

Parallelism

Extracting more parallelism is arguably main goal
of most systems research for the last decade.

(Pervasive Parallelism Lab @ Stanford)

- ① THE MAIN TRENDS (from 5 years ago...) (Slides)
- ② WHAT A MODERN TRAINING SETUP LOOKS LIKE
- ③ DATA, MODEL & PIPELINE PARALLEL

Goal: KEEP IT HIGH LEVEL, GIVE YOU A TASTE

Disclaimer: I WAS ACTIVE IN THIS AREA IN RESEARCH

→ BIASED SLIDES

HAVE COFOUNDED, FUNDED COMPANIES IN
SPACE

⇒ TO THE SLIDES!

WHAT DOES A MODERN LARGE TRAINING SETUP LOOK LIKE? (TEAMS)



Single GPU, TPU, APU WHATEVER

5 YEARS AGO, WE WERE EXCITED ABOUT
125 TFLOP → 1000 TFLOP (1 PF)
(this year, HALF precision)

Largest supercomputer IN CANADA IS 5.6 PF

Germany IS 44 PF

NOTE: NOT comparable (Built on 5yo ten,

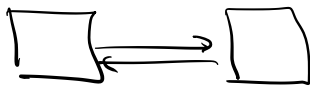
⇒ SPECIALIZATION ≠ RAPID GROWTH (LARGE #s)

⇒ ALSO CREATING: MANY GPUS CONTAIN MULTIPLE CHIPS INSIDE (SOON)

⇒ CONTAIN LOTS OF HIGH BANDWIDTH MEMORY (16, 32, ...) (HBM) (2016 AND)

MANY GPUS → Box. (8 or 16) NVIDIA cards DLX.

→ PLUG INTO PCI-E



↑ ~60 GB/s (4 TB/s with HBM)

PER LAKE.

16 LANES ~ 1 TB/s (WHOLE DEVICE)

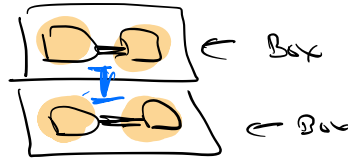
NVLink 200 G/s
L1 P2P network.

Boxes connected INTO RACKS. (4 PER RACK)

RACKS connected INTO "BAYS" → Datacenter depend.

TERMS

CHIP } many (8-16x)
Box } many (4x)
RAIC } many
POD



Key CHALLENGES

- ① Compute (How much can you use?) ~ 10-15% util!
- ② Memory (Does your model fit 175B → Spill out)
- ③ NETWORK (loosely at all levels)

WE'VE CODESIGNED MODELS TO EXTRACT PERFORMANCE.

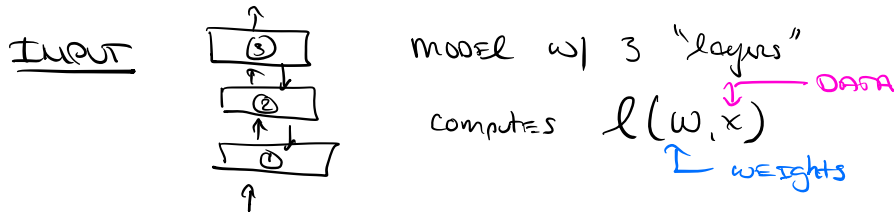
EXAMPLE Matrix multiply is highly parallel
⇒ CAN get 80% utilization or more!
→ TEAMS of folks took YEARS for attention.
"typical code" runs at 10-15% out of the box.

MAJOR TECHNIQUES

→ DATA Parallel

→ MODEL \neq Pipeline Parallel

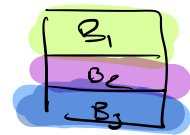
DATA PARALLEL



AND 3 NODES (CHIPS, BOXES, WHATEVER) that can for compute

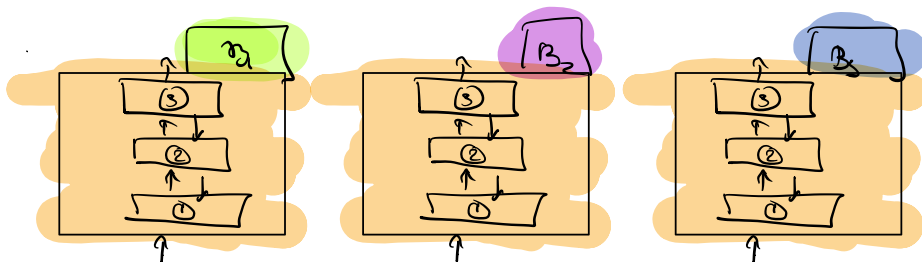


Goal
Compute a Batch of data $B = \sum_{x \in B} l(w, x)$



AND IT'S GRADIENT wrt w .

$$\frac{\partial L(w, B)}{\partial w} = \sum_{x \in B} \frac{\partial l(w, x)}{\partial w}$$



Straightforward! $L(w, B) = L(w, B_1) + L(w, B_2) + L(w, B_3)$

$$\frac{\partial L}{\partial w}(w, B) = \frac{\partial L}{\partial w}(w, B_1) + \frac{\partial L}{\partial w}(w, B_2) + \frac{\partial L}{\partial w}(w, B_3)$$

1. SEND BATCHES
2. EACH COMPUTE [IN PARALLEL]
3. SEND BACK GRADIENT (to ONE NODE OR HOST)
4. UPDATE MODEL WEIGHTS (take a step)
5. BROADCAST WEIGHTS BACK TO EACH WORKER

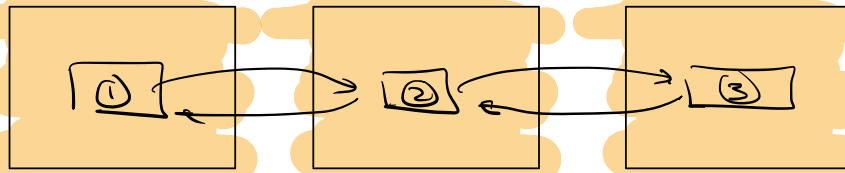
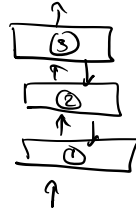
Utilization very high if NODE IS (AND) DOING WORK IN BATCHES

Communication

Entire model sent 3 times (+ 3 times for gradient)

Requirement: Model fits in a node. GDS4?! NO CHANCE!

Model "Parallelism"



Only each layer AND IT'S WORKING MEAN need to fit.

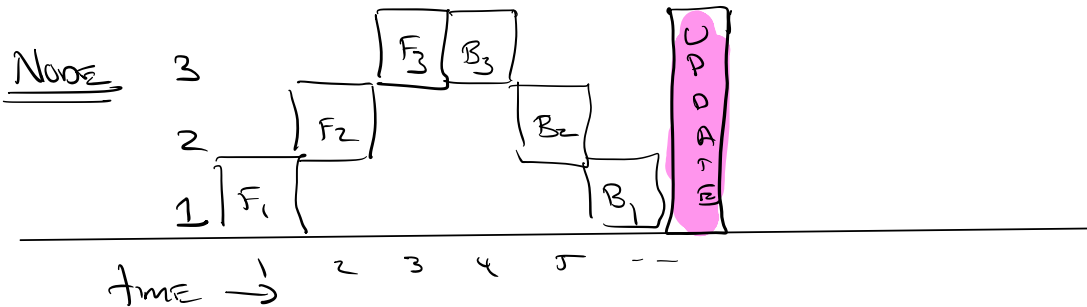
Recall "Forward and Back"

$$z_1 = f_1(w_1, x)$$

$$z_2 = f_2(w_2, z_1) \quad \text{FORWARD PASS} \quad \underline{\text{EASY to compute}}$$

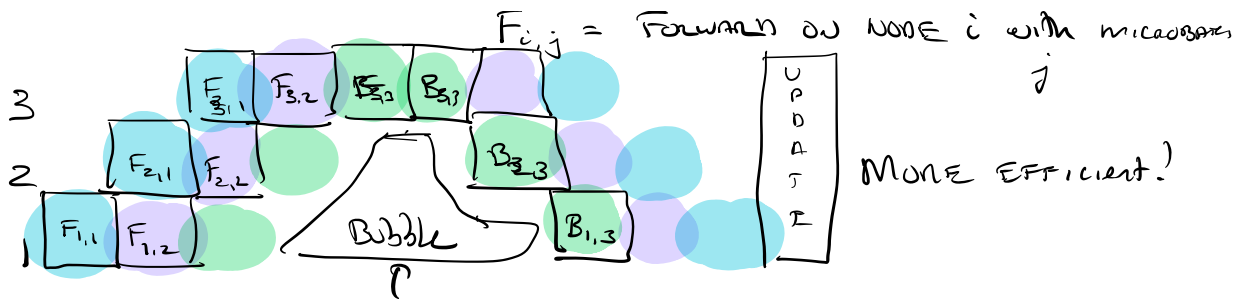
$$z_3 = f_3(w_3, z_2)$$

Backward $\frac{\partial z_3}{\partial w_2} = \frac{\partial f_3(w_3, z_2)}{\partial z_2} \cdot \frac{\partial z_2}{\partial w_2}$
 (example) \hookrightarrow NOT COMPUTED UNTIL FORWARD PASS



NOT PARALLEL AT ALL!

IDEA MICROBATCHES Smaller batches to overlap (PIPE)



WHITE SPACE IS WASTED 12 WHITE SQUARES

$3 \cdot 3 \cdot 2$ 18 SQUARES w/ WORK

$\Rightarrow \frac{18}{30}$ PER NODE UTILIZATION

$\approx 60\%$

\Rightarrow PIPELINE, PIPELINE SLOW HOW TO IMPROVE.

Key IDEA: ONE-STATE, HIGHER!

Stu Optimizations

→ COMMUNICATIONAL

→ REDUCED PRECISION