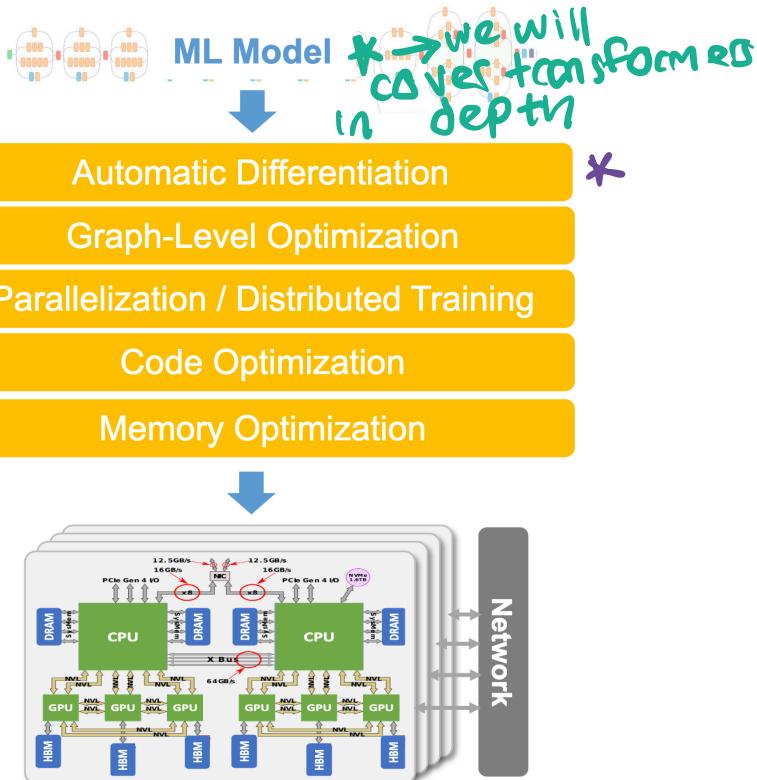


- \* Agenda: class
- Overview of Topics we will cover in the
  - Data Parallelism
    - ↳ Parameter servers
    - ↳ All Reduce Algorithms

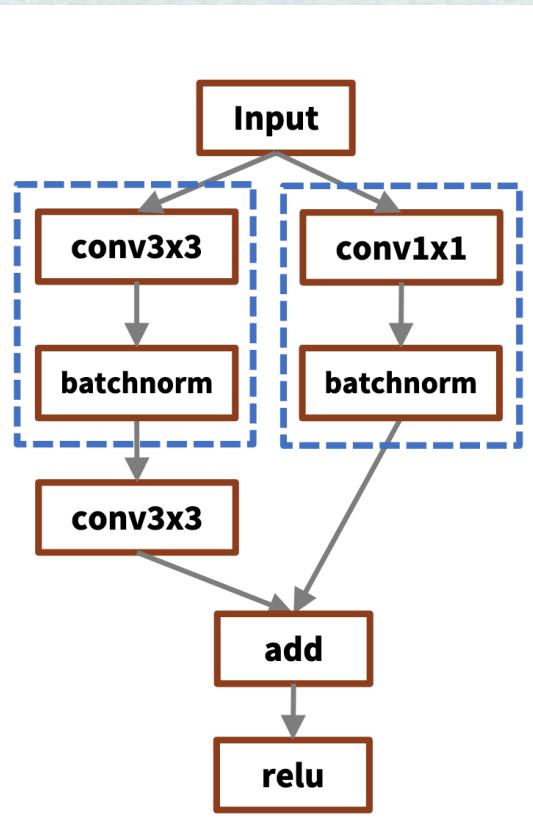
## Recap: Deep Learning Systems



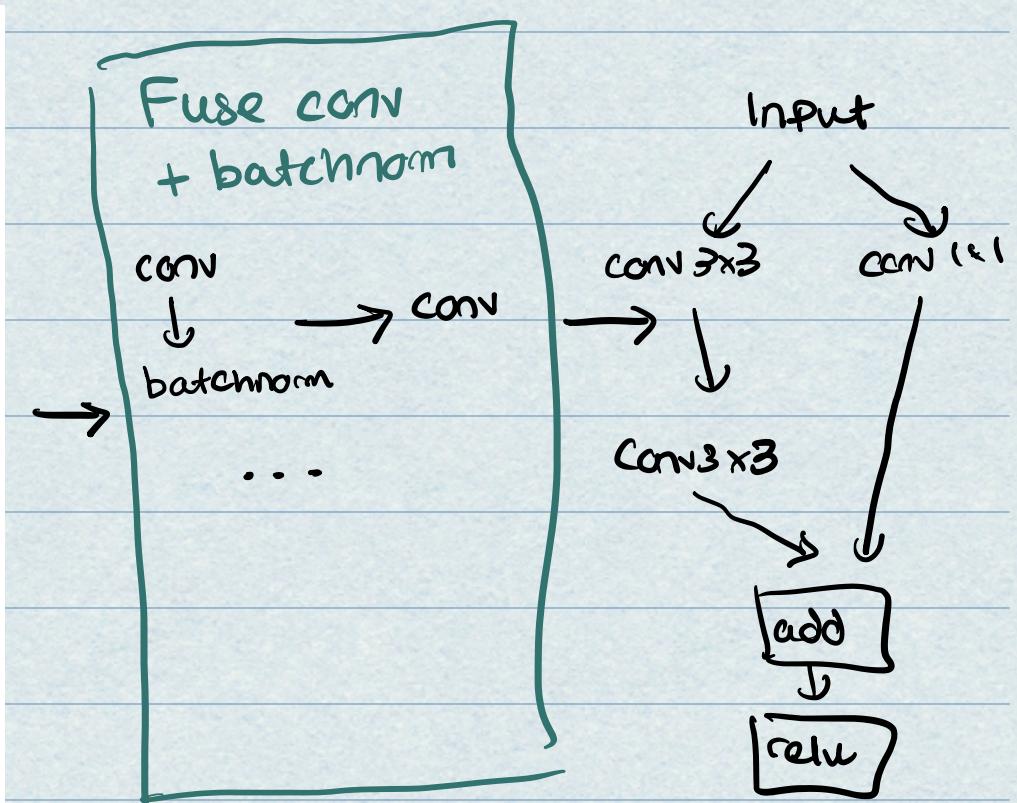
Auto Diff (from last time):

- automatically construct FW and BW computation GRAPH
- implication: 1 unified FW/BW pass to discover weight updates (vs. separate for each weight)

## \* Graph-Level - Optimizations

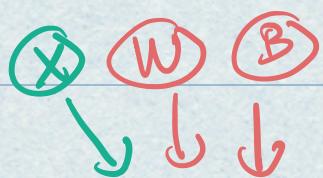


Input  
Computational  
Graph

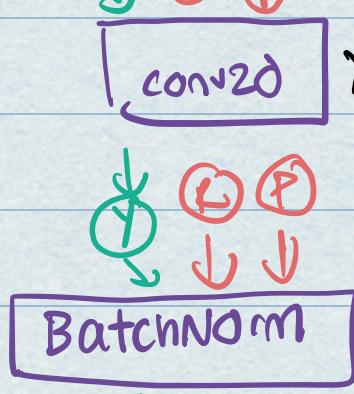


Potential graph  
transformations

## \* Example: Fusing Conv and Batch Norm



- $W, B, R, P$ : constant pre-trained weights
- $X, Y, Z$  involve learned params (dnn)



$$Y(n, c, h, w) = \sum_{d, u, v} (X(n, d, h+u, w+v) \cdot W(c, d, u, v)) + B(n, c, h, w)$$

$\rightarrow$  not in sum

$$Z(n, c, h, w) = Y(n, c, h, w) \cdot R(c) + P(c)$$

\* We can Fuse conv and Batch norm



$$W_a(n, c, h, w) = W(n, c, h, w) \cdot R(c)$$

conv2d



$$B_2(n, c, h, w) = B(n, c, h, w) \cdot R(c) + P(c)$$

$$(\text{intuition}) = ((x \cdot w) + b) \cdot R + P$$

$$= (x \cdot w \cdot R) + (b \cdot R + P)$$

$$Z(n, c, h, w) = \left( \sum_{d, u, v} X(n, d, h+u, w+v) \cdot W_2(c, d, u, v) \right) + B_2(n, c, h, w)$$

- Discussion Q: Why is fusion a good idea?

where is it applied in other areas of systems?

\* How do these rules get applied?

- tensorflow has ~200 rules

(53,000 LOC!)

- includes:

- Fuse conv + relu
- Fuse conv + batchnorm
- <sup>multi</sup>Fuse n conv

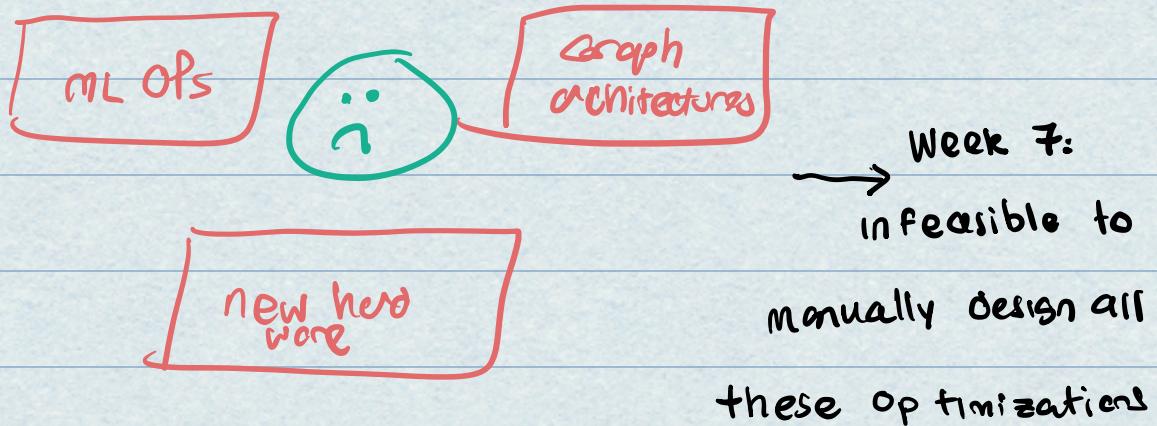
\* What are limitations of rule based optimizers?

1) Robustness: heuristics may not always make things faster

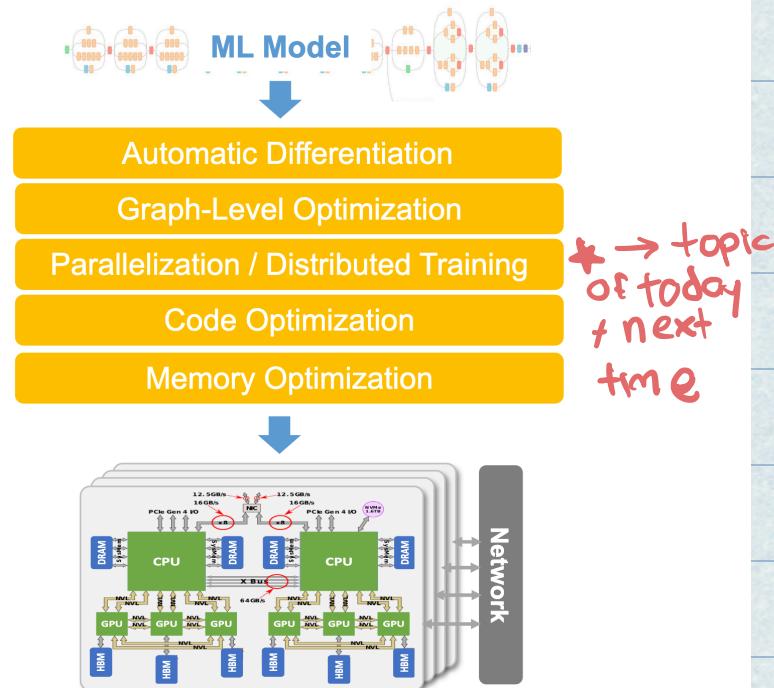
2) Scalability of implementation:

- tensorflow uses ~4K LOC to optimize convs

3) Still can miss some subtle optimizations



## Recap: Deep Learning Systems



Code optimization:

- given a hardware  
how do we map a  
specific operator  
onto it?

$$C = \sum_{k=0}^K A(k, y) \bullet B(k, x)$$

Simple option:

```
for y in range(1024):  
    for x in range(1024):  
        C[y][x] = 0  
        for k in range(1024):  
            C[y][x] += A[k][y] * B[k][x]
```

\* Now one X and  
Y laid out  
in memory  
on the device?

- How is this handled in systems today?
  - HW/SW engineers manually write

## operator libraries

- cuDNN, cuSparse, CURAND, cuBLAS

for GPUs

↳ "  
"cudnn convolutionfw"  
"cUBLAS SGEMM"

- **Problem:** cannot provide support immediately for new ops
- **Problem:** what if HW changes?

## - Week 6 and Project 3:

- how do we program HW?

- HW for ML

- Frameworks to automatically

generate HW code

## - overview of what's to come:

- week 1-2: parallelism

- week 2-5: Thinking about  
HW & HW-aware algorithms w/  
a focus on attention, revisit parallelism

- week 6: GPUs

- Week 7-8: ML frameworks,

# Quantization, sparsity

- After spring break: TBD topics in research
- what are general themes/tradeoffs?

## 1) Memory vs. Computation

- can redesign procedures to avoid recomputations, but comes at cost of memory

## 2) Memory vs. communication

- can scale up, but this will add communication overhead

## 3) Accuracy vs. memory vs. computation:

- can quantize → how much perf overhead?

## 4) How do we restructure methods to perform better on HW?

- Level of alg it self:

- efficient HW-aware

- attention

- using OJ paging to improve attention

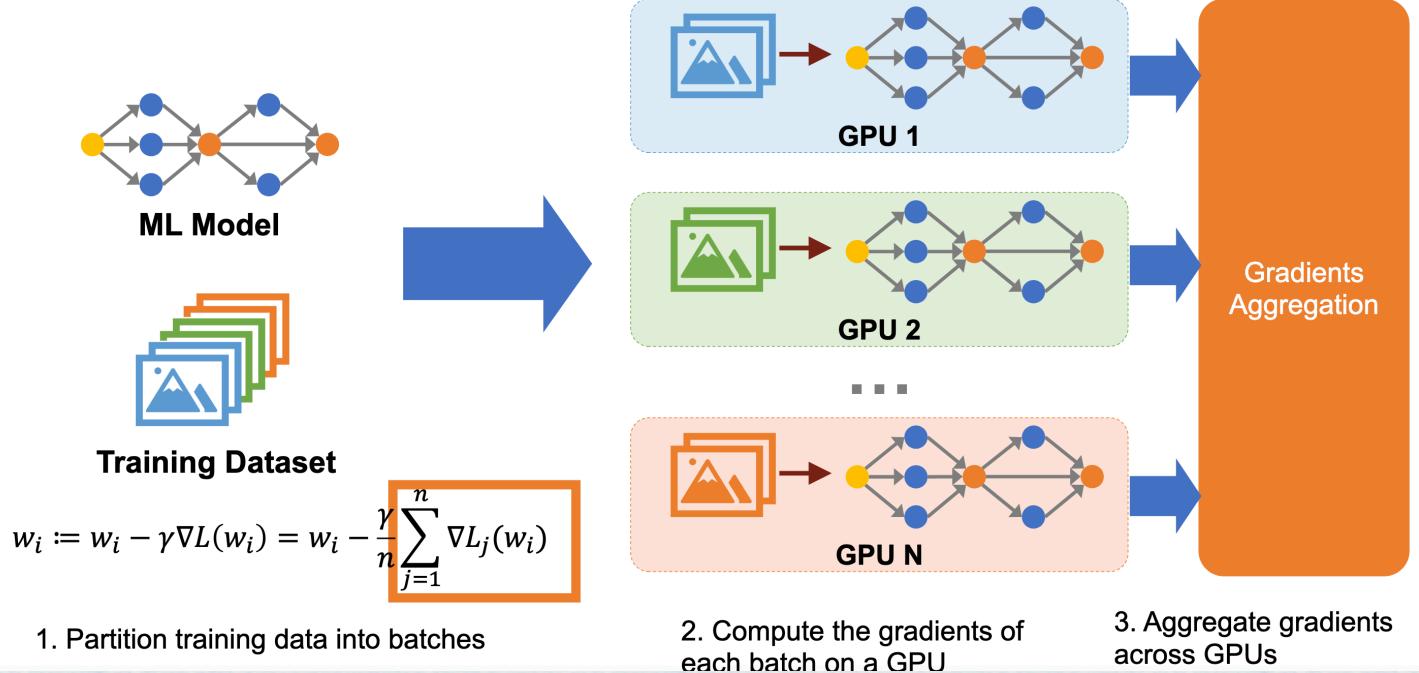
- Level of computation graph / code optimizations

# PARALLELISM

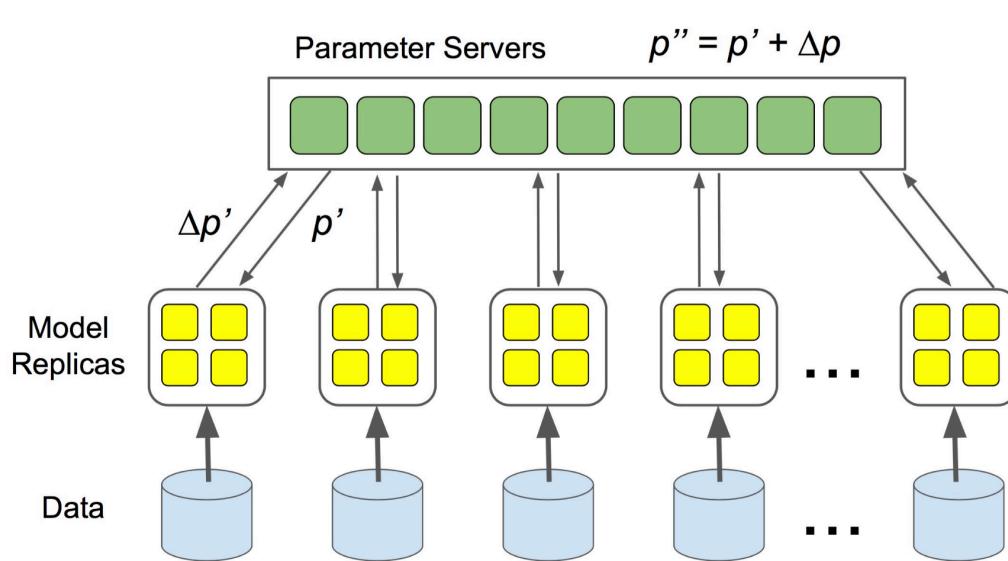
- Data Parallelism

what does  
orange box look  
like?

## Data Parallelism



## - option 1: parameter server



• workers

push gradient  
updates to  
centralized  
server and pulls  
updates  
back

- each worker must have 1 wire to PS, PS must have N wires

## Cons:

- centralized communication →

all workers communicate w/ param.

server for update; cannot scale

to large numbers of workers

## Pros:

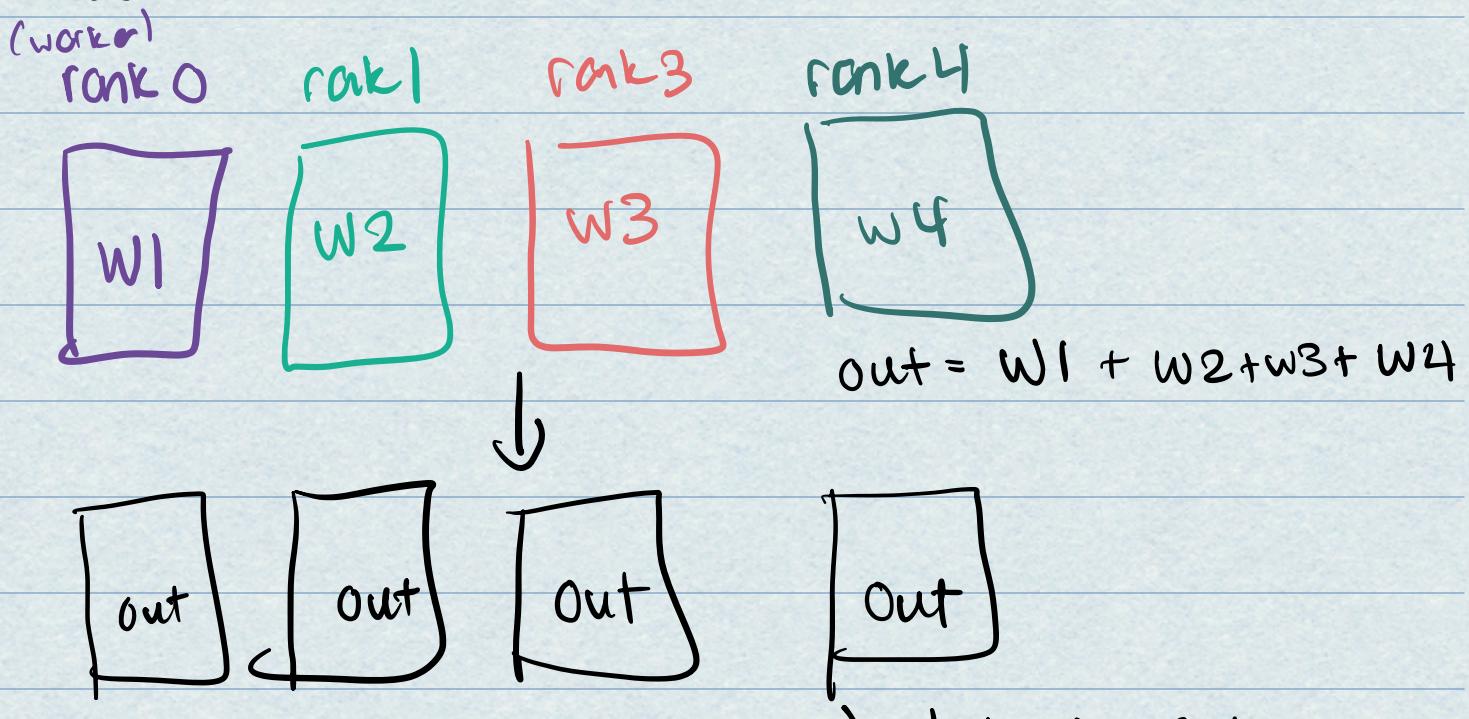
- can potentially scale # of

param. servers by splitting

which param each is responsible for

- Can we DECENTRALIZE communication?

· solution: All-reduce



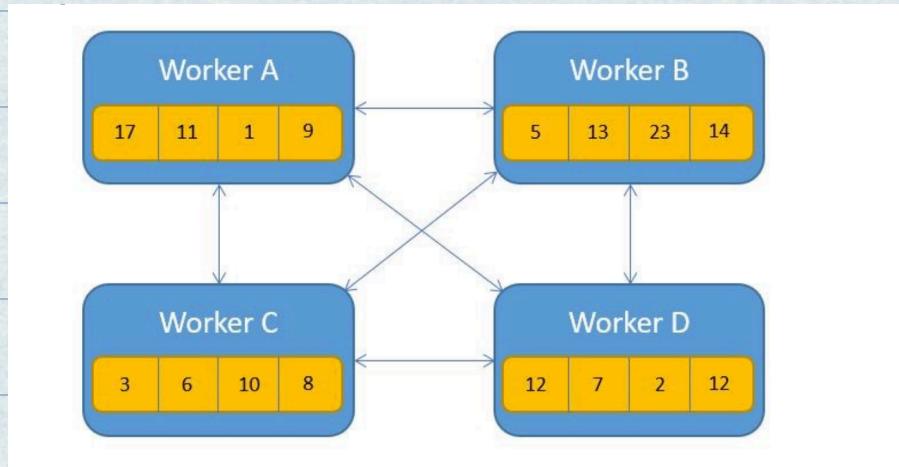
- can also do min, max, prod, bitwise ops

-Q: How do we actually do all reduce?

## 1) Naive All-reduce

-each worker sends its local state to all other workers

- consider N workers and M params



→ each worker must have  $N - 1$  wires of bandwidth B

→ communication:  $O(N \cdot (N-1) \cdot M)$

→ each worker communicates pt-to-pt

w/ every other worker:

- imagine param.server w/o server
- some scalability issue

## 2) Ring all-reduce

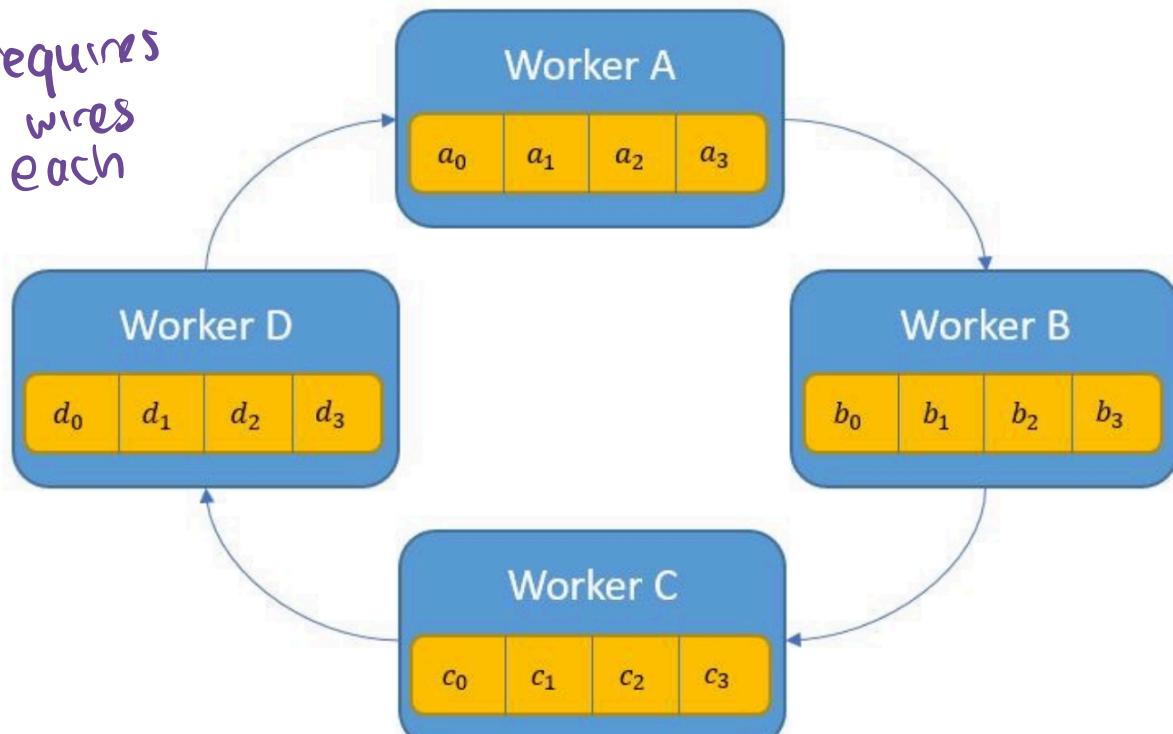
- construct ring of N workers & divide

M params into N slices

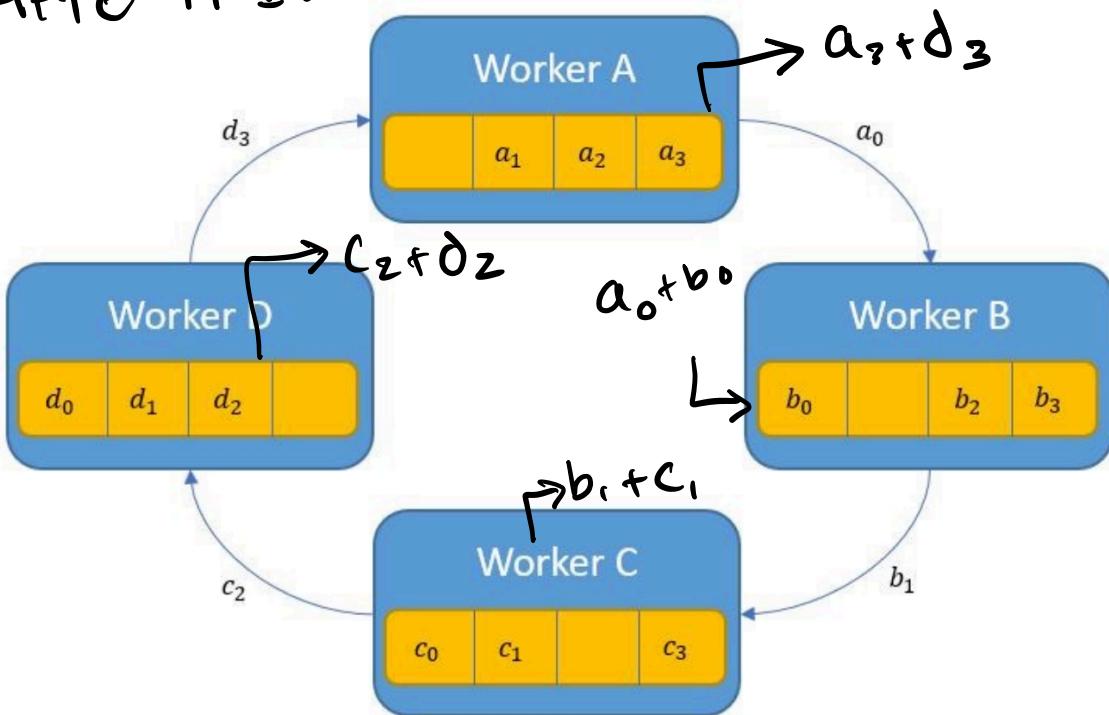
→ step 1: Aggregation: each worker sends a single slice:  $(M/N)$  params: to next worker on ring; repeat  $N - 1$  times

→ goal of step 1: end up w/ sum of each slice on 1 of the workers

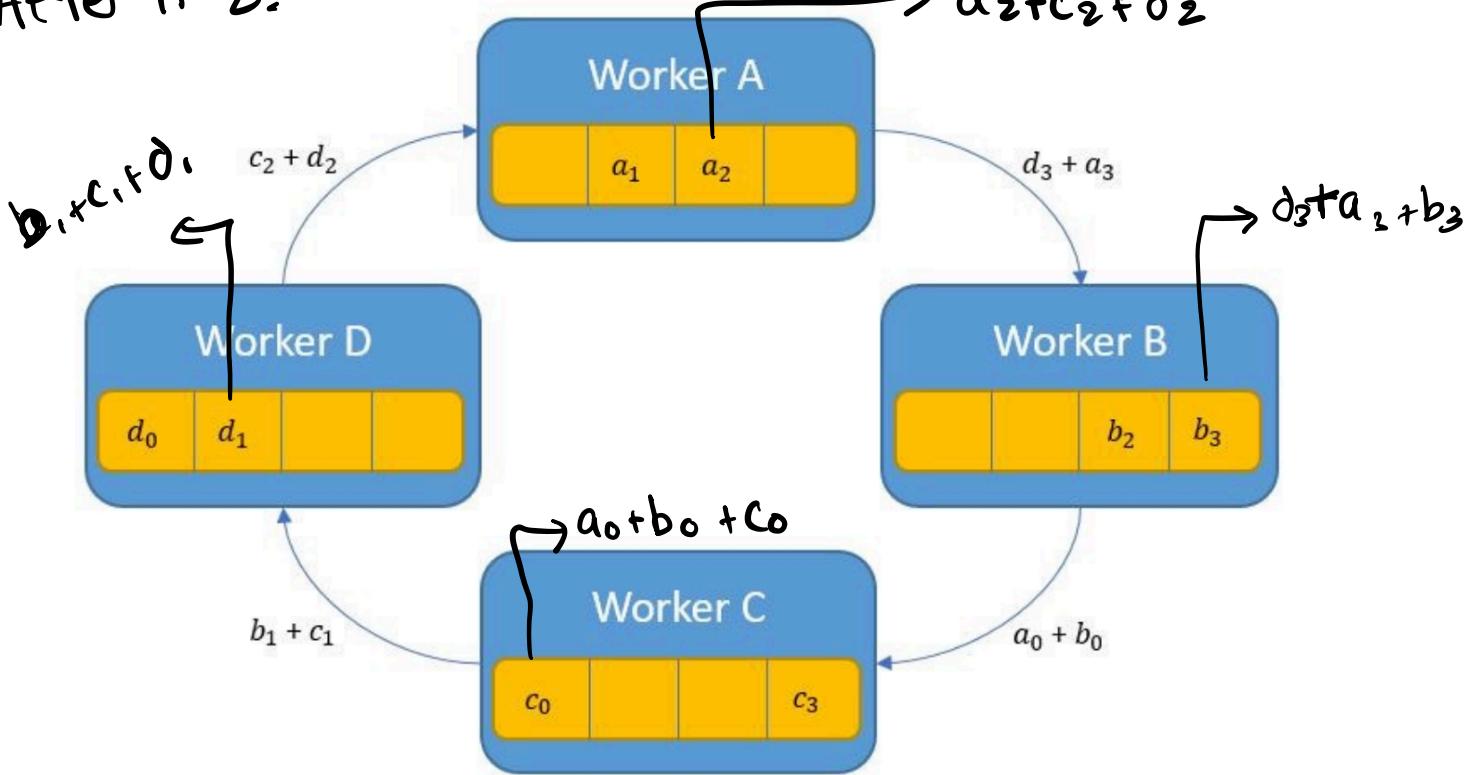
\* requires  
2 wires  
each



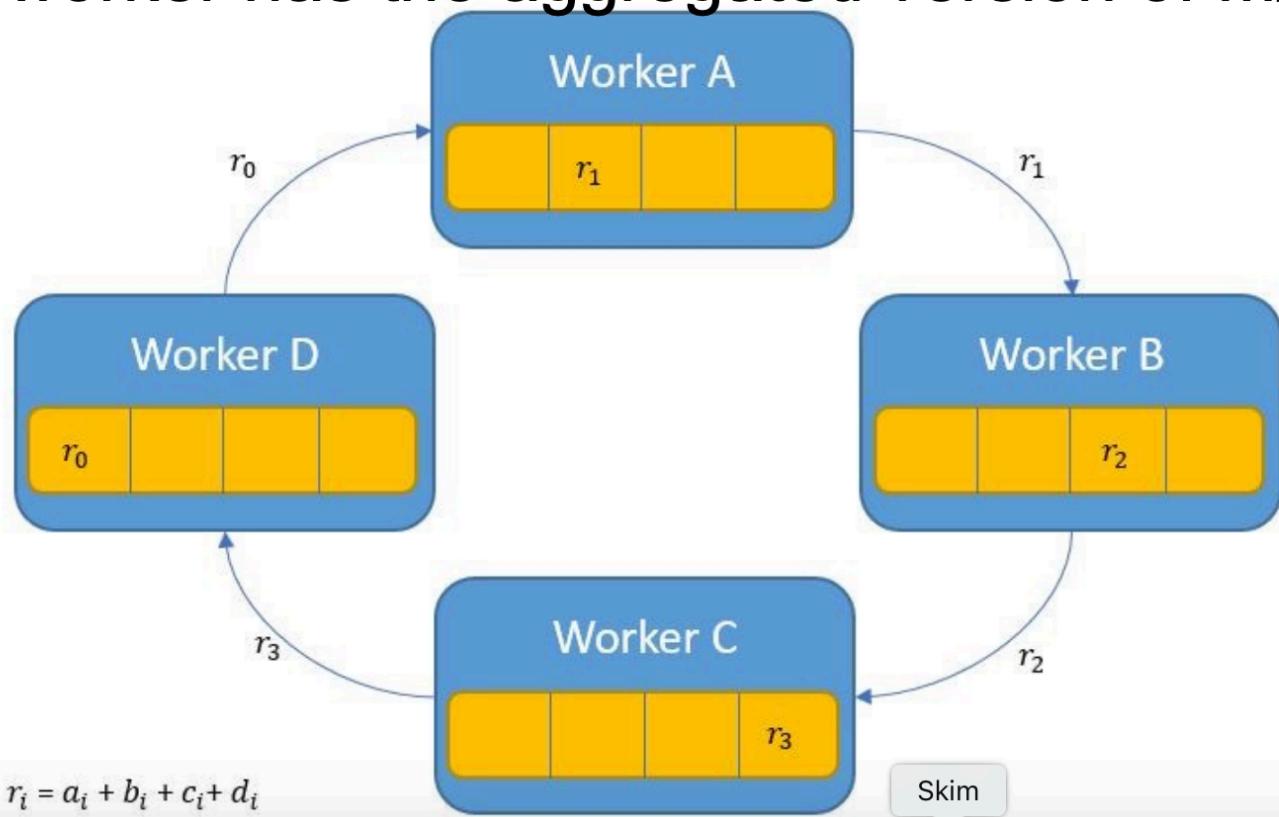
After It 1:



After It 2:

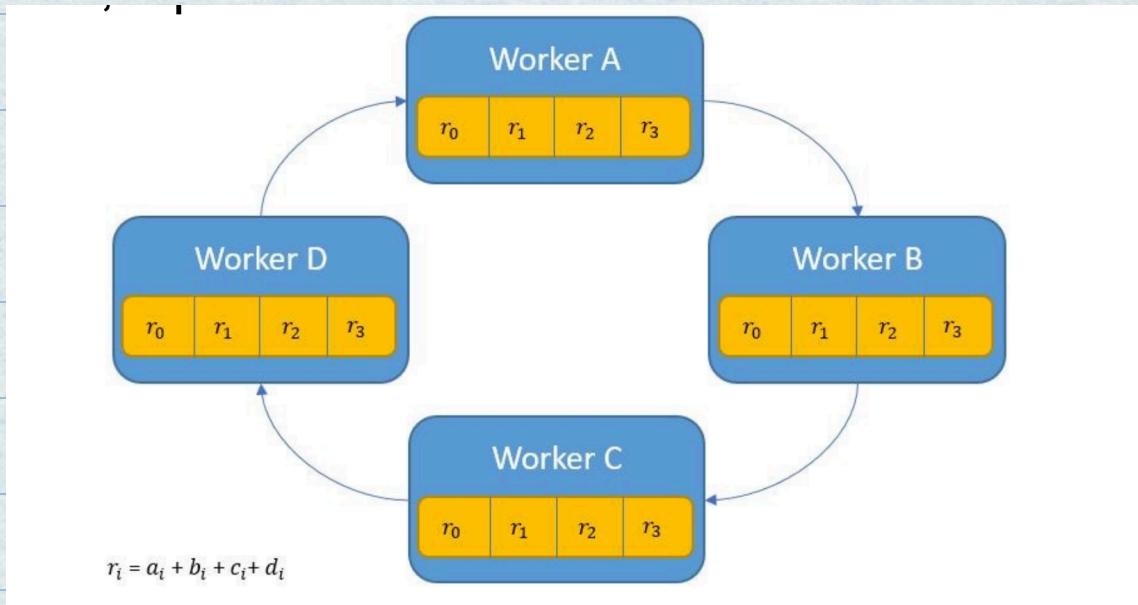


• After It 3 sum for each partition is  
on one of the workers



• Phase 2 of ring all reduce:

- Broadcast → each worker sends one slice of aggregated params to next worker; repeat  $N-1$  times



• what is total communication amount:

- discuss w/ neighbor

$$\rightarrow O(2 \cdot (N-1) \cdot \frac{M}{N} \cdot N)$$

→ 2 times  $N-1$  iterations

→ each it, a worker sends  $M/N$  data

\* TREE Allreduce:

- construct tree of  $N$  workers

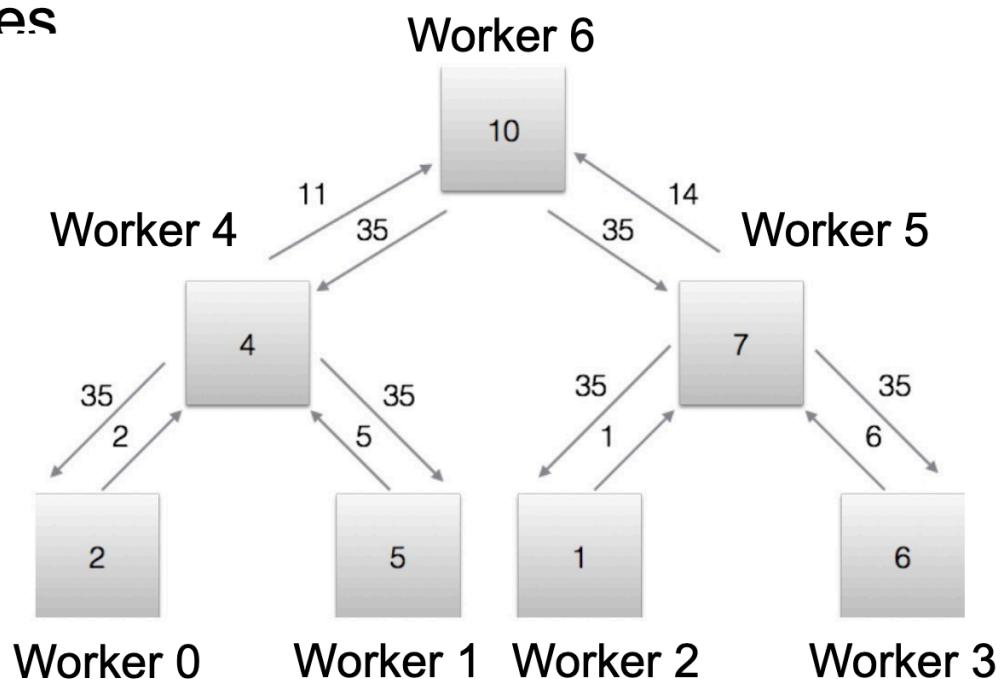
- step 1: aggregation:

→ each worker sends  $M$  params to its parent; repeat  $\log(N)$  times

- step 2: broadcast:

→ each worker sends  $M$  params to its children, repeat  $\log(N)$  times

18s



total comm.

→ amount

$2 \cdot M \cdot N$  → each worker sends

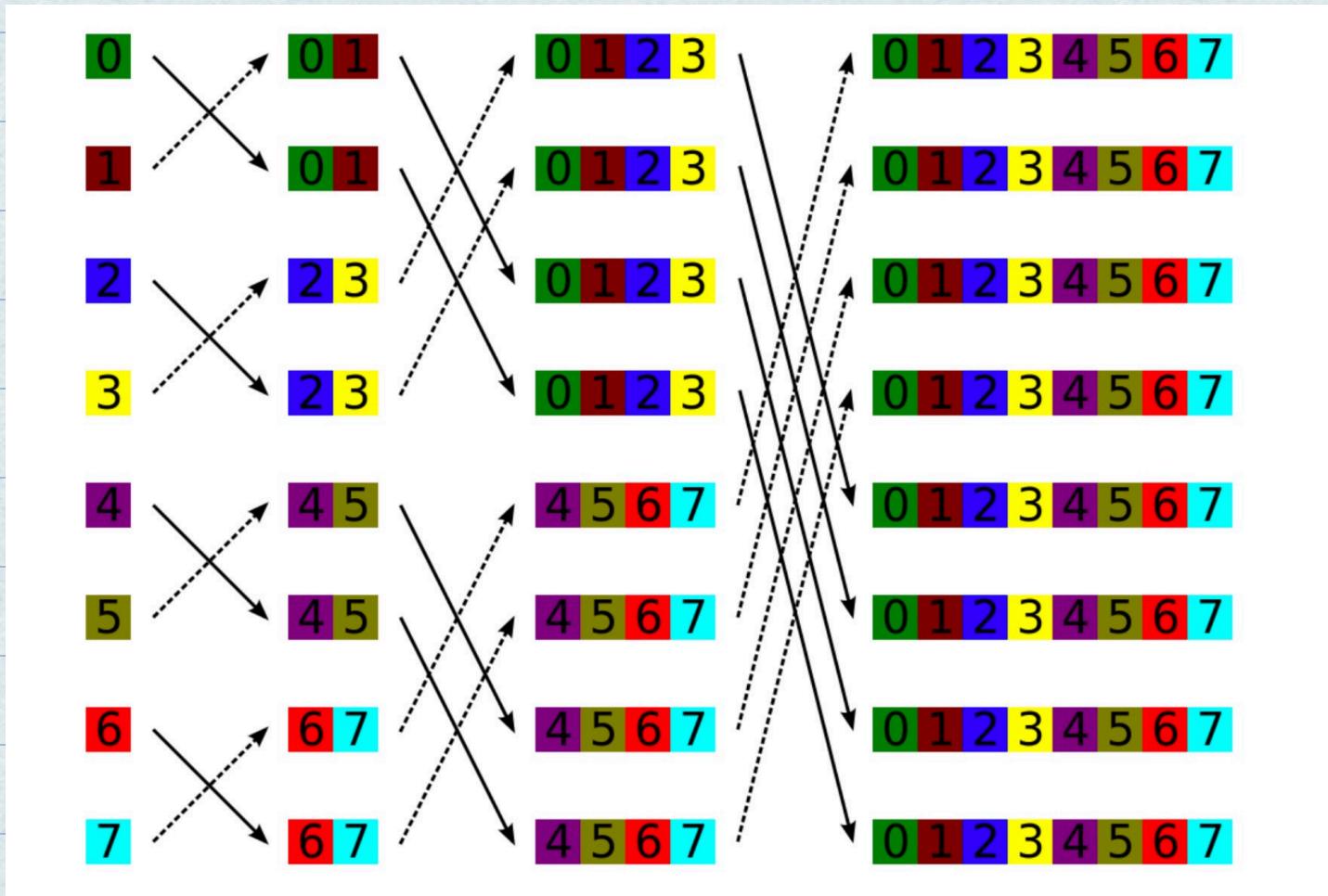
→ Agg:  $M \cdot N$   $N_{\text{cnt}}$  of data, once up once down

→ Broadcast:  $M \cdot N$

→ each worker requires ~3 wires

\* **Butterfly All-reduce** (I promise its the last one):

Construct butterfly network



→ Repeat  $\log(N)$  times:

- each worker sends  $M$  params to its target node in butterfly network
- each worker aggregates gradients locally

→ comm. amount:  $O(M \cdot N \cdot \log(N))$

→ each worker requires  $\sim \log N$  wires

Method	Param server	Ring	Tree	Butterfly
# wires of bond width $B$	<ul style="list-style-type: none"> <li>- each worker has <math>1</math> wires to PS</li> <li>- PS has <math>N</math> wires to each worker of <math>B</math></li> </ul>	- each worker has $2$ wires of BW $B$	- each worker has $3$ wires of BW $B$	- each worker has $\log N$ wires of BW $B$
# rounds/steps	$2$ (to PS and back)	$2 \cdot (N-1)$	$2 \cdot \log N$	$\log N$
amt of data each worker sends in a step	$M$ (each worker sends $1$ )	$\frac{M}{N}$	→ not all workers sending at all steps	$M$
Time per step	$\frac{M}{B}$ (for broadcast cost, assume PS can broadcast in parallel)	$\frac{M}{N \cdot B}$	$\frac{M}{B}$ (can broadcast or args. in parallel)	$\frac{M}{B}$
Agg data sent: multi-ply row $\exists$ by # of workers	$2 \cdot M \cdot N$	$2(N-1) \cdot \frac{M}{N} \cdot N$ $= 2(N-1) \cdot M$	$2 \cdot M \cdot N$ (up and down tree, each worker sends $N$ data)	$(MN \log N)$
Total Time =	$\frac{2 \cdot M}{B}$	$\frac{2(N-1) \cdot M}{N \cdot B}$	$\frac{2 \log N \cdot M}{B}$	$\frac{\log(M)}{B}$

# of "rounds"  
times time  
per round

Q: does tree or ring scale better? (both send  $\tilde{n}$  same amt of ags. data):

- ring appears to scale w/  $(N-1)/N$
- while tree w/  $\log N$
- in reality:  $B$  is really a function of message size (due to fixed overhead):

$$B = f_B(\text{messagelength})$$

so time per step:

$$\frac{\text{ring } F_B\left(\frac{m}{N}\right)}{\frac{m}{N \cdot f_B(m/N)}} \ll \frac{\text{tree } f_B(m)}{f_B(N)}$$

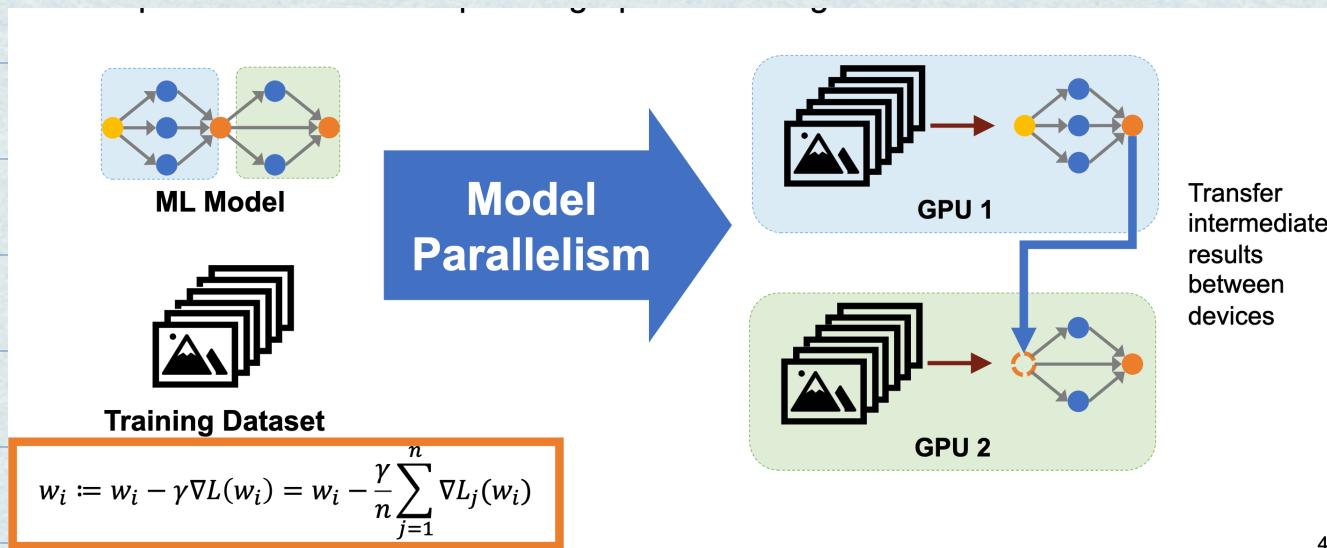
so for large  $N$ :

$$\rightarrow \text{total time} \frac{2(N-1) \cdot m}{N \cdot f_B(m/N)} \gg \frac{2 \cdot \log N \cdot m}{f_B(m)}$$

(time per step · # steps)

↳ hence, tree is preferred

# Model Parallelism



4!

- You will implement layer-by-layer model parallelism (partition the model by layers)