



# Workshop on Systems for LLMs

Systems for Large Language Models (LLMs) Workshop with Hands-On Exploring the Infrastructure Behind Intelligent Systems

#### **Sessions Overview:**

LLM building blocks
GPUs for LLM
Inference serving with vLLM
Distributed Training

Tools: Hugging Face, PyTorch



#### **Details**

Date: 6<sup>th</sup>,7<sup>th</sup> September

Time: 9:30 am-5pm

Register @

https://forms.gle/a8rQQ835jJG42i3eA

Supported by: IBM-IITB Academic Research Partnership





# Distributed Training

## Recap: The Train Loop

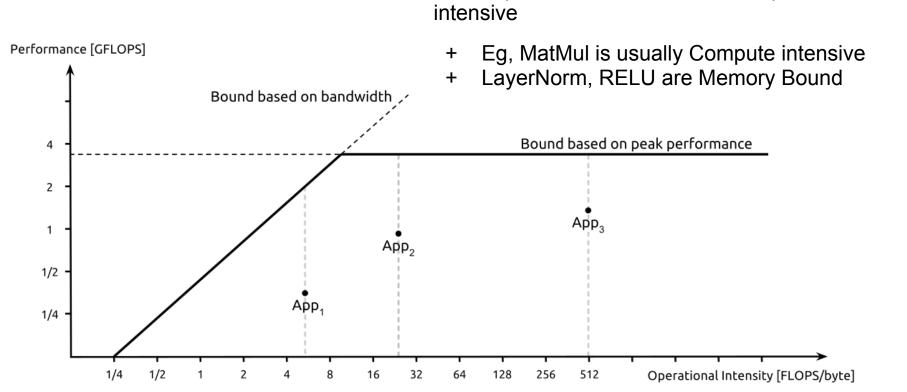
```
For e in epochs:
For b in batches:
loss = model.forward(b)
loss.backward()

optimizer.step()
```

#### **GPU Resources**

- + Compute
  - + The number of SMs in the GPU, eg. 128 SMs on an A100
  - + Capacity of the GPU, totally represented as TFLOPs
- + Memory
  - + HBM capacity: eg. 80 GB on an A100
  - + To a lesser extent: Caches and per SM local memory capacity

## Compute Bound vs Memory Bound



Each Op can be classified into Compute vs Memory

## Memory requirements for training

Simple envelope calculation: assuming 8B model in fp32 datatype with AdamW optimizer

Model weights: 8B \* 4 bytes per param = ~30GB

Gradients: 4 bytes per param =  $\sim$ 30GB

Optimizer state: 12 bytes per param = ~90GB

+ Activations = f(batch size, seq len)

### Flow of memory

Model and Optimizer are flat.

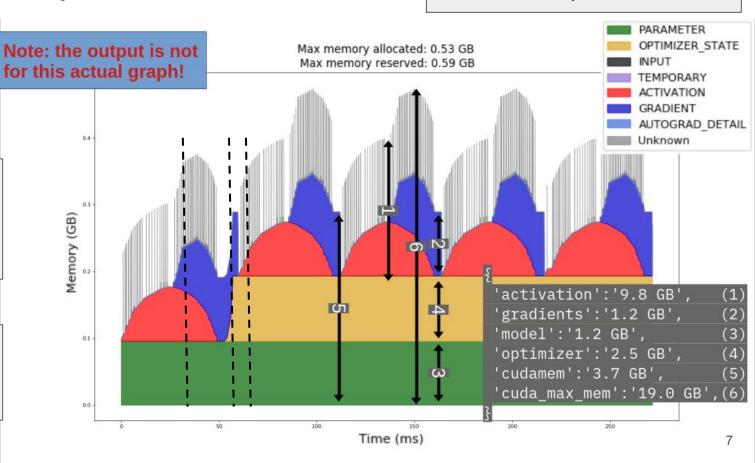
Forward:

Activations increase

Backward:

Activations are freed, gradients are created

Optimizer Step: Gradients are freed.



# Memory requirements for training various model sizes

Model	Weights	Gradients	Optimizer	Sum (w/o activations)
Granite 2b	7.5 GB	7.5 GB	22 GB	37 GB
Llama 8b	30 GB	30 GB	90 GB	150 GB
Qwen 32B	120 GB	120 GB	360 GB	600 GB

# Memory Availability of commonly used GPUs

Vendor	Model	Memory
Nvidia	A100	80 GB
Nvidia	V100	32GB
Nvidia	A6000	23 GB
Nvidia	L40s	48 GB
AMD	Radeon RX 7900 XTX	24GB

### Extras: Impact of Optimizer choice

Instead of using Adam, earlier models used to use SGD optimizer,

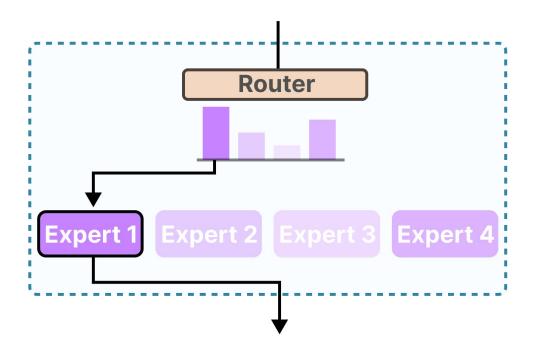
Which has: 8 bytes per param

But Adam/AdamW is better from an algorithmic point of view

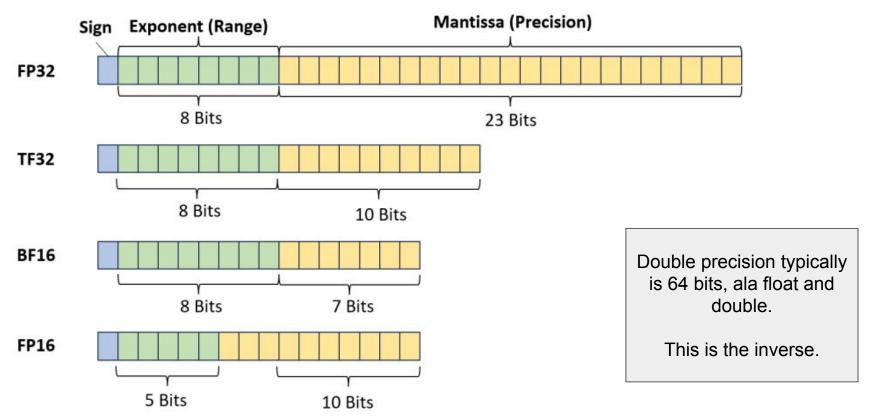
## Extras: Impact of MoE

MoE architecture activates a fewer set of params during inferences, but the full model has to be loaded anyway

In case of training, it makes no difference



## Other precision Types: half precision



#### **TF32**

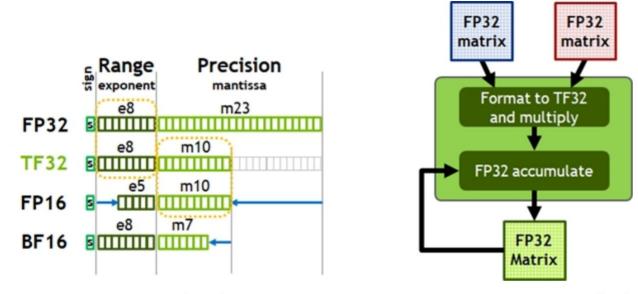


Figure 7. TensorFloat-32 (TF32) provides the range of FP32 with the precision of FP16 (left). A100 accelerates tensor math with TF32 while supporting FP32 input and output data (right), enabling easy integration into DL and HPC programs and automatic acceleration of DL frameworks.

# Hardware support for other precisions

A100 80GB PCIe	A100 80GB SXM	
9.7 TFLOPS		
19.5 TFLOPS		
19.5 TFLOPS		
156 TFLOPS   312 TFLOPS*		
312 TFLOPS   624 TFLOPS*		
312 TFLOPS   624 TFLOPS*		
624 TOPS   1248 TOPS*		
80GB HBM2e	80GB HBM2e	
1,935 GB/s	2,039 GB/s	
	9.7 TF  19.5 TI  19.5 TI  19.5 TI  156 TFLOPS    312 TFLOPS    312 TFLOPS    624 TOPS	

# H100 specs

Technical Specifications				
	H100 SXM	H100 NVL		
FP64	34 teraFLOPS	30 teraFLOPS		
FP64 Tensor Core	67 teraFLOPS	60 teraFLOPS		
FP32	67 teraFLOPS	60 teraFLOPS		
TF32 Tensor Core*	989 teraFLOPS	835 teraFLOPS		
BFLOAT16 Tensor Core*	1,979 teraFLOPS	1,671 teraFLOPS		
FP16 Tensor Core*	1,979 teraFLOPS	1,671 teraFLOPS		
FP8 Tensor Core*	3,958 teraFLOPS	3,341 teraFLOPS		
INT8 Tensor Core*	3,958 TOPS	3,341 TOPS		
GPU Memory	80GB	94GB		
<b>GPU Memory Bandwidth</b>	3.35TB/s	3.9TB/s		

#### 1 SM in the A100





#### **Automatic Mixed Precision**

Model weights: half - 16 bits

Gradients: half - 16 bits

Optimizer: full - 32 bits (4 bytes) \* 12

Actually faster to run and takes lesser memory

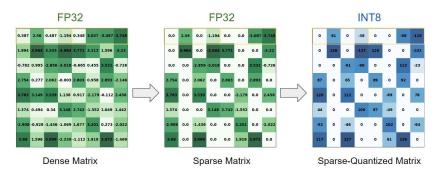
#### What more from here?

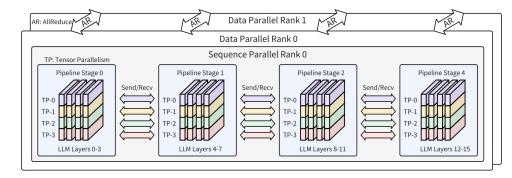
#### Fit model within this GPU

- Quantization
- + Sparsity

Throw more resources at this problem

**Distributed Training!** 





## Distributed Training Algorithms vs Systems

- + Data Parallelism (DDP)
- + Tensor Parallelism (TP)
- + Pipeline Parallelism (PP)

- Full Sharded Data Parallelism (FSDP)
  - + With variants like Full/Hybrid Shard

- + Pytorch in-built DDP
- + Hugging Face Accelerate
- + Nvidia Megatron
- + Pytorch in-built FSDP
- Microsoft DeepSpeed ZeRO family

- + 3D parallelism
- + Context/Sequence Parallelism
- + Expert Parallelism (EP)

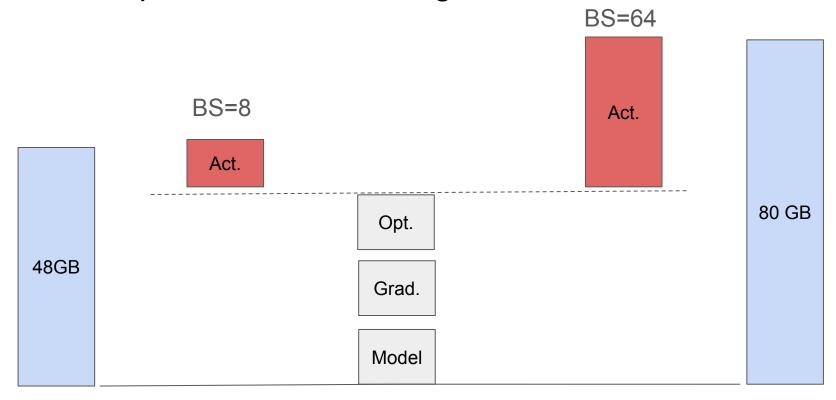
# **Basic Tech**

## Does your model fit into on a single GPU?

- Yes we go with DDP.
- 2. No consider other options.

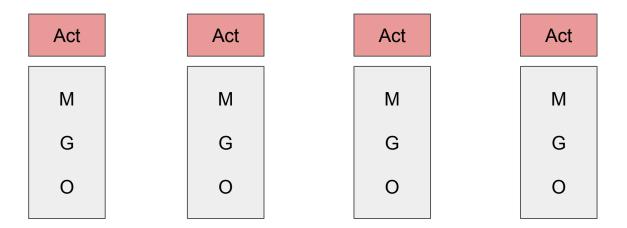
Here: model => (model + optimizer + gradients), the remaining goes to activations

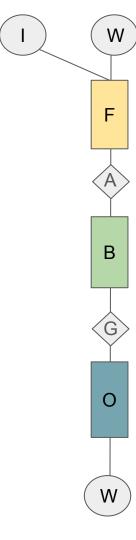
# Limited Scope of Vertical Scaling



#### DDP: Data Parallelism

Can we have multiple gpus loading the same model and run training on different subsets of the data?

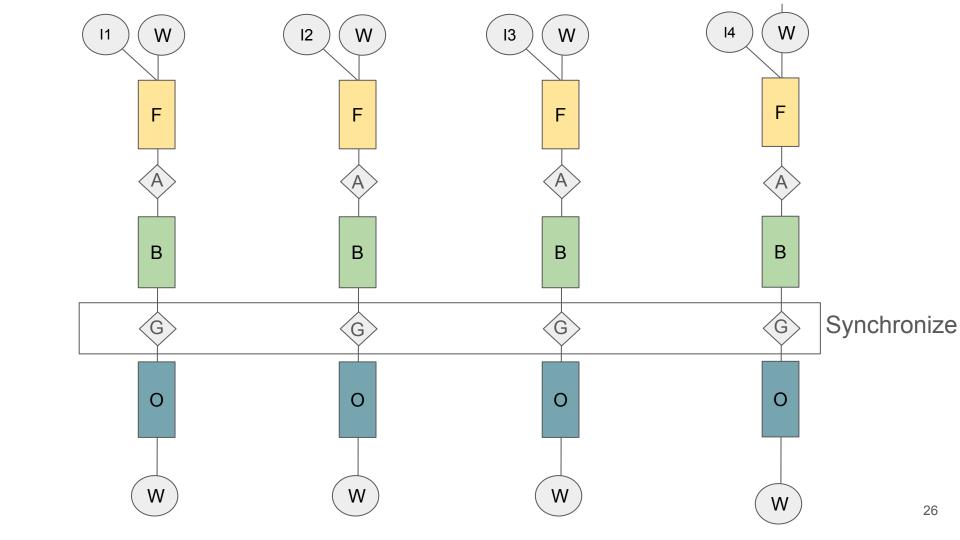




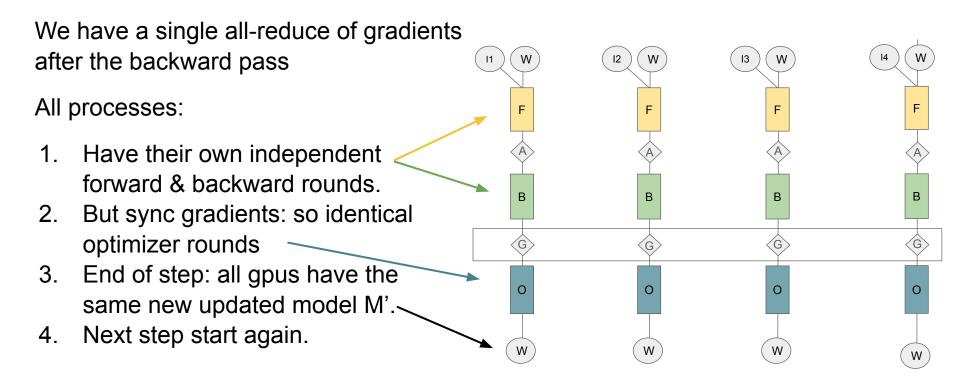
## How to do simple Horizontal Scaling?

- 1. Divide your input data
- 2. Load identical training setup on each GPU same model, optimizer setup
- 3. Run forward pass on all gpus:
  - a. We will have intermediate activations created on each gpu
- Run backward pass on all gpus:
  - a. Now each process has its own gradients
- 5. Run optimizer on all gpus

This is a problem!



#### **DDP** Details



## What does "Synchronize" mean?

It means, we just add up the gradients.

Why does this work?

Gradients are always accumulate-able

Same mechanism is used during Gradient Accumulation.

### How to synchronize gradients?

Each process has its own gradients: Num params \* 2/4.

For a 8B model, this can mean 15-25 GB of data.

Gradients are automatically created by Pytorch during the backward pass - all of these "tensors" are in GPU memory.

The all-reduce NCCL primitive









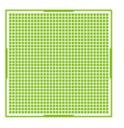
**Download NCCL** 

Documentation

**Developer Guide** 

GitHub

**Watch GTC Webinar** 





1 GPU

multi-GPU, multi-node

#### **Performance**

NCCL conveniently removes the need for developers to optimize their applications for specific machines. NCCL provides fast collectives over multiple GPUs both within and across nodes.

#### **Ease of Programming**

NCCL uses a simple C API, which can be easily accessed from a variety of programming languages.NCCL closely follows the popular collectives API defined by MPI (Message Passing Interface).

#### Compatibility

NCCL is compatible with virtually any multi-GPU parallelization model, such as: single-threaded, multi-threaded (using one thread per GPU) and multi-process (MPI combined with multi-threaded operation on GPUs).

#### CCL: honorable mentions

NCCL is not the only CCL.

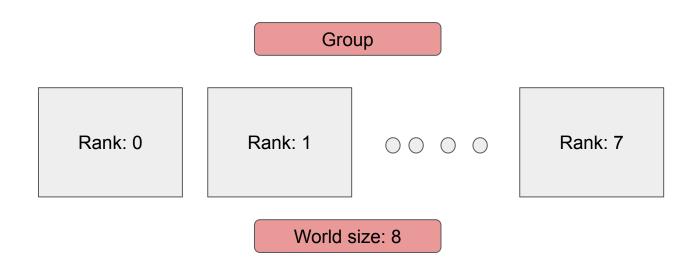
Communication between CPU cores (if using a CPU only training): Gloo

Other specialized hardware: Google TPU, or Intel XPU etc.

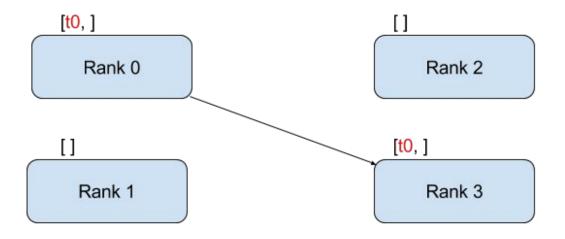
Some research work on optimizing CCLs: eg, UCCL from Berkeley

## Collective setup - the higher level api

- Each process can use the CCL functions as needed, using ranks to refer to others.
- API exposed by NCCL to Pytorch to higher layers.
  - a. At pytorch called `torch.distributed`. Abstracts the underlying CCL lib used.

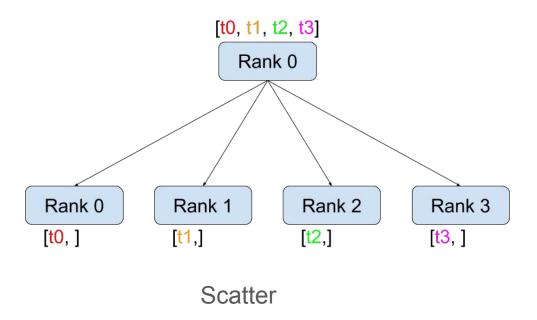


#### P2P communication

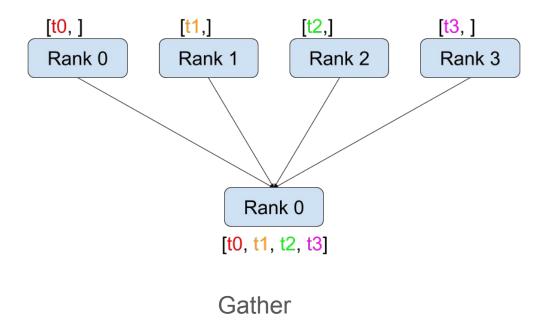


Point to Point Communication: Send and recv

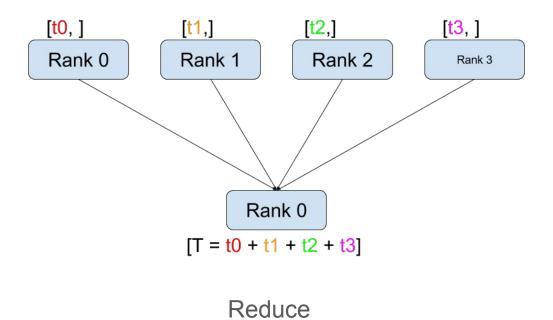
#### **Collective Communication**



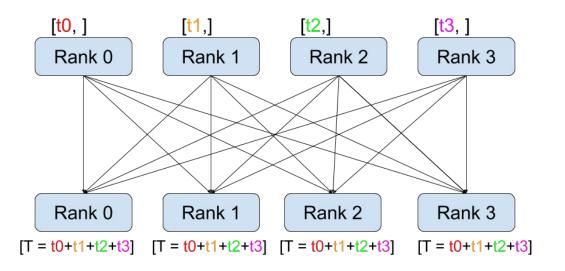
#### **Collective Communication**



#### **Collective Communication**

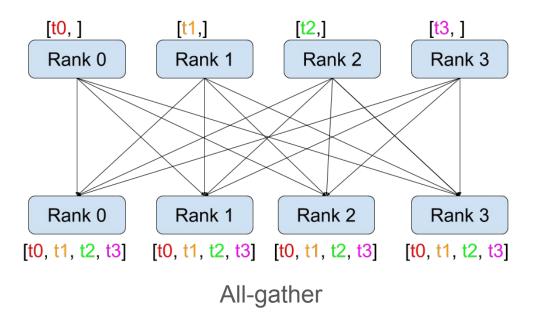


### **Collective Communication**



All-Reduce

### **Collective Communication**



#### How does NCCL work?

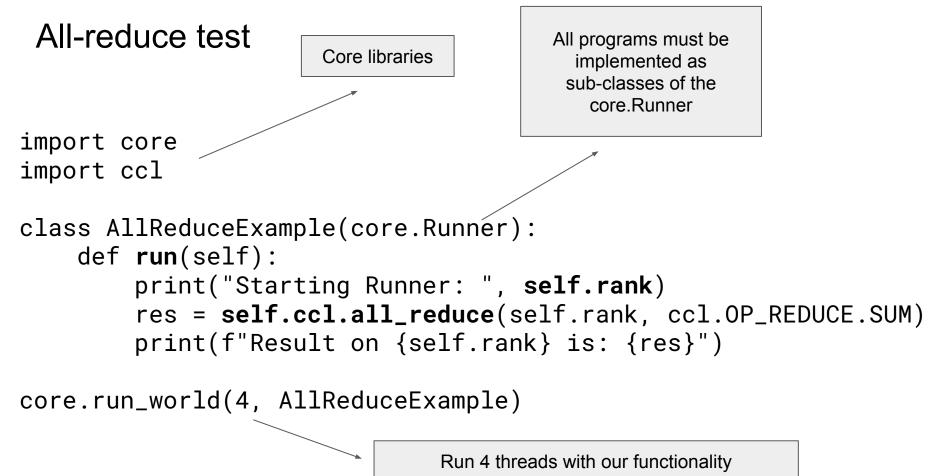
- Network paths:
  - a. For GPUs on the same node: we have the NVLink.
  - For GPUs across the node: we have NVSwitch and other higher level switches.
- 2. Topology management: ring vs tree etc
- 3. Chunking of tensors -> chunks -> packets for the lower level transport

# Hands-on (CCL primitives: in the toy framework)

## How to do a simple hands-on?

We will use a simple Python emulator where:

- 1. Each "gpu" process is run as a thread.
- CCL is completely simplified into a single global object with thread IPC



## Test the setup

Add random sleep to each process.

Write a program where all even numbered ranks send a random number to the next process, which will print it out.

# The Simplified Model (single.py)

```
import core
m = core.Model(4, 2, 6)
# create some samples for input
i, o = core.gen_inp(10, 2)
                                                                   x4
print("Output from simple single gpu run: ")
out = m.forward(i)
loss = out - o
                                                           10 samples of size 2 each
m.backward(loss)
                                  Dummy outputs to
print("Output: ", out)
                                    manage loss
print("Loss: ", loss)
                                     gradients
m.print_grads()
                                                               No optimizer step
```

## Look at the Model code in core.py

```
self.n = num_layers
self.in_dim = in_dim
self.hidden_dim = hidden_dim

self.blocks = []
self.inputs = []
self.grads = []
```

```
def forward(self, x):
    # Input is of the form: seqlen * in_dim

for i in range(self.n):
    self.inputs[i]["a"] = x # save input
    x = x @ self.blocks[i]["a"]
    self.inputs[i]["b"] = x # save input
    x = x @ self.blocks[i]["b"]

return x
```

# The backward pass for simple matmul

$$X = \begin{pmatrix} x_{1,1} & x_{1,2} \\ x_{2,1} & x_{2,2} \end{pmatrix} \qquad W = \begin{pmatrix} w_{1,1} & w_{1,2} & w_{1,3} \\ w_{2,1} & w_{2,2} & w_{2,3} \end{pmatrix} \tag{1}$$

$$Y = XW (2)$$

$$= \begin{pmatrix} x_{1,1}w_{1,1} + x_{1,2}w_{2,1} & x_{1,1}w_{1,2} + x_{1,2}w_{2,2} & x_{1,1}w_{1,3} + x_{1,2}w_{2,3} \\ x_{2,1}w_{1,1} + x_{2,2}w_{2,1} & x_{2,1}w_{1,2} + x_{2,2}w_{2,2} & x_{2,1}w_{1,3} + x_{2,2}w_{2,3} \end{pmatrix}$$
(3)

https://cs231n.stanford.edu/handouts/linear-backprop.pdf

Since L is a scalar and Y is a matrix of shape  $N \times M$ , the gradient  $\frac{\partial L}{\partial Y}$  will be a matrix with the same shape as Y, where each element of  $\frac{\partial L}{\partial Y}$  gives the derivative of the loss L with respect to one element of Y:

$$\frac{\partial L}{\partial Y} = \begin{pmatrix} \frac{\partial L}{\partial y_{1,1}} & \frac{\partial L}{\partial y_{1,2}} & \frac{\partial L}{\partial y_{1,3}} \\ \frac{\partial L}{\partial y_{2,1}} & \frac{\partial L}{\partial y_{2,2}} & \frac{\partial L}{\partial y_{2,3}} \end{pmatrix}$$
(4)

By the chain rule, we know that:

$$\frac{\partial L}{\partial X} = \frac{\partial L}{\partial Y} \frac{\partial Y}{\partial X} \qquad \qquad \frac{\partial L}{\partial W} = \frac{\partial L}{\partial Y} \frac{\partial Y}{\partial W} \tag{5}$$

The terms  $\frac{\partial Y}{\partial X}$  and  $\frac{\partial Y}{\partial W}$  in Equation 5 are *Jacobian matrices* containing the partial derivative of each element of Y with respect to each element of the inputs X and W.

$$\frac{\partial L}{\partial X} = \begin{pmatrix} \frac{\partial L}{\partial y_{1,1}} w_{1,1} + \frac{\partial L}{\partial y_{1,2}} w_{1,2} + \frac{\partial L}{\partial y_{1,3}} w_{1,3} & \frac{\partial L}{\partial y_{1,1}} w_{2,1} + \frac{\partial L}{\partial y_{1,2}} w_{2,2} + \frac{\partial L}{\partial y_{1,3}} w_{2,3} \\ \frac{\partial L}{\partial y_{2,1}} w_{1,1} + \frac{\partial L}{\partial y_{2,2}} w_{1,2} + \frac{\partial L}{\partial y_{2,3}} w_{1,3} & \frac{\partial L}{\partial y_{2,1}} w_{2,1} + \frac{\partial L}{\partial y_{2,2}} w_{2,2} + \frac{\partial L}{\partial y_{2,3}} w_{2,3} \end{pmatrix}$$
(22)

$$= \begin{pmatrix} \frac{\partial L}{\partial y_{1,1}} & \frac{\partial L}{\partial y_{1,2}} & \frac{\partial L}{\partial y_{1,3}} \\ \frac{\partial L}{\partial y_{2,1}} & \frac{\partial L}{\partial y_{2,2}} & \frac{\partial L}{\partial y_{2,3}} \end{pmatrix} \begin{pmatrix} w_{1,1} & w_{2,1} \\ w_{1,2} & w_{2,2} \\ w_{1,3} & w_{2,3} \end{pmatrix}$$
(23)

$$= \boxed{\frac{\partial L}{\partial Y} W^T} \tag{24}$$

Using the same strategy of thinking about components one at a time, you can derive a similarly simple equation to compute  $\frac{\partial L}{\partial W}$  without explicitly forming the Jacobian  $\frac{\partial Y}{\partial W}$ :

$$\frac{\partial L}{\partial W} = X^T \frac{\partial L}{\partial Y} \tag{25}$$

```
def backward(self, loss_grad):
    # single var used to track grad propogation across layers
    # starts out with the input loss_grad
    act grad = loss grad
    for i in range(self.n-1, 0-1, -1):
        self.grads[i]["b"] = self.inputs[i]["b"].T @ act_grad
        act_grad = act_grad @ self.blocks[i]["b"].T
        self.grads[i]["a"] = self.inputs[i]["a"].T @ act_grad
        act_grad = act_grad @ self.blocks[i]["a"].T
```



# Implementing DDP

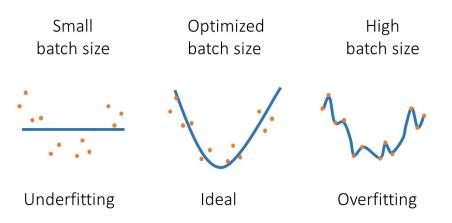
```
i, o = \dots
class DDPRunner(core.Runner):
    def run(self):
        # make a local copy of the model
        m = copy.deepcopy(m)
        _i = ... # get input shard of relevance to us
        ... # implement forward pass here
        loss = out - o
        ... # implement backward pass here
        ... # any final synchronization here
        if self.rank == 0:
            global m_final
            m final = m
core.run world(2, DDPRunner)
# single model processing here
# ...
m.compare(m_final)
# This test should return true
```

### Problems with DDP

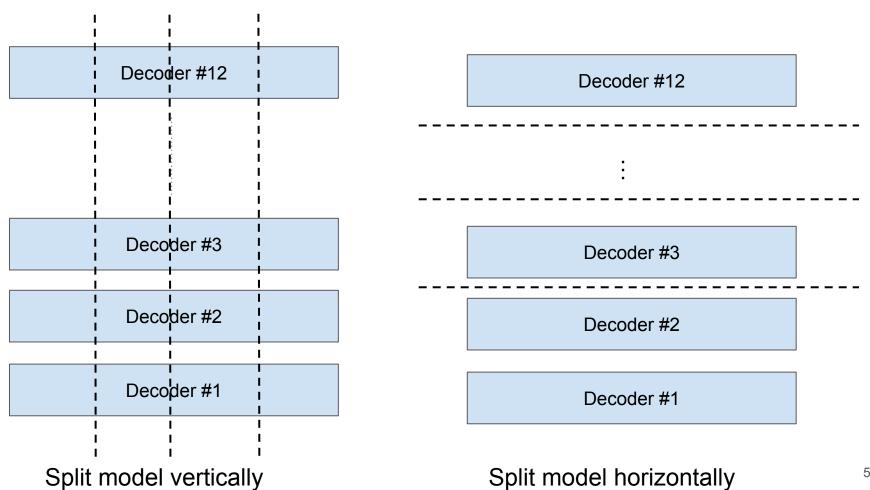
Can only scale by batch size.

Having too large of a batch size is bad for the DL learning

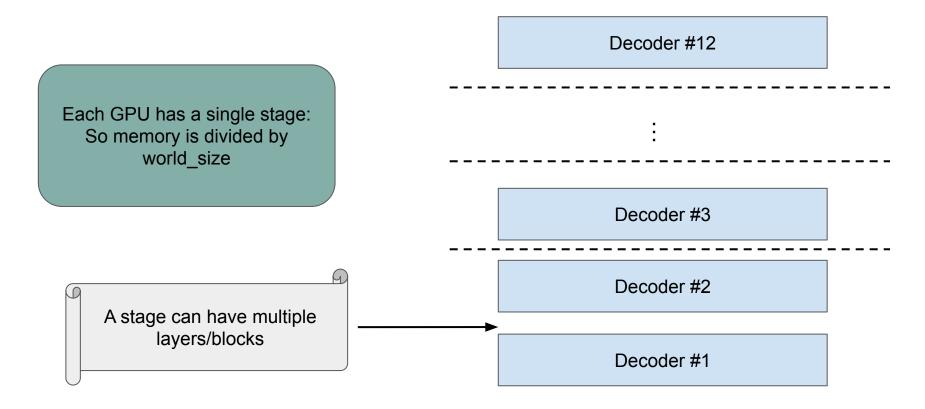
Doesn't work if your model is too big to fit on a single gpu



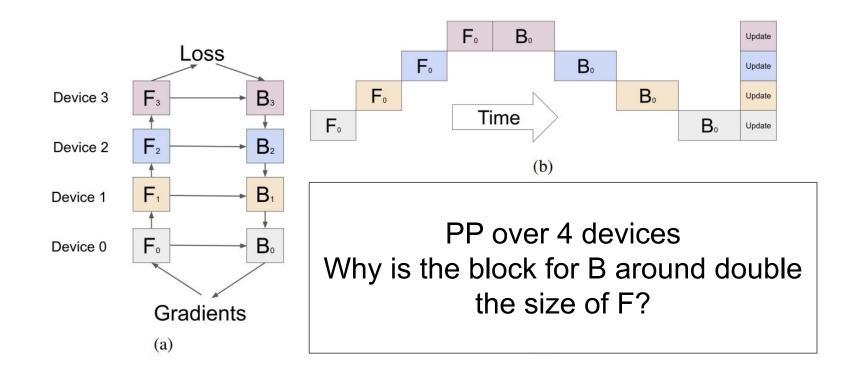




## Pipeline Parallelism

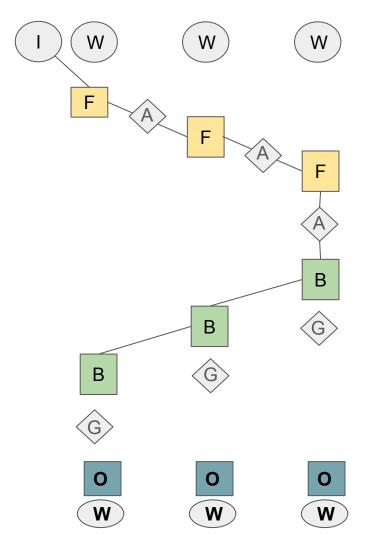


## Pipeline Parallelism



#### What is needed?

- 1. Split up the model loading only needed layers on each process
- 2. Within stages:
  - a. Run forward/backward normally
- 3. At stage boundaries:
  - a. Use p2p send/recv CCL calls
- 4. At end of forward: flow in reverse direction for backward pass
- 5. Each optimizer is only focused on it's section of the model



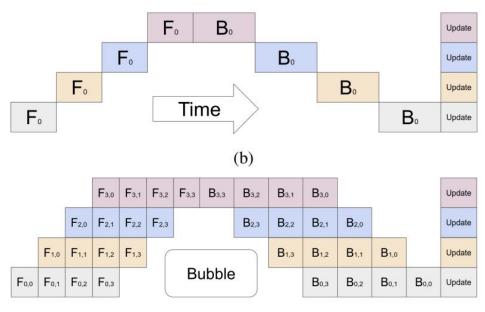
```
\mathbf{m} = \dots
i, o = \dots
class PPRunner(core.Runner):
    def run(self):
        # make a local copy of the model
        m = copy.deepcopy(m)
        ... # identify what subset of the model we will use for this rank
        ... # implement forward pass depending on our rank,
        # including any sync before and after
        loss = out - o # only on the final rank
        ... # implement backward pass here, depending on rank
        ... # any final synchronization here
        # purely test code
        global m_final
        ... # copy gradients of our portion into a single copy
core.run world(4, DDPRunner)
# single model processing here
# ...
m.compare(m_final)
# This test should return true
```

#### Problem with PP

**Bubbles** 

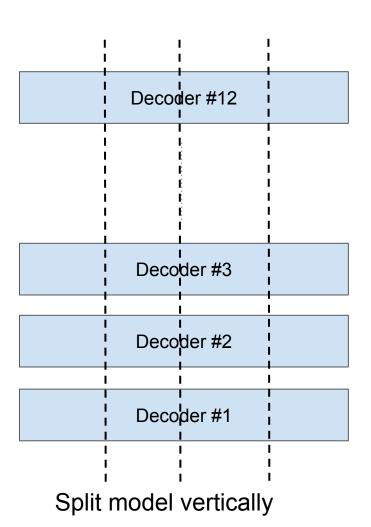
Lot of prior art on bubble management:

- 1. Mini/micro batches
- 2. Reuse of existing bubbles for secon
- 3. Alternative schedules
- 4. Fusion (of pipelines)



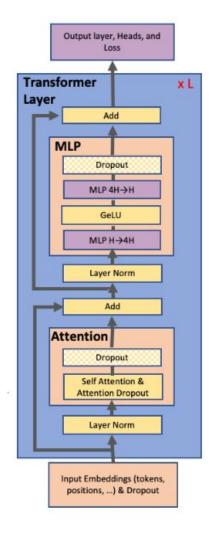
(c)

# Can we split up the model in a way that does not lead to bubbles?



What is the challenge with this?

There is **compute balance** between the processes, but how does this process now look like?



Let's focus on a single MLP block in 1 Decoder layer:

Linear layers A and B, say input is X

X = XA

X = GELU(X)

X = XB

X = DROPOUT(X)

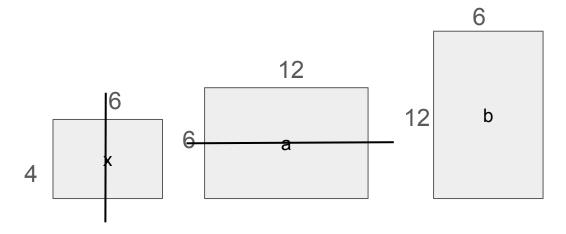
#### import numpy as np

```
np.random.seed(10)
x = np.random.rand(4, 6)
a = np.random.rand(6, 12)
b = np.random.rand(12, 6)
x = x @ a
                                       a
                              6
                                      a[0:2] \rightarrow along rows
                 12
                                      a[:, 0:2] -> along
    6
                                       columns
          6
   Χ
                                      # check and verify the
                                       shapes
```

# Let's try one option

One option to parallelize the GEMM is to split the weight matrix A along its rows and input X along its columns as:

$$X = [X_1, X_2], A = \begin{bmatrix} A_1 \\ A_2 \end{bmatrix}.$$
 (2)



## Solution

```
x @ a
x1 = x[:, 0:3]
x2 = x[:, 3:6]
a1 = a[0:3]
a2 = a[3:6]
(x1@a1) + (x2@a2)
```

# How to continue the chain of computation?

X = XA

X = GELU(X)

X = XB

X = DROPOUT(X)

Let's simplify for the moment and work with:

Y = XA

Z = YB

Y1 = X1A1, Y2 = X2A2

How will you partition B?

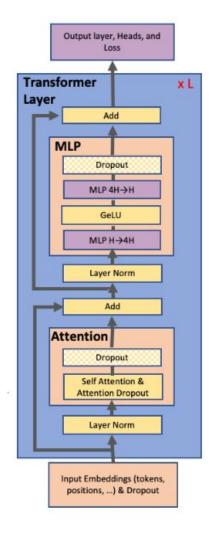
How will you arrange the computation so that the 2 processes can proceed independently?

## Solution

```
x @ a
x1 = x[:, 0:3]
x2 = x[:, 3:6]
a1 = a[0:3]
a2 = a[3:6]
y1 = x1@a1
y2 = x2@a2
```

$$z1 = y1 @ b$$
  
 $z2 = y2 @ b$ 

$$z == z1 + z2$$



Let's focus on a single MLP block in 1 Decoder layer:

Linear laters A and B, say input is X

X = XA

X = GELU(X)

X = XB

X = DROPOUT(X)

## The problem

So, we can safely synchronize once at the end of the MLP block.

But, this doesn't work because of GELU:

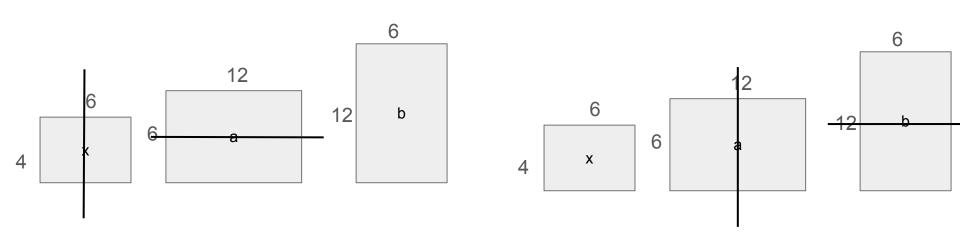
$$GELU_{tanh}(x) = 0.5x \left( 1 + \tanh\left(\sqrt{\frac{2}{\pi}} \left(x + 0.044715x^3\right)\right) \right)$$

$$G(x) + G(y) != G(x + y)$$

This means we need to synchronize before the GELU is applied...

# Let's try another option

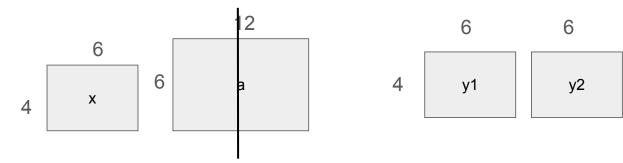
Split the a matrix vertically!



## Solution

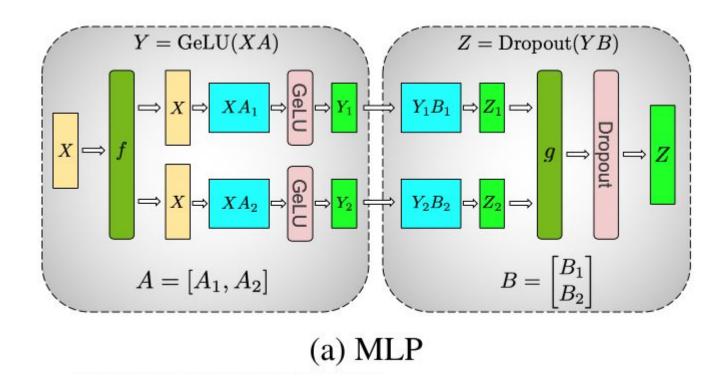
```
y = x @ a
                                      b1 = b[0:6]
# keep x as-is
                                      b2 = b[6:12]
X
                                      z1 = y1@b1
a1 = a[:,0:6]
                                      z2 = y2@b2
a2 = a[:,6:12]
                                      z == z1 + z2
y1 = x@a1
y2 = x@a2
```

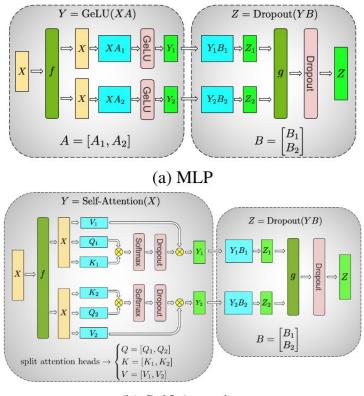
# Why does this work?



Each element of y is fully-ready At this point

### Final logic for the MLP block - with 1 all-reduce for fwd



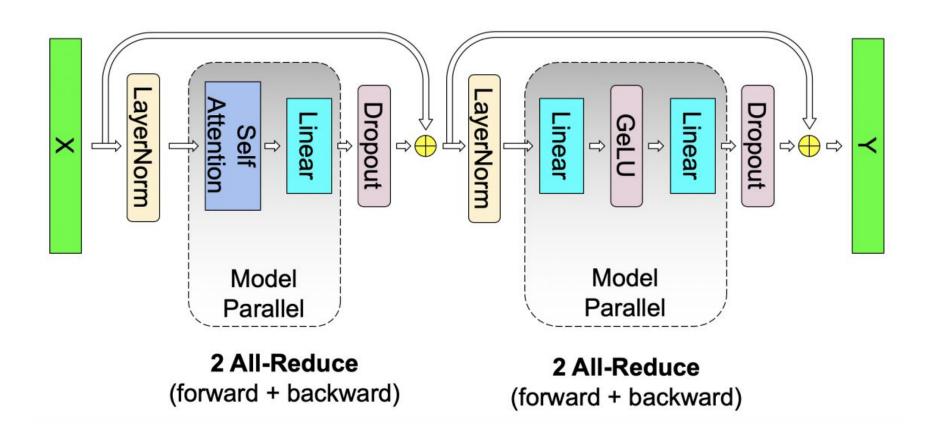


(b) Self-Attention

Figure 3. Blocks of Transformer with Model Parallelism. f and g are conjugate. f is an identity operator in the forward pass and all reduce in the backward pass while g is an all reduce in the forward pass and identity in the backward pass.

Some of the less intensive layers are simply duplicated

Even that has been improved with other advanced techniques



#### Complexities

Refer to:

https://danielvegamyhre.github.io/ml/performance/2025/03/30/illustrated-megatron.html for a fuller explanation of all combinations of TP

# Implementing TP

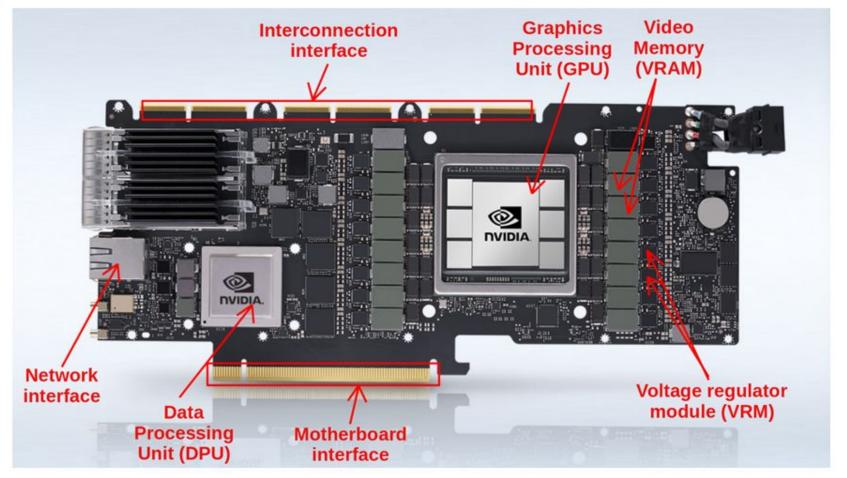
#### The Problem with TP

Too many network connections: 4 all-reduces per block

All-reduce is very expensive across Node boundaries



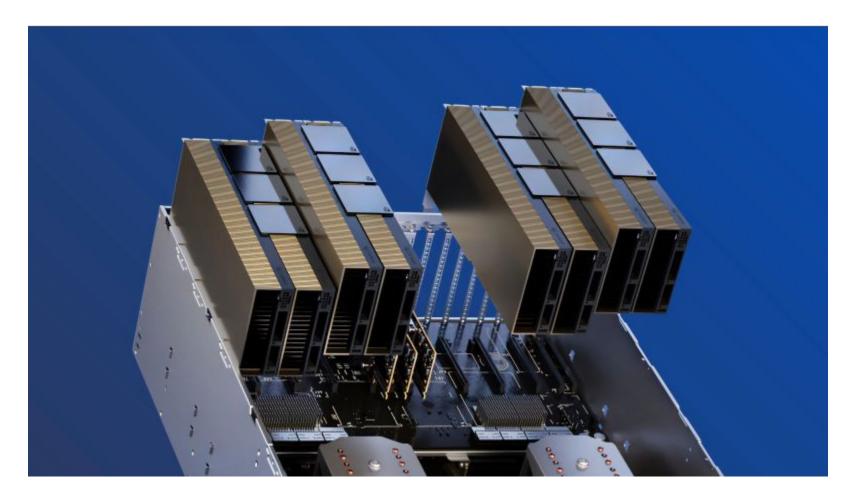


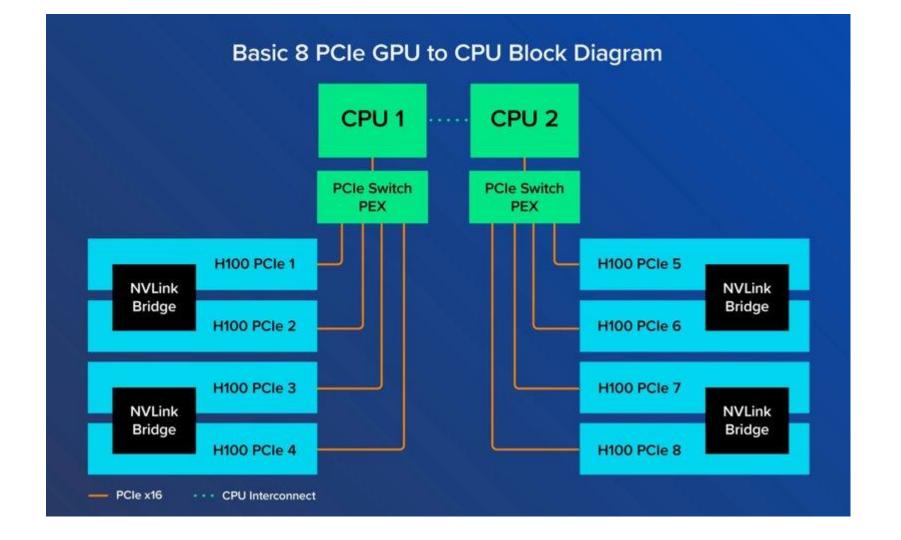




https://www.youtube.com/watch?v=5daQpuK0KYU

#### 8 Nvidia H100 PCie devices connected into a node



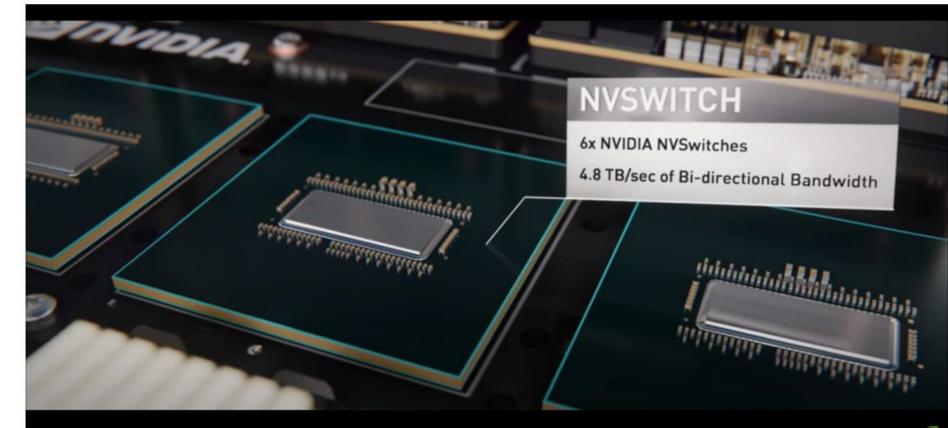


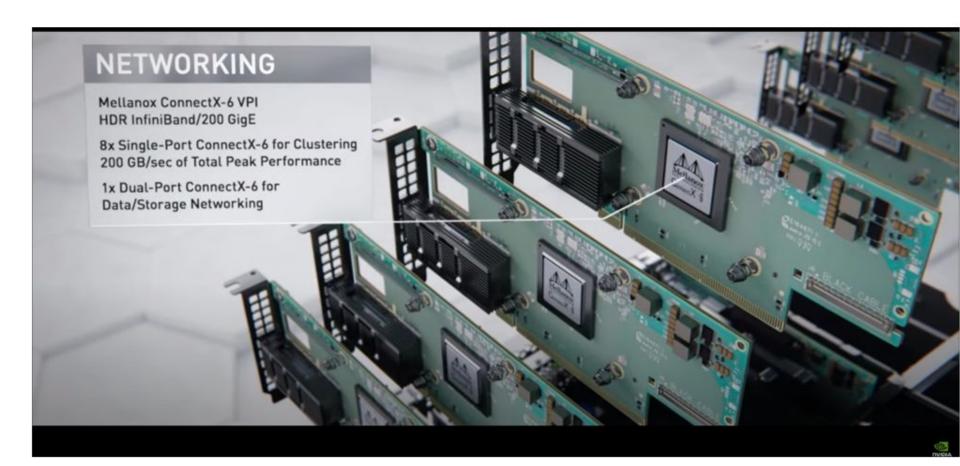
## Nvidia DGX A100 "node"

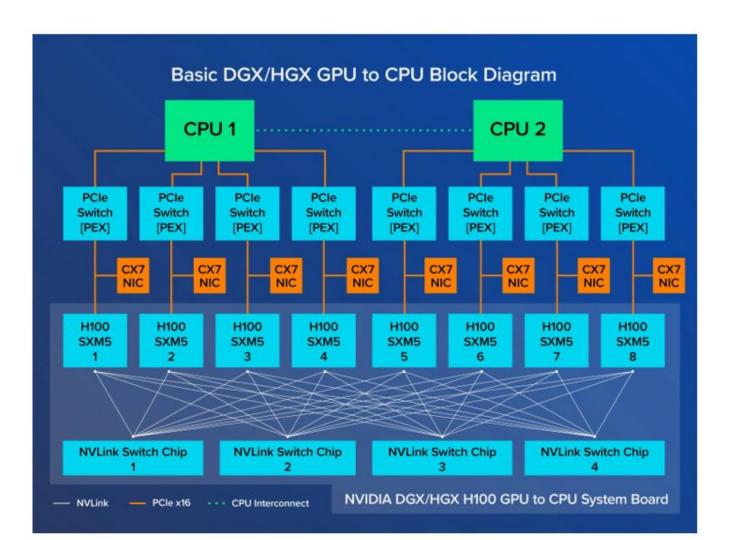


#### 8 A100's slotted onto a DGX Base









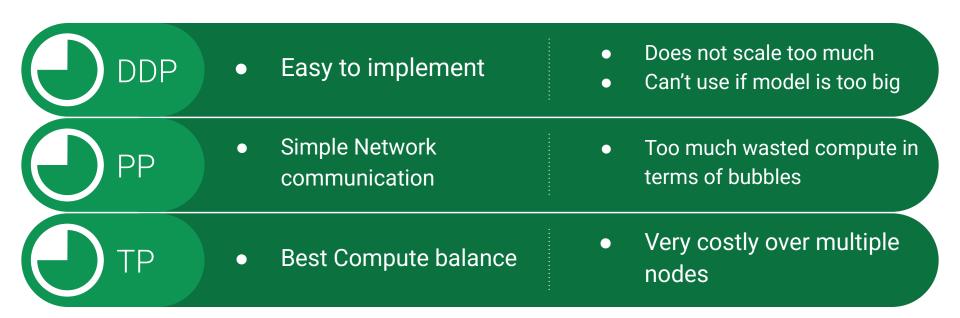
#### Now what?

All-reduce within a single node is very fast,

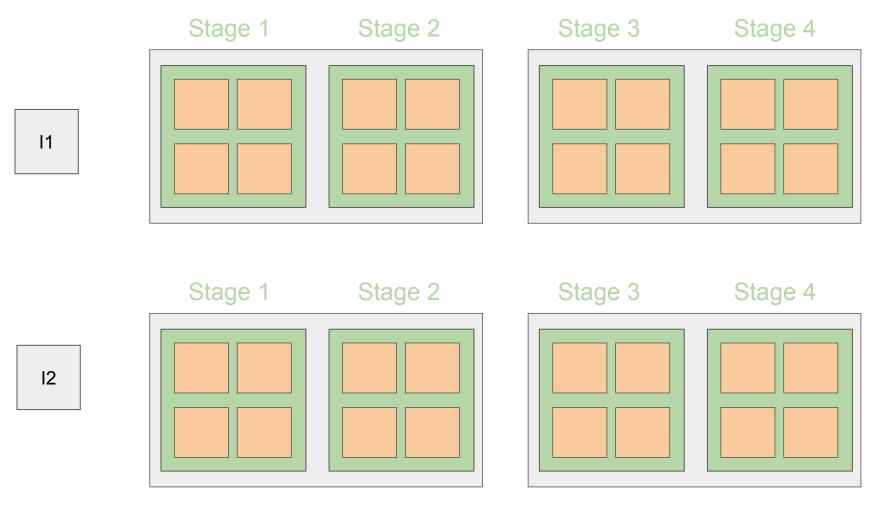
But the moment you have to go to multiple Nodes, > 8 gpus, hit is going to be very high..

What can you do about this?

# 3 techniques compared



# 3D Parallelism: d-t-p



### 3d-parallelism

Tries to use best of all 3 techniques

Optimal schedule, it d-p-t depends on a lot of factors:

- Model size
- 2. Data size
- 3. Number of gpus
- 4. GPU networking interconnects
- 5. Task-characteristics: inference vs training
- 6. etc

# FSDP: Fully sharded Data Parallel

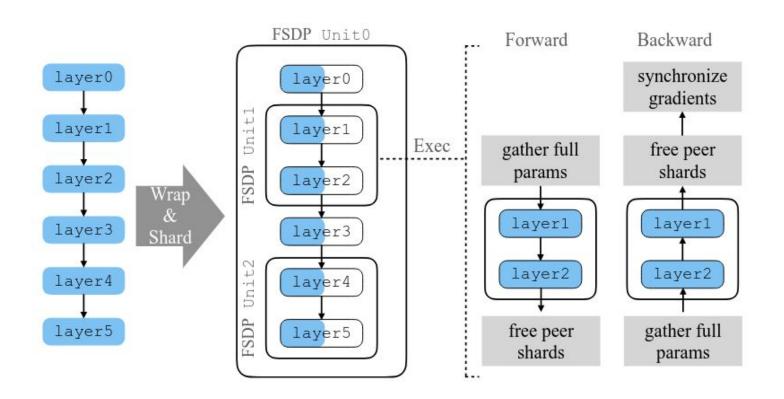


Figure 1: FSDP Algorithm Overview

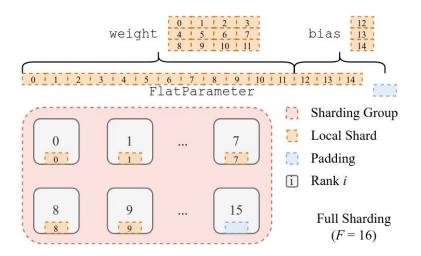


Figure 3: Full Sharding Across 16 GPUs

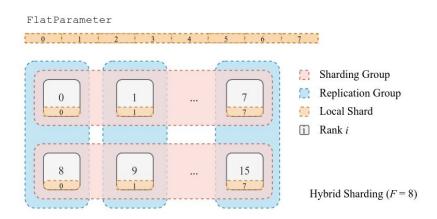


Figure 4: Hybrid Sharding on 16 GPUs: GPUs are configured into 2 sharding groups and 8 replication groups

# Implementing FSDP

#### Challenges

Computation-communication overlap

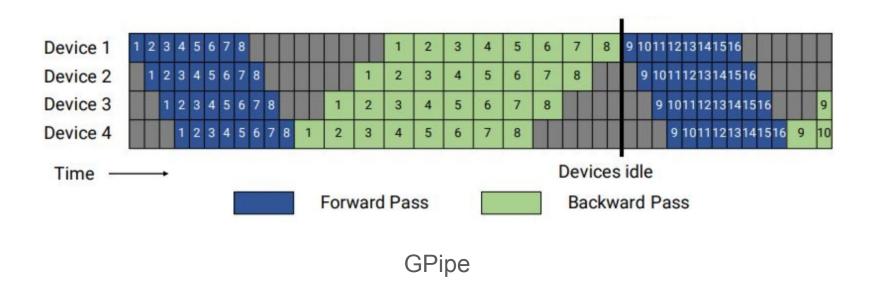
Network speeds dominate computation speeds: leading to around 50% wasted compute on large clusters

Designing large clusters of gpus to manage communications

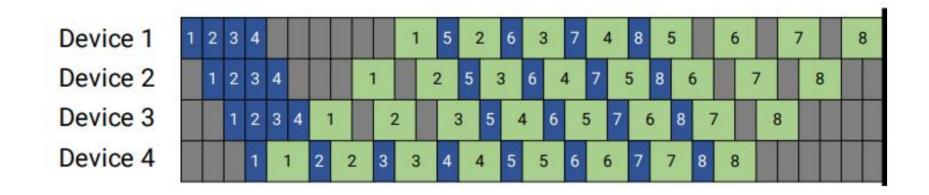
Congestion control due to bursty nature of training communication

# Some Prior Art

## Pipeline Parallelism advancements - GPipe



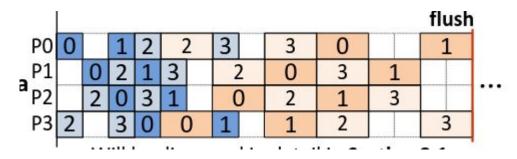
## Pipeline Parallelism advancements - PipeDream



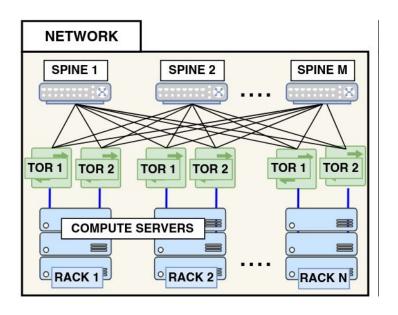
#### Chimera

Table 2: Comparison between different pipeline schemes.

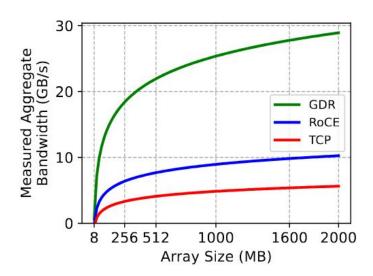
Pipeline Schemes	Bubble Ratio	Weights Memory	Activations Memory	Convergence Friendly
PipeDream [38] PipeDream-2BW [39]	≈ 0 (Å)(Å ≈ 0 (Å)(Å	$[M_{\theta}, D * M_{\theta}]^{1} \blacksquare \qquad \qquad 2M_{\theta}  \triangle$	$[M_a, D * M_a]^1$ $(C)$ $[M_a, D * M_a]^1$ $(C)$	Asynchronous 📭
GPipe [26] GEMS [28]	$(D-1)/(N+D-1) \qquad \qquad \qquad \qquad \approx (D-1)/(D+\frac{1}{2}) \qquad \qquad$	$M_{ heta}$ $\mathcal{C}$ $2M_{ heta}$ $\mathcal{C}$	$N*M_a$ $M_a$	Synchronous 🖒
DAPPLE [16] Chimera (this work)	(D-1)/(N+D-1) $(D-2)/(2N+D-2)$	$M_{ heta}$ $\mathcal{C}$ $2M_{ heta}$ $\mathcal{C}$	$[M_a, D * M_a]^1$ $(D/2 + 1)M_a, D * M_a]^1 (D+1)$	



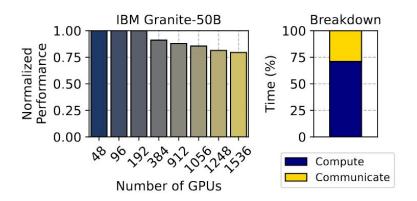
# IBM's Compute Setup - Vela (ASPLOS'25)



#### Vela



(a) 1024 GPUs AllReduce performance.



**Figure 1.** Training performance of IBM Granite-50B model on Vela at different scales normalized to ideal throughput.

# Alibaba HPN (SIGCOMM '24)

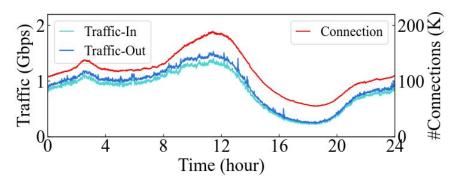


Figure 1: Traditional cloud computing traffic pattern.

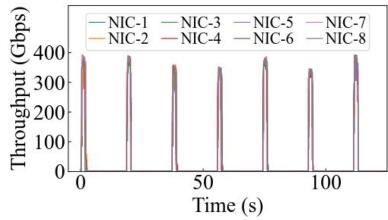


Figure 2: NIC egress traffic pattern during production model training.

## SuperBench (Microsoft), ATC'24

**Cloud Environment** Cloud infrastructure environments can introduce additional incidents, as they differ from vendors' qualification environments in terms of power, temperature, and other factors. For example, data centers in tropical areas experience more incidents due to higher temperatures. We observed a 35× increase in defective InfiniBand links with  $> 10^{-12}$  bit error rate in data centers in tropical areas compared to data centers in higher latitudes, leading to significantly degraded performance for training and inference. Another example is GPU throttling. Even within the same data center, different racks or locations within the same rack can exhibit varying temperatures. However, all GPUs are designed with identical cooling and heatsinks by vendors, resulting in GPUs located in warmer locations potentially experiencing thermal throttling if they cannot receive more cool air.

# Astral (Tencent) (SIGCOMM '25)

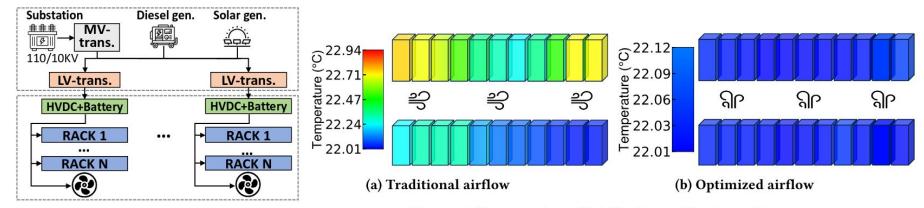


Figure 4: Hierarchy of the distributed HVDC power system.

Figure 5: Temperature distribution with air cooling.

## Recap

-> How multiple GPUs are used in concert

-> What are the challenges, research areas?

### Distributed Training Algorithms vs Systems

- + Data Parallelism (DDP)
- + Tensor Parallelism (TP)
- + Pipeline Parallelism (PP)

- Full Sharded Data Parallelism (FSDP)
  - + With variants like Full/Hybrid Shard

- + Pytorch in-built DDP
- + Hugging Face Accelerate
- + Nvidia Megatron
- + Pytorch in-built FSDP
- + Microsoft DeepSpeed ZeRO family

- + 3D parallelism
- + Context/Sequence Parallelism
- + Expert Parallelism (EP)

Systems

3 - Inference

4 - Dist. Training

Mechanics

1 - LLM Core

2 - Pytorch/GPUS

Math