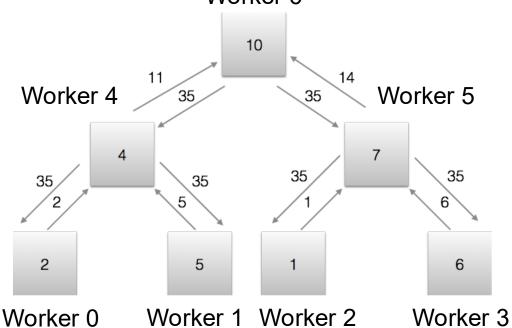
- Construct a ring of N workers, divide M parameters into N slices
- Step 1 (Aggregation): each worker send one slice (M/N parameters) to the next worker on the ring; repeat N times
- Step 2 (Broadcast): each worker send one slice of aggregated parameters to the next worker; repeat N times
- Overall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

Tree AllReduce

- 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

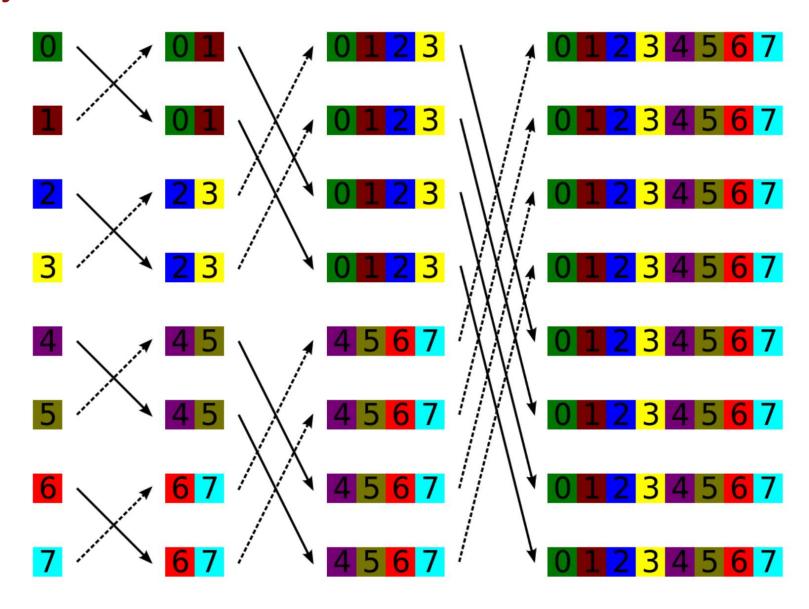
Worker 6



Tree AllReduce

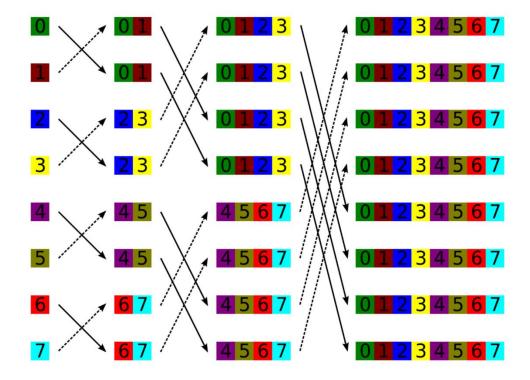
- Construct a tree of N workers;
- Step 1 (Aggregation): each worker sends M parameters to its parent; repeat log(N) times
- Step 2 (Broadcast): each worker sends M parameters to its children; repeat log(N) times
- Overall communication: 2 * N * M parameters
 - Aggregation: M * N parameters
 - Broadcast: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

Butterfly Network



Butterfly AllReduce

- 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 AllReduce Methods

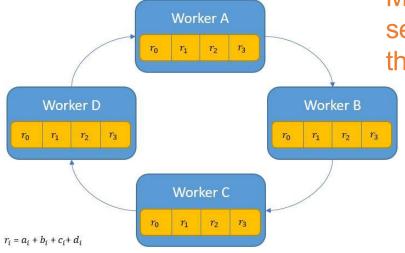
	Parameter Server		Ring AllReduce	Tree AllReduce	Butterfly AllReduce
Overall communication	$2 \times N \times M$	$N^2 \times M$	$2 \times N \times M$	$2 \times N \times M$	$N \times M$ $\times \log N$

Question: Ring AllReduce is more efficient and scalable then Tree AllReduce and Parameter Server, why? Ring AllReduce v.s. Tree AllReduce v.s. Parameter Server

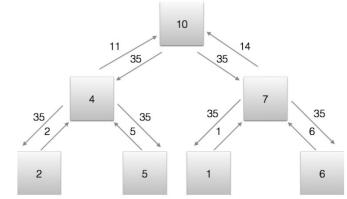
Ring AllReduce:

- Best latency
- Balanced workload across workers

 More scalable since each worker sends 2*M parameters (independent to the number of workers)



Each worker sends M/N parameters per iteration; repeat for 2*N iterations Latency: M/N * (2*N) / bandwidth



Each worker sends M parameters per iteration; repeat for 2*log(N) iterations Latency: M * 2 * log(N) / bandwidth

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

 $p'' = p' + \Delta p$

Latency: M * N / bandwidth

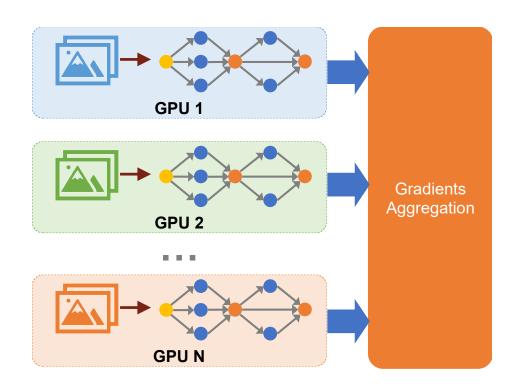
Parameter Servers

Model Replicas

Data

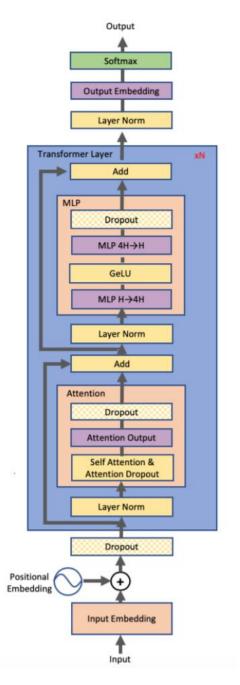
An Issue with Data Parallelism

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



Large Model Training Challenges

	Bert-		Turing	
	Large	GPT-2	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	54x	547x
Memory Footprint	5.12GB	24GB	275GB	2800GB

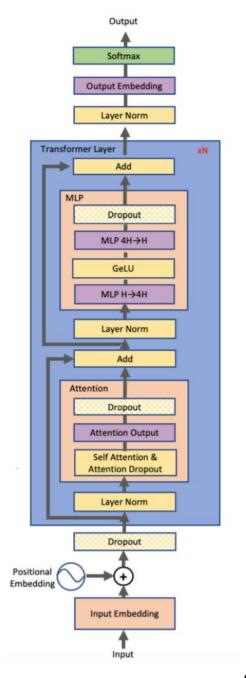


Large Model Training Challenges

	Bert- Large	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	54x	547x
Memory Footprint	5.12GB	24GB	275GB	2800GB

NVIDIA V100 GPU memory capacity: 16G/32G NVIDIA A100 GPU memory capacity: 40G/80G

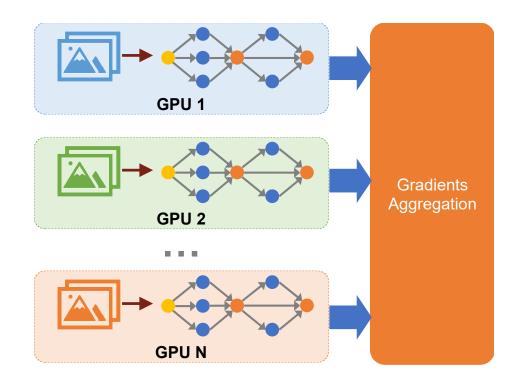
Out of Memory







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



Revisit: Stocastic Gradient Descent

```
For t = 1 to T
\Delta w = \eta \times \frac{1}{b} \sum_{i=1}^{b} \nabla \left( loss(f_w(x_i, y_i)) \right) // compute derivative and update
w -= \Delta w // apply update
End
```

Adaptive Learning Rates (Adam)

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

$$\Delta w = adam(g)$$

$$w -= \Delta w // apply update$$
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 ω^j

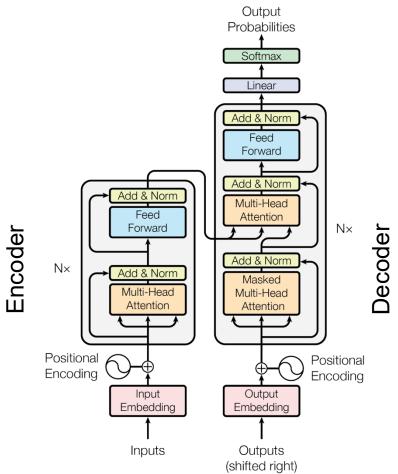
 ν_t : Exponential Average of gradients along ω_j

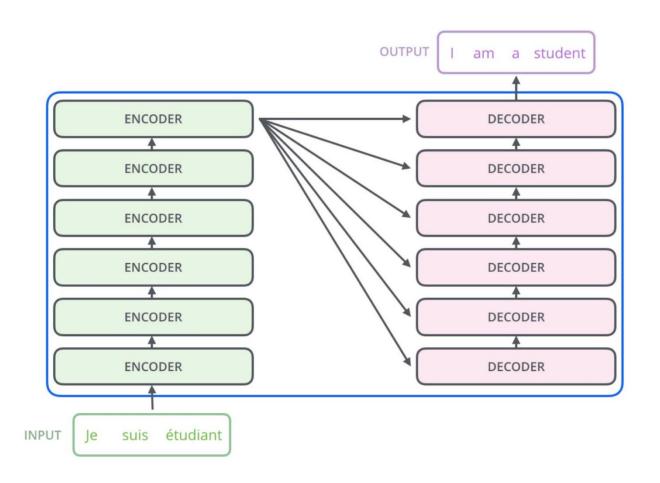
 s_t : Exponential Average of squares of gradients along ω_j

 $\beta_1, \beta_2: Hyperparameters$

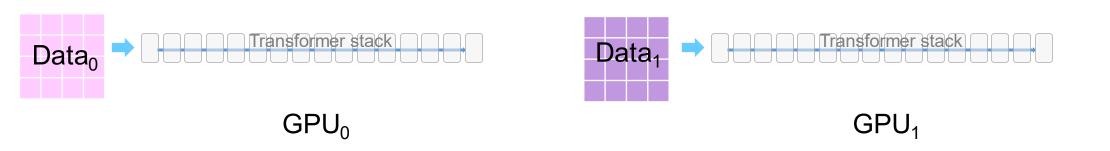
[1] Kingma and Ba, "Adam: A Method for Stochastic Optimization", 2014, https://arxiv.org/abs/1412.6980

Transformer for Language Models

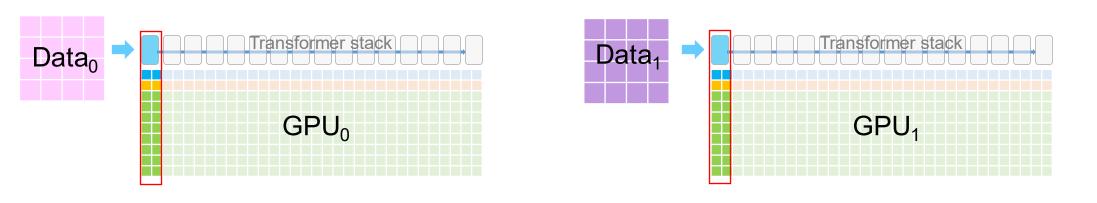




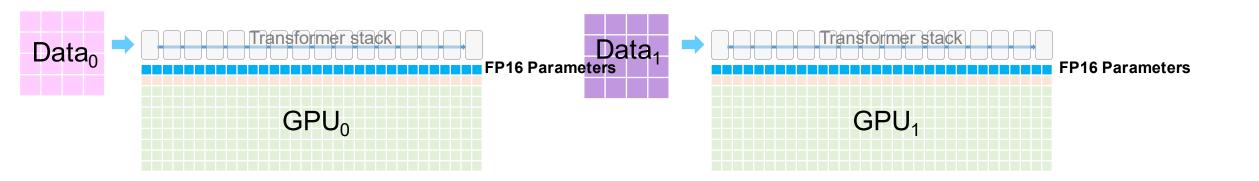
Ashish Vaswani et. al. Attention is all you need.



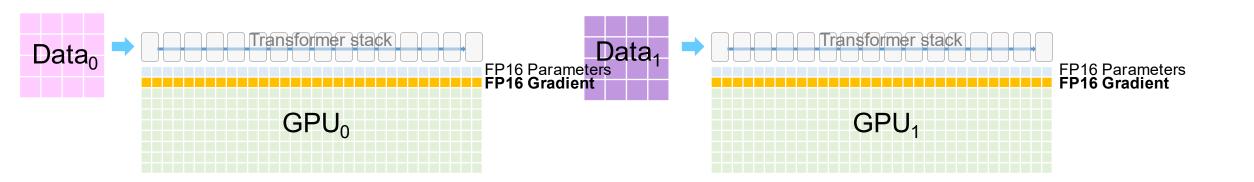
A 16-layer transformer model = 1 layer



Each cell represents GPU memory used by its corresponding transformer layer

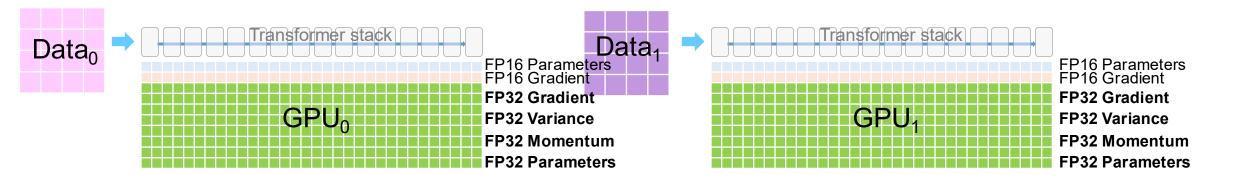


FP16 parameter



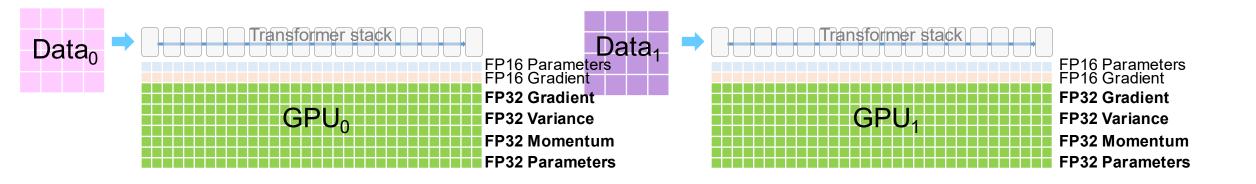
- FP16 parameter
- FP16 Gradients

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

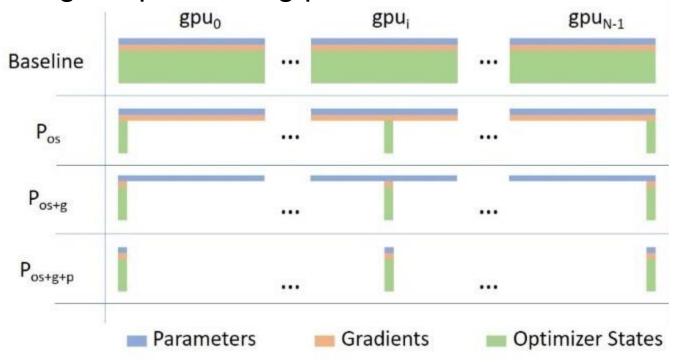
Example 1B parameter model -> 20GB/GPU

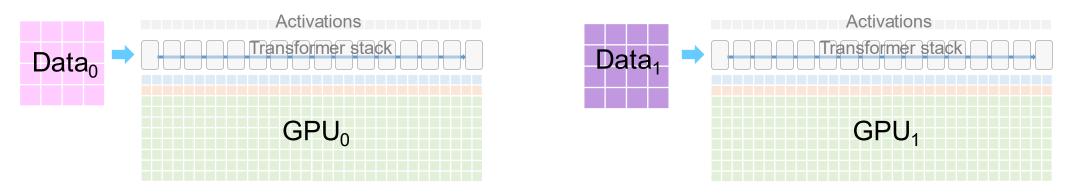
Memory consumption doesn't include:

Input batch + activations

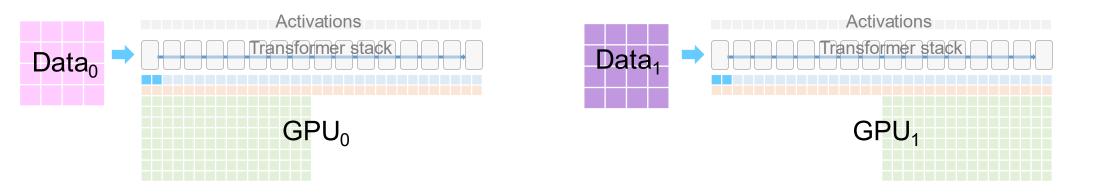
ZeRO-DP: ZeRO powered Data Parallelism

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

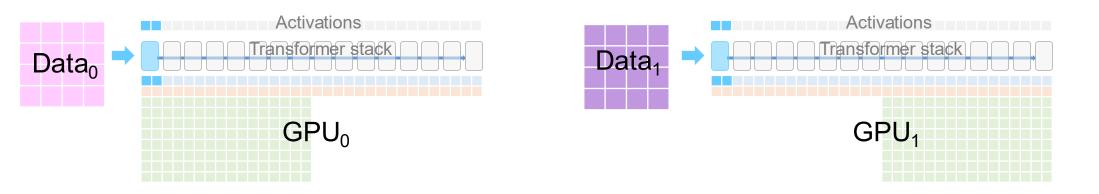




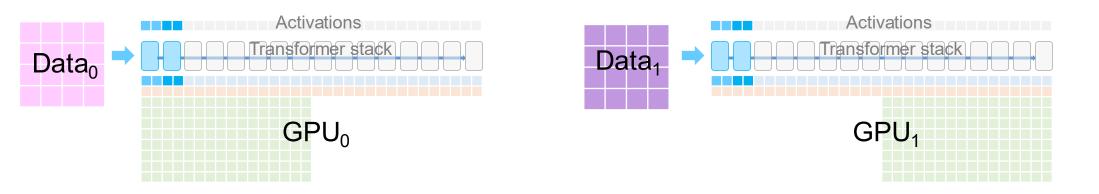
ZeRO Stage 1



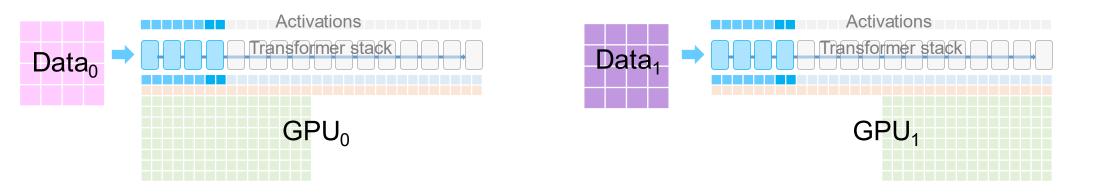
- ZeRO Stage 1
- Partitions optimizer states across GPUs



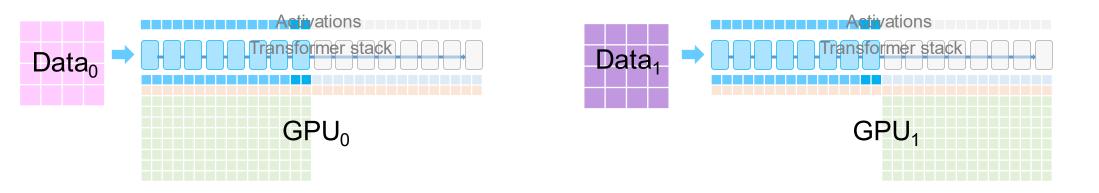
- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks



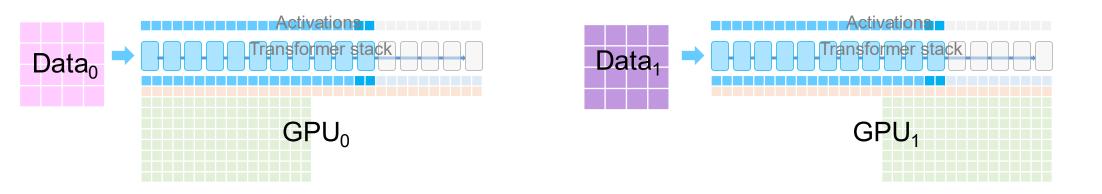
- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks



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



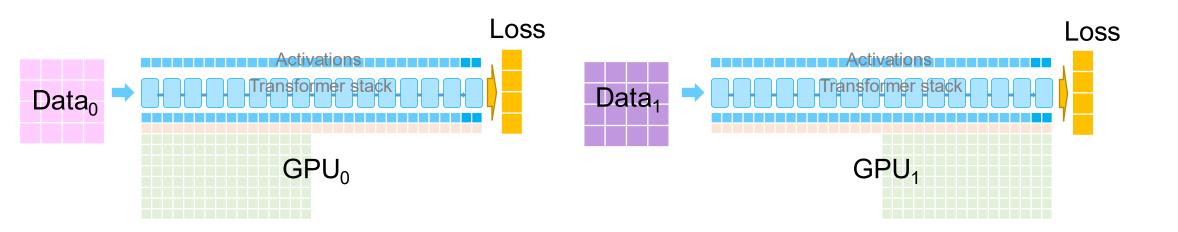
- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks



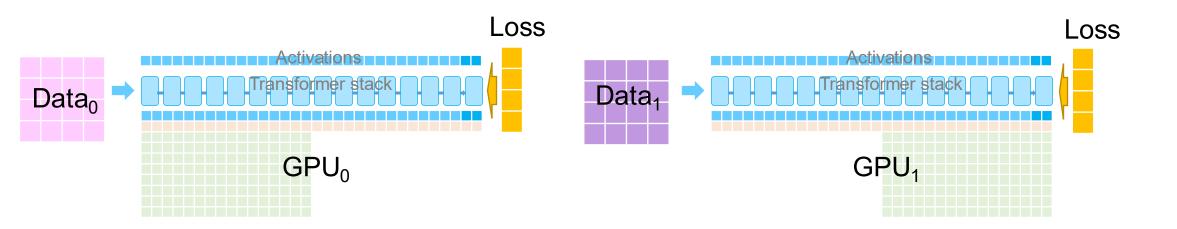
- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks



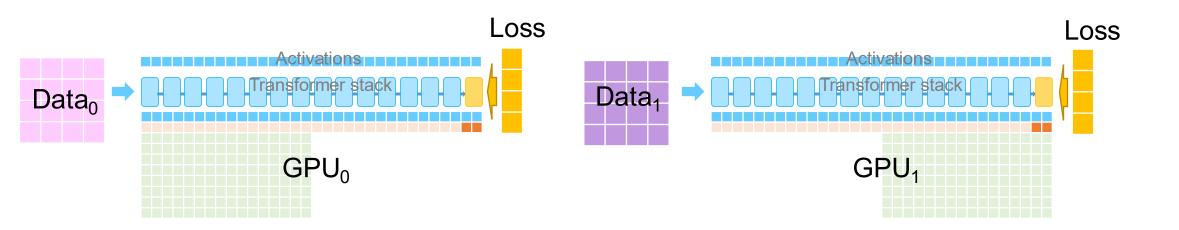
- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks



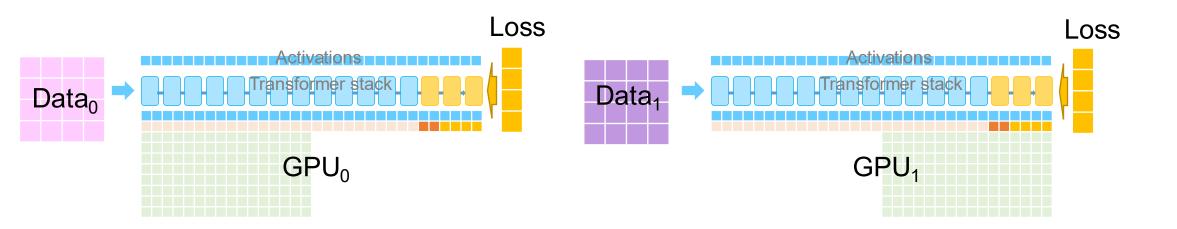
- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks



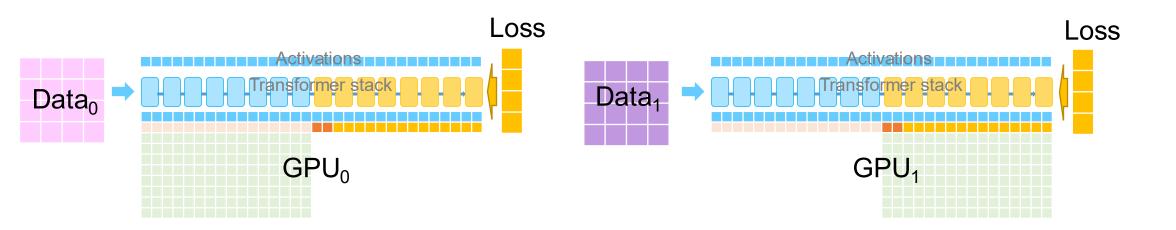
- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks
- Backward propagation to generate FP16 gradients



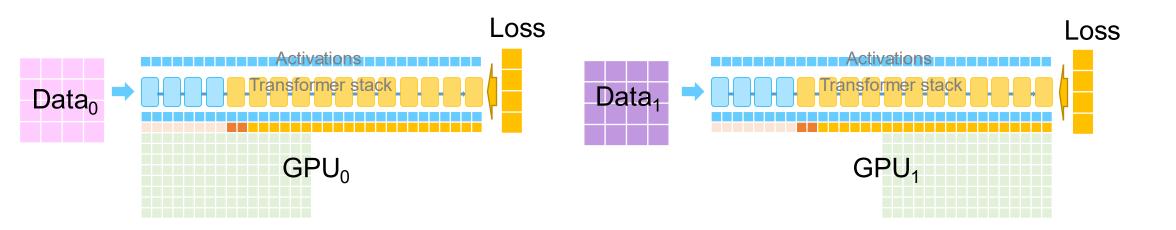
- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks
- Backward propagation to generate FP16 gradients



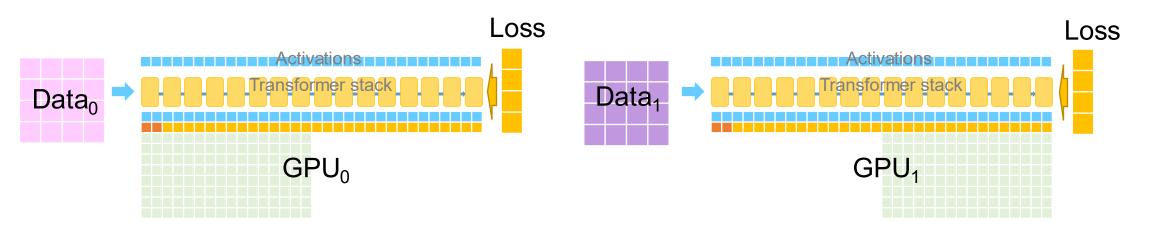
- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks
- Backward propagation to generate FP16 gradients



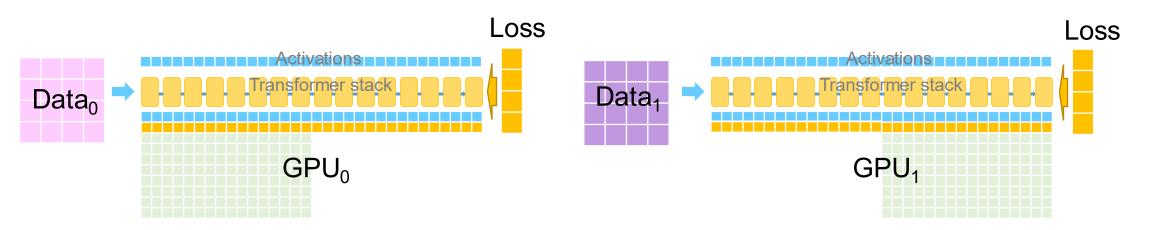
- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks
- Backward propagation to generate FP16 gradients



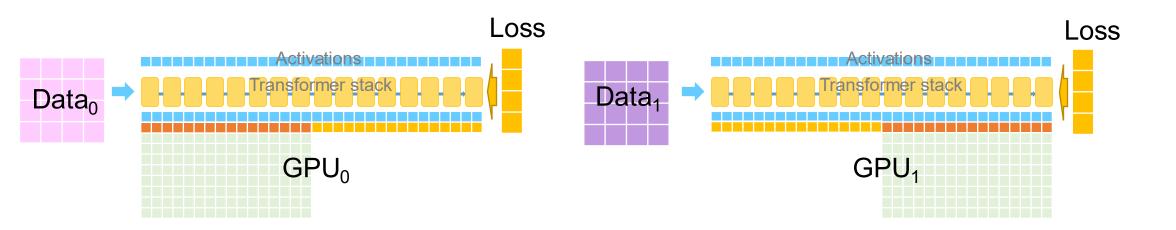
- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks
- Backward propagation to generate FP16 gradients



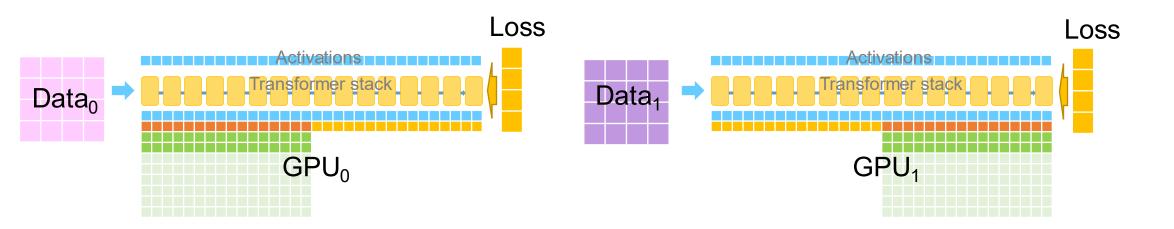
- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks
- Backward propagation to generate FP16 gradients



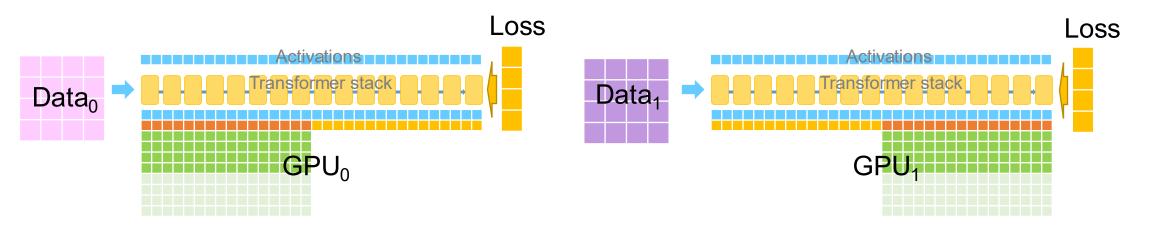
- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks
- Backward propagation to generate FP16 gradients and AllReduce to average



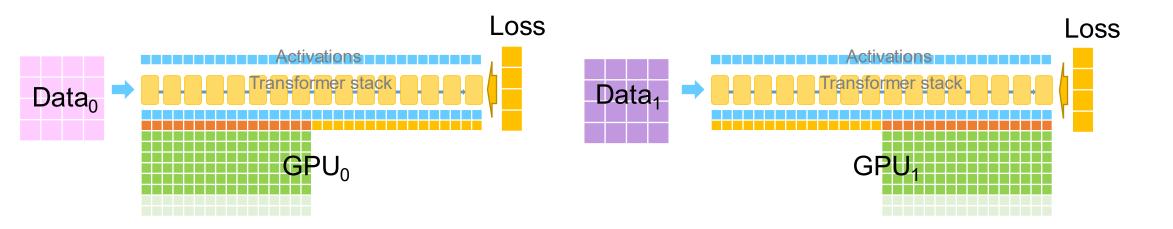
- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks
- Backward propagation to generate FP16 gradients and AllReduce to average



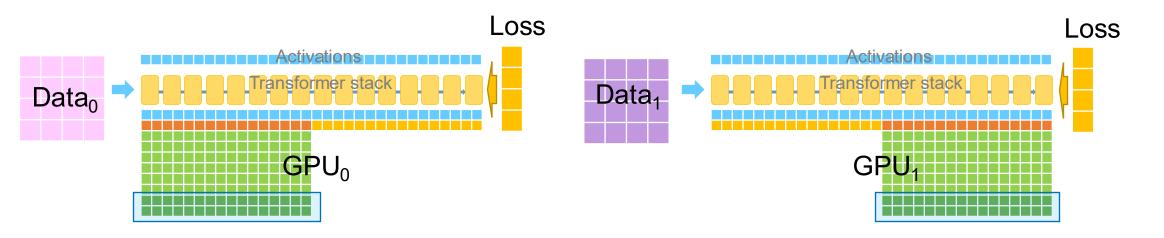
- ZeRO Stage 1
- Partitions optimizer states across GPUs
- Run Forward across the transformer blocks
- Backward propagation to generate FP16 gradients and AllReduce to average
- Update the FP32 weights with ADAM optimizer



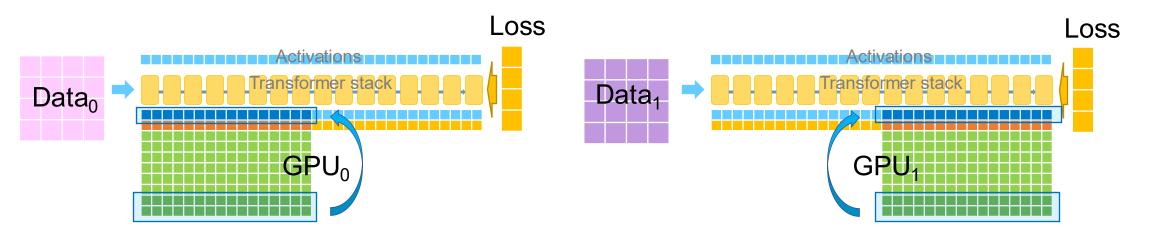
- Run Forward across the transformer blocks
- Backward propagation to generate FP16 gradients and AllReduce to average
- Update the FP32 weights with ADAM optimizer



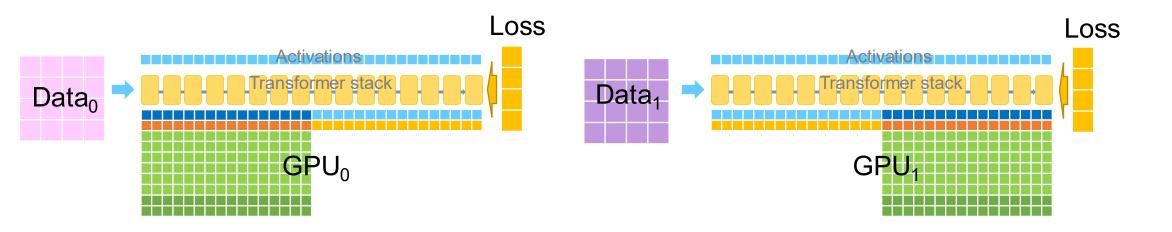
- Run Forward across the transformer blocks
- Backward propagation to generate FP16 gradients and AllReduce to average
- Update the FP32 weights with ADAM optimizer



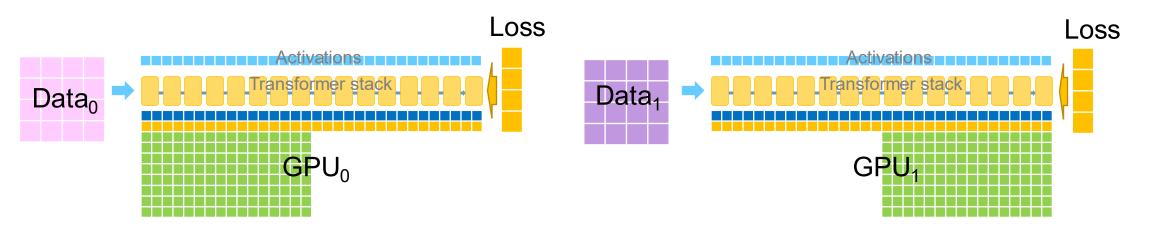
- Run Forward across the transformer blocks
- Backward propagation to generate FP16 gradients and AllReduce to average
- Update the FP32 weights with ADAM optimizer



- Run Forward across the transformer blocks
- Backward propagation to generate FP16 gradients and AllReduce to average
- Update the FP32 weights with ADAM optimizer
- Update the FP16 weights



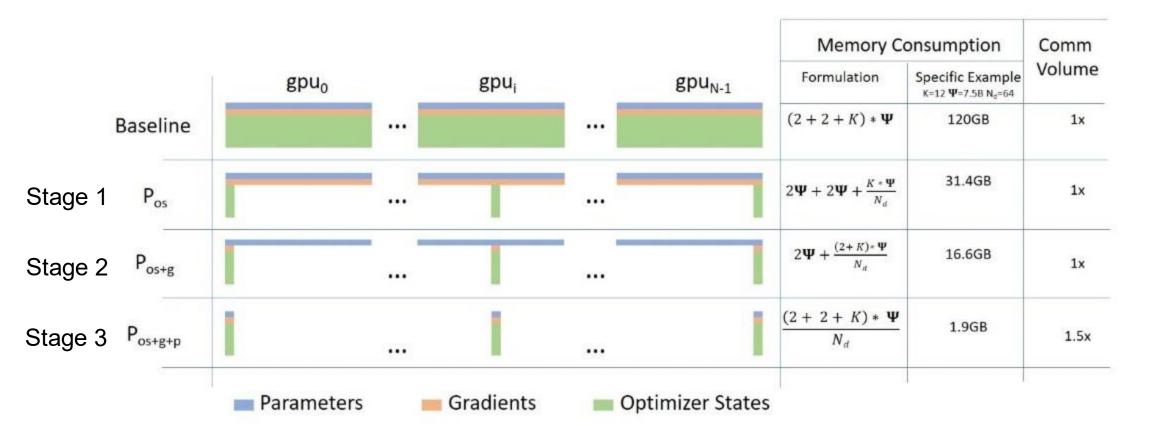
- Run Forward across the transformer blocks
- Backward propagation to generate FP16 gradients and AllReduce to average
- Update the FP32 weights with ADAM optimizer
- Update the FP16 weights
- All Gather the FP16 weights to complete the iteration

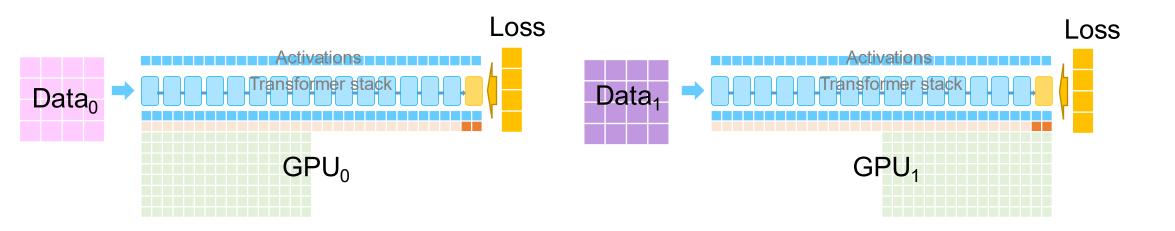


- Run Forward across the transformer blocks
- Backward propagation to generate FP16 gradients and AllReduce to average
- Update the FP32 weights with ADAM optimizer
- Update the FP16 weights
- All Gather the FP16 weights to complete the iteration

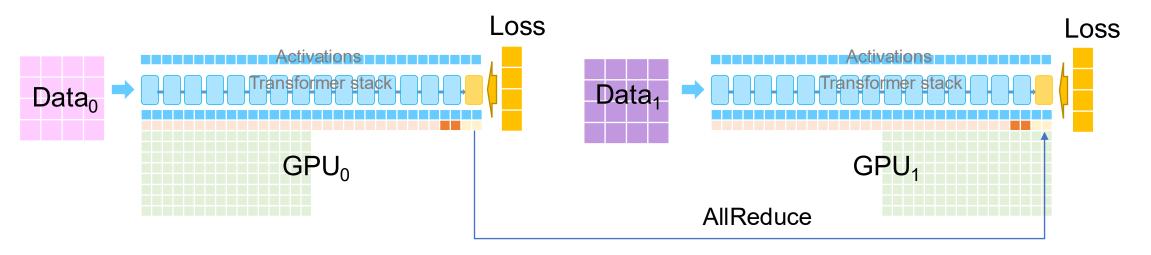
ZeRO: Zero Redundancy Optimizer

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

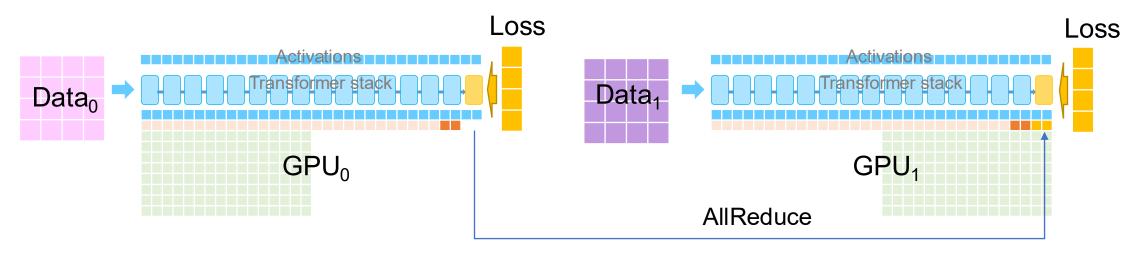




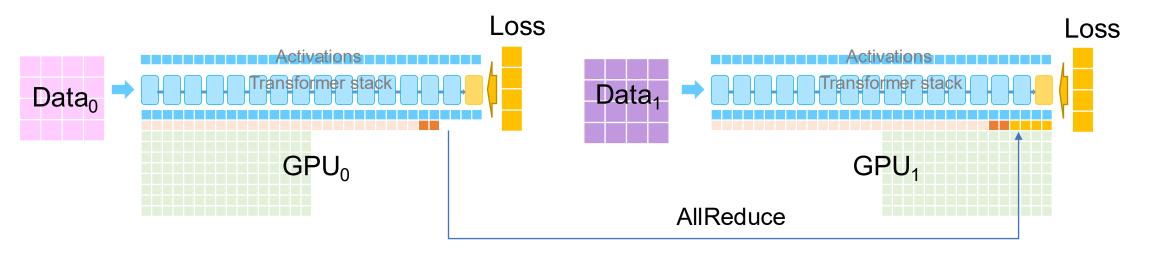
- Partitioning gradients across GPUs
- The forward process remains the same as stage 1



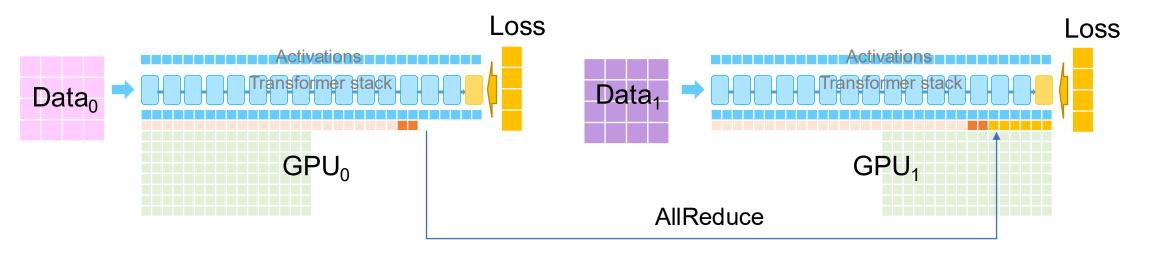
- Partitioning gradients across GPUs
- Perform AllReduce right after back propagation of each layer



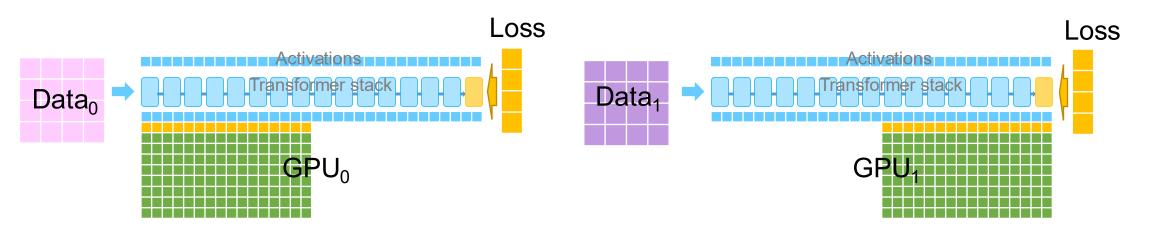
- Partitioning gradients across GPUs
- Only one GPU keeps the gradients after AllReduce



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



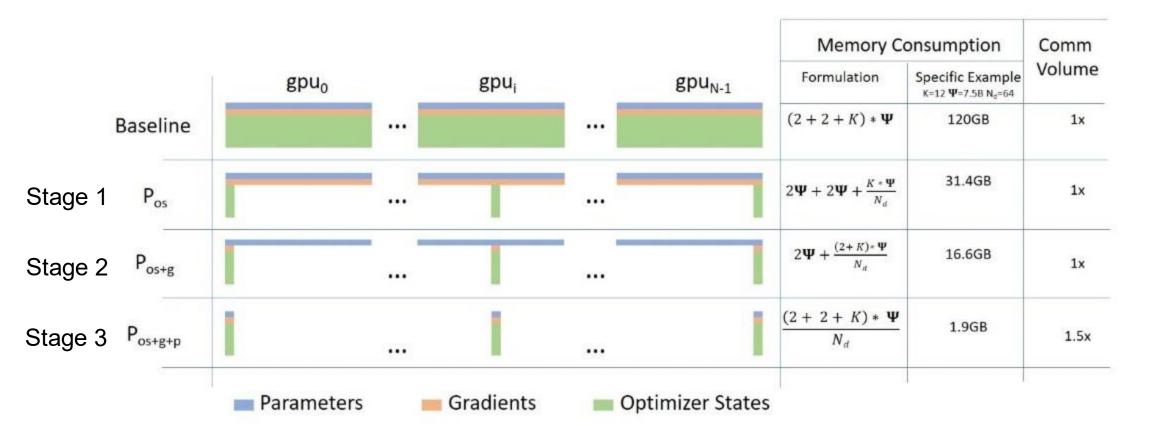
- Partitioning gradients across GPUs
- Reduce gradients on GPUs responsible for updating parameters



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

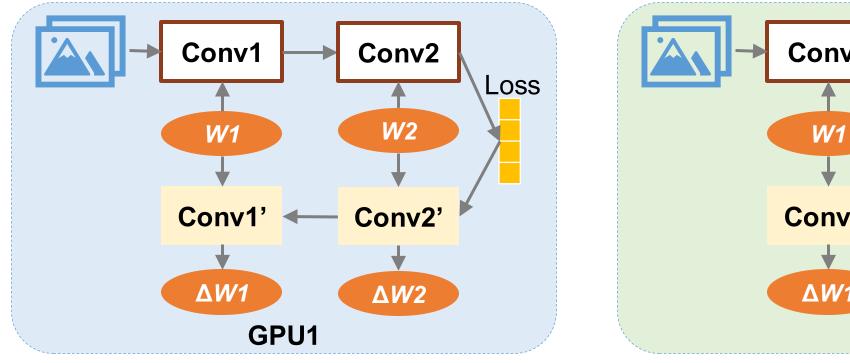
ZeRO: Zero Redundancy Optimizer

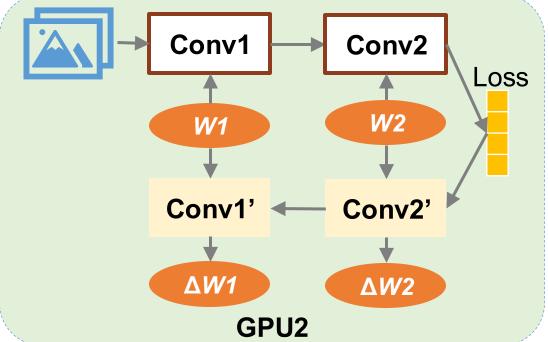
- Progressive memory savings and communication volume
- Turning NLR 17.2B is powered by Stage 1 and Megatron



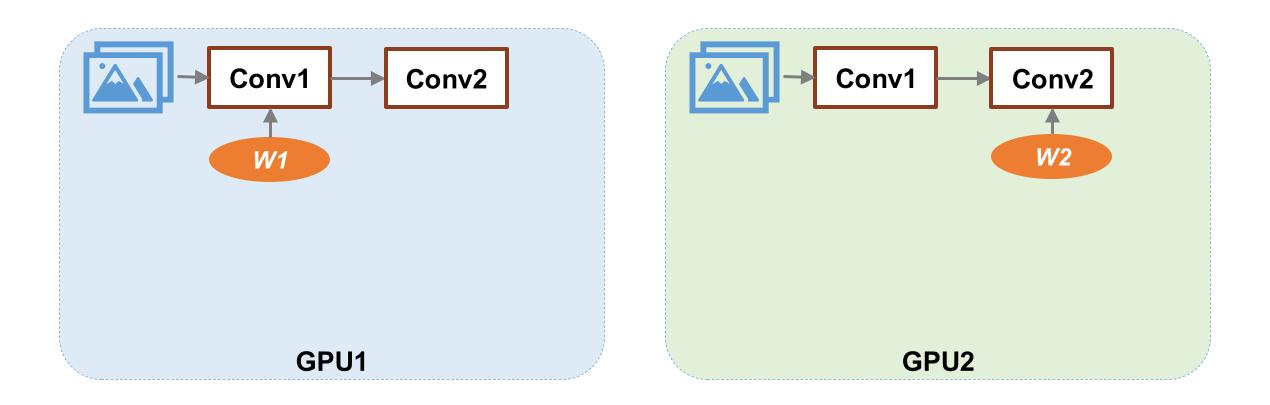
ZeRO Stage 3: Partitioning Parameters

• In data parallel training, all GPUs keep all parameters during training

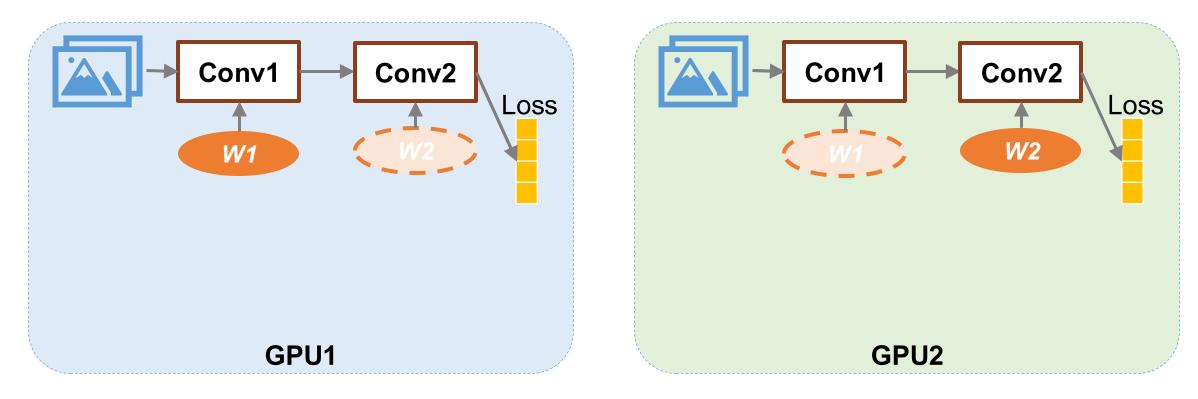




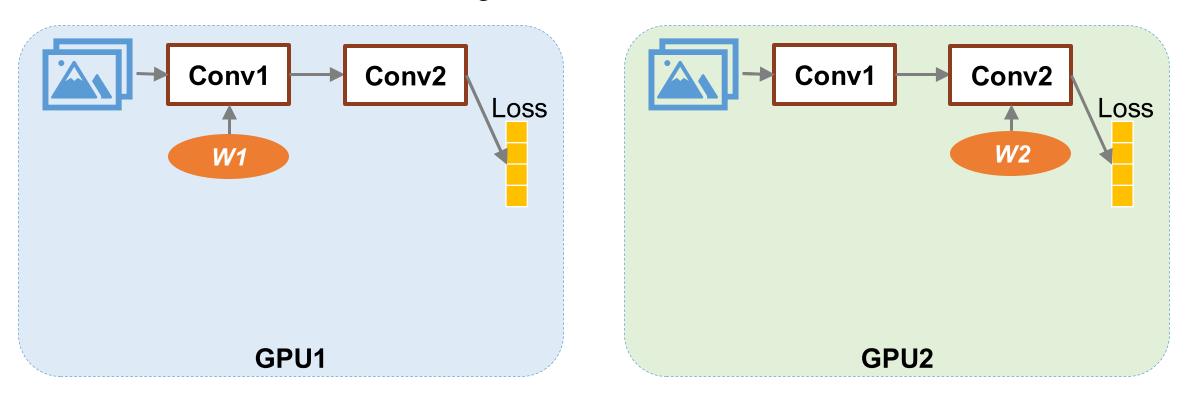
• In ZeRO, model parameters are partitioned across GPUs



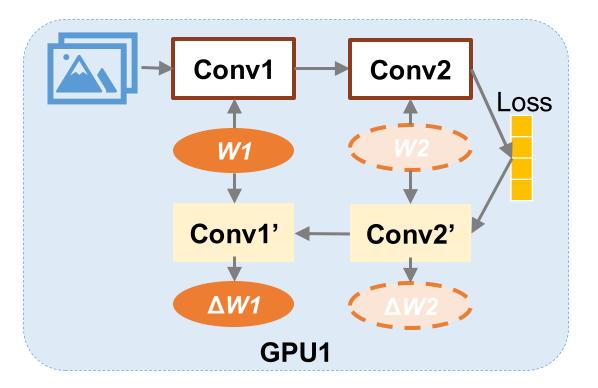
- In ZeRO, model parameters are partitioned across GPUs
- GPUs broadcast their parameters during forward

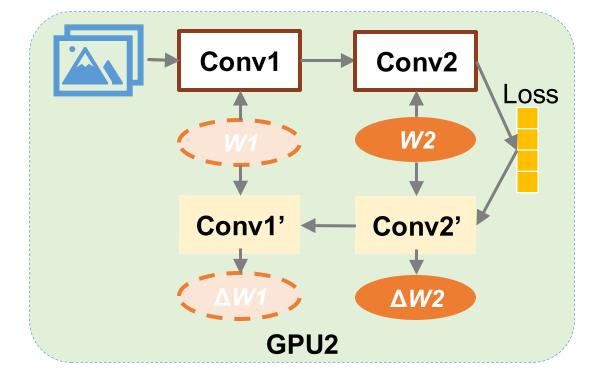


- In ZeRO, model parameters are partitioned across GPUs
- Parameters are discarded right after use



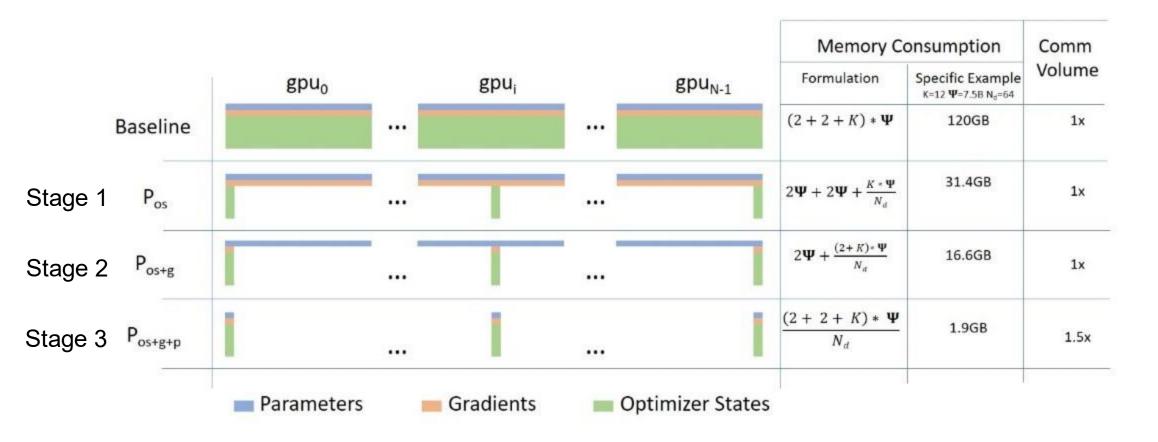
- In ZeRO, model parameters are partitioned across GPUs
- GPUs broadcast their parameters again during backward

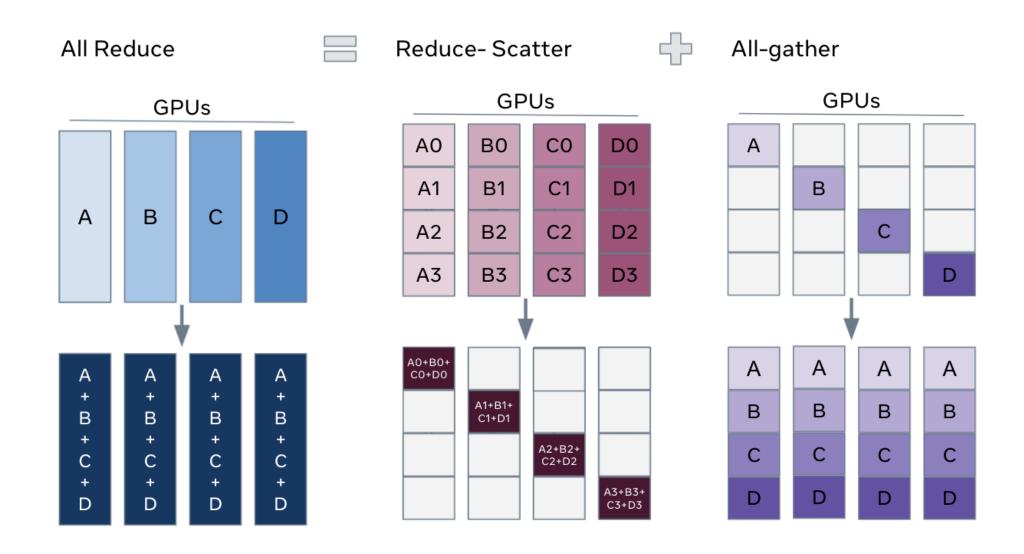




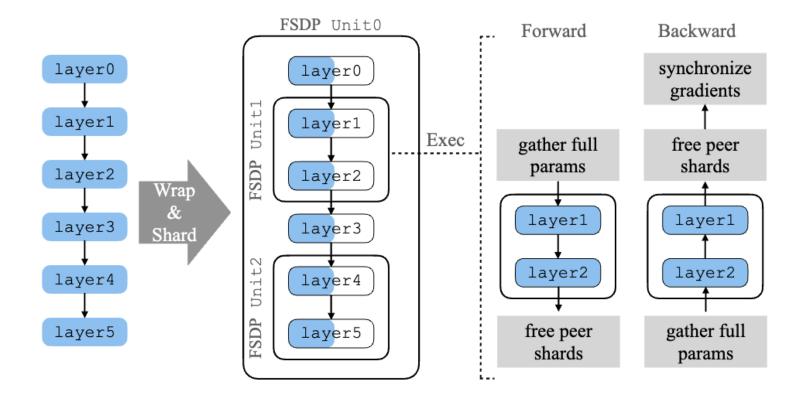
ZeRO: Zero Redundancy Optimizer

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

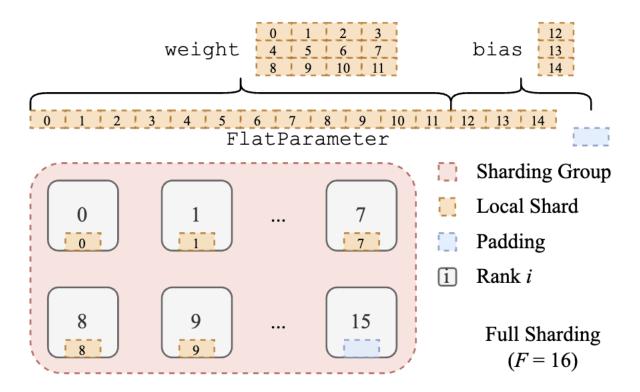




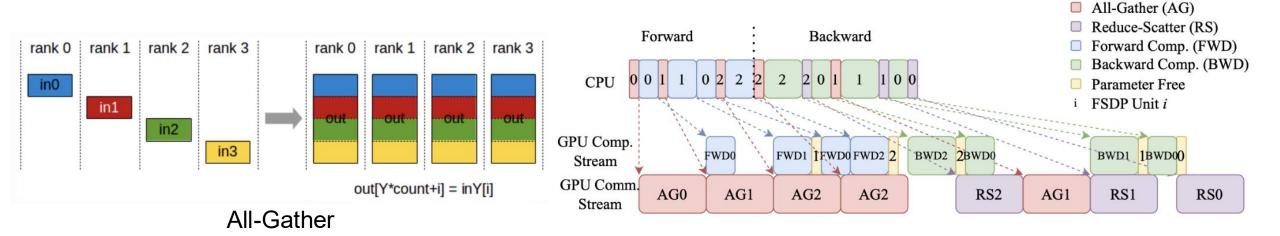
1. Divide model parameters into FSDP units



- 1. Divide model parameters into FSDP units
- 2. Share each unit across multiple GPUs

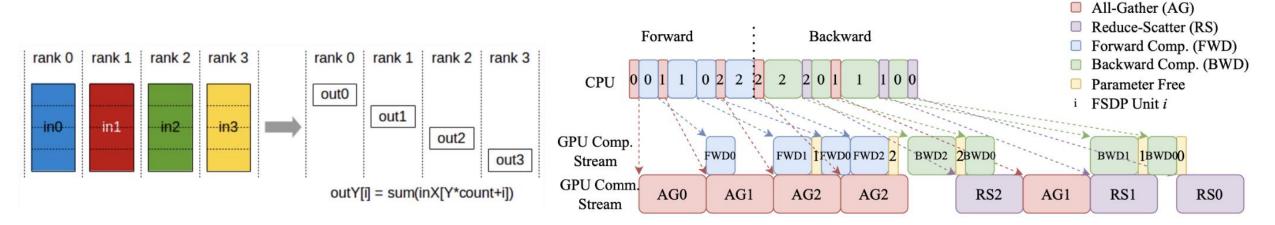


- 1. Divide model parameters into FSDP units
- 2. Share each unit across multiple GPUs
- 3. Run forward pass;
 - 1. Perform all-gather so that each GPU get all parameters of a unit
 - 2. Run forward pass & discard parameter shards



- 1. Divide model parameters into FSDP units
- 2. Share each unit across multiple GPUs
- 3. Run forward pass
- 4. Run backward pass

- 1. Perform all-gather again to get all parameters of a unit
- 2. Each GPU computes gradients for all parameters
- 3. Perform reduce-scatter to aggregate full gradients



Reduce-Scatter

Programming in PyTorch FSDP

```
from torch.distributed.fsdp import fully shard, FSDPModule
                                  from torch.distributed.tensor import Dtensor
                                  model = Transformer()
Sharding individual
                                  for layer in model.layers:
layers and entire model
                                    fully_shard(layer)
                                  fully shard(model)
Parameters are of type
                                  for param in model.parameters():
DTensor after sharding
                                    assert isinstance(param, DTensor)
Optimizer will be sharded
                                  optim = torch.optim.Adam(model.parameters(), Ir=1e-2)
automatically
                                  for _ in range(epochs):
                                    x = torch.randint(0, vocab_size, (batch_size, seq_len))
Normal training forward
                                    loss = model(x).sum()
& backward as before
                                    loss.backward()
                                    optim.step()
                                    optim.zero grad()
```

Summary

- Data-parallel training
 - Parameter server
 - Ring AllReduce
 - Tree AllReduce
 - Butterfly AllReduce
- ZeRO: zero redundancy optimizer
- PyTorch FSDP: Fully Shared Data Parallelism