

Agenda:

- Limitations of Data Parallel and ZeRO

- model parallel

- Tensor model parallel

- split at layer boundaries

- Pipeline model Parallel

* ISSUE WITH DP: MEMORY

- each GPU saves a replica of the entire model

- can't train large models that exceed

GPU memory! (including BOTH model weights + saved activations)

size of hidden layers

	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

doesn't fit on A100 & V100

- NVIDIA V100: 16GB or 32GB

- NVIDIA A100: 40GB or 80GB

* ZeRO (on which FSDP is based): zero-redundancy

Optimizers

- where is redundancy coming from in model training?

- review: SGD / Adam

SGD:

for $t = 1 \dots T$:

$$\Delta w = \eta \times \frac{1}{B} \sum_{i=1}^B \nabla (\text{loss}(f_w(x_i, y_i)))$$

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

learning rate (points to η)
backward pass (points to ∇)
forward pass (points to f_w)

ADAM:

for $t = 1 \dots T$:

$$grad = \frac{1}{B} \sum_{i=1}^B \nabla (\text{loss}(f_w(x_i, y_i)))$$
$$\Delta w = \text{ADAM}(grad) \rightarrow$$

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

Adam might involve:

$$v_t = \beta_1 * v_{t-1} - (1 - \beta_1) * g_t$$

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

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

g_t : Gradient at time t along ω^j

v_t : Exponential Average of gradients along ω_j

s_t : Exponential Average of squares of gradients along ω_j

β_1, β_2 : Hyperparameters

optimizer

needs to store:

- params themselves
(to make final update)

- momentum (v_t)

- variance (s_t)

- * where does REDUNDANCY come from?

• in DP - each worker stores a copy of:

→ optimizer states (even though...
model updates end up being same)

→ gradients (which are the same after
all-reduce)

→ params

↳ could be 2 for fp16,
4 for fp32

→ params = $P \cdot \text{sizeof}(\text{precision params})$

→ gradients = $P \cdot \text{sizeof}(\text{precision gradient})$

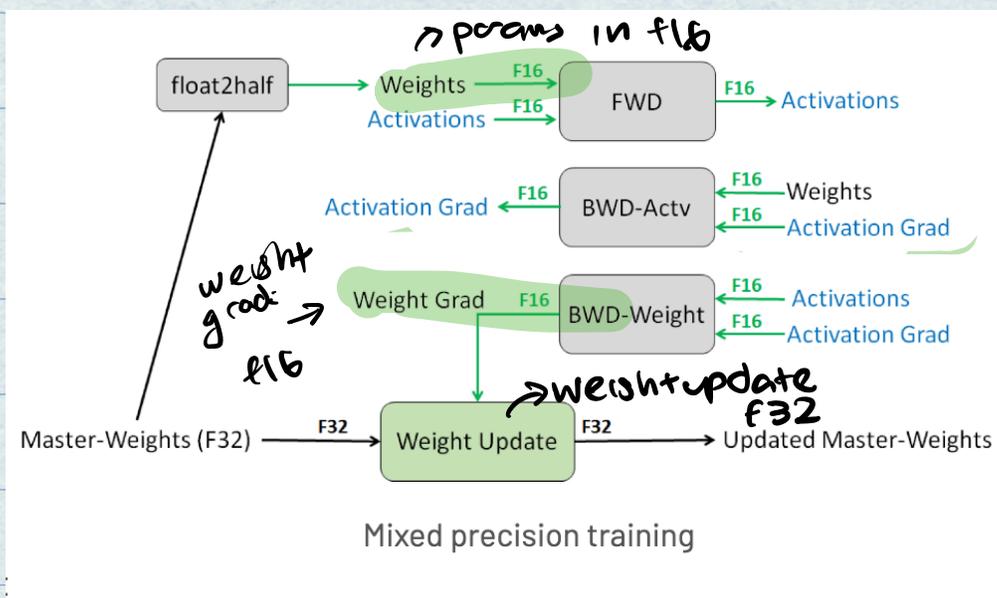
→ OS state = $K \cdot P \cdot \text{sizeof}(\text{precision optimizer})$

↳ # of variables optimizer class

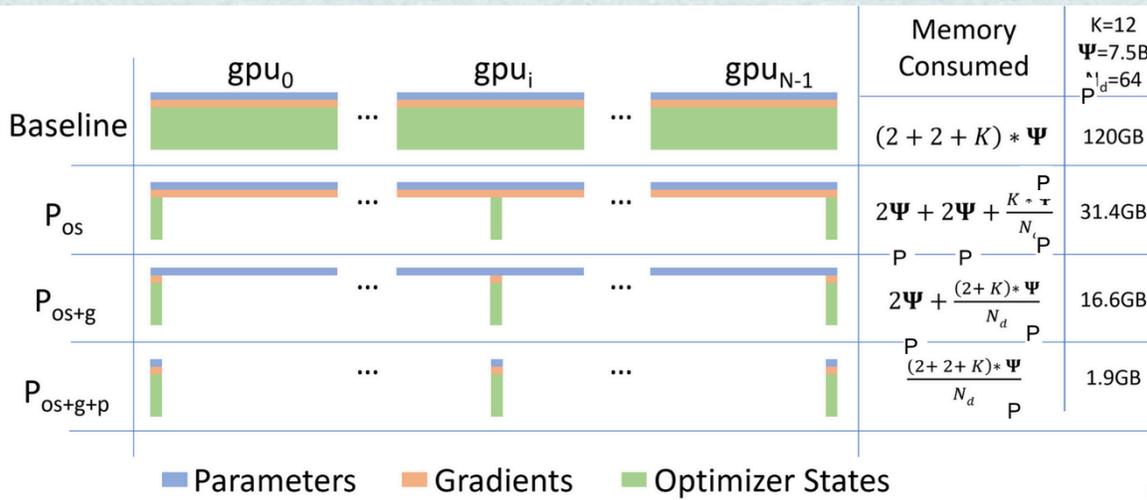
holds PER param (1 for SGD,

3 for Adam)

* Example of MIXED-PRECISION TRAINING



* ZeRO: shard all the things!



overhead per-GPU
for MP training of
model w/
7.5 Billion params
64 workers,
Adam optimizer

* ZERO STAGE 1: OPTIMIZER-STATE PARTITIONING

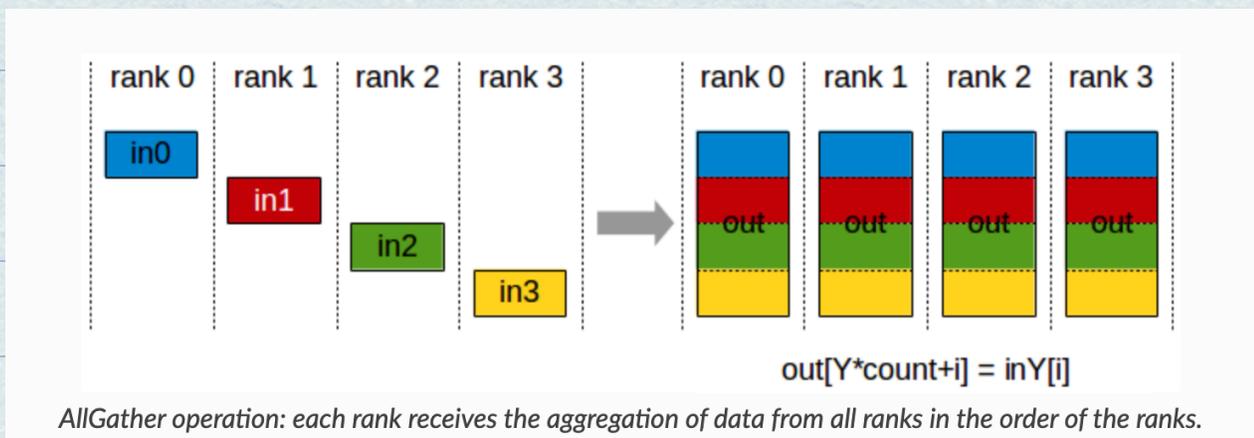
- Group optimizer into N_d equal partitions
- Each worker only updates state for its partition
- How does each worker actually

get the updated param from optimizer?

- use All-GATHER operation:

- gathers $\left(\begin{matrix} \text{optstate} \\ \xrightarrow{N_d} \end{matrix} \right)$ data from each worker

- distributes result to all workers



when local process is DONE w/ all-gather,

can update local params and free

everything except its partition

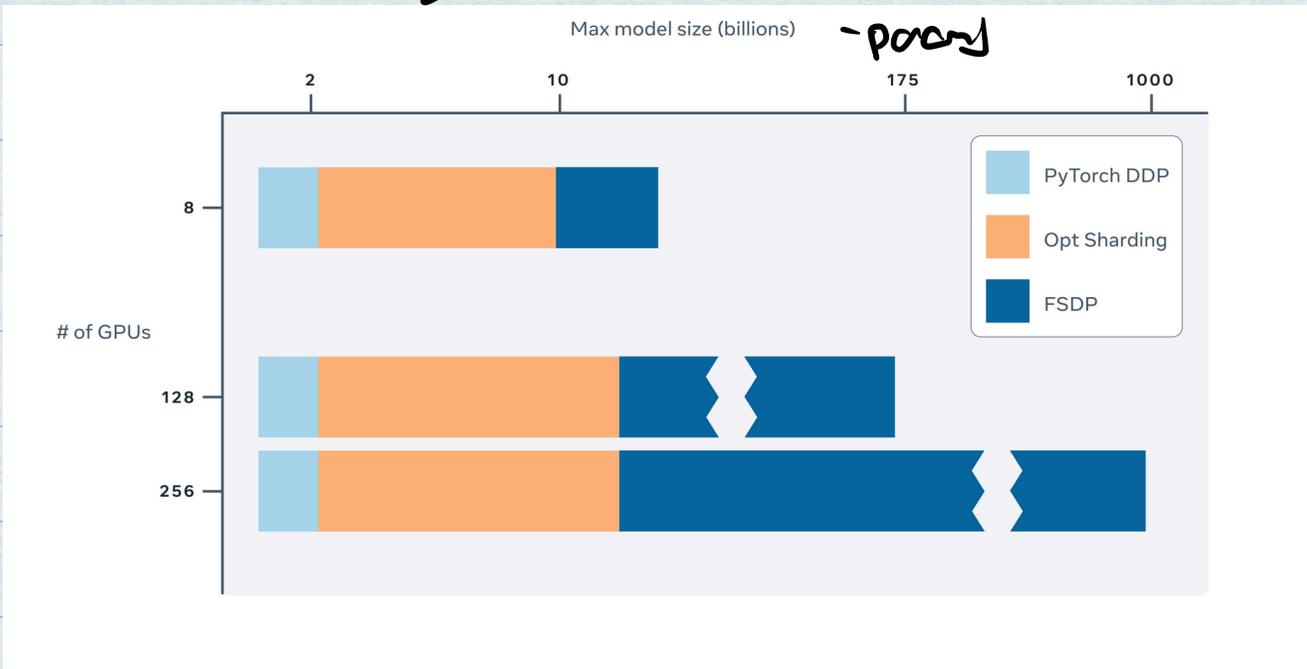
* ZERO stage 2: Partition Gradients Also

- Because local optimizers are only responsible for updating a partition of params, don't need to keep all gradients on all workers - need to only store gradients for its partition

* ZERO stage 3: Add partitioning of PARAMS

- instead of storing all params on all layers - fetch params as needed for FW and BW pass
- increases comm. to 1.5x compared to DP (see Zero-paper), but reduces memory proportionally by # of workers

Meta's experiences w/ FSDP (pytorch implementation of zero)



* Intuition for gradient sharding

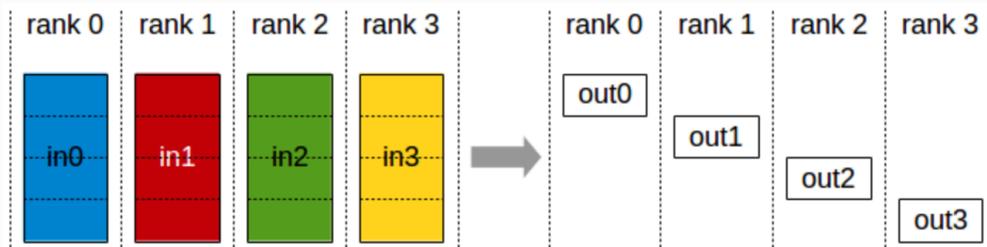
- q: how do we do grad update if grad is distributed?

- combine REDUCE-SCATTER w/ All-gather:

ReduceScatter

The ReduceScatter operation performs the same operation as Reduce, except that the result is scattered in equal-sized blocks between ranks, each rank getting a chunk of data based on its rank index.

The ReduceScatter operation is impacted by a different rank to device mapping since the ranks determine the data layout.

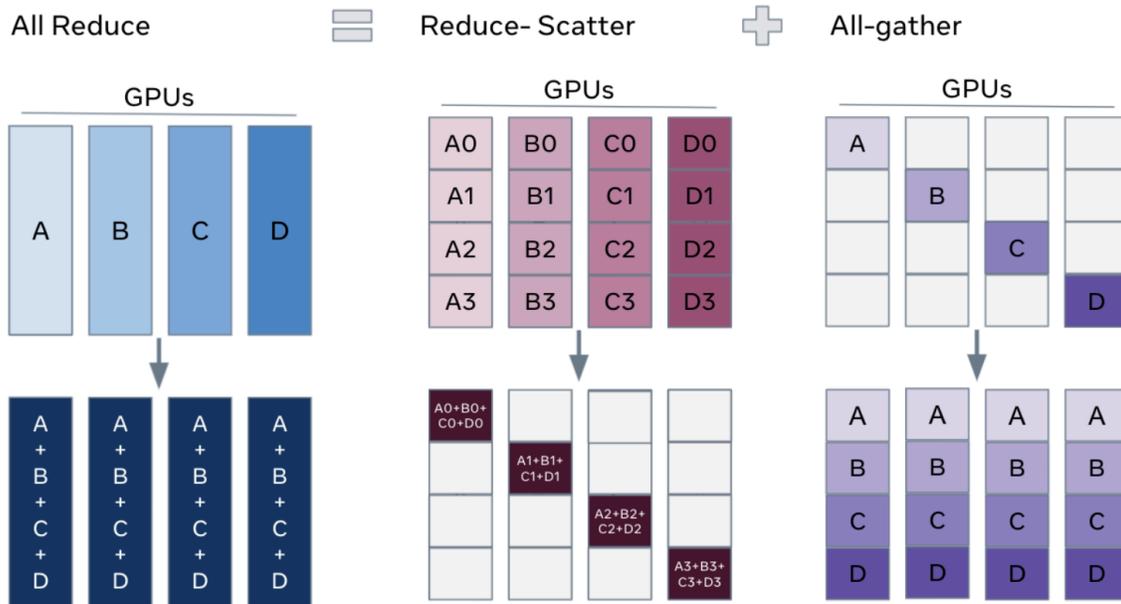


$$\text{outY}[i] = \text{sum}(\text{inX}[\text{Y*count}+i])$$

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

- shard gradient updates per param into 4 (0,1,2,3,4)

+ T000: instead sharding a bit better



All-reduce as a combination of reduce-scatter and all-gather. The standard all-reduce operation to aggregate gradients can be decomposed into two separate phases: reduce-scatter and all-gather. During the reduce-scatter phase, the gradients are summed in equal blocks among ranks on each GPU based on their rank index. During the all-gather phase, the sharded portion of aggregated gradients available on each GPU are made available to all GPUs (see here for details on those operators).

Part 2: Model Parallelism

- in DP we replicated model, but sharded data

- in MP: replicate data but partition model

- Two general approaches:

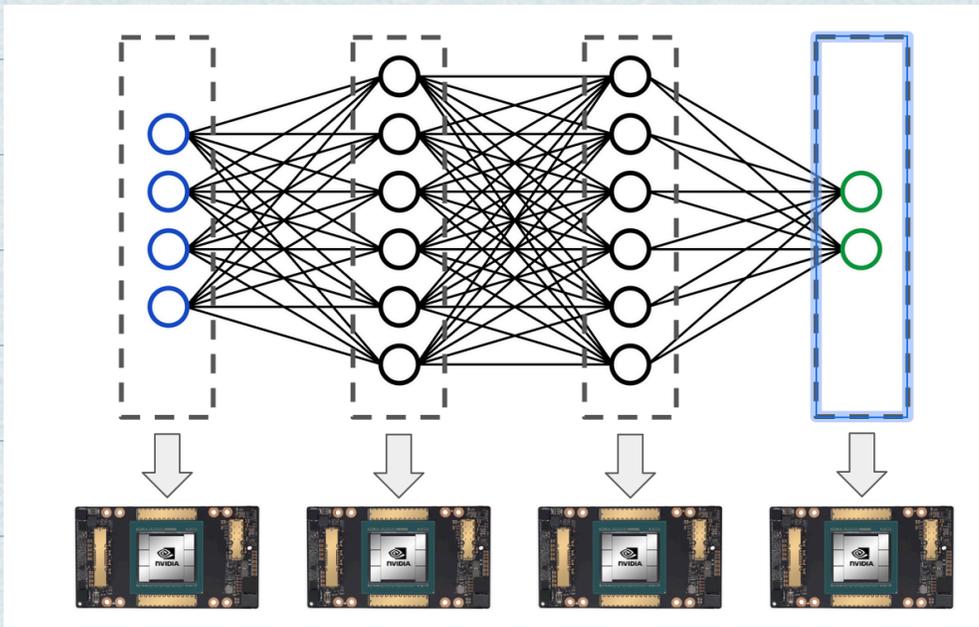
1) slice model "vertically": place dif.

subsets of dif. layers on each rank

2) slice model horizontally: shard model

weights across dif. GPUs

- * vertical slicing MP:

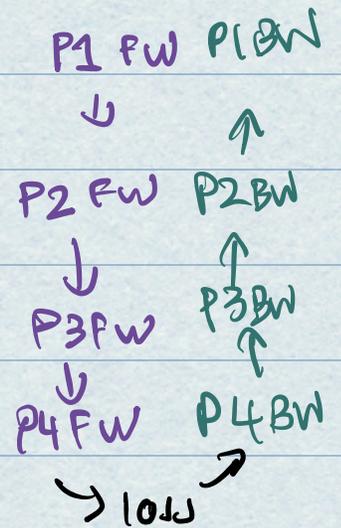
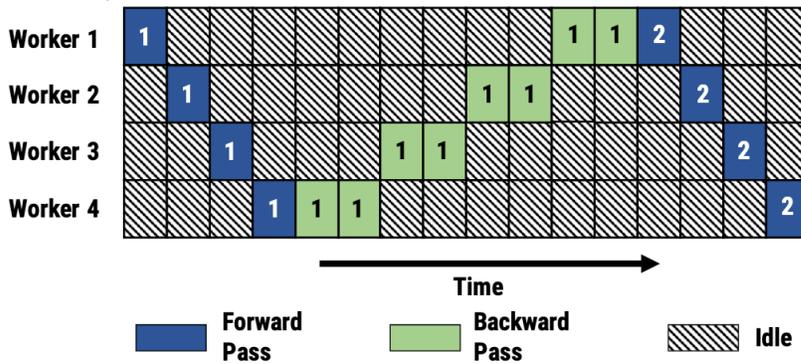


1

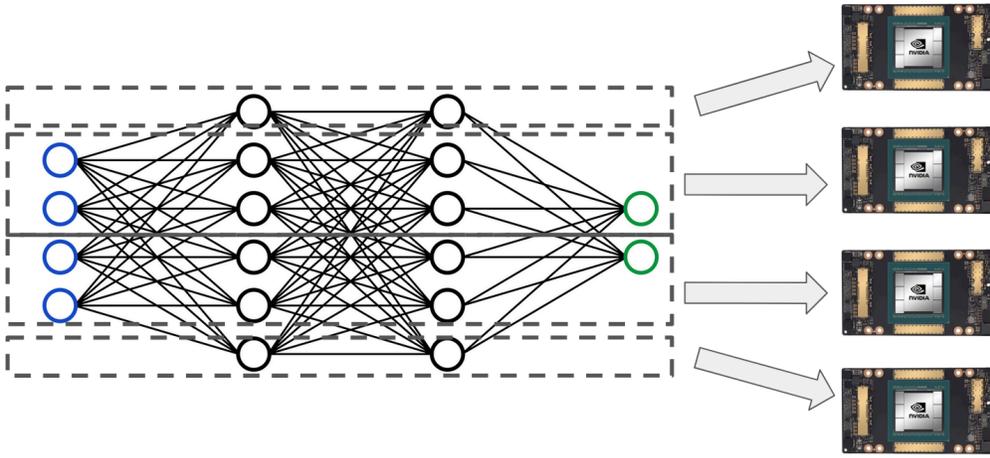
leads to:

- under-utilization of compute resources
(only 1 worker being used at a time)
- low overall throughput due to low utilization

(from Pipedream paper)



• Tensor-model parallelism: slice layers horizontally



- how does this work?

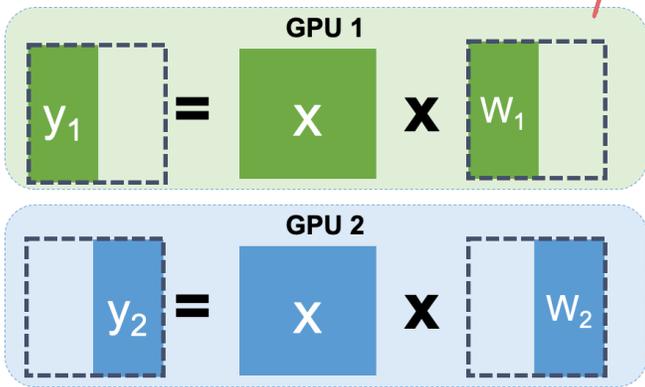
Tensor Model Parallelism

$$B \begin{Bmatrix} C_{out} \\ m \end{Bmatrix} y = B \begin{Bmatrix} C_{in} \\ m \end{Bmatrix} x \times W \begin{Bmatrix} C_{out} \\ m \end{Bmatrix}$$

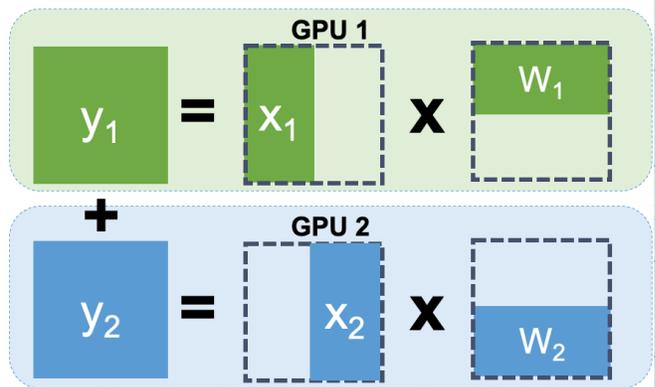
output
input
parameters

- Partition parameters/gradients *within* a layer

2 parts: partition + reduce



Tensor Model Parallelism (partition output)



Tensor Model Parallelism (reduce output)
y = y1 + y2

* what is the comm. overhead of tensor model parallelism?

- lets start by looking at data:

DATA PARALLEL

Tensor model parallelism

FW processing

partition: $O(B \times C_{in})$

reduce: $O(B \times C_{out})$

BW processing

partition: $O(B \times C_{out})$

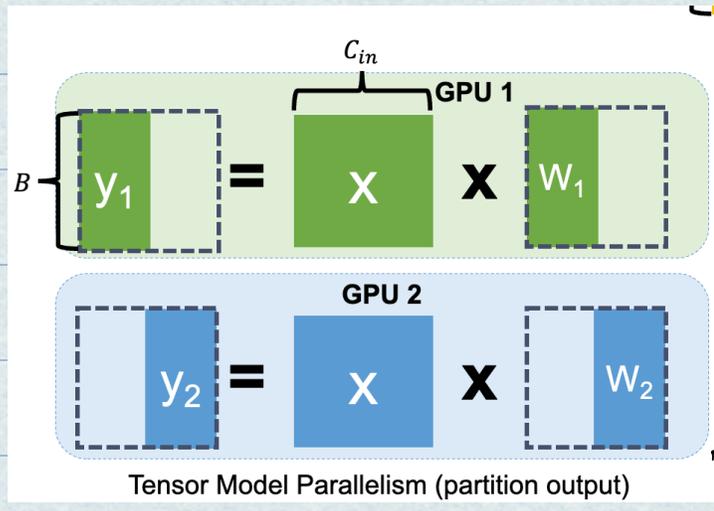
reduce: $O(B \times C_{out})$

Gradient sync

$O(C_{out} \times C_{in})$

→ TMP require sync after each fw/bw pass to collect output. But no syncs of gradients

SHARD W BY COLUMN (duplicate X):



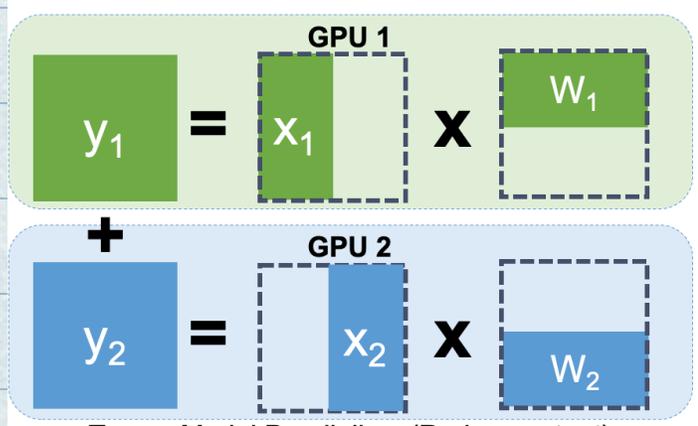
option 1: replicate X and split W by column

Forward Processing	Backward Propagation	Gradients Sync
$O(B * C_{in})$	$O(B * C_{in})$	0

Communication Cost of Tensor Model Parallelism

$y = (y_1, y_2)$

SHARD W BY ROW (split X by column to work):



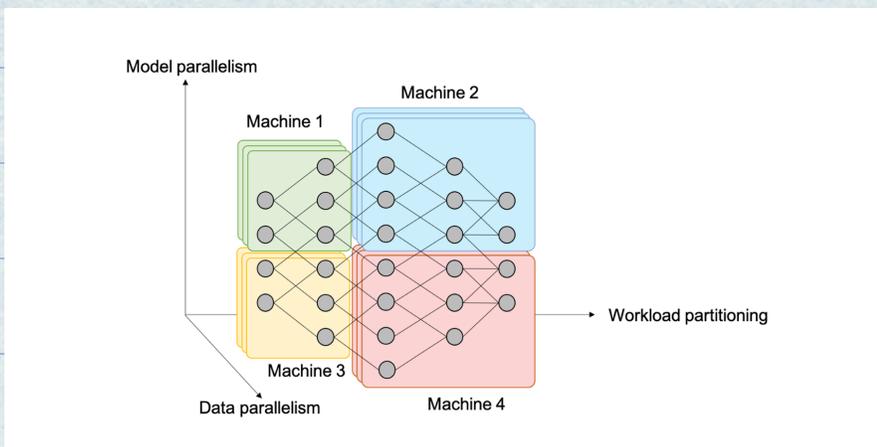
option 2: split X by col and W by row

Forward Processing	Backward Propagation	Gradients Sync
$O(B * C_{out})$	$O(B * C_{out})$	0

Communication Cost of Tensor Model Parallelism

$y = y_1 + y_2$ → if there's a non-linear op after, need to save outputs

- could even combine model and data parallel:



microbatch size

- caveat: larger minibatch sizes leads to accuracy loss

- smaller microbatch sizes reduces utilization

• Play w/ p (number of stages):

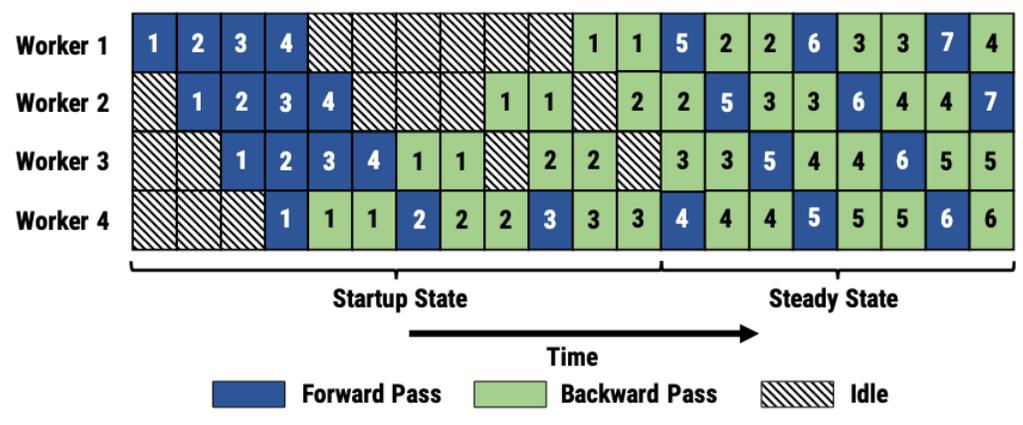
- Decrease pipeline size

- caveat: memory limits

• Another issue: ACTIVE memory requirement

- for EACH microbatch - need to keep output of intermediate layers - for later BW pass

• Idea: solve this w/ "IFIB" schedule:



- this way, each worker can delete

State for each microbatch after

BW paid

... see paper for more variations!

SUMMARY: when are dif. parallelism approaches effective?

- Data: if model weights and activation fit into GPU memory

- Tensor MP: model weights / activation DO NOT fit in GPU mem, but we have fast networking (in a single box)

- pipeline MP: effective if model weights/activation DON'T fit into GPU mem, but we have slow networking

- either multiple machines

- OR 1 box w/ slow interconnect.

* next time: intro to transformers!

- we will revisit tensor model parallelism

References:

Model size table: CMU 15-442 Lecture on ML Parallelization Part 1 led by Zhihao Jia and Tianqi Chen
Adam optimizer method screenshot: Kingma and Ba, "Adam: A Method for Stochastic Optimization", 2014, <https://arxiv.org/abs/1412.6980>.
Mixed precision training: Mixed Precision Training. Narang, et al. <https://arxiv.org/pdf/1710.03740.pdf>
General Information about ZeRO: ZeRO: Memory Optimizations Toward Training Trillion Parameter Models. Rajbhandari et al. <https://arxiv.org/pdf/1910.02054.pdf>
ZeRO diagram: From ZeRO paper (above), taken from Azalia's slides
All-gather image: NCCL documentation
Meta FSDP Article: <https://engineering.fb.com/2021/07/15/open-source/fsdp/>
Reduce-scatter image: NCCL documentation
All-reduce = reduce-scatter + All-gather: Meta Article
Model-Parallelism Diagram: Azalia's slides
Tensor-Parallelism Diagram: Azalia's slides
Breakdown of tensor-model parallelism of partition and reduce: CMU 15-442 Lecture on ML Parallelism Part 2
Combined Data and Model Parallel: CMU 15-442 Lecture on ML Parallelism Part 2
Vertical Model Parallel Utilization Diagram: Pipedream paper
Gpipe Utilization diagram: Pipedream paper
Pipedream 1F1B diagram: Pipedream paper
Megatron LM: <https://arxiv.org/pdf/2104.04473>

General Lecture Flow:

CMU 15-442 Lecture on ML Parallelization Part 1

CS 229s 2023, Lecture on Parallelism Fundamentals, Given By Azalia Mirhoseini