

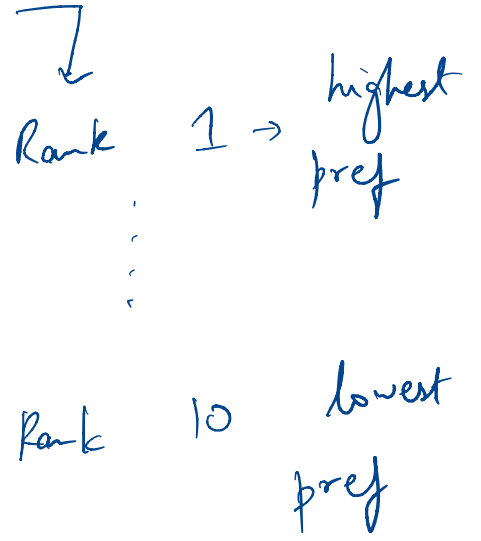
hello!

# CS 744: PIPEDREAM

Shivaram Venkataraman

Fall 2021

# ADMINISTRIVIA

- Assignment 2 is due Wednesday AM! *Tuesday evening!*
- Course project groups due Oct 11, Monday! 

Rank 1 → highest pref

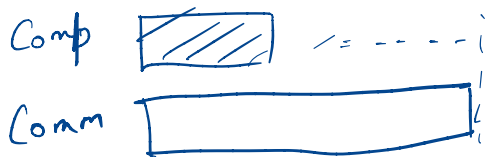
...

Rank 10 → lowest pref
- Project proposal aka Introduction (10/25)

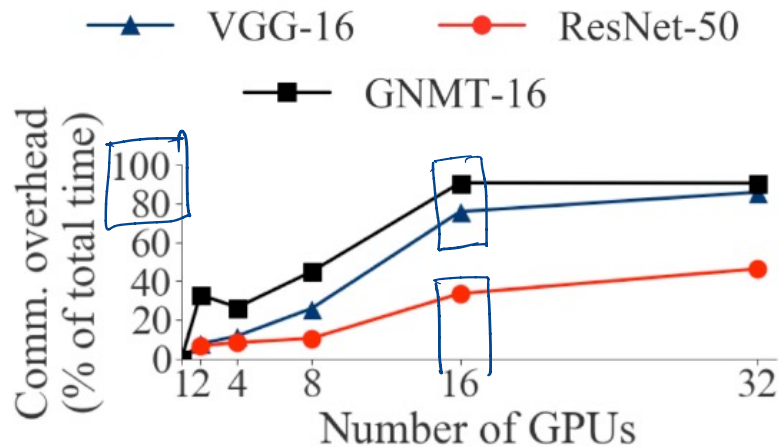
# LIMITATIONS OF DATA PARALLEL

① Overhead of comm. is high with data parallelism

② Under utilization



③ Overhead / Scaling is different for GNMT vs. Resnet

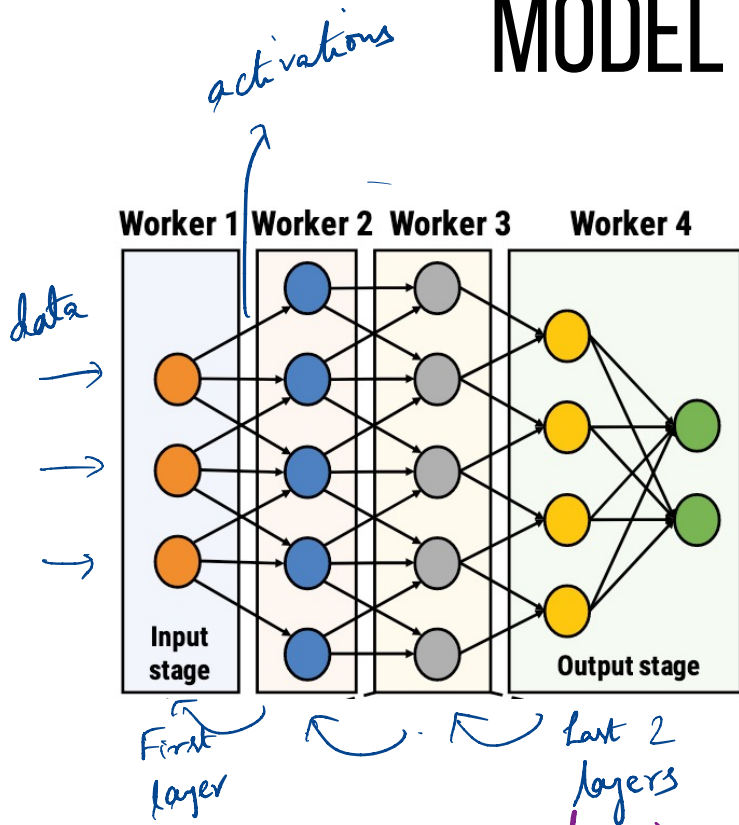


8xV100s with NVLink (AWS)  
PyTorch + NCCL 2.4

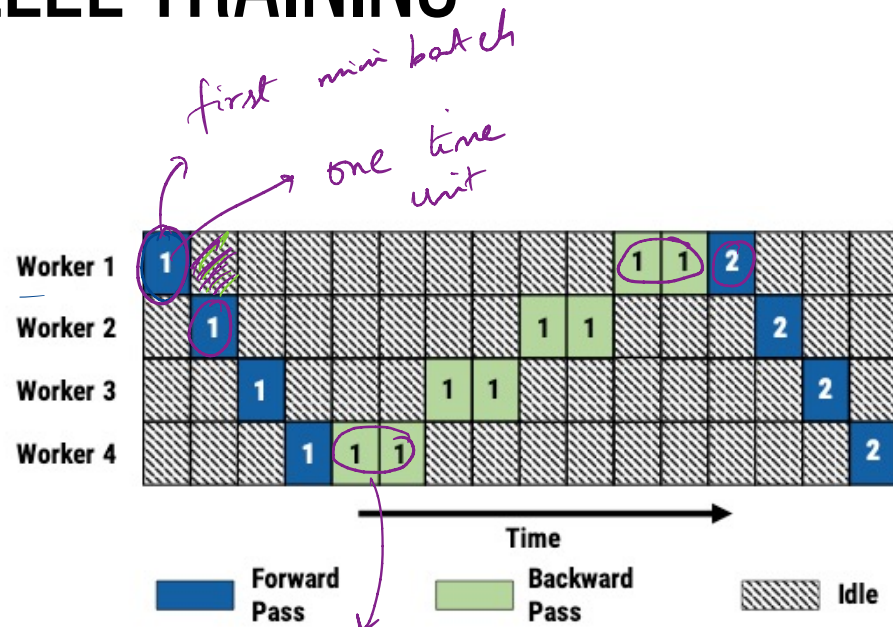
5

“fraction of training time spent  
in communication stalls”

# MODEL PARALLEL TRAINING



$bs = 128$   
examples  
CIFAR-10



- Memory overhead is lower
- Comm. per worker is only related to activation size

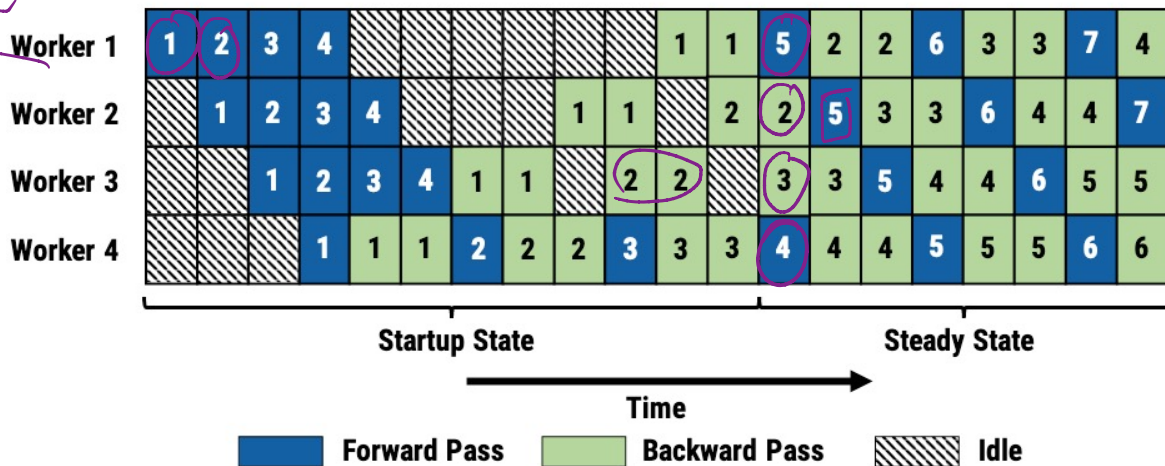
backwards pass takes 2 time unit

Low UTILIZATION

# PIPELINE PARALLEL

→ Instead of having  
one minibatch in flight  
you can have four! →

- ① Partitioning
- ② Scheduling
- ③ Learning =



## Advantages?

- High utilization after startup
- The comm. per worker is proportional to activation size

# CHALLENGE 1: WORK PARTITIONING

Stage: Collection of layers at a given worker

Goal: Balanced stages in the pipeline. Why?

Steady state throughput is the throughput of the slowest stage

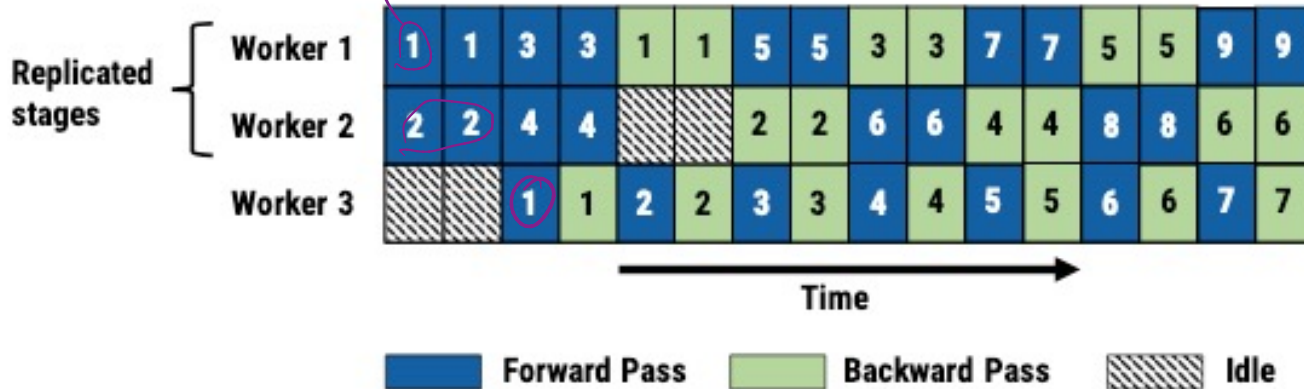
□ □ □  
↳ bad utilization if unbalanced

Stages can be replicated!

mini batch

$bs = 128$

□<sub>W<sub>1</sub></sub> □<sub>W<sub>2</sub></sub> □<sub>W<sub>3</sub></sub> → layers are fast



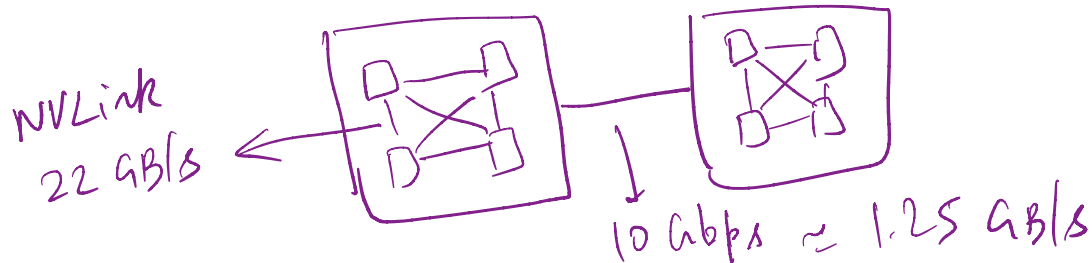
# WORK PARTITIONING

*Given a model*

Profiler: computation time for forward, backward *for each layer*  
size of output activations, gradients (network transfer)  
size of parameters (memory)

Dynamic programming algorithm

Intuition: Find optimal partitions within a server,  
Then find best split across servers using that



# CHALLENGE 2: WORK SCHEDULING

Traditional data parallel

forward iter(i)

backward iter(i)

forward iter(i+1)

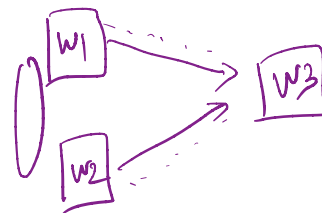
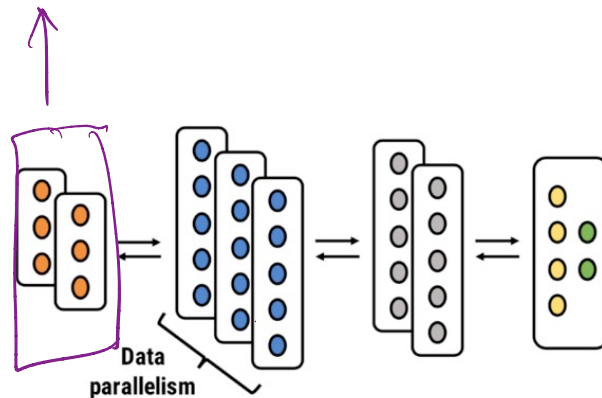
...

Pipeline parallel: Worker can

Forward pass to push to downstream

Backward pass to push to upstream

handling forward & backward pass for layer 1





# CHALLENGE 2: WORK SCHEDULING

Num active batches  $\sim$  num\_workers / num\_replicas\_input

Schedule one-forward-one-backward (IFIB)  $\rightarrow$

Round-robin for replicated stages  $\rightarrow$

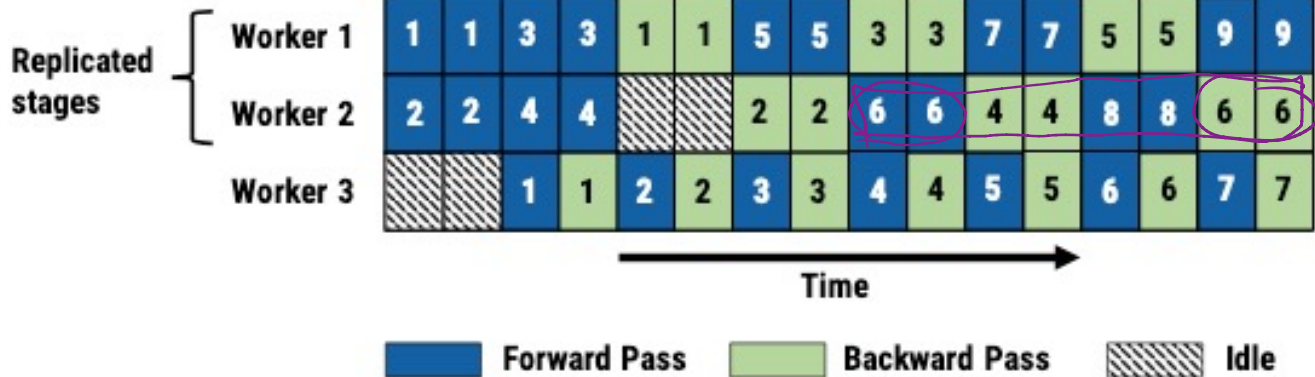
same worker for fwd, backward

$\rightarrow$  mini batch

1  $\rightarrow$  W1  
2  $\rightarrow$  W2  
3  $\rightarrow$  W1  
...

Worker 2

FW (6)  
BW (4)  
FW (8)  
BW (6)

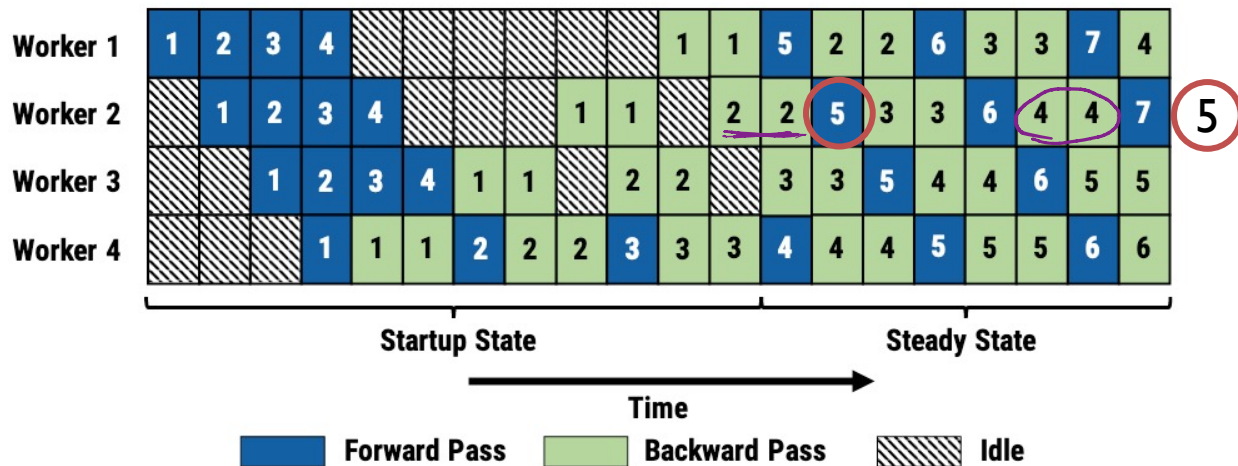


# CHALLENGE 3: EFFECTIVE LEARNING

Naïve pipelining

Different model versions forward and backward

W2 : fwd pass of mini batch 5  
↳ model version 2  
backward pass for mini batch 5  
↳ model version 4



# CHALLENGE 3: EFFECTIVE LEARNING

## Weight stashing

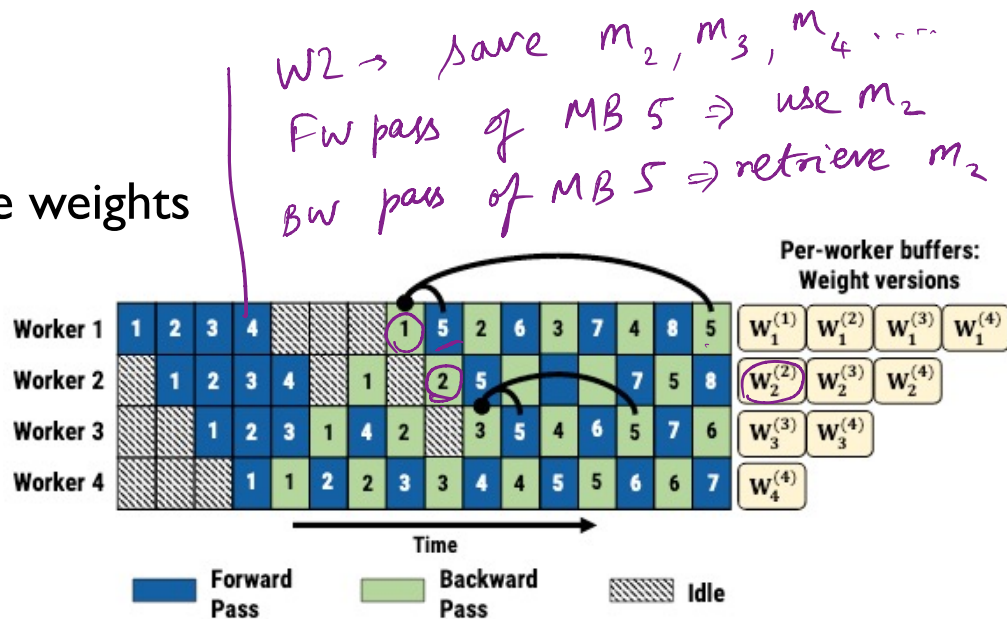
Maintain multiple versions of the weights

One per active mini-batch  $\approx 4$

Use latest version for forward pass.

Retrieve for backward

No guarantees across stages!



# STALENESS, MEMORY OVERHEAD

How to avoid staleness:

Vertical sync

Memory overhead

Similar to data parallel?

if you have  $k$  mini batches in flight

$\Rightarrow$  same as  $k \times b$  batchsize DP with

When you first process mini batch  
note model version  $\leftarrow$  pass it along

$w_1, m_0 \rightarrow MB_1$

$w_2, m_0 \rightarrow MB_1$

$w_4, m_0 \rightarrow MB_1$

$w_4, m_0 \rightarrow MB_1$

$w_1, m_0 \rightarrow MB_1$

# SUMMARY

Pipeline parallelism: Combine inter-batch and intra-batch

Partitioning: Replication, dynamic programming

Scheduling: IFIB

Weight management: Stashing, vertical sync

# DISCUSSION

<https://forms.gle/j2GCDyqCejBH8DaCA>

List two takeaways from the following table

Model Name	Model Size	GPUs (#Servers x #GPUs/Server)	PipeDream Config	Speedup over DataParallel (Epoch Time)
Resnet-50	97MB	4x4 2x8	16 16 <i>replication</i>	1x 1x <i>same as Data Parallel</i>
VGG-16 <i>Diff hardware topology → diff speedups</i>	528MB	4x4 <u>2x8</u>	15-1 <u>15-1</u>	<u>5.28x</u> 2.98x
GNMT-8	1.1GB	3x4 2x8	Straight 16	2.95x 1x



What are some other workload scenarios (e.g. things we discussed for MapReduce or Spark) that could use similar ideas of pipelined parallelism? Develop such one example and its execution

Sync between stages  $\rightarrow$  can be pipelined

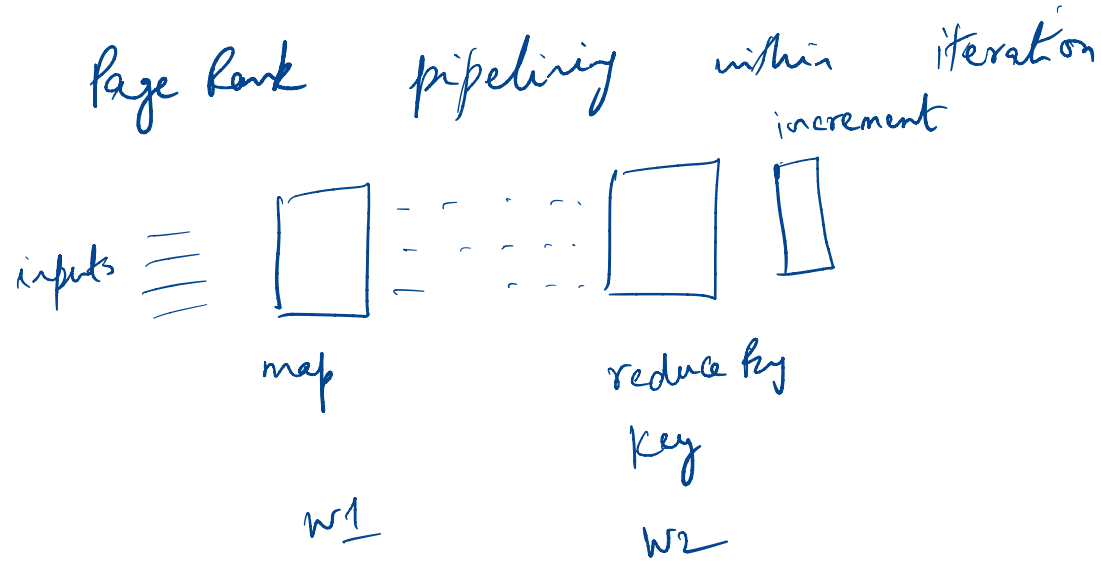
MR  $\rightarrow$  if we don't have a dep. on all map tasks,  
can we start reducers early

PageRank\*  $\rightarrow$  where output  $(i-2)$  is input to iter  $(i)$   
 $(i-1)$   $\dots$   $(i+1)$

can we run  
two of these  
at the same time



What are some other workload scenarios (e.g. things we discussed for MapReduce or Spark) that could use similar ideas of pipelined parallelism? Develop such one example and its execution



# NEXT STEPS

Next class: Ray

Assignment 2 is due soon!

Course project: Oct 11 (Monday) Submit titles, groups