Lecture 11: Large-Scale Distributed Training

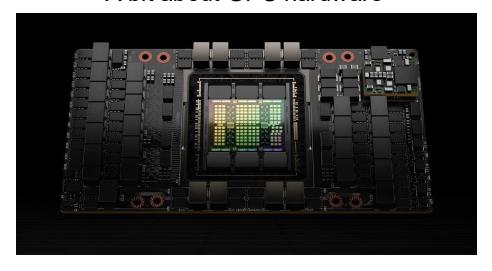
Today: Large-Scale Distributed Training

Running Example: Llama3-405B

- GPT4 kicked off a trend of not sharing any model details:
 "Given both the competitive landscape and the safety implications of large-scale models like GPT-4, this report contains no further details about the architecture (including model size), hardware, training compute, dataset construction, training method, or similar."
- Llama3: Open-source LLM released by Meta in April 2024; paper shares many model and training details
- Llama4: Released initial models April 2025, but no paper yet

GPUs and How to Train On Them

A bit about GPU hardware

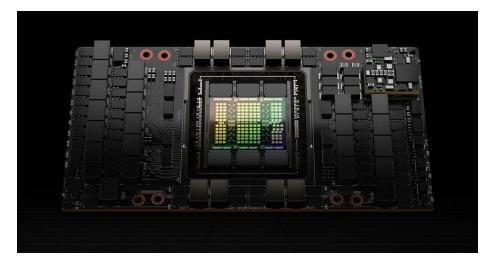


How to train on lots of GPUs



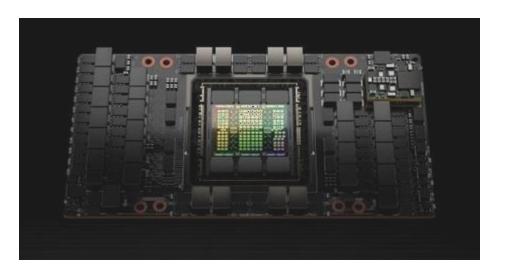
GPUs and How to Train On Them

A bit about GPU hardware

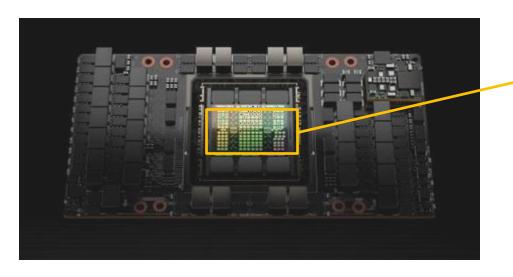


How to train on lots of GPUs



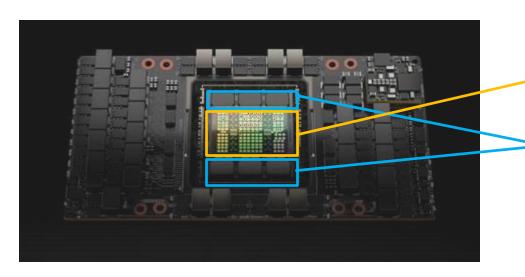


GPU = Graphics Processing UnitOriginally for graphics Now a general parallel processor



GPU = Graphics Processing UnitOriginally for graphics Now a general parallel processor

Compute Cores



GPU = Graphics Processing Unit

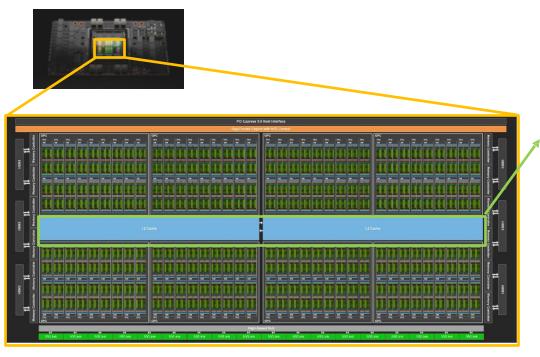
Originally for graphics Now a general parallel processor

Compute Cores

80 GB of HBM Memory
3352 GB/sec bandwidth to cores



H100 Compute Cores



H100 Compute Cores

50MB of L2 Cache



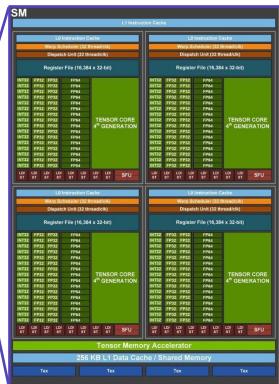
H100 Compute Cores

50MB of L2 Cache

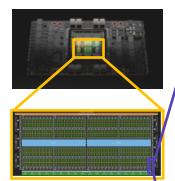
132 Streaming Multiprocessors (SMs) These are independent parallel cores (Actually 144 here; only 132 are enabled due to yield)



Sort of like a CPU core with vector instructions



H100 Streaming Multiprocessor

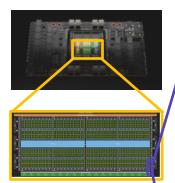


Sort of like a CPU core with vector instructions



H100 Streaming Multiprocessor

256 KB L1 cache, 256 KB registers



Sort of like a CPU core with vector instructions

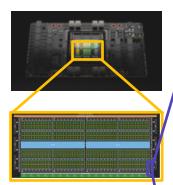


H100 Streaming Multiprocessor

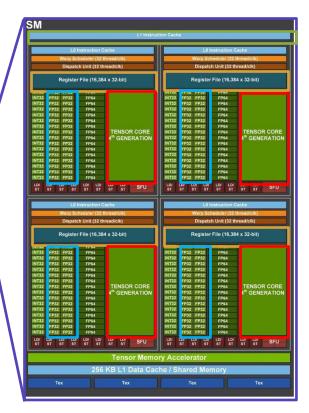
256 KB L1 cache, 256 KB registers

128 **FP32 Cores**

Computes a*x + b per clock cycle 2 FLOPs = Floating Point Operations 256 FLOP/cycle per SM



Sort of like a CPU core with vector instructions



H100 Streaming Multiprocessor

256 KB L1 cache, 256 KB registers

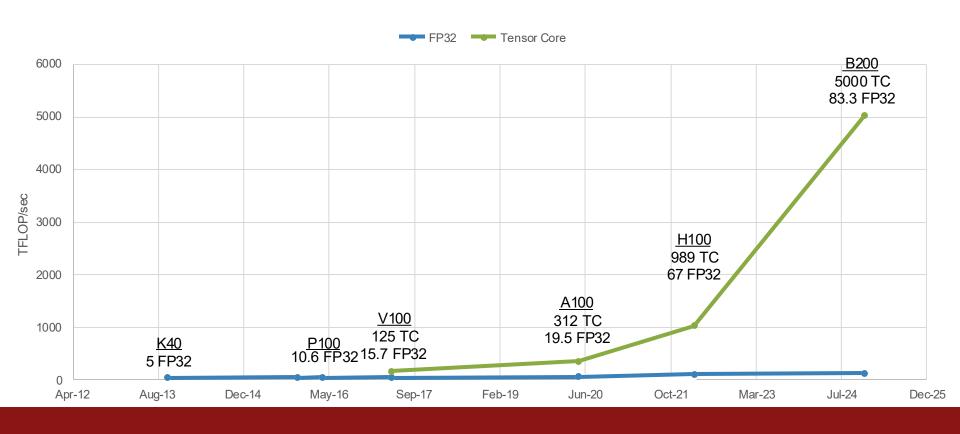
128 **FP32 Cores**

Computes a*x + b per clock cycle 2 FLOPs = Floating Point Operations 256 FLOP/cycle per SM

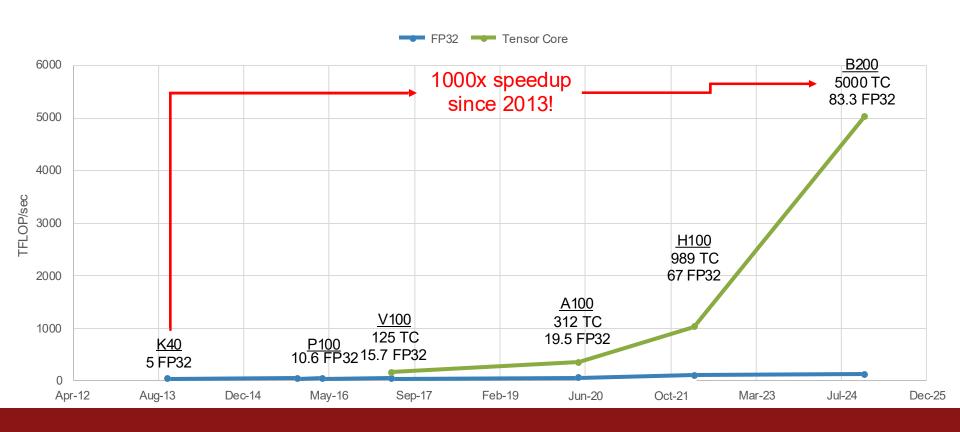
4 Tensor Cores

Computes AX + B per clock cycle
Matrix operation: [16x4][4x8] + [16x8]
16*4*8*2 = 1024 FLOPs
4096 FLOP/cycle per SM
Mixed precision: 16-bit / 32-bit

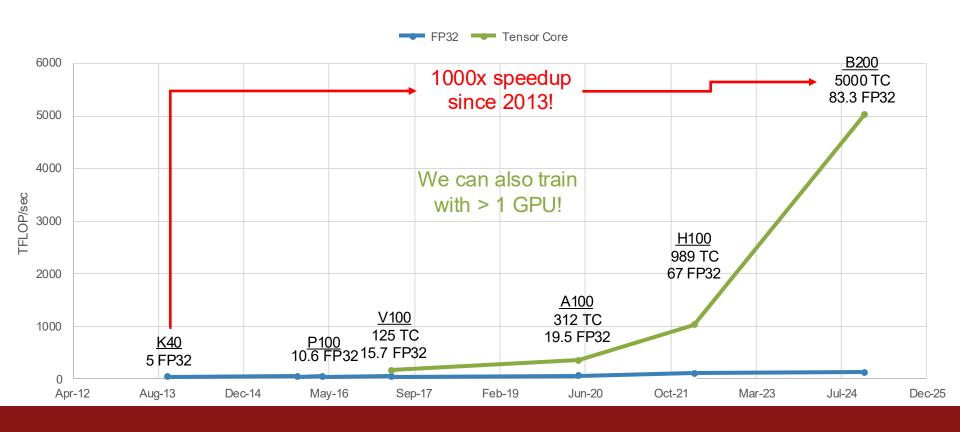
GPUs Have Gotten Much Faster!



GPUs Have Gotten Much Faster!



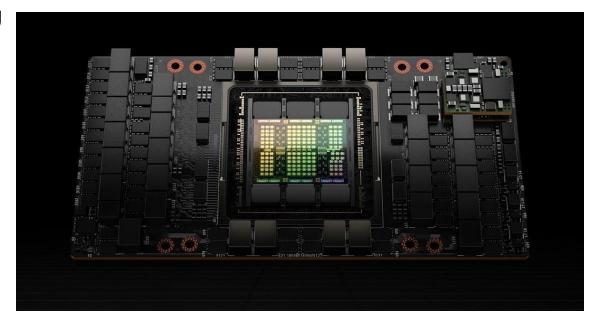
GPUs Have Gotten Much Faster!



NVIDIAH100 GPU

H100 GPU

3352 GB/sec inside the GPU



NVIDIAH100 GPU

H100 GPU 3352 GB/sec inside the GPU

Server = 8x GPU 900 GB/sec between GPUs

GPU Server

















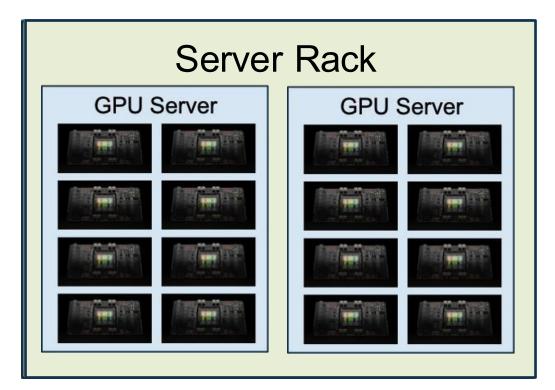
Case Study: Meta's Llama3 Cluster

H100 GPU

3352 GB/sec inside the GPU

Server = 8x GPU 900 GB/sec between GPUs

Rack = 2 Servers = 16x GPU



Case Study: Meta's Llama3 Cluster

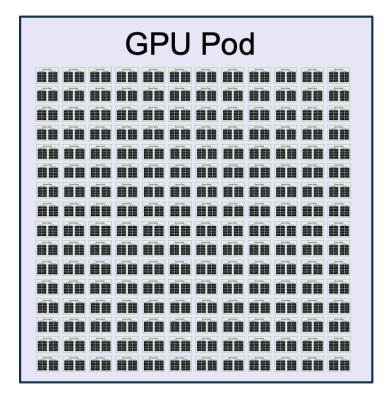
H100 GPU

3352 GB/sec inside the GPU

Server = 8x GPU 900 GB/sec between GPUs

Rack = 2 Servers = 16x GPU

Pod = 192 Racks = 3072 GPUs 50 GB/sec between GPUs



Llama Team, "The Llama 3 Herd of Models", https://arxiv.org/abs/2407.21783

Case Study: Meta's Llama3 Cluster

H100 GPU

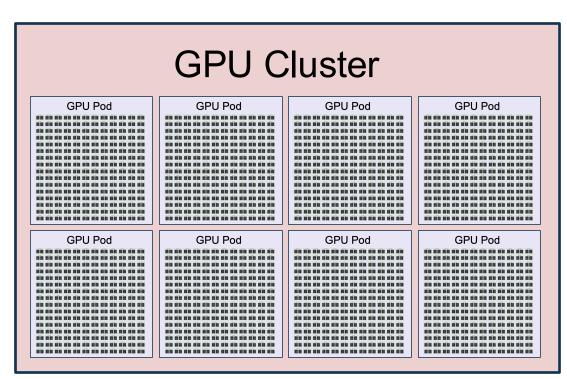
3352 GB/sec inside the GPU

Server = 8x GPU 900 GB/sec between GPUs

Rack = 2 Servers = 16x GPU

Pod = 192 Racks = 3072 GPUs 50 GB/sec between GPUs

Cluster = 8 Pods = 24,576 GPUs < 50GB/sec between GPUs

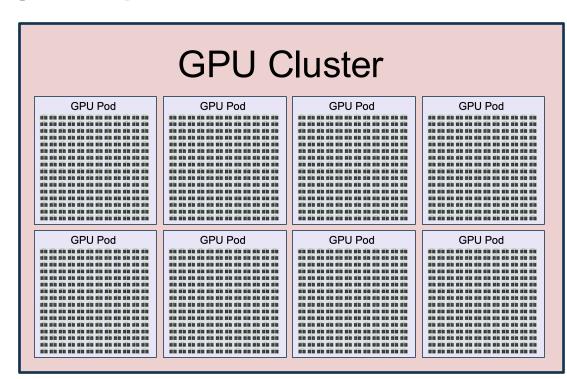


GPU Cluster = One Big Computer

Total Cluster Stats

24,576 GPUs
1.875 PB of GPU memory
415M FP32 cores
13M Tensor Cores
24.3 FFI OP/sec = 24.3 x 10¹⁸

Goal: Train one giant neural network on this cluster



Google: Tensor Processing Units (TPUs)

Custom chips designed by Google

TPU v5p: 459 TFLOP/sec BF16 per chip 95GB of memory per chip Arranged in pods of 8960 chips

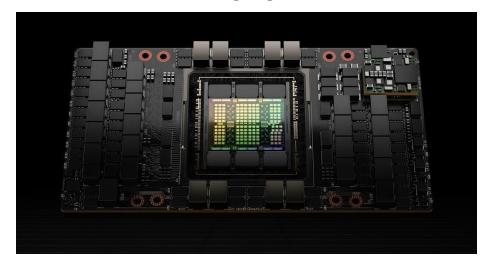


Other Training Chips

AMD MI325X 1300 TFLOP/sec BF16 256GB memory AWS Trainium2
667 TFLOP/sec BF16
96GB memory
Packed in UltraServers with 64 chips

Today: GPUs and How to Train On Them

A bit about GPU hardware

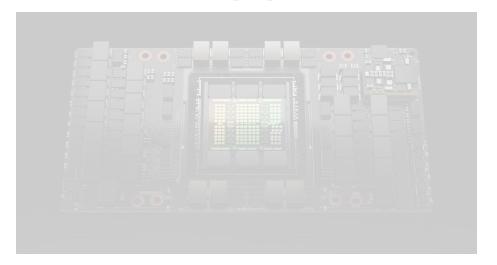


How to train on lots of GPUs

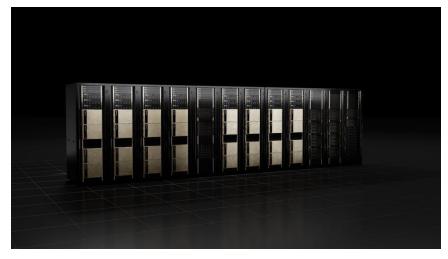


Today: GPUs and How to Train On Them

A bit about GPU hardware



How to train on lots of GPUs



How to train on lots of GPUs

A model with L layers operates on tensors of shape (Batch, Sequence, Dim)

Data Parallelism (DP)

Split on Batch dimension

Context Parallelism (CP)

Split on Sequence dimension

Pipeline Parallelism (PP)

Split on L dimension

Tensor Parallelism (TP)

Split on Dim dimension

How to train on lots of GPUs

A model with L layers operates on tensors of shape (Batch, Sequence, Dim)

Data Parallelism (DP)
Split on Batch dimension

Context Parallelism (CP)
Split on Sequence dimension

Pipeline Parallelism (PP)
Split on L dimension

Tensor Parallelism (TP)Split on Dim dimension

Recall: Loss is usually averaged over a minibatch of N samples

Recall: Loss is usually averaged over a minibatch of N samples

<u>Idea</u>: Use minibatch of MN samples, split over M GPUs

Recall: Loss is usually averaged over a minibatch of N samples

<u>Idea</u>: Use minibatch of MN samples, split over M GPUs

Gradients are linear, so each GPU computes its own gradient:

$$L = \frac{1}{MN} \sum_{i=1}^{M} \sum_{j=1}^{N} \ell(x_{i,j}, W)$$
$$\frac{\partial L}{\partial W} = \frac{1}{M} \sum_{i=1}^{M} \left(\frac{1}{N} \sum_{j=1}^{N} \frac{\partial}{\partial W} \ell(x_{i,j}, W) \right)$$

Recall: Loss is usually averaged over a minibatch of N samples

Idea: Use minibatch of MN samples, split over M GPUs

Gradients are linear, so each GPU computes its own gradient:

$$L = \frac{1}{MN} \sum_{i=1}^{M} \sum_{j=1}^{N} \ell(x_{i,j}, W)$$
 Each GPU computes gradient on N examples
$$\frac{\partial L}{\partial W} = \frac{1}{M} \sum_{i=1}^{M} \left(\frac{1}{N} \sum_{j=1}^{N} \frac{\partial}{\partial W} \ell(x_{i,j}, W) \right)$$

Each GPU

Recall: Loss is usually averaged over a minibatch of N samples

<u>Idea</u>: Use minibatch of MN samples, split over M GPUs

Gradients are linear, so each GPU computes its own gradient:

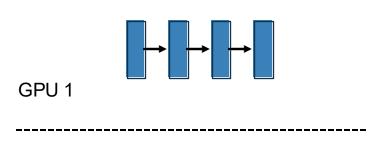
$$L = \frac{1}{MN} \sum_{i=1}^{M} \sum_{j=1}^{N} \ell(x_{i,j}, W)$$
 Each GPU computes gradient on N examples
$$\frac{\partial L}{\partial W} = \frac{1}{M} \sum_{i=1}^{M} \left(\frac{1}{N} \sum_{j=1}^{N} \frac{\partial}{\partial W} \ell(x_{i,j}, W) \right)$$
 Average gradients across M GPUs

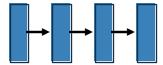
Recall: Loss is usually averaged over a minibatch of N samples

<u>Idea</u>: Use minibatch of MN samples, split over M GPUs

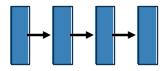
Gradients are linear, so each GPU computes its own gradient:

1. Each GPU has it's own copy of model + optimizer





GPU 2



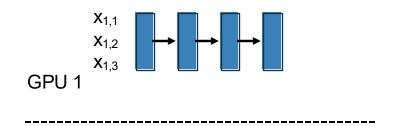
GPU 3

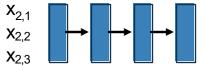
Recall: Loss is usually averaged over a minibatch of N samples

<u>Idea</u>: Use minibatch of MN samples, split over M GPUs

Gradients are linear, so each GPU computes its own gradient:

- 1. Each GPU has it's own copy of model + optimizer
- 2. Each GPU loads its own batch of data





GPU 2

X_{3,1} X_{3,2} X_{3,3}

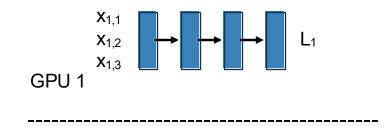
GPU 3

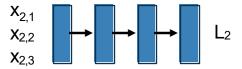
Recall: Loss is usually averaged over a minibatch of N samples

<u>Idea</u>: Use minibatch of MN samples, split over M GPUs

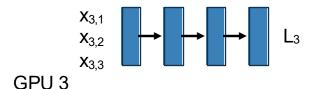
Gradients are linear, so each GPU computes its own gradient:

- 1. Each GPU has it's own copy of model + optimizer
- 2. Each GPU loads its own batch of data
- 3. Each GPU runs forward to compute loss





GPU 2

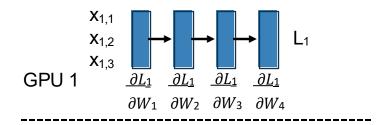


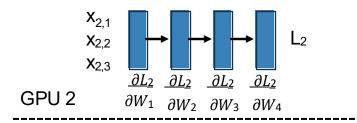
Recall: Loss is usually averaged over a minibatch of N samples

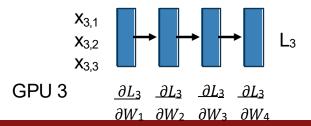
<u>Idea</u>: Use minibatch of MN samples, split over M GPUs

Gradients are linear, so each GPU computes its own gradient:

- 1. Each GPU has it's own copy of model + optimizer
- Each GPU loads its own batch of data
- 3. Each GPU runs forward to compute loss
- 4. Each GPU runs backward to compute gradients





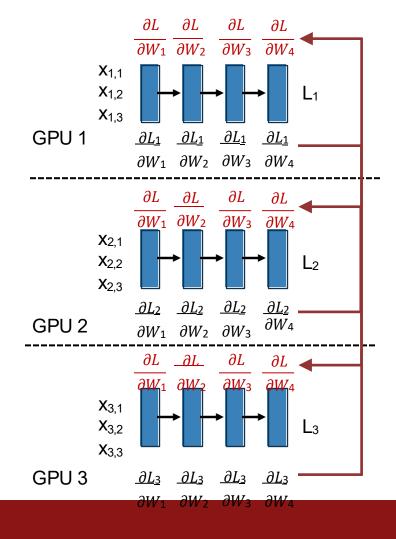


Recall: Loss is usually averaged over a minibatch of N samples

<u>Idea</u>: Use minibatch of MN samples, split over M GPUs

Gradients are linear, so each GPU computes its own gradient:

- 1. Each GPU has it's own copy of model + optimizer
- Each GPU loads its own batch of data
- 3. Each GPU runs forward to compute loss
- 4. Each GPU runs backward to compute gradients
- 5. Average gradients across all GPUs

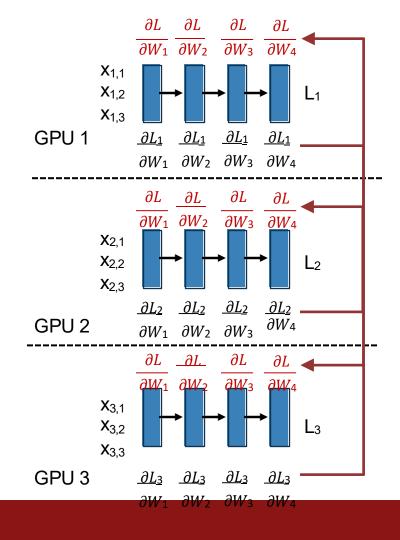


Recall: Loss is usually averaged over a minibatch of N samples

<u>Idea</u>: Use minibatch of MN samples, split over M GPUs

Gradients are linear, so each GPU computes its own gradient:

- 1. Each GPU has it's own copy of model + optimizer
- Each GPU loads its own batch of data
- 3. Each GPU runs forward to compute loss
- 4. Each GPU runs backward to compute gradients
- 5. Average gradients across all GPUs
- 6. Each GPU updates its own weights



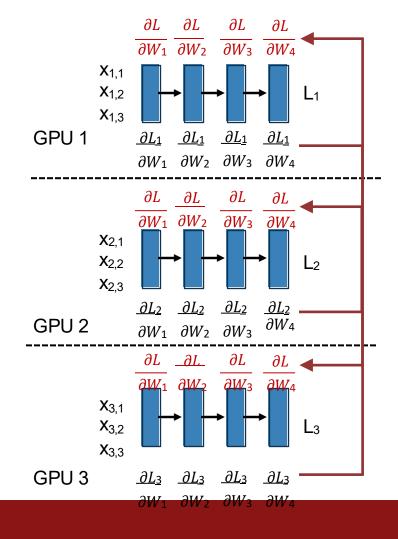
Recall: Loss is usually averaged over a minibatch of N samples

<u>Idea</u>: Use minibatch of MN samples, split over M GPUs

Gradients are linear, so each GPU computes its own gradient:

- 1. Each GPU has it's own copy of model + optimizer
- Each GPU loads its own batch of data
- 3. Each GPU runs forward to compute loss
- 4. Each GPU runs backward to compute gradients
- 5. Average gradients across all GPUs
- 6. Each GPU updates its own weights

(4) and (5) can run in parallel!



Recall: Loss is usually averaged over a minibatch of N samples

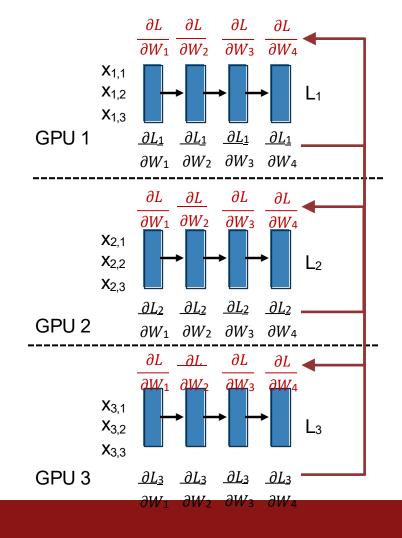
<u>Idea</u>: Use minibatch of MN samples, split over M GPUs

Gradients are linear, so each GPU computes its own gradient:

<u>Problem</u>: Model size constrained by GPU memory.

Each weight needs 4 numbers (weight, grad, Adam β_1 , β_2). Each number needs 2 bytes.

1B params takes 8GB; 10B params fills up 80GB GPU.



Recall: Loss is usually averaged over a minibatch of N samples

<u>Idea</u>: Use minibatch of MN samples, split over M GPUs

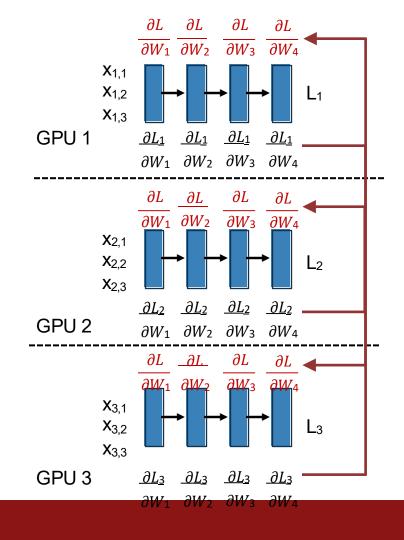
Gradients are linear, so each GPU computes its own gradient:

<u>Problem</u>: Model size constrained by GPU memory.

Each weight needs 4 numbers (weight, grad, Adam β_1 , β_2). Each number needs 2 bytes.

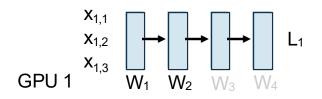
1B params takes 8GB; 10B params fills up 80GB GPU.

Solution: Split model weights across GPUs



Split model weights across GPUs

Each weight W_iis owned by one GPU, which also holds its grads and optim states



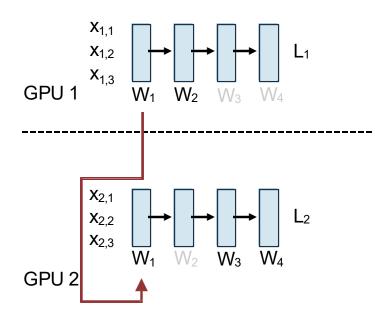
 $X_{2,1}$ $X_{2,2}$ $X_{2,3}$ W_1 W_2 W_3 W_4

GPU₂

Split model weights across GPUs

Each weight W_iis owned by one GPU, which also holds its grads and optim states

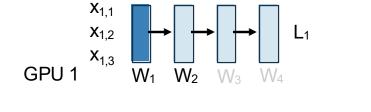
1. Before forward for layer i, the GPU that owns W_i broadcasts it to all GPUs

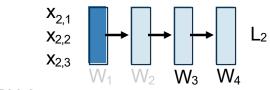


Split model weights across GPUs

Each weight W_iis owned by one GPU, which also holds its grads and optim states

- 1. Before forward for layer i, the GPU that owns W_i broadcasts it to all GPUs
- 2. All GPUs run forward for layer i, then delete their local copy of W_i



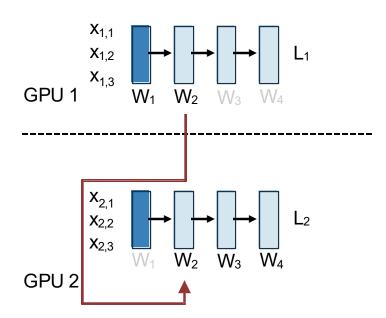


GPU 2

Split model weights across GPUs

Each weight W_iis owned by one GPU, which also holds its grads and optim states

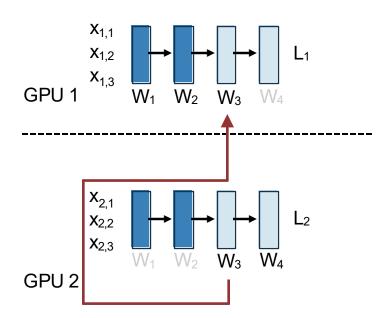
- 1. Before forward for layer i, the GPU that owns W_i broadcasts it to all GPUs
- 2. All GPUs run forward for layer i, then delete their local copy of W_i



Split model weights across GPUs

Each weight W_iis owned by one GPU, which also holds its grads and optim states

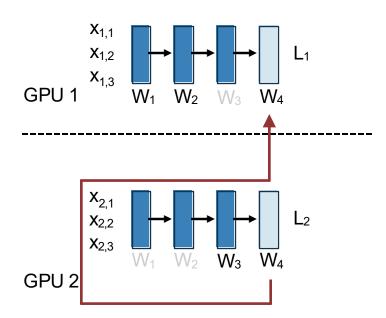
- 1. Before forward for layer i, the GPU that owns W_i broadcasts it to all GPUs
- 2. All GPUs run forward for layer i, then delete their local copy of W_i



Split model weights across GPUs

Each weight W_iis owned by one GPU, which also holds its grads and optim states

- 1. Before forward for layer i, the GPU that owns W_i broadcasts it to all GPUs
- 2. All GPUs run forward for layer i, then delete their local copy of W_i

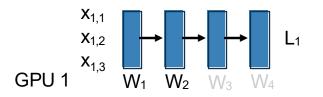


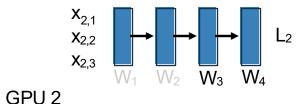
Split model weights across GPUs

Each weight W_i is owned by one GPU, which also holds its grads and optim states

- 1. Before forward for layer i, the GPU that owns W_i broadcasts it to all GPUs
- 2. All GPUs run forward for layer i, then delete their local copy of W;



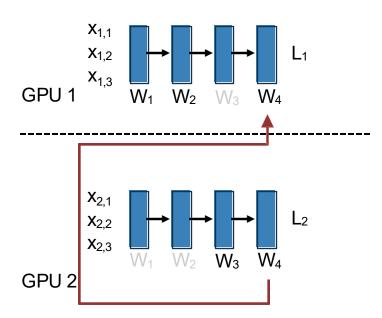




Split model weights across GPUs

Each weight W_iis owned by one GPU, which also holds its grads and optim states

- 1. Before forward for layer i, the GPU that owns W_i broadcasts it to all GPUs
- 2. All GPUs run forward for layer i, then delete their local copy of W_i
- Before backward for layer i, owner broadcasts W_i to all GPUs



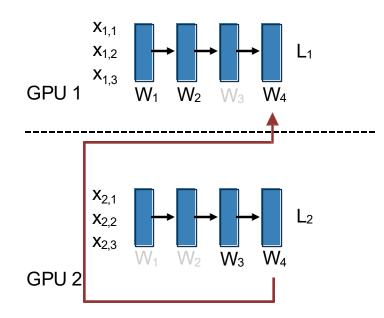
Optimization: don't delete last weight at end of forward to avoid immediately resending it

Fully Sharded Data Parallelism (FSPD)

Split model weights across GPUs

Each weight W_iis owned by one GPU, which also holds its grads and optim states

- Before forward for layer i, the GPU that owns W_i broadcasts it to all GPUs
- 2. All GPUs run forward for layer i, then delete their local copy of W_i
- 3. Before backward for layer i, owner broadcasts W_i to all GPUs

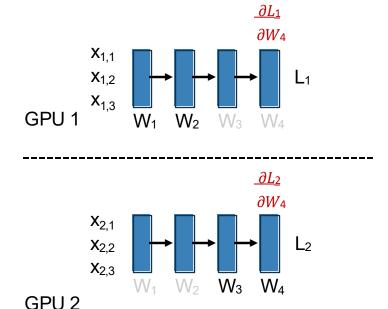


Optimization: don't delete last weight at end of forward to avoid immediately resending it

Split model weights across GPUs

Each weight W_i is owned by one GPU, which also holds its grads and optim states

- 1. Before forward for layer i, the GPU that owns W_i broadcasts it to all GPUs
- 2. All GPUs run forward for layer i, then delete their local copy of W_i
- 3. Before backward for layer i, owner broadcasts W_i to all GPUs
- All GPUs run backward for layer i to compute local dL/dW_i and delete W_i

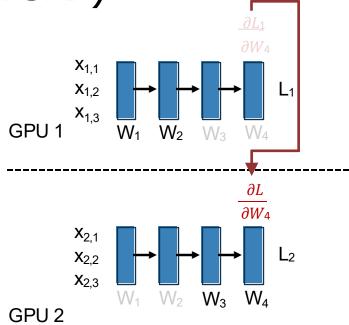


Split model weights across GPUs

Each weight Wis owned by one GPU, which also holds its grads and optim states

- Before forward for layer i, the GPU that owns Wi broadcasts it to all GPUs
- All GPUs run forward for layer i, then delete their local copy of W_i
- Before backward for layer i, owner 3. broadcasts Wi to all GPUs
- All GPUs run backward for layer i to compute local dL/dWi and delete Wi
- After backward for layer i, all GPUs send local dL/dWi to owning GPU and delete local dL/dWi

Fetch Wi+1 while computing forward with Wi



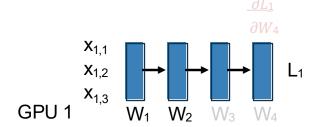
Optimization: don't delete last weight at end of forward to avoid immediately resending it

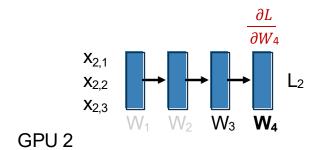
Optimization: don't delete last weight at end of forward to avoid immediately resending it

Split model weights across GPUs

Each weight W_iis owned by one GPU, which also holds its grads and optim states

- Before forward for layer i, the GPU that owns W_i broadcasts it to all GPUs
- 2. All GPUs run forward for layer i, then delete their local copy of W_i
- Before backward for layer i, owner broadcasts W_i to all GPUs
- All GPUs run backward for layer i to compute local dL/dW_i and delete W_i
- 5. After backward for layer i, all GPUs send local dL/dW_i to owning GPU and delete local dL/dW_i
- 6. Owner of Wi makes gradient update





Optimization: don't delete last weight at end of forward to avoid immediately resending it

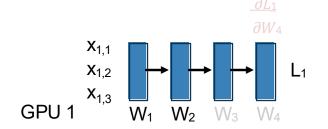
Split model weights across GPUs

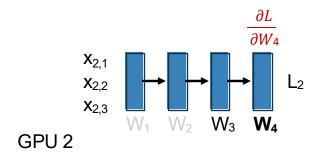
Each weight W_i is owned by one GPU, which also holds its grads and optim states

- Before forward for layer i, the GPU that owns W_i broadcasts it to all GPUs
- 2. All GPUs run forward for layer i, then delete their local copy of W_i
- Before backward for layer i, owner broadcasts W_i to all GPUs
- All GPUs run backward for layer i to compute local dL/dW_i and delete W_i
- After backward for layer i, all GPUs send local dL/dWi to owning GPU and delete local dL/dWi
- 6. Owner of W_i makes gradient update

Fetch W_{i+1} while computing forward with W_i

Fetch W_i while computing with W_{i+1}; send dL/dW_i and update W_i while computing with W_{i-1}





Optimization: don't delete last weight at end of forward to avoid immediately resending it

Fully Sharded Data Parallelism (FSPD)

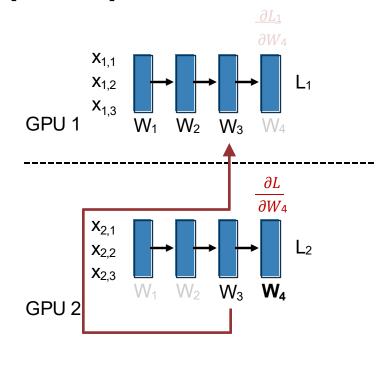
Split model weights across GPUs

Each weight W_iis owned by one GPU, which also holds its grads and optim states

- Before forward for layer i, the GPU that owns W_i broadcasts it to all GPUs
- 2. All GPUs run forward for layer i, then delete their local copy of W_i
- Before backward for layer i, owner broadcasts W_i to all GPUs
- All GPUs run backward for layer i to compute local dL/dW_i and delete W_i
- After backward for layer i, all GPUs send local dL/dWi to owning GPU and delete local dL/dWi
- 6. Owner of W_i makes gradient update

Fetch W_{i+1} while computing forward with W_i

Fetch W_i while computing with W_{i+1}; send dL/dW_i and update W_i while computing with W_{i-1}



Optimization: don't delete last weight at end of forward to avoid immediately resending it

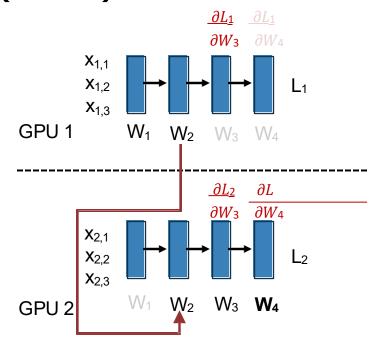
Split model weights across GPUs

Each weight W_i is owned by one GPU, which also holds its grads and optim states

- Before forward for layer i, the GPU that owns W_i broadcasts it to all GPUs
- 2. All GPUs run forward for layer i, then delete their local copy of W_i
- Before backward for layer i, owner broadcasts W_i to all GPUs
- All GPUs run backward for layer i to compute local dL/dW_i and delete W_i
- After backward for layer i, all GPUs send local dL/dW_i to owning GPU and delete local dL/dW_i
- 6. Owner of W_i makes gradient update

Fetch W_{i+1} while computing forward with W_i

Fetch W_i while computing with W_{i+1}; send dL/dW_i and update W_i while computing with W_{i-1}



Split model weights across GPUs

Each weight W_iis owned by one GPU, which also holds its grads and optim states

- Before forward for layer i, the GPU that owns W_i broadcasts it to all GPUs
- 2. All GPUs run forward for layer i, then delete their local copy of W_i
- Before backward for layer i, owner broadcasts W_i to all GPUs
- All GPUs run backward for layer i to compute local dL/dW_i and delete W_i
- 5. After backward for layer i, all GPUs send local dL/dW_i to owning GPU and delete local dL/dW_i
- 6. Owner of W_i makes gradient update

Fetch W_{i+1} while computing forward with W_i

Fetch W_i while computing with W_{i+1}; send dL/dW_i and update W_i while computing with W_{i-1}

 ∂L_1 ∂W_2 $X_{1.1}$ $X_{1,2}$ $X_{1,3}$ GPU 1 ∂L_2 ∂L ∂W_A $X_{2,1}$ L_2 $X_{2,2}$ $X_{2.3}$ GPU₂ All at the same time: Send grads and update W₃ Run backward with W₂ Fetch W₁

Optimization: don't delete last weight at end of forward to avoid immediately resending it

Rajbhandrari et al, "ZeRO: Memory Optimizations Toward Training Trillion Parameter Models", arXiv 2019

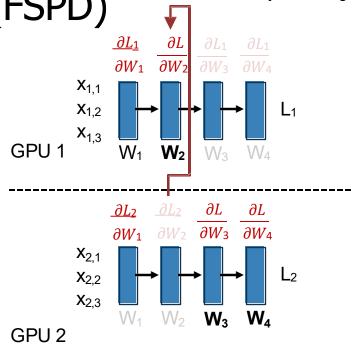
Split model weights across GPUs

Each weight W_i is owned by one GPU, which also holds its grads and optim states

- Before forward for layer i, the GPU that owns W_i broadcasts it to all GPUs
- 2. All GPUs run forward for layer i, then delete their local copy of W_i
- Before backward for layer i, owner broadcasts W_i to all GPUs
- All GPUs run backward for layer i to compute local dL/dW_i and delete W_i
- After backward for layer i, all GPUs send local dL/dWi to owning GPU and delete local dL/dWi
- 6. Owner of W_i makes gradient update

Fetch W_{i+1} while computing forward with W_i

Fetch W_i while computing with W_{i+1}; send dL/dW_i and update W_i while computing with W_{i-1}



Optimization: don't delete last weight at end of forward to avoid immediately resending it

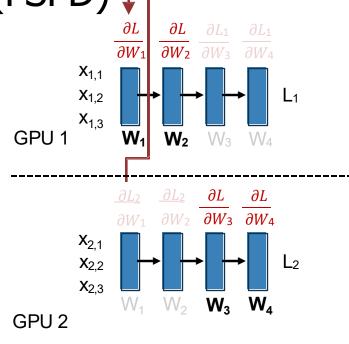
Split model weights across GPUs

Each weight W_i is owned by one GPU, which also holds its grads and optim states

- 1. Before forward for layer i, the GPU that owns W_i broadcasts it to all GPUs
- 2. All GPUs run forward for layer i, then delete their local copy of W_i
- Before backward for layer i, owner broadcasts W_i to all GPUs
- All GPUs run backward for layer i to compute local dL/dW_i and delete W_i
- After backward for layer i, all GPUs send local dL/dWi to owning GPU and delete local dL/dWi
- 6. Owner of W_i makes gradient update

Fetch W_{i+1} while computing forward with W_i

Fetch W_i while computing with W_{i+1}; send dL/dW_i and update W_i while computing with W_{i-1}



Optimization: don't delete last weight at end of forward to avoid immediately resending it

Optimization: don't delete last weight at end of forward to avoid immediately resending it

Fully Sharded Data Parallelism (FSPD)

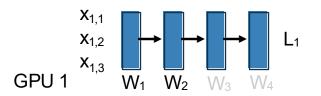
Split model weights across GPUs

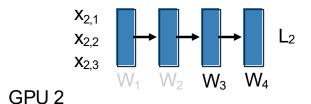
Each weight W_iis owned by one GPU, which also holds its grads and optim states

- 1. Before forward for layer i, the GPU that owns W_i broadcasts it to all GPUs
- 2. All GPUs run forward for layer i, then delete their local copy of W_i
- Before backward for layer i, owner broadcasts W_i to all GPUs
- All GPUs run backward for layer i to compute local dL/dW_i and delete W_i
- 5. After backward for layer i, all GPUs send local dL/dW_i to owning GPU and delete local dL/dW_i
- 6. Owner of W_i makes gradient update

Fetch W_{i+1} while computing forward with W_i

Fetch W_i while computing with W_{i+1}; send dL/dW_i and update W_i while computing with W_{i-1}





Repeat with next batch of data Data was being pre-fetched during forward+backward

Hybrid Sharded Data Parallel (HSDP)

Split N = M*K GPUs into M groups of K

Each group of K GPUs does FSDP, splits model weights across all K GPUs. K can be O(100) GPUs.

Do DP across the M groups.

Example: HSDP with M=2 groups of K=4 GPUs GPU (1, 1) W_3 W_4

Hybrid Sharded Data Parallel (HSDP)

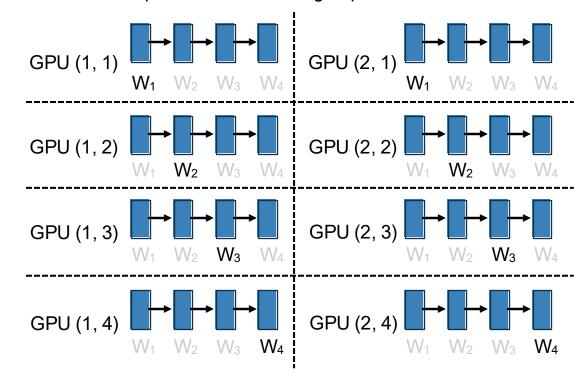
Split N = M*K GPUs into M groups of K

Each group of K GPUs does FSDP, splits model weights across all K GPUs. K can be O(100) GPUs.

Do DP across the M groups.

Multidimensional parallelism: Use different parallelism strategies at the same time! Organize GPUs in a 2D grid

Example: HSDP with M=2 groups of K=4 GPUs



Hybrid Sharded Data Parallel (HSDP)

Split N = M*K GPUs into M groups of K

Each group of K GPUs does FSDP, splits model weights across all K GPUs. K can be O(100) GPUs.

Do DP across the M groups.

Multidimensional parallelism: Use different parallelism strategies at the same time! Organize GPUs in a 2D grid

3x communication inside each group of K: W in forward, W + dL/dW in backward. Keep them in the same node / pod.

1x communication across the M groups: dL/dW in backward. Can use slower communication.

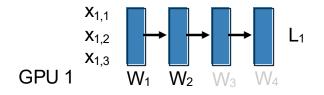
Example: HSDP with M=2 groups of K=4 GPUs GPU (1, 1) W_4

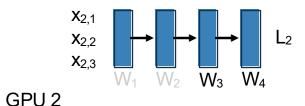
Data Parallelism (DP, FSPD, HSDP)

Split data and model weights across GPUs

Can now scale up to big models that don't fit in a single GPU!

A model with 100B params needs 4 numbers per param (param, grad, Adam β_1 , β_1); 2 bytes per number takes 800GB; splitting over 80 GPUs is just 10GB per GPU!





Data Parallelism (DP, FSPD, HSDP)

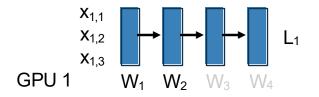
Split data and model weights across GPUs

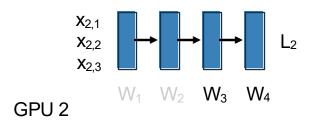
Can now scale up to big models that don't fit in a single GPU!

A model with 100B params needs 4 numbers per param (param, grad, Adam β_1 , β_1); 2 bytes per number takes 800GB; splitting over 80 GPUs is just 10GB per GPU!

Problem: Model activations can fill up memory. Llama3-405B Transformer has 126 layers, D=16,384, seq length 4096. Just FFN hidden activations need 2*126*(4*16384)*4096 bytes

= 63GB; plus need other activations.





Data Parallelism (DP, FSPD, HSDP)

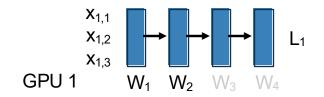
Split data and model weights across GPUs

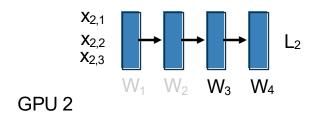
Can now scale up to big models that don't fit in a single GPU!

A model with 100B params needs 4 numbers per param (param, grad, Adam β_1 , β_1); 2 bytes per number takes 800GB; splitting over 80 GPUs is just 10GB per GPU!

Problem: Model activations can fill up memory. Llama3-405B Transformer has 126 layers, D=16,384, seq length 4096. Just FFN hidden activations need 2*126*(4*16384)*4096 bytes = 63GB; plus need other activations.

<u>Solution</u>: Don't keep all activations in memory; recompute them on the fly!



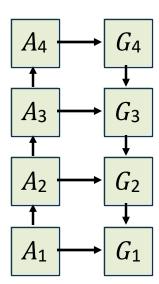


Activation Checkpointing

Each layer in the network is two functions:

Forward: Compute next-layer activations $A_{i+1} = F \xrightarrow{\rightarrow} A_i$

Backward: Compute prev-layer gradients $G_i = F_i^{\leftarrow}(A_i, G_{i+1})$



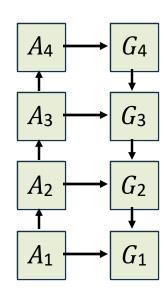
Activation Checkpointing

Each layer in the network is two functions:

Forward: Compute next-layer activations $A_{i+1} = F \xrightarrow{i} A_i$

Backward: Compute prev-layer gradients $G_i = F_i^{\leftarrow}(A_i, G_{i+1})$

Q: How much compute and memory does this take? Assume each F_i^{\rightarrow} and F_i^{\leftarrow} is O(1) compute and memory.



Activation Checkpointing

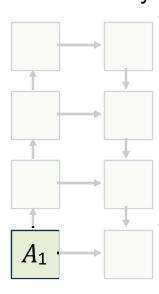
Each layer in the network is two functions:

Forward: Compute next-layer activations $A_{i+1} = F \xrightarrow{}_{i} A_{i}$

Backward: Compute prev-layer gradients $G_i = F_i^{\leftarrow}(A_i, G_{i+1})$

Q: How much compute and memory does this take? Assume each F_i^{\rightarrow} and F_i^{\leftarrow} is O(1) compute and memory.

Compute: 1
Current Memory: 1
Peak Memory: 1



Each layer in the network is two functions:

Forward: Compute next-layer activations $A_{i+1} = F \xrightarrow{}_{i} A_{i}$

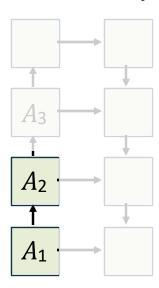
Backward: Compute prev-layer gradients $G_i = F_i^{\leftarrow}(A_i, G_{i+1})$

Q: How much compute and memory does this take? Assume each F_i^{\rightarrow} and F_i^{\leftarrow} is O(1) compute and memory.

Compute: 2

Current Memory: 2

Peak Memory: 2



Each layer in the network is two functions:

Forward: Compute next-layer activations $A_{i+1} = F \xrightarrow{}_{i} A_{i}$

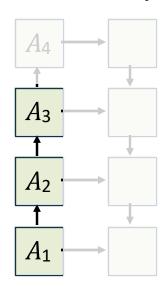
Backward: Compute prev-layer gradients $G_i = F_i^{\leftarrow}(A_i, G_{i+1})$

Q: How much compute and memory does this take? Assume each F_i^{\rightarrow} and F_i^{\leftarrow} is O(1) compute and memory.

Compute: 3

Current Memory: 3

Peak Memory: 3



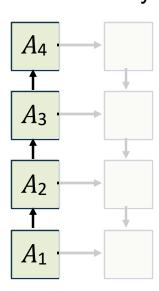
Each layer in the network is two functions:

Forward: Compute next-layer activations $A_{i+1} = F \xrightarrow{} A_i$

Backward: Compute prev-layer gradients $G_i = F_i^{\leftarrow}(A_i, G_{i+1})$

Q: How much compute and memory does this take? Assume each F_i^{\rightarrow} and F_i^{\leftarrow} is O(1) compute and memory.

Compute: 4
Current Memory: 4
Peak Memory: 4



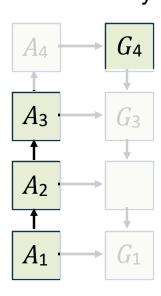
Each layer in the network is two functions:

Forward: Compute next-layer activations $A_{i+1} = F \xrightarrow{} A_i$

Backward: Compute prev-layer gradients $G_i = F_i^{\leftarrow}(A_i, G_{i+1})$

Q: How much compute and memory does this take? Assume each F_i^{\rightarrow} and F_i^{\leftarrow} is O(1) compute and memory.

Compute: 5
Current Memory: 4
Peak Memory: 4



Each layer in the network is two functions:

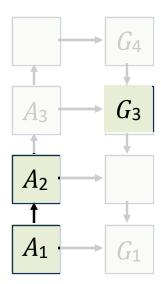
Forward: Compute next-layer activations $A_{i+1} = F \xrightarrow{}_{i} A_{i}$

Backward: Compute prev-layer gradients $G_i = F_i^{\leftarrow}(A_i, G_{i+1})$

Q: How much compute and memory does this take? Assume each F_i^{\rightarrow} and F_i^{\leftarrow} is O(1) compute and memory.

Compute: 6
Current Memory: 3

Peak Memory: 4



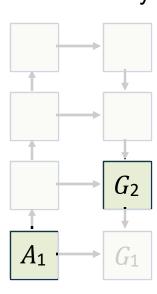
Each layer in the network is two functions:

Forward: Compute next-layer activations $A_{i+1} = F \xrightarrow{}_{i} A_{i}$

Backward: Compute prev-layer gradients $G_i = F_i^{\leftarrow}(A_i, G_{i+1})$

Q: How much compute and memory does this take? Assume each F_i^{\rightarrow} and F_i^{\leftarrow} is O(1) compute and memory.

Compute: 7
Current Memory: 2
Peak Memory: 4



Each layer in the network is two functions:

Forward: Compute next-layer activations $A_{i+1} = F \xrightarrow{}_{i} A_{i}$

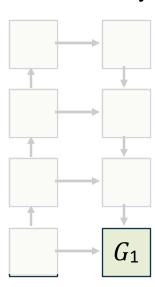
Backward: Compute prev-layer gradients $G_i = F_i^{\leftarrow}(A_i, G_{i+1})$

Q: How much compute and memory does this take? Assume each F_i^{\rightarrow} and F_i^{\leftarrow} is O(1) compute and memory.

Compute: 8

Current Memory: 1

Peak Memory: 4



Each layer in the network is two functions:

Forward: Compute next-layer activations $A_{i+1} = F \xrightarrow{}_{i} A_{i}$

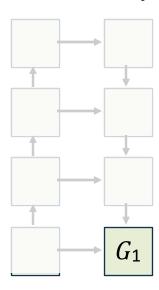
Backward: Compute prev-layer gradients $G_i = F_i^{\leftarrow}(A_i, G_{i+1})$

Forward+backward: O(N) compute, O(N) memory

Compute: 8

Current Memory: 1

Peak Memory: 4



Each layer in the network is two functions:

Forward: Compute next-layer activations $A_{i+1} = F \xrightarrow{}_{i} A_{i}$

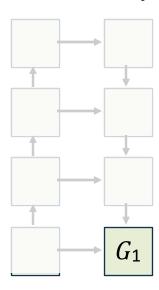
Backward: Compute prev-layer gradients $G_i = F_i^{\leftarrow}(A_i, G_{i+1})$

Forward+backward: O(N) compute, O(N) memory

Compute: 8

Current Memory: 1

Peak Memory: 4



Each layer in the network is two functions:

Forward: Compute next-layer activations $A_{i+1} = F \xrightarrow{i} A_i$

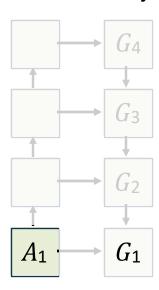
Backward: Compute prev-layer gradients $G_i = F_i^{\leftarrow}(A_i, G_{i+1})$

Forward+backward: O(N) compute, O(N) memory

Compute: 1

Current Memory: 1

Peak Memory: 1



Each layer in the network is two functions:

Forward: Compute next-layer activations $A_{i+1} = F \xrightarrow{i} A_i$

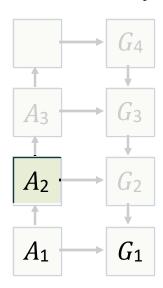
Backward: Compute prev-layer gradients $G_i = F_i^{\leftarrow}(A_i, G_{i+1})$

Forward+backward: O(N) compute, O(N) memory

Compute: 2

Current Memory: 1

Peak Memory: 1



Each layer in the network is two functions:

Forward: Compute next-layer activations $A_{i+1} = F \xrightarrow{i} A_i$

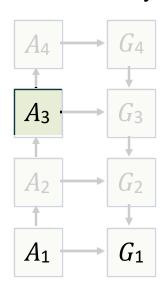
Backward: Compute prev-layer gradients $G_i = F_i^{\leftarrow}(A_i, G_{i+1})$

Forward+backward: O(N) compute, O(N) memory

Compute: 3

Current Memory: 1

Peak Memory: 1



Each layer in the network is two functions:

Forward: Compute next-layer activations $A_{i+1} = F \xrightarrow{i} A_i$

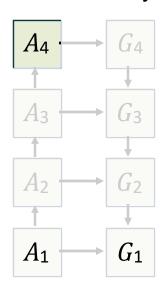
Backward: Compute prev-layer gradients $G_i = F_i^{\leftarrow}(A_i, G_{i+1})$

Forward+backward: O(N) compute, O(N) memory

Compute: 4

Current Memory: 1

Peak Memory: 1



Each layer in the network is two functions:

Forward: Compute next-layer activations $A_{i+1} = F \xrightarrow{} A_i$

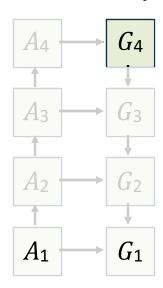
Backward: Compute prev-layer gradients $G_i = F_i^{\leftarrow}(A_i, G_{i+1})$

Forward+backward: O(N) compute, O(N) memory

Compute: 5

Current Memory: 1

Peak Memory: 1



Each layer in the network is two functions:

Forward: Compute next-layer activations $A_{i+1} = F \xrightarrow{} A_i$

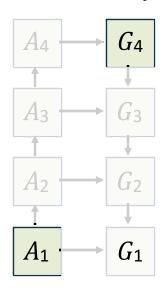
Backward: Compute prev-layer gradients $G_i = F_i^{\leftarrow}(A_i, G_{i+1})$

Forward+backward: O(N) compute, O(N) memory

Compute: 6

Current Memory: 2

Peak Memory: 2



Each layer in the network is two functions:

Forward: Compute next-layer activations $A_{i+1} = F \xrightarrow{} A_i$

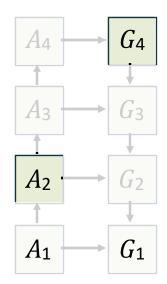
Backward: Compute prev-layer gradients $G_i = F_i^{\leftarrow}(A_i, G_{i+1})$

Forward+backward: O(N) compute, O(N) memory

Compute: 7

Current Memory: 2

Peak Memory: 2



Each layer in the network is two functions:

Forward: Compute next-layer activations $A_{i+1} = F \xrightarrow{} A_i$

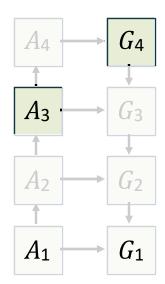
Backward: Compute prev-layer gradients $G_i = F_i^{\leftarrow}(A_i, G_{i+1})$

Forward+backward: O(N) compute, O(N) memory

Compute: 8

Current Memory: 2

Peak Memory: 2



Each layer in the network is two functions:

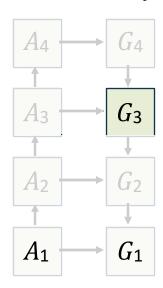
Forward: Compute next-layer activations $A_{i+1} = F \xrightarrow{} A_i$

Backward: Compute prev-layer gradients $G_i = F_i^{\leftarrow}(A_i, G_{i+1})$

Forward+backward: O(N) compute, O(N) memory

Compute: 9

Current Memory: 1 Peak Memory: 2



Each layer in the network is two functions:

Forward: Compute next-layer activations $A_{i+1} = F \xrightarrow{i} A_i$

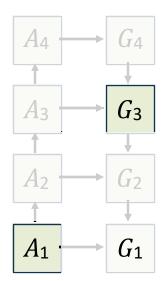
Backward: Compute prev-layer gradients $G_i = F_i^{\leftarrow}(A_i, G_{i+1})$

Forward+backward: O(N) compute, O(N) memory

Compute: 10

Current Memory: 2

Peak Memory: 2



Each layer in the network is two functions:

Forward: Compute next-layer activations $A_{i+1} = F \xrightarrow{i} A_i$

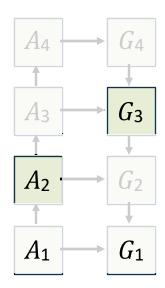
Backward: Compute prev-layer gradients $G_i = F_i^{\leftarrow}(A_i, G_{i+1})$

Forward+backward: O(N) compute, O(N) memory

Compute: 11

Current Memory: 2

Peak Memory: 2



Each layer in the network is two functions:

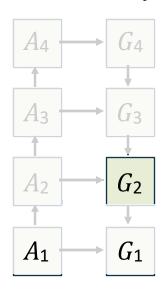
Forward: Compute next-layer activations $A_{i+1} = F \xrightarrow{i} A_i$

Backward: Compute prev-layer gradients $G_i = F_i^{\leftarrow}(A_i, G_{i+1})$

Forward+backward: O(N) compute, O(N) memory

Compute: 12

Current Memory: 1 Peak Memory: 2



Each layer in the network is two functions:

Forward: Compute next-layer activations $A_{i+1} = F \xrightarrow{i} A_i$

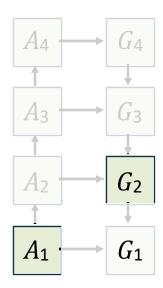
Backward: Compute prev-layer gradients $G_i = F_i^{\leftarrow}(A_i, G_{i+1})$

Forward+backward: O(N) compute, O(N) memory

Compute: 13

Current Memory: 2

Peak Memory: 2



Each layer in the network is two functions:

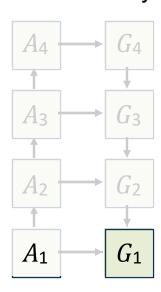
Forward: Compute next-layer activations $A_{i+1} = F \xrightarrow{i} A_i$

Backward: Compute prev-layer gradients $G_i = F_i^{\leftarrow}(A_i, G_{i+1})$

Forward+backward: O(N) compute, O(N) memory

Compute: 14

Current Memory: 1 Peak Memory: 2



Each layer in the network is two functions:

Forward: Compute next-layer activations $A_{i+1} = F \xrightarrow{}_{i} A_{i}$

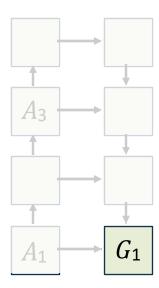
Backward: Compute prev-layer gradients $G_i = F_i^{\leftarrow}(A_i, G_{i+1})$

Forward+backward: O(N) compute, O(N) memory Full Recomputation: O(N²) compute, O(1) memory

Compute: 14

Current Memory: 1

Peak Memory: 2



Each layer in the network is two functions:

Forward: Compute next-layer activations $A_{i+1} = F \xrightarrow{i} A_i$

Backward: Compute prev-layer gradients $G_i = F_i^{\leftarrow}(A_i, G_{i+1})$

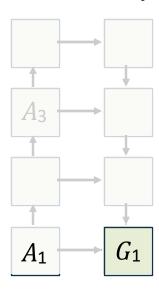
Forward+backward: O(N) compute, O(N) memory Full Recomputation: O(N2) compute, O(1) memory

Problem: N² compute is bad!

Compute: 14

Current Memory: 2

Peak Memory: 2



Each layer in the network is two functions:

Forward: Compute next-layer activations $A_{i+1} = F \xrightarrow{}_{i} A_{i}$

Backward: Compute prev-layer gradients $G_i = F_i^{\leftarrow}(A_i, G_{i+1})$

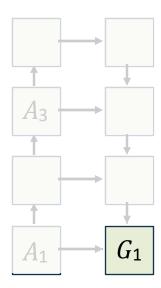
Forward+backward: O(N) compute, O(N) memory Full Recomputation: O(N²) compute, O(1) memory

<u>Problem</u>: N² compute is bad!

Compute: 14

Current Memory: 2

Peak Memory: 2



Idea: Don't recompute everything; save a checkpoint every C layers

Each layer in the network is two functions:

Forward: Compute next-layer activations $A_{i+1} = F \xrightarrow{i} A_i$

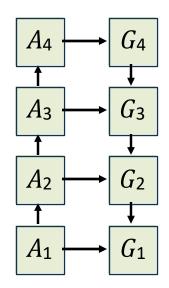
Backward: Compute prev-layer gradients $G_i = F_i^{\leftarrow}(A_i, G_{i+1})$

Forward+backward: O(N) compute, O(N) memory Full Recomputation: O(N²) compute, O(1) memory C checkpoints: O(N²/C) compute, O(C) memory

Compute: 14

Current Memory: 2

Peak Memory: 2



Idea: Don't recompute everything; save a checkpoint every C layers

Each layer in the network is two functions:

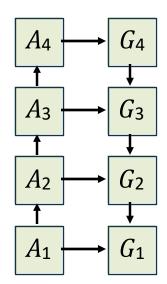
Forward: Compute next-layer activations $A_{i+1} = F \xrightarrow{i} A_i$

Backward: Compute prev-layer gradients $G_i = F_i^{\leftarrow}(A_i, G_{i+1})$

Forward+backward: O(N) compute, O(N) memory Full Recomputation: $O(N^2)$ compute, O(1) memory C checkpoints: $O(N^2/C)$ compute, O(C) memory VN checkpoints: O(N VN) compute, O(VN) memory Compute: 14

Current Memory: 2

Peak Memory: 2



Idea: Don't recompute everything; save a checkpoint every C layers

How to train on lots of GPUs

HSDP + Activation checkpointing can take you a long way!

Scaling recipe:

- 1. Use data parallelism up to ~128 GPUs, models with ~1B params
- 2. Always set per-GPU batch size to max out GPU memory
- 3. If your model is >1B params, consider **FSDP**
- 4. Add activation checkpointing to fit larger batches per GPU
- 5. If you have >256 GPUs, consider **HSDP**
- 6. If you have >1K GPUs, models >50B params, or sequence lengths > 16K then use more advanced strategies (CP, PP, TP)

Problem: Lots of knobs to tune! How should we set them?

Solution: Maximize Model Flops Utilization (MFU)

Recall: H100 can theoretically do 989.4 TFLOP/sec of 16-bit matrix multiplies on Tensor Cores

Question: How much throughput can we see in practice?

Recall: H100 can theoretically do 989.4 TFLOP/sec of 16-bit matrix multiplies on Tensor Cores

Question: How much throughput can we see in practice?

Hardware FLOPs Utilization (HFU): The fraction of theoretical matmul performance we actually achieve

Recall: H100 can theoretically do 989.4 TFLOP/sec of 16-bit matrix multiplies on Tensor Cores

Question: How much throughput can we see in practice?

Hardware FLOPs Utilization (HFU): The fraction of theoretical matmul performance we actually achieve

Benchmark for the best-case scenario: only matrix multiply

```
h100 tflop per sec = 989.4
sizes = [512, 1024, 2048, 4096,
         8192, 16_384, 32_768]
for N in sizes:
    x = torch.randn(N, N, device="cuda",
                    dtvpe=torch.bfloat16)
    flops = 2 * N * N * N
    times = []
    for i in range(12):
        t0 = time.time()
        y = x @ x
        if i > 2: times.append(time.time() - t0)
    sec = np.mean(times)
    tflops_per_sec = flops / sec / 10 ** 12
    hfu = 100 * tflops per sec / h100 tflop per sec
    print(f"N: {N}, "
          f"TFLOP/sec: {tflops_per_sec:.2f}, "
          f"HFU: {hfu:.2f}%")
```

Run this with CUDA_LAUNCH_BLOCKING=1, otherwise GPU kernels launch async and measurements are wrong

Recall: H100 can theoretically do 989.4 TFLOP/sec of 16-bit matrix multiplies on Tensor Cores

Question: How much throughput can we see in practice?

Hardware FLOPs Utilization (HFU):

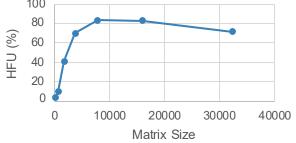
The fraction of theoretical matmul performance we actually achieve

Benchmark for the best-case scenario: only matrix multiply

```
h100 tflop per sec = 989.4
sizes = [512, 1024, 2048, 4096,
         8192, 16_384, 32_768]
for N in sizes:
    x = torch.randn(N, N, device="cuda",
                    dtvpe=torch.bfloat16)
    flops = 2 * N * N * N
    times = []
    for i in range(12):
        t0 = time.time()
        y = x @ x
        if i > 2: times.append(time.time() - t0)
    sec = np.mean(times)
    tflops_per_sec = flops / sec / 10 ** 12
    hfu = 100 * tflops per sec / h100 tflop per sec
    print(f"N: {N}, "
          f"TFLOP/sec: {tflops_per_sec:.2f}, "
          f"HFU: {hfu:.2f}%")
```

Large matrix 100

Large matrix multiply gets ~80% HFU on H100



Matmul HFU on H100

Recall: H100 can theoretically do 989.4 TFLOP/sec of 16-bit matrix multiplies on Tensor Cores

Question: How much throughput can we see in practice?

Hardware FLOPs Utilization (HFU):

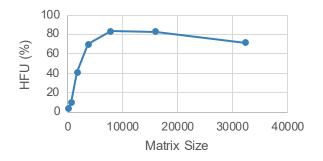
The fraction of theoretical matmul performance we actually achieve

<u>Problem</u>: HFU does not account for activation checkpointing or "helper" computation like data augmentation, optimizer, preprocessing

Benchmark for the best-case scenario: only matrix multiply h100_tflop_per_sec = 989.4 sizes = [512, 1024, 2048, 4096, 8192, 16_384, 32_768] for N in sizes: x = torch.randn(N, N, device="cuda", dtvpe=torch.bfloat16) flops = 2 * N * N * Ntimes = [] for i in range(12): t0 = time.time() y = x @ xif i > 2: times.append(time.time() - t0) sec = np.mean(times) tflops_per_sec = flops / sec / 10 ** 12 hfu = 100 * tflops per sec / h100 tflop per sec print(f"N: {N}, " f"TFLOP/sec: {tflops_per_sec:.2f}, " f"HFU: {hfu:.2f}%")

Matmul HFU on H100

Large matrix multiply gets ~80% HFU on H100



Model FLOPs Utilization (MFU)

<u>Idea</u>: What fraction of the GPU's theoretical peak FLOPs is being used for "useful" model computation?

- Compute FLOP_{theoretical} = total number of matrix multiply FLOPs in the forward + backward pass (can approximate backward = 2x forward) (Ignore nonlinearities, normalization, elementwise ops like residuals. They will run on FP32 cores)
- 2. Look up FLOP/sec_{theoretical} = theoretical max throughput of your device (H100: 989 TFLOP/sec)
- **3** Compute t_{theoretical} = FLOP_{theoretical} / FLOP/sec_{theoretical}
- 4. Measure t_{actual} = Actual time for a full iteration of data loading, forward, backward, optimizer step
- 5. MFU = $t_{theoretical}$ / t_{actual}

Model FLOPs Utilization (MFU)

Idea: What fraction of the GPU's theoretical peak FLOPs is being used for "useful" model computation?

- Compute FLOP_{theoretical} = total number of matrix multiply FLOPs in the forward + backward pass (can approximate backward = 2x forward) (Ignore nonlinearities, normalization, elementwise ops like residuals. They will run on FP32 cores)
- Look up FLOP/sec_{theoretical} = theoretical maxthroughput of your device (H100: 989 TFLOP/sec)
- 3 Compute t_{theoretical} = FLOP_{theoretical} / FLOP/sec_{theoretical}
- 4. Measure t_{actual} = Actual time for a full iteration of data loading, forward, backward, optimizer step
- 5. MFU = $t_{theoretical} / t_{actual}$

```
L, D, N = 8, 8192, 8192
flop fwd = N * L * 2 * D * D
flop_bwd = 2 * flop_fwd
flop theoretical = flop fwd + flop bwd
t theoretical = flop theoretical / (989.4 * 10 ** 12)
layers = []
for _ in range(L):
    layers += [torch.nn.Linear(D, D), torch.nn.ReLU()]
model = torch.nn.Sequential(*layers).cuda()
optimizer = torch.optim.Adam(model.parameters(), lr=1e-4)
for _ in range(20):
    torch.cuda.synchronize()
                              Example: Wide
    t0 = time.time()
   x = torch.randn(
                              MLP with big
       N, D, device="cuda",
                              batch size gets
       dtype=torch.float32
                              ~49% MFU on
   with torch.autocast(
       device type="cuda",
                              H100
       dtype=torch.bfloat16,
       y = model(x)
    loss = ((x - y) ** 2.0).sum()
    loss.backward()
   optimizer.step()
   optimizer.zero_grad()
    torch.cuda.synchronize()
    t actual = time.time() - t0
   mfu = t_theoretical / t_actual
   print(f"MFU: {100*mfu:.2f}%")
```

<u>Idea</u>: What fraction of the GPU's theoretical peak FLOPs is being used for "useful" model computation?

- Compute FLOP_{theoretical} = total number of matrix multiply FLOPs in the forward + backward pass (can approximate backward = 2x forward) (Ignore nonlinearities, normalization, elementwise ops like residuals. They will run on FP32 cores)
- Look up FLOP/sec_{theoretical} = theoretical max throughput of your device (H100: 989 TFLOP/sec)
- **3** Compute t_{theoretical} = FLOP_{theoretical} / FLOP/sec_{theoretical}
- 4. Measure t_{actual} = Actual time for a full iteration of data loading, forward, backward, optimizer step
- 5. MFU = $t_{theoretical}$ / t_{actual}

Optimize distributed training setup to maximize MFU!

```
L, D, N = 8, 8192, 8192
flop fwd = N * L * 2 * D * D
flop_bwd = 2 * flop_fwd
flop theoretical = flop fwd + flop bwd
t theoretical = flop theoretical / (989.4 * 10 ** 12)
layers = []
for _ in range(L):
    layers += [torch.nn.Linear(D, D), torch.nn.ReLU()]
model = torch.nn.Sequential(*layers).cuda()
optimizer = torch.optim.Adam(model.parameters(), lr=1e-4)
for _ in range(20):
    torch.cuda.synchronize()
                              Example: Wide
    t0 = time.time()
   x = torch.randn(
                              MLP with big
       N, D, device="cuda",
                              batch size gets
       dtype=torch.float32
                              ~49% MFU on
   with torch.autocast(
       device type="cuda",
                              H100
       dtype=torch.bfloat16,
       y = model(x)
    loss = ((x - y) ** 2.0).sum()
    loss.backward()
   optimizer.step()
   optimizer.zero_grad()
    torch.cuda.synchronize()
    t actual = time.time() - t0
   mfu = t_theoretical / t_actual
    print(f"MFU: {100*mfu:.2f}%")
```

<u>Idea</u>: What fraction of the GPU's theoretical peak FLOPs is being used for "useful" model computation?

MFU >30% is good, >40% is excellent

<u>Idea</u>: What fraction of the GPU's theoretical peak FLOPs is being used for "useful" model computation?

MFU >30% is good, >40% is excellent

Model	# of Parameters (in billions)	Accelerator chips	Model FLOPS utilization	
GPT-3	175B	V100	21.3%	
Gopher	280B	4096 TPU v3	32.5%	
Megatron-Turing NLG	530B	2240 A100	30.2%	
PaLM	540B	6144 TPU v4	46.2%	

Chowdhery et al, "PaLM: Scaling Language Modeling with Pathways", arXiv 2022

Example: Llama3-405B training on H100 GPUs

GPUs	TFLOPs/GPU	BF16 MFU
8,192	430	43%
16,384	400	41%
16,384	380	38%

Llama Team, "The Llama3 Herd of Models", arXiv 2024

<u>Idea</u>: What fraction of the GPU's theoretical peak FLOPs is being used for "useful" model computation?

MFU >30% is good, >40% is excellent

Model	# of Parameters (in billions) Accelerator chips		Model FLOPS utilization	
GPT-3	175B	V100	21.3%	
Gopher	280B	4096 TPU v3	32.5%	
Megatron-Turing NLG	530B	2240 A100	30.2%	
PaLM	540B	6144 TPU v4	46.2%	

Chowdhery et al, "PaLM: Scaling Language Modeling with Pathways", arXiv 2022

Example: Llama3-405B training on H100 GPUs

GPUs	TFLOPs/GPU	BF16 MFU
8,192	430	43%
$16,\!384$	400	41%
16,384	380	38%

Llama Team, "The Llama3 Herd of Models", arXiv 2024

More recent devices sometimes get worse MFU since their peak FLOPs increases much faster than their memory bandwidth

A100 => H100: 3.1x FLOPs

2.1x memory bandwidth

How to train on lots of GPUs

A model with L layers operates on tensors of shape (Batch, Sequence, Dim)

Data Parallelism (DP)
Split on Batch dimension

Context Parallelism (CP)
Split on Sequence dimension

Pipeline Parallelism (PP)
Split on L dimension

Tensor Parallelism (TP)Split on Dim dimension

How to train on lots of GPUs

A model with L layers operates on tensors of shape (Batch, Sequence, Dim)

Data Parallelism (DP)

Split on Batch dimension

Context Parallelism (CP)

Split on Sequence dimension

Pipeline Parallelism (PP)

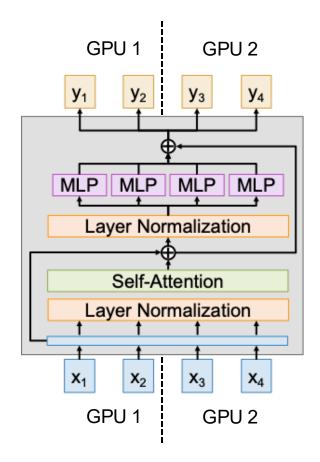
Split on L dimension

Tensor Parallelism (TP)

Split on Dim dimension

(Usually for Transformers)

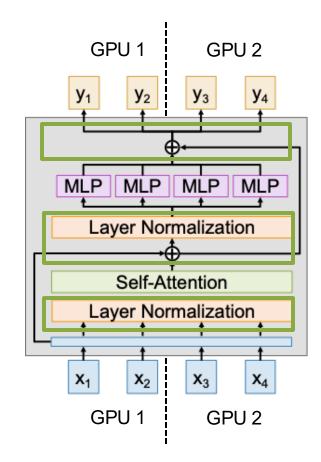
<u>Idea</u>: Transformers operate on L-length sequences. Use multiple GPUs to process a single long sequence



(Usually for Transformers)

<u>Idea</u>: Transformers operate on L-length sequences. Use multiple GPUs to process a single long sequence

Normalization, residual connections: Easy, they have no weights and trivially parallelizable

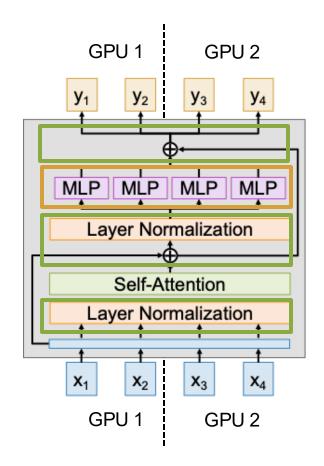


(Usually for Transformers)

<u>Idea</u>: Transformers operate on L-length sequences. Use multiple GPUs to process a single long sequence

Normalization, residual connections: Easy, they have no weights and trivially parallelizable

MLP: Trivially parallelizable, but has weights. Each GPU keeps a copy of the weights and communicates gradients like in DP



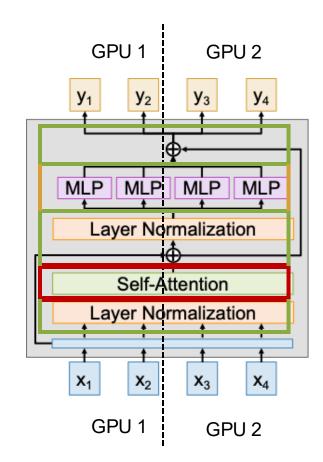
(Usually for Transformers)

<u>Idea</u>: Transformers operate on L-length sequences. Use multiple GPUs to process a single long sequence

Normalization, residual connections: Easy, they have no weights and trivially parallelizable

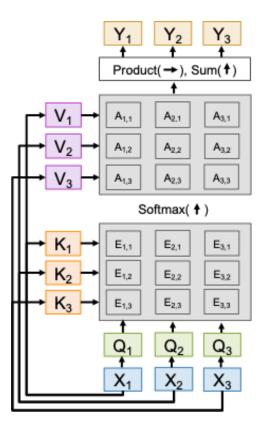
MLP: Trivially parallelizable, but has weights. Each GPU keeps a copy of the weights and communicates gradients like in DP

Attention: More complex, need to dig in



(Usually for Transformers)

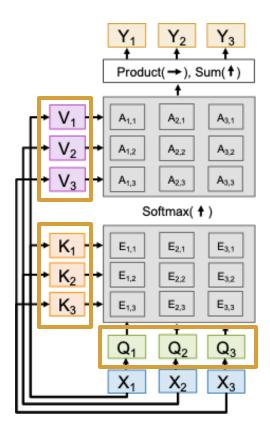
<u>Idea</u>: Transformers operate on L-length sequences. Use multiple GPUs to process a single long sequence



(Usually for Transformers)

<u>Idea</u>: Transformers operate on L-length sequences. Use multiple GPUs to process a single long sequence

QKV Projection: Same as MLP, parallelize over the sequence and sync gradients as in DP

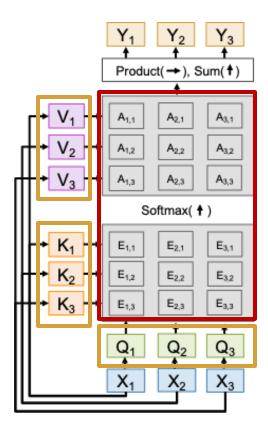


(Usually for Transformers)

<u>Idea</u>: Transformers operate on S-length sequences. Use multiple GPUs to process a single long sequence

QKV Projection: Same as MLP, parallelize over the sequence and sync gradients as in DP

Attention operator: Hardest to parallelize



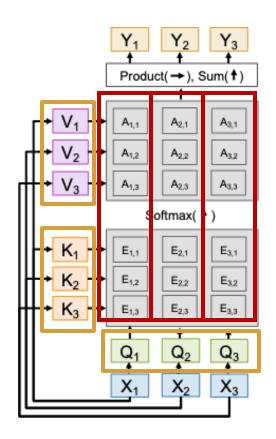
(Usually for Transformers)

<u>Idea</u>: Transformers operate on S-length sequences. Use multiple GPUs to process a single long sequence

QKV Projection: Same as MLP, parallelize over the sequence and sync gradients as in DP

Attention operator: Hardest to parallelize

(Option 1) Ring Attention: Divide into blocks and distribute over GPUs. Inner loop over keys/values, outer loop over queries. Complex to implement but can scale to very long sequences.



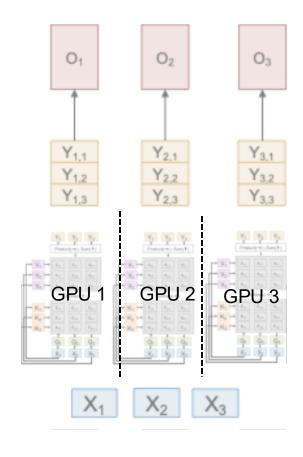
(Usually for Transformers)

<u>Idea</u>: Transformers operate on S-length sequences. Use multiple GPUs to process a single long sequence

QKV Projection: Same as MLP, parallelize over the sequence and sync gradients as in DP

Attention operator: Hardest to parallelize

(Option 2) Ulysses: Don't try to distribute attention matrix, instead parallelize over heads in multihead attention. Simpler, but max parallelism = number of heads



(Usually for Transformers)

<u>Idea</u>: Transformers operate on S-length sequences. Use multiple GPUs to process a single long sequence

Often used for long-sequence finetuning.

Example: Llama3-405B training:

- Stage 1: S=8192, no context-parallelism
- Stage 2: S=131,072, 16-way context-parallelism (8192 per GPU)

How to train on lots of GPUs

A model with L layers operates on tensors of shape (Batch, Sequence, Dim)

Data Parallelism (DP)

Split on Batch dimension

Context Parallelism (CP)

Split on Sequence dimension

Pipeline Parallelism (PP)

Split on L dimension

Tensor Parallelism (TP)

Split on Dim dimension

How to train on lots of GPUs

A model with L layers operates on tensors of shape (Batch, Sequence, Dim)

Data Parallelism (DP)

Split on Batch dimension

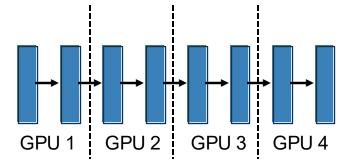
Context Parallelism (CP)

Split on Sequence dimension

Pipeline Parallelism (PP)
Split on L dimension

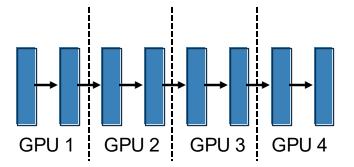
Tensor Parallelism (TP)Split on Dim dimension

<u>Idea</u>: Split the layers of the model across GPUs. Copy activations between layers at GPU boundaries.



<u>Idea</u>: Split the layers of the model across GPUs. Copy activations between layers at GPU boundaries.

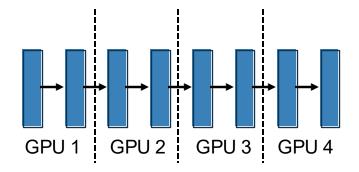
<u>Problem</u>: Sequential dependencies; GPUs are mostly sitting idle.



Idea: Split the layers of the model across GPUs. Copy activations between layers at GPU boundaries.

Problem: Sequential dependencies; GPUs are mostly sitting idle.

Max MFU with N-way PP is 1/N

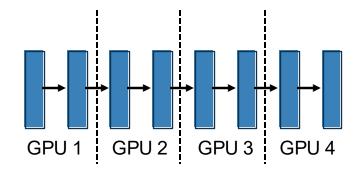


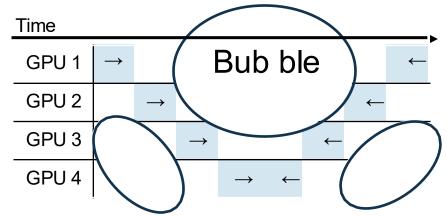
_	Time									
	GPU 1	\rightarrow							←	
	GPU 2		\rightarrow					\downarrow		
	GPU 3			\rightarrow						
	GPU 4				\rightarrow	←				

<u>Idea</u>: Split the layers of the model across GPUs. Copy activations between layers at GPU boundaries.

Problem: Sequential dependencies; GPUs are mostly sitting idle.

Max MFU with N-way PP is 1/N





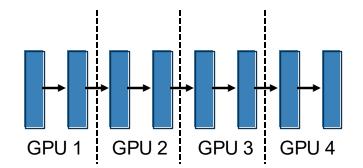
Huang et al, "Gpipe: Efficient Training of Giant Neural Networks using Pipeline Parallelism", arXiv 2018

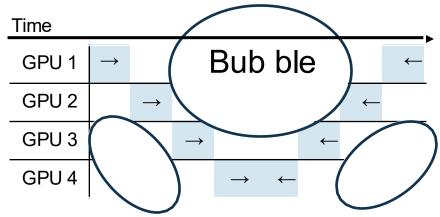
Idea: Split the layers of the model across GPUs. Copy activations between layers at GPU boundaries.

Problem: Sequential dependencies; GPUs are mostly sitting idle.

Max MFU with N-way PP is 1/N

Solution: Run multiple **microbatches** at the same time, pipeline them through the GPUs

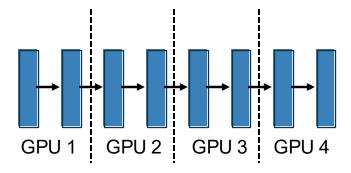




Huang et al, "Gpipe: Efficient Training of Giant Neural Networks using Pipeline Parallelism", arXiv 2018

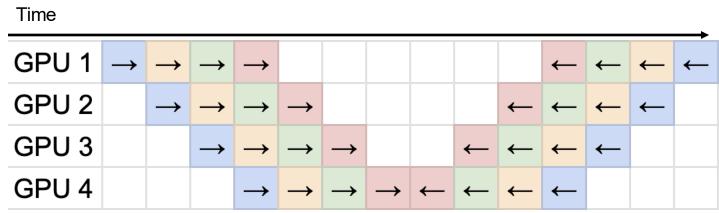
Pipeline Parallelism (PP) - Microbatches

<u>Idea</u>: Split the layers of the model across GPUs. Copy activations between layers at GPU boundaries.



Example: 4-way PP with 4 microbatches.

Max MFU increases from 1/4 = 25% to $16/28 \approx 57.1\%$



Huang et al, "Gpipe: Efficient Training of Giant Neural Networks using Pipeline Parallelism", arXiv 2018

How to train on lots of GPUs

A model with L layers operates on tensors of shape (Batch, Sequence, Dim)

Data Parallelism (DP)

Split on Batch dimension

Context Parallelism (CP)

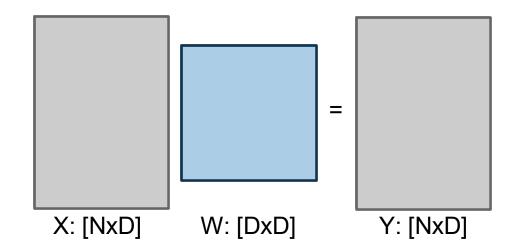
Split on Sequence dimension

Pipeline Parallelism (PP)
Split on L dimension

Tensor Parallelism (TP)Split on Dim dimension

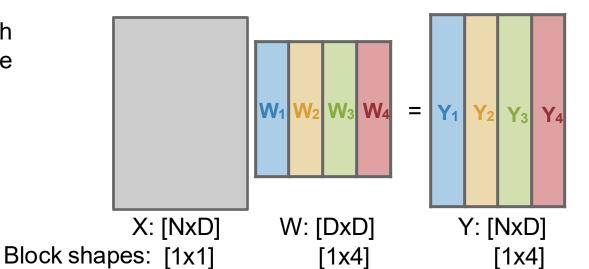
<u>Idea</u>: Split the weights of each linear layer across GPUs, use block matrix multiply

$$XW = Y (1 GPU)$$



XW = Y (4-way TP)

Idea: Split the weights of each linear layer across GPUs, use block matrix multiply

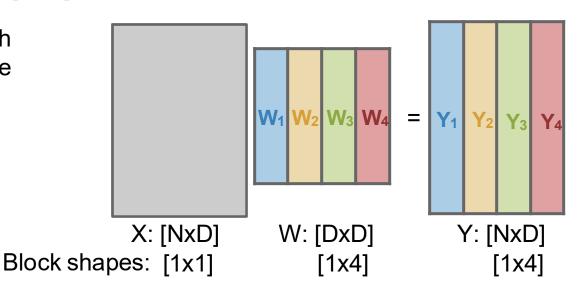


GPU i computes XW_i = Y_i

XW = Y (4-way TP)

Idea: Split the weights of each linear layer across GPUs, use block matrix multiply

<u>Problem</u>: Need to gather parts of Y after forward, can't overlap with communication



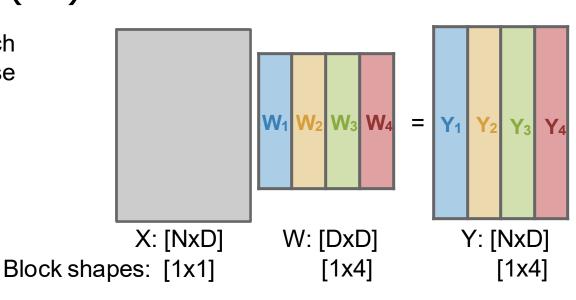
GPU i computes
$$XW_i = Y_i$$

XW = Y (4-way TP)

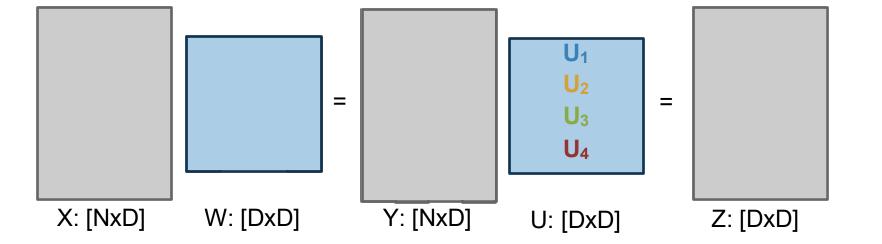
Idea: Split the weights of each linear layer across GPUs, use block matrix multiply

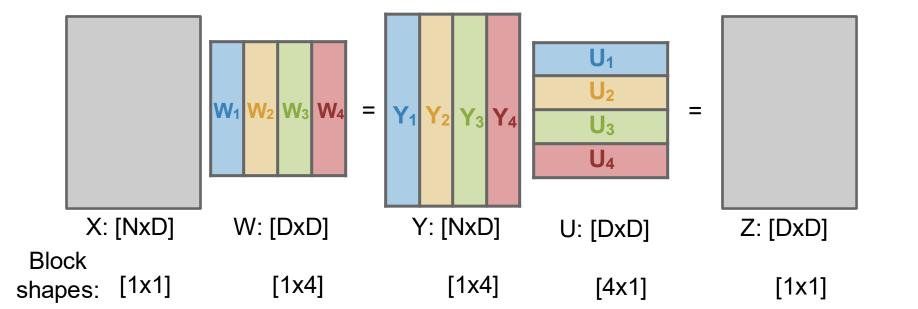
<u>Problem</u>: Need to gather parts of Y after forward, can't overlap with communication

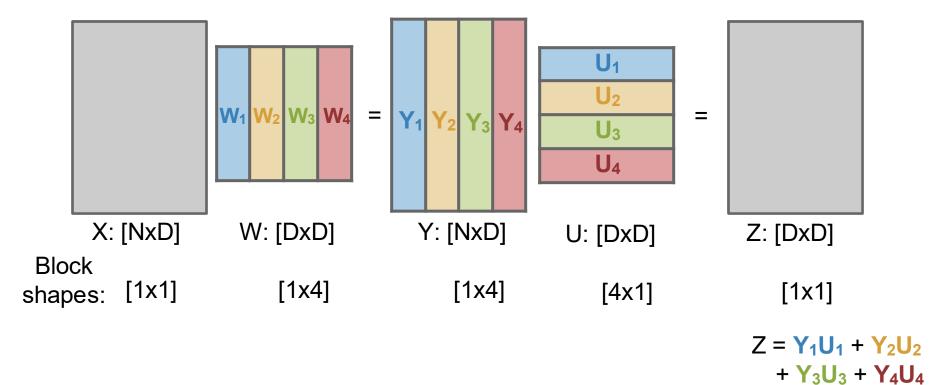
Trick: With 2 consecutive TP layers, shard first over row and second over column to avoid communication

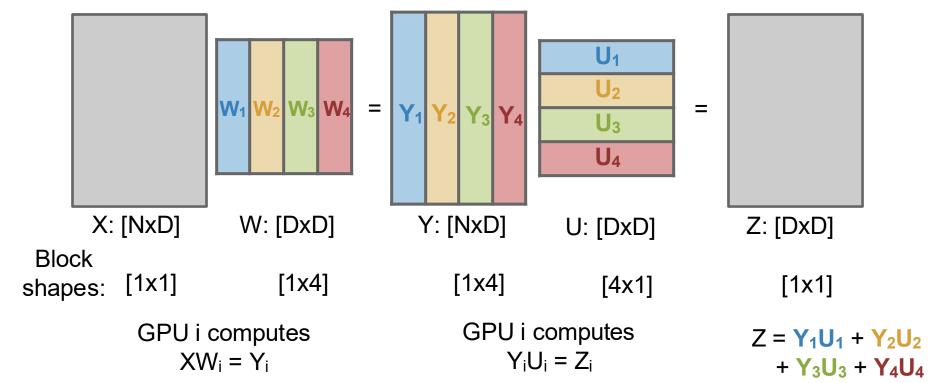


GPU i computes
$$XW_i = Y_i$$

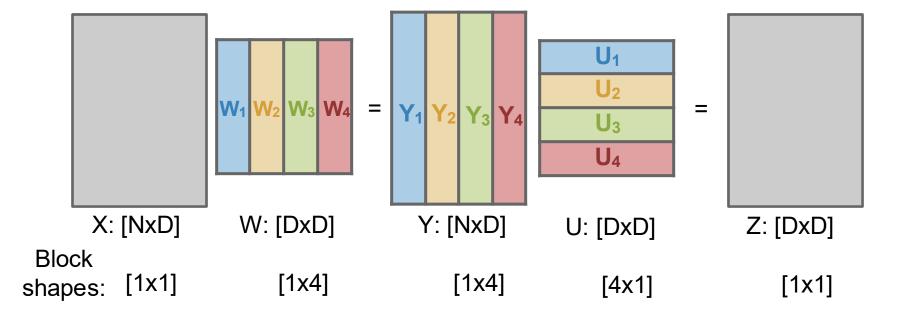








(4-way TP) XW = Y (layer 1) YU = Z (layer 2)



No need for communication after XW=Y! Each GPU computes one term of Z, then broadcasts to all other GPUs

$$Z = Y_1U_1 + Y_2U_2 + Y_3U_3 + Y_4U_4$$

How to train on lots of GPUs

A model with L layers operates on tensors of shape (Batch, Sequence, Dim)

Data Parallelism (DP)

Split on Batch dimension

Context Parallelism (CP)

Split on Sequence dimension

Pipeline Parallelism (PP)
Split on L dimension

Tensor Parallelism (TP)Split on Dim dimension

How to train on lots of GPUs

A model with L layers operates on tensors of shape (Batch, Sequence, Dim)

Data Parallelism (DP)

Split on Batch dimension

Context Parallelism (CP)

Split on Sequence dimension

Q: Which to use for largest models?

A: All of them!

Pipeline Parallelism (PP)

Split on L dimension

Tensor Parallelism (TP)

Split on Dim dimension

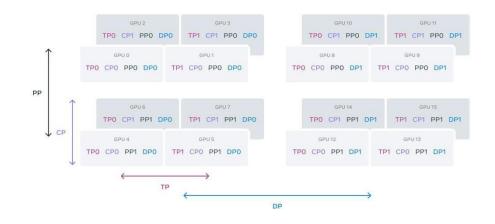
ND Parallelism

Use TP, CP, PP, and DP all at the same time!

Arrange GPUs in a 4D grid

GPUs index in the grid gives its rank along each parallelism dimension

Optimize setup to maximize MFU

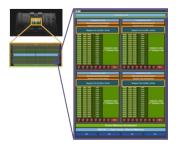


Example: LLama3-405B

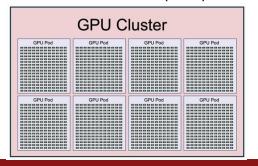
GPUs	TP	СР	PP	DP	Seq. Len.	Batch size/DP	Tokens/Batch	TFLOPs/GPU	BF16 MFU
8,192	8	1	16	64	8,192	32	16M	430	43%
16,384	8	1	16	128	8,192	16	16M	400	41%
16,384	8	16	16	8	131,072	16	16M	380	38%

Summary: Large-Scale Distributed Training

A GPU is a parallel processor with hundreds of cores



A GPU cluster has O(10K) GPUs



Split up the computation along different axes Consider a model with many Layers, operating on tensors of shape (Batch, Seq, Dim)

- Data Parallel (DP): Split on Batch
- Context Parallel (CP): Split on Seq
- Pipeline Parallel (PP): Split on Layers
- Tensor Parallel (TP): Split on Dim

Activation Checkpointing saves memory by recomputing during backward

Tune parallelism recipe to maximize **Model Flops Utilization (MFU)**