

MULTI-GPU FINE-TUNING

WITH DDP AND FSDP

Trelis Research

OVERVIEW

1. WHAT KIND OF MULTI-GPU TRAINING SHOULD I DO?
2. FIGURING OUT VRAM REQUIREMENTS.
3. TRAINING TYPES:
 - A. MODEL PARALLEL (MP)
 - B. DISTRIBUTED DATA PARALLEL (DDP)
 - C. FULLY SHARDED DATA PARALLEL (FSDP)
4. CODE MODIFICATIONS FOR EACH APPROACH.
5. DEMO FOR A) TINYLLAMA IN DDP AND B) CODE LLAMA 34B IN FSDP



WHAT GPU TRAINING TO DO?

1. WHAT IS MY MODEL SIZE?
2. WHAT ACCURACY DO I WANT?

- A. BEST - FULL-FINE TUNING
- B. GOOD - LORA
- C. OK - QUANTIZED LORA

VRAM REQUIREMENTS

3. DO I WANT TO:
 - A. MINIMIZE NUMBER OF GPUS / COST
 - B. MAXIMISE TRAINING SPEED

GPU SETUP
(MP, DDP, FSDP)



EFFECT OF *MODEL SIZE* ON VRAM REQUIREMENTS!

8 BITS = 1 BYTE

16 BIT NUMBER => 2 BYTES

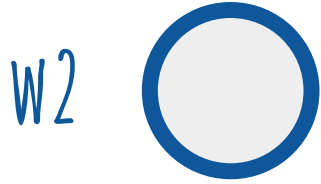
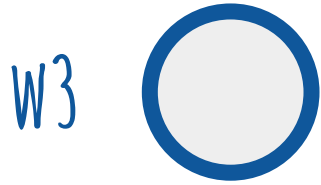


VRAM REQUIREMENTS!

= MODEL PARAMETERS (E.G. 7B)
+ GRADIENTS
+ OPTIMIZER STATES
+ ACTIVATIONS

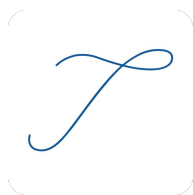


TOY MODEL (FORWARD)

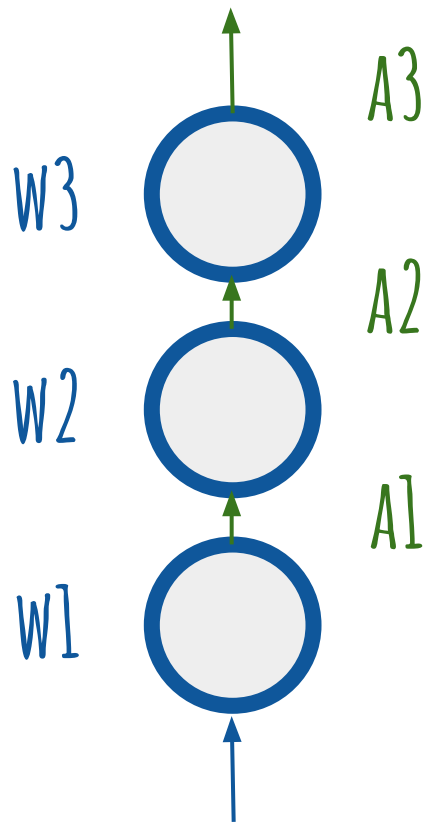


THREE "LAYERS"

ONE WEIGHT PER LAYER!



TOY MODEL (FORWARD)



THREE "LAYERS"

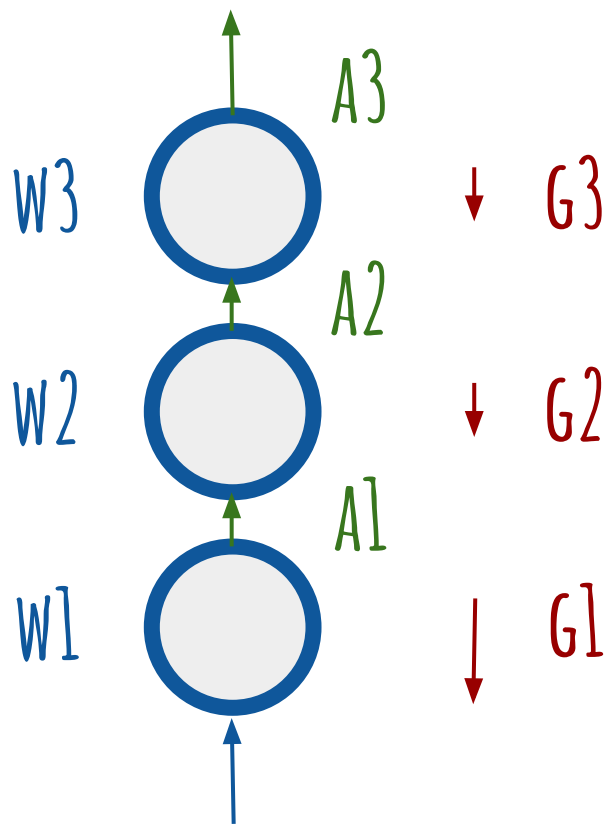
ONE WEIGHT PER LAYER!

OUTPUT/ACTIVATION = WEIGHT * INPUT

LOSS (a_3) = OUTPUT - REFERENCE

T

TOY MODEL (BACKWARD)



THREE "LAYERS"

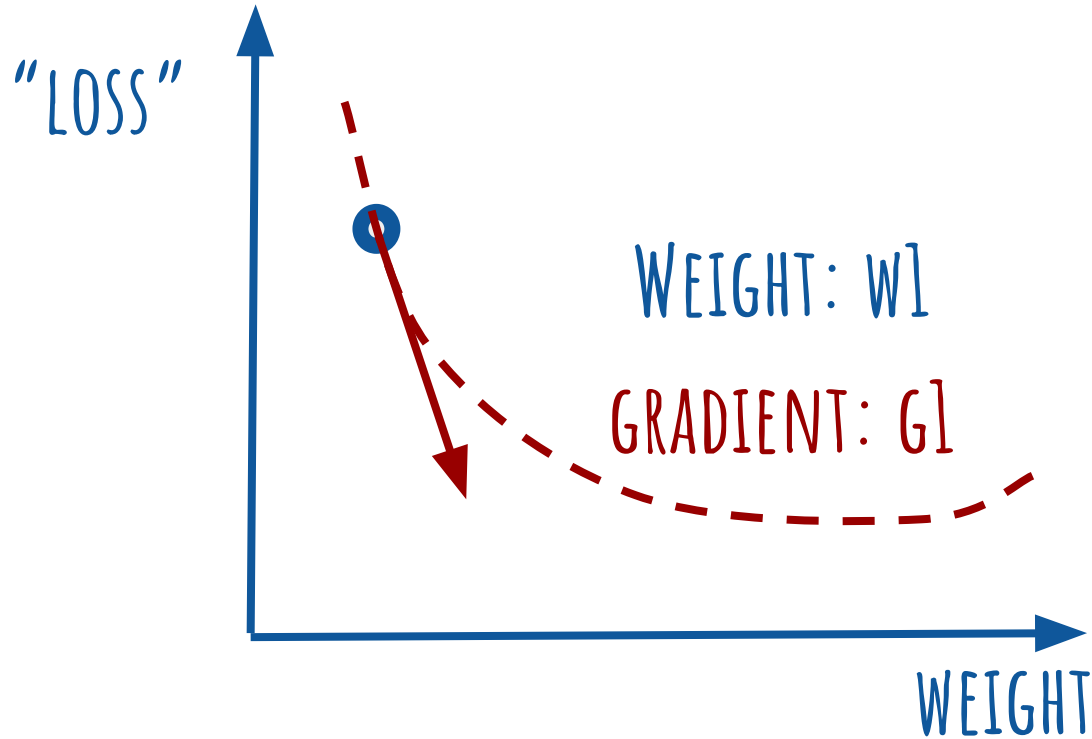
ONE WEIGHT PER LAYER!

ACTIVATION = WEIGHT * INPUT

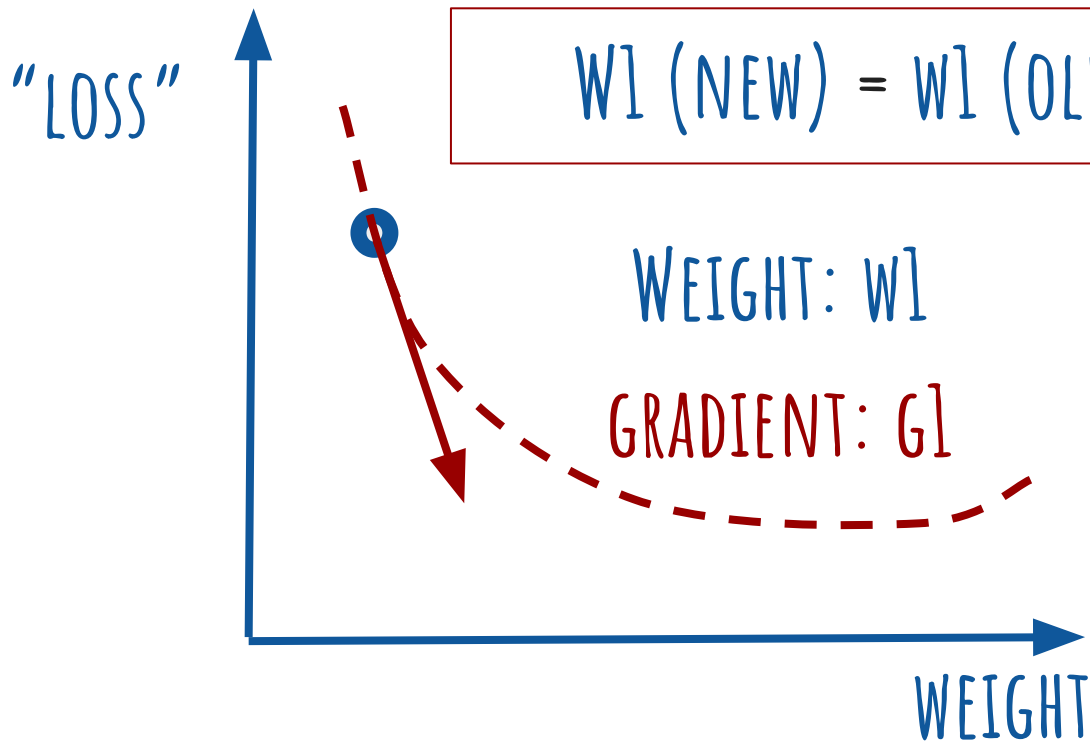
GRADIENTS (CHAIN RULE)

T

TOY MODEL

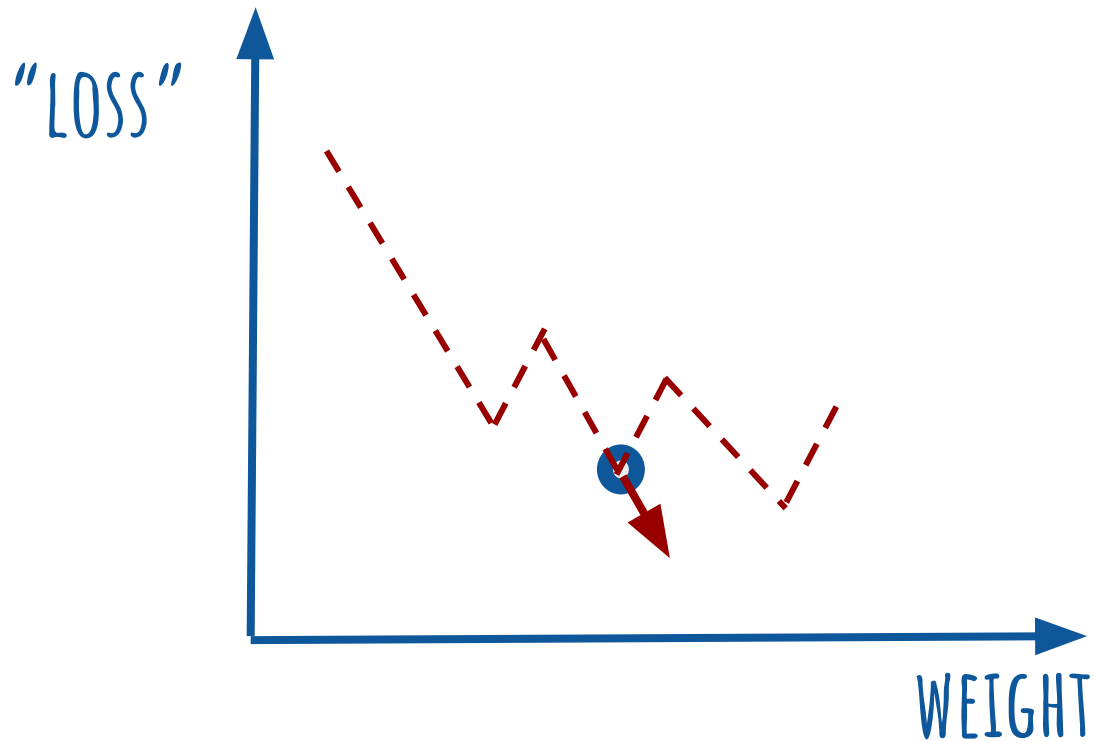


TOY OPTIMIZER



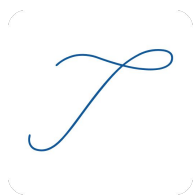
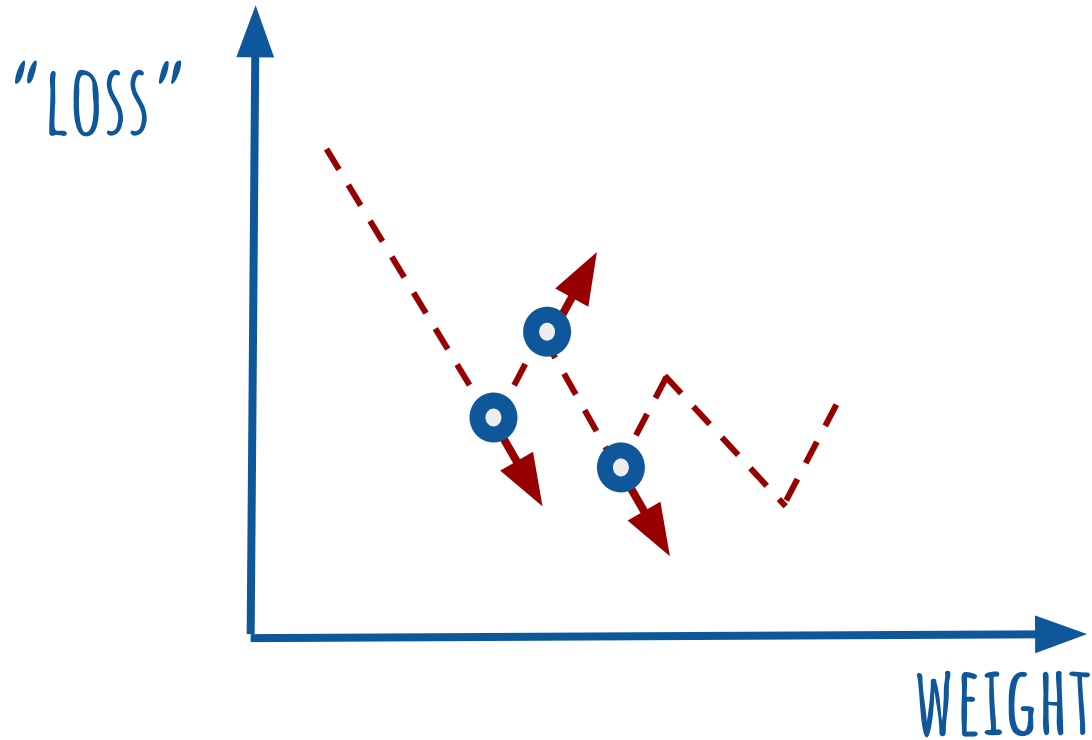
J

DON'T JUST USE THE GRADIENT



T

USE A HISTORICAL AVERAGE

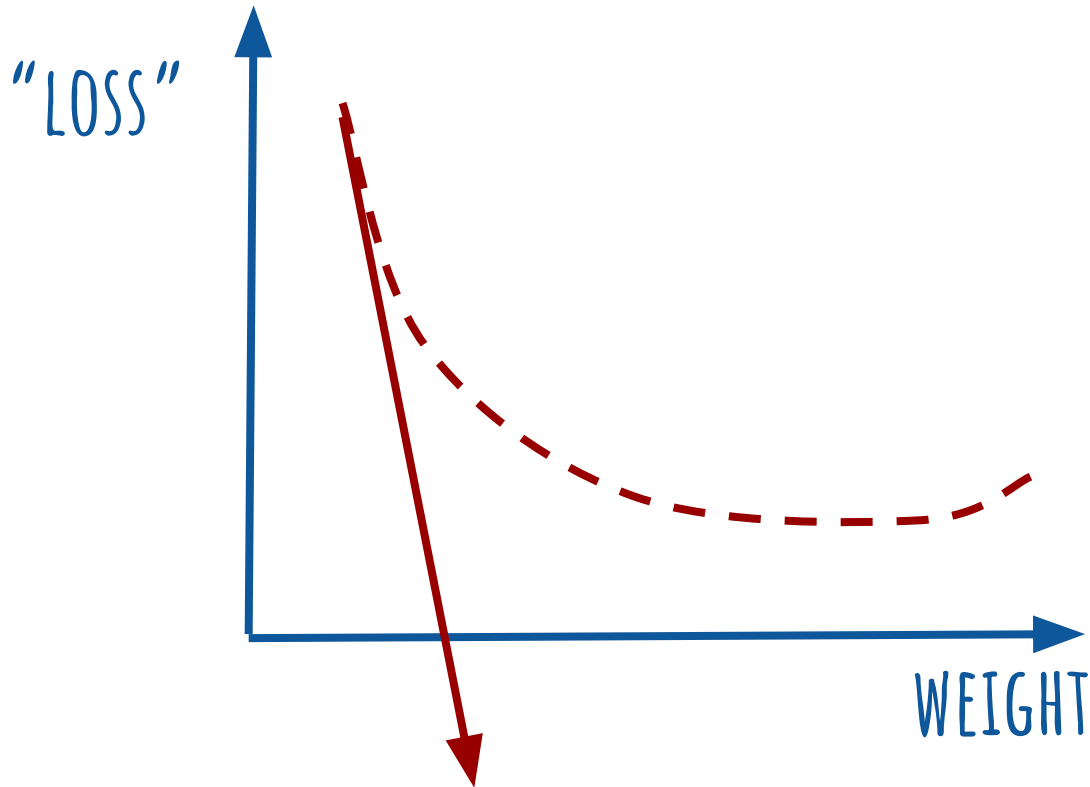


CONSIDER A HISTORICAL AVERAGE

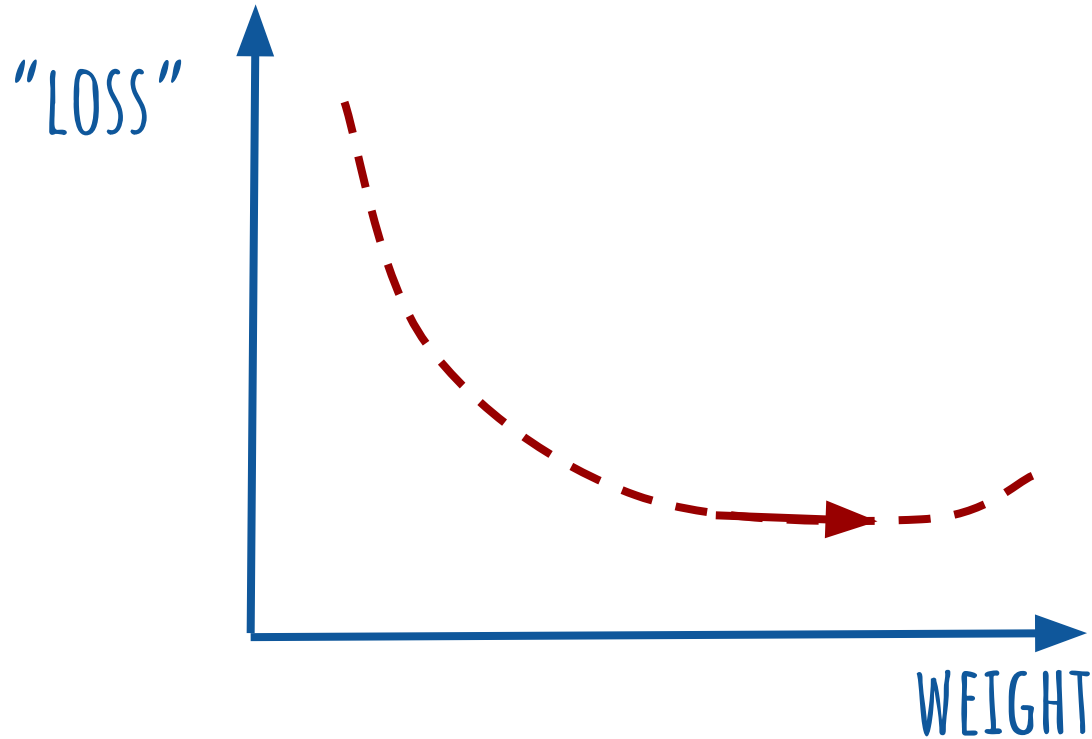
$$G = 1\% \cdot G(\text{NEW}) + 99\% \cdot G(\text{OLD})$$



DAMPEN LARGE GRADIENTS



AMPLIFY SMALL GRADIENTS



T

AMPLIFY SMALL, DAMPEN BIG

CONSIDER THE HISTORICAL VARIANCE

$$G^2 = 0.01\% \cdot G^2(\text{NEW}) + 99.9\% \cdot G^2(\text{OLD})$$

BIG VARIANCE \rightarrow DAMPEN

SMALL VARIANCE \rightarrow ACCELERATE



ACTUAL OPTIMIZER (ADAM)

1. AVERAGE THE GRADIENTS

$$G_M = 1\% \cdot G(\text{NEW}) + 99\% \cdot G(\text{OLD}).$$

2. USE THE VARIANCE TO DAMPEN/AMPLIFY

$$G_V^2 = 0.01\% \cdot G^2(\text{NEW}) + 99.9\% \cdot G^2(\text{OLD}).$$



ACTUAL OPTIMIZER (ADAM)

FOR EVERY TRAINABLE MODEL PARAMETER:

- "GRADIENT HISTORY" (MOMENTUM TERM)
- "GRADIENT VARIANCE HISTORY" (ADAPTIVE TERM)



MODEL + GRADS + OPTIMIZER = 16 BYTES

FOR EVERY TRAINABLE MODEL PARAMETER:

- PARAMETER
- GRADIENT
- MOMENTUM TERM
- ADAPTIVE TERM



MODEL + GRADS + OPTIMIZER = 16 BYTES

FOR EVERY TRAINABLE MODEL PARAMETER:

- PARAMETER [16 BITS = 2 BYTES]
- GRADIENT [16 BITS = 2 BYTES]
- MOMENTUM TERM [32 BITS = 4 BYTES]
- ADAPTIVE TERM [32 BITS = 4 BYTES]
- PARAMETER IN 32 BITS [4 BYTES]



MODEL (16BIT)

PARAMETER [16 BITS = 2 BYTES]

+ GRADIENT [16 BITS = 2 BYTES] $\times 1/64$

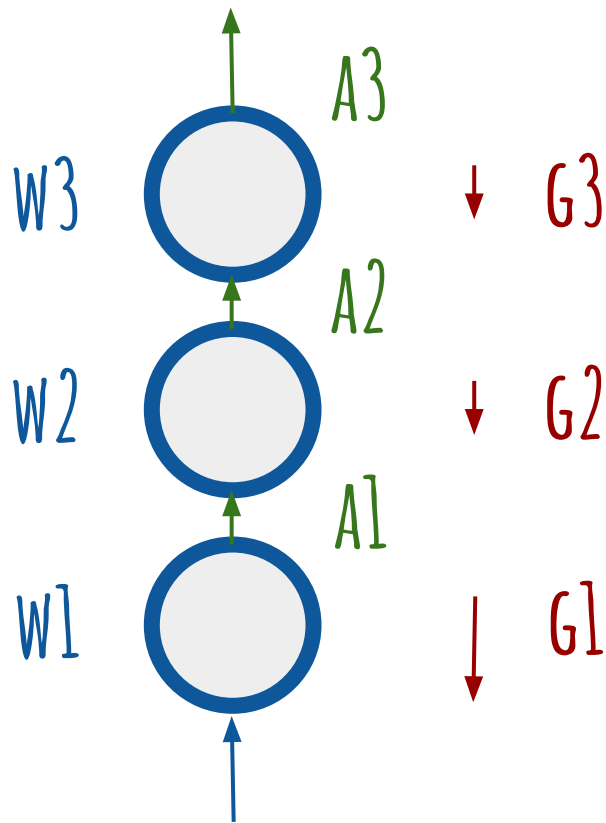
+ MOMENTUM TERM [32 BITS = 4 BYTES] $\times 1/64$

+ ADAPTIVE TERM [32 BITS = 4 BYTES] $\times 1/64$

+ PARAMETER IN 32 BITS [4 BYTES] $\times 1/64$

+ ACTIVATIONS!!! (SCALES WITH CONTEXT AND BATCH SIZE) 

ACTIVATIONS SIZE...



$$\text{ACTIVATION} = \text{WEIGHT} * \text{INPUT}$$

$$(\text{INPUT} = \text{BATCH SIZE} * \text{SEQ LENGTH})$$

T

GPU SETUP PRINCIPLE

YOUR GPUS MUST FIT THE MODEL + GRADS + OPTIMIZER
+ AT LEAST ONE BATCH!

TIP:

1. START WITH ENOUGH GPUS TO FIT MODEL + GRADS + OPTIMIZER
2. START WITH BATCH = 1 AND INCREASE TIL OOM.



WHAT GPU TRAINING TO DO?

1. WHAT IS MY MODEL SIZE?
2. WHAT ACCURACY DO I WANT?

- A. BEST - FULL-FINE TUNING
- B. GOOD - LORA
- C. OK - QUANTIZED LORA

VRAM REQUIREMENTS

3. DO I WANT TO:
 - A. MINIMIZE NUMBER OF GPUS / COST
 - B. MAXIMISE TRAINING SPEED

GPU SETUP
(MP, DDP, FSDP)



MODEL (4BIT)

PARAMETER [4 BITS = 0.5 BYTES]

+ GRADIENT [16 BITS = 2 BYTES] $\times 1/64$

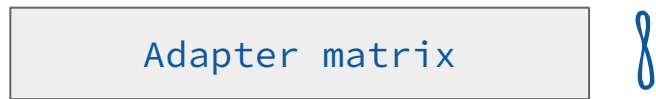
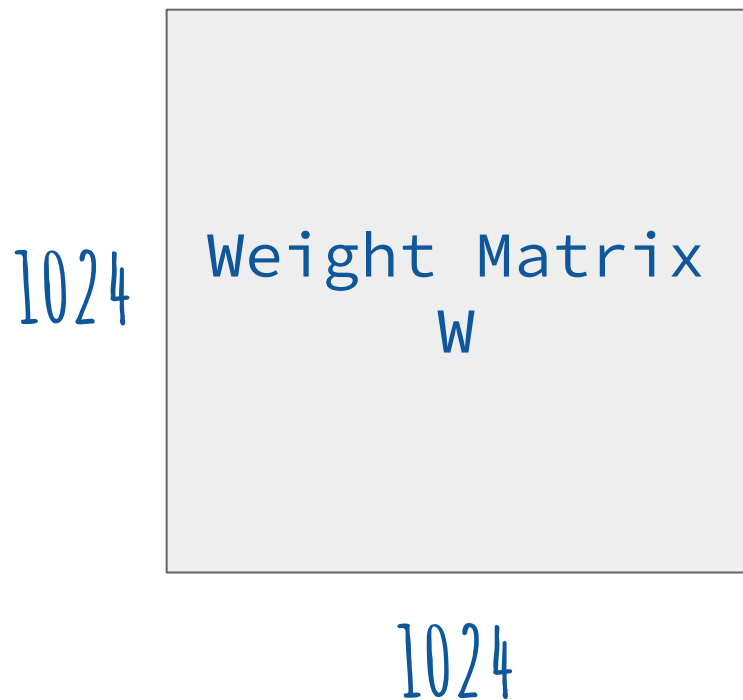
+ MOMENTUM TERM [32 BITS = 4 BYTES] $\times 1/64$

+ ADAPTIVE TERM [32 BITS = 4 BYTES] $\times 1/64$

+ PARAMETER IN 32 BITS [4 BYTES] $\times 1/64$

+ ACTIVATIONS!!! (SCALES WITH CONTEXT AND BATCH SIZE) 

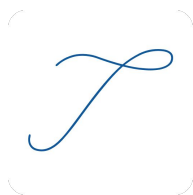
FREEZING THE CORE MODEL + TRAINING LOW RANK ADAPTERS (LoRA)



$$1024 \times 8 = 8192$$



1024



LORA SAVINGS

$$(2 * 8 * 1024) / (1024 * 1024)$$

= 1/64 PARAMETERS TO TRAIN!!!



MODEL + GRADS + OPTIMIZER W/ LORA

PARAMETER [16 BITS = 2 BYTES]

+ GRADIENT [16 BITS = 2 BYTES] $\times 1/64$

+ MOMENTUM TERM [32 BITS = 4 BYTES] $\times 1/64$

+ ADAPTIVE TERM [32 BITS = 4 BYTES] $\times 1/64$

+ PARAMETER IN 32 BITS [4 BYTES] $\times 1/64$

+ ACTIVATIONS



MODEL + OPTIMIZER W/ LORA = 2 BYTES

PARAMETER [16 BITS = 2 BYTES]

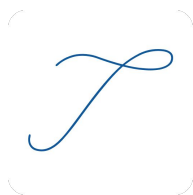
.

.

.

.

+ ACTIVATIONS



MODEL(4BIT) + OPTIM. + LORA = 0.5 BYTES

PARAMETER [5 BITS = 0.5 BYTES]

.

.

.

.

+ ACTIVATIONS

T

WHAT GPU TRAINING TO DO?

1. WHAT IS MY MODEL SIZE?
2. WHAT ACCURACY DO I WANT?

- A. BEST - FULL-FINE TUNING
- B. GOOD - LORA
- C. OK - QUANTIZED LORA

VRAM REQUIREMENTS

3. DO I WANT TO:
 - A. MINIMIZE NUMBER OF GPUS / COST
 - B. MAXIMISE TRAINING SPEED

GPU SETUP
(MP, DDP, FSDP)



WHAT GPU TRAINING TO DO?

1. WHAT IS MY MODEL SIZE?
2. WHAT ACCURACY DO I WANT?
 - A. BEST - FULL-FINE TUNING
 - B. GOOD - LORA
 - C. OK - QUANTIZED LORA

VRAM REQUIREMENTS

```
graph LR; Q1[1. WHAT IS MY MODEL SIZE?] --> VRAM[VRAM REQUIREMENTS]; Q2[2. WHAT ACCURACY DO I WANT?] --> VRAM; VRAM --> GPU[GPU SETUP (MP, DDP, FSDP)]; Q3[3. GPU SETUP: DO I WANT TO:] --> GPU;
```

3. GPU SETUP: DO I WANT TO:
 - A. MINIMIZE NUMBER OF GPUS / COST
 - B. MAXIMISE TRAINING SPEED

GPU SETUP
(MP, DDP, FSDP)

T

GPU SETUP

1. YOU CAN ALWAYS DO *MODEL PARALLEL (MP)*

2. DOES MODEL FIT ON A SINGLE GPU?

YES (SINGLE GPU)

MORE SPEED WITH

DISTRIBUTED DATA PARALLEL (DDP)

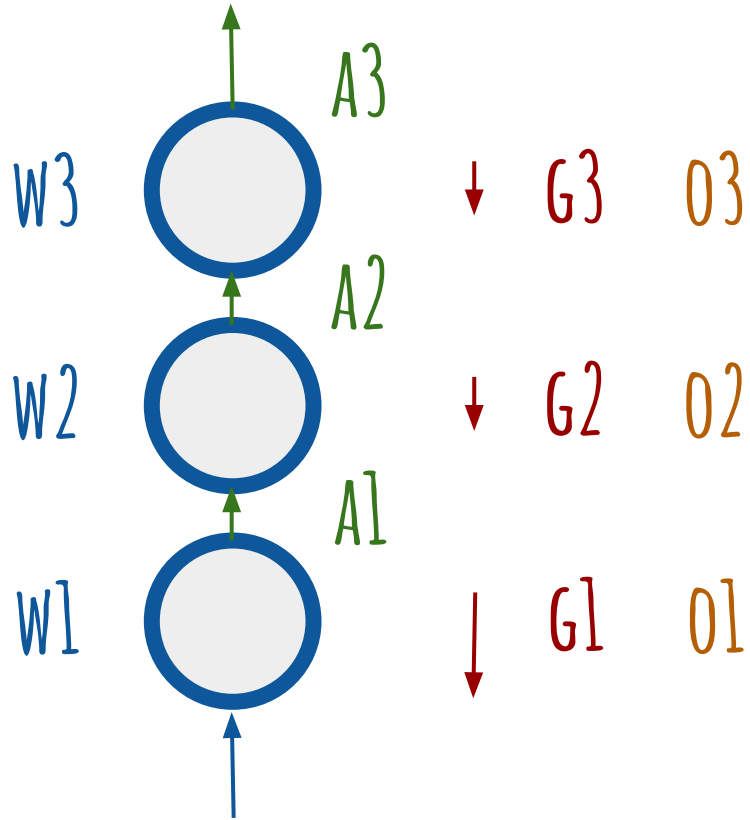
NO (NEED MULTIPLE GPUS)

MORE SPEED WITH

FULLY SHARDED DATA PARALLEL (FSDP)



TOY MODEL



WEIGHTS

GRADIENTS

OPTIMIZER STATES



IF MODEL FITS ON SINGLE GPU

GPU 0

W3, G3, O3

W2, G2, O2

W1, G1, O1

↑ BATCH = 1

WEIGHTS

GRADIENTS

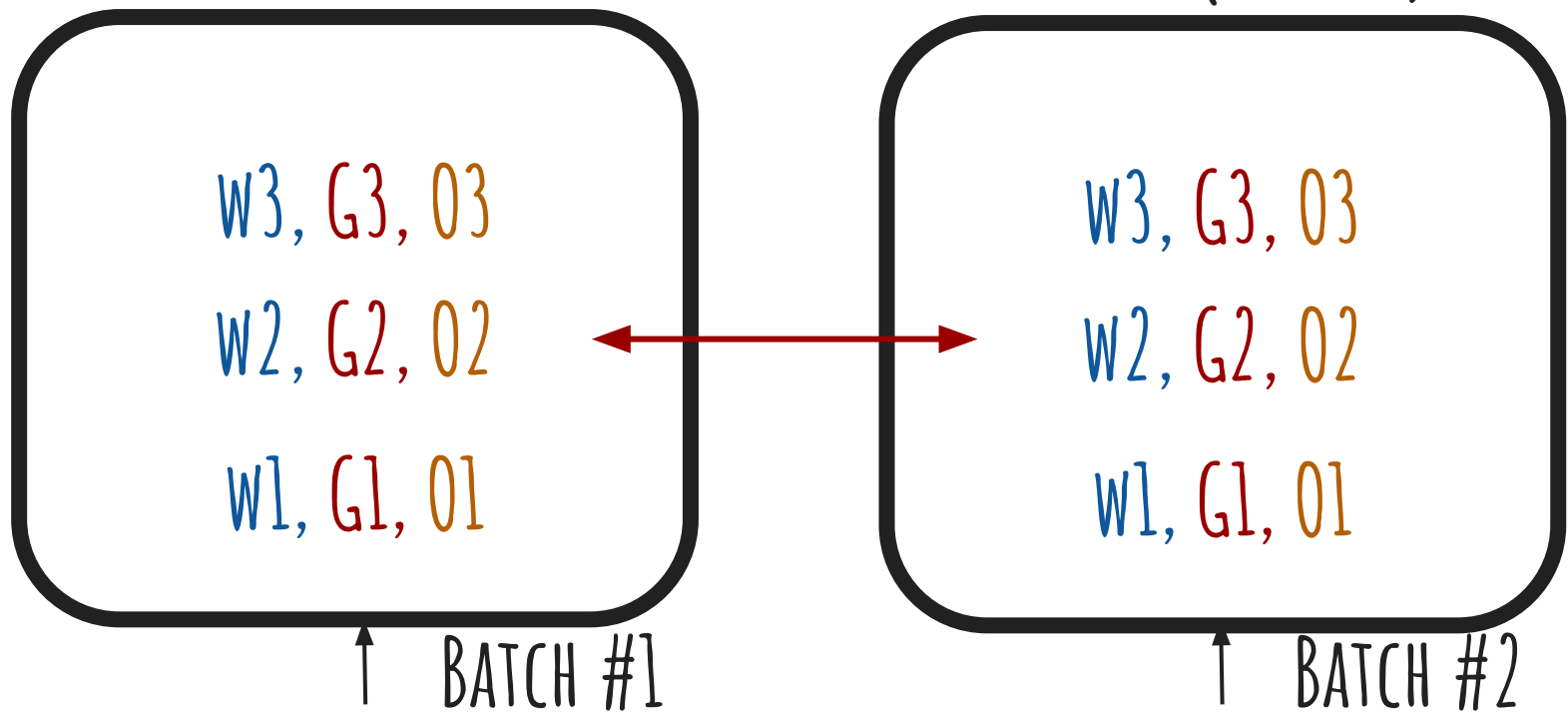
OPTIMIZER STATES

T

SINGLE GPU -> DATA PARALLEL

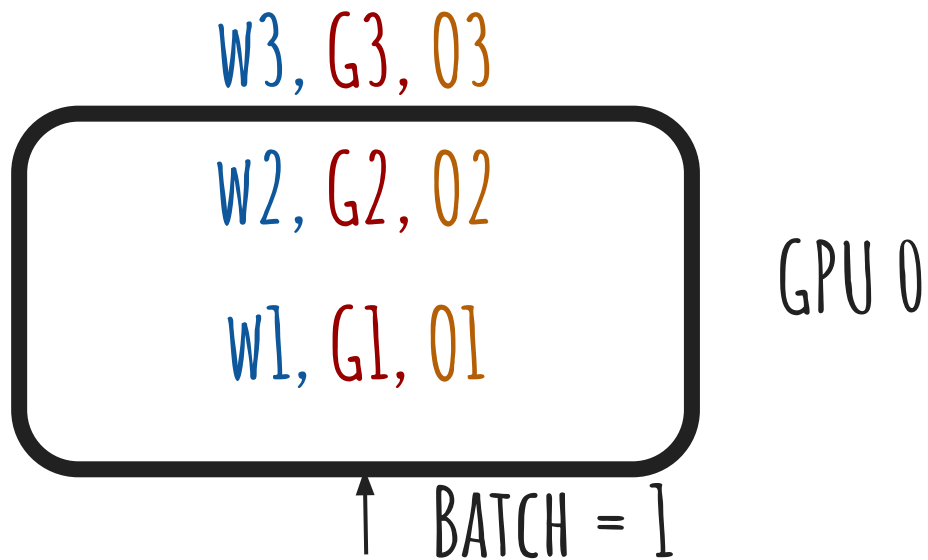
GPU 0

GPU 1 (REPLICA)



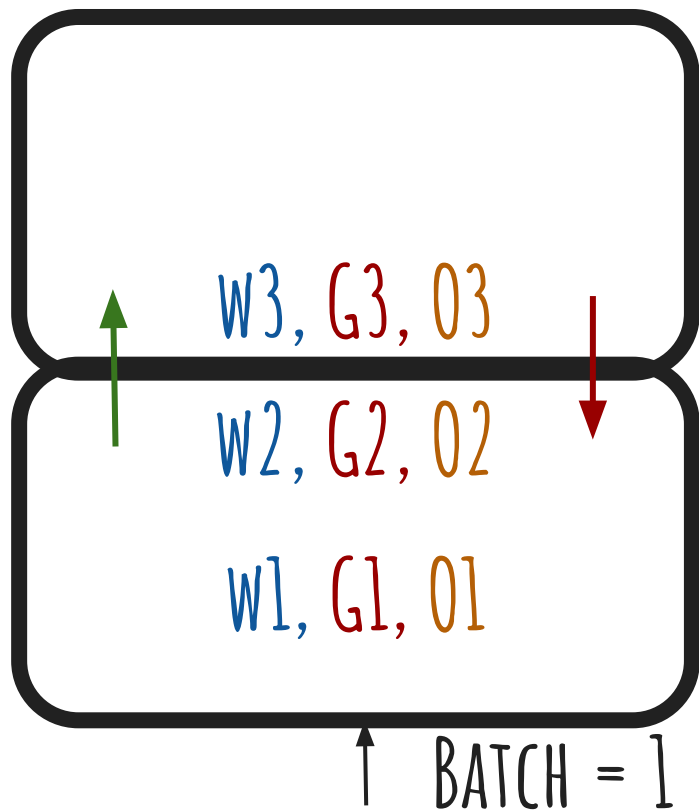
T

MODEL DOESN'T FIT



T

NAIVE MODEL PARALLEL



GPU 1

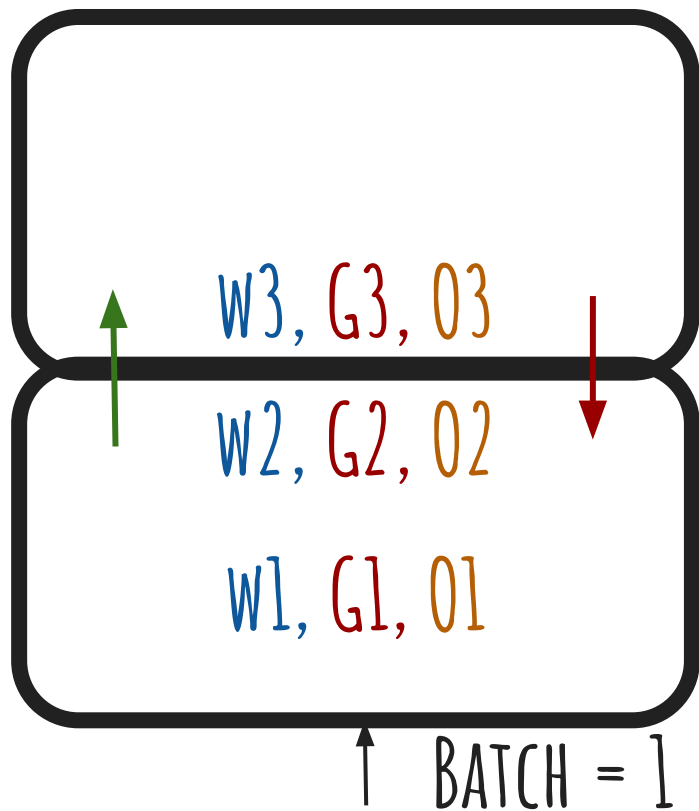
FORWARD PASS:
PASS ACTIVATIONS

GPU 0

BACK-PROP:
PASS GRADIENT INFO

T

NAIVE MODEL PARALLEL



PROBLEM:

~1 GPU ACTIVE AT A TIME
BAD "UTILIZATION"



FORWARD PASS: FSDP

GPU 0 (SHARD 0)

$w_3(0)$, $G_3(0)$, $o_3(0)$

$w_2(0)$, $G_2(0)$, $o_2(0)$

$w_1(0)$, $G_1(0)$, $o_1(0)$

↑
BATCH #1

GPU 1 (SHARD 1)

$w_3(1)$, $G_3(1)$, $o_3(1)$

$w_2(1)$, $G_2(1)$, $o_2(1)$

$w_1(1)$, $G_1(1)$, $o_1(1)$

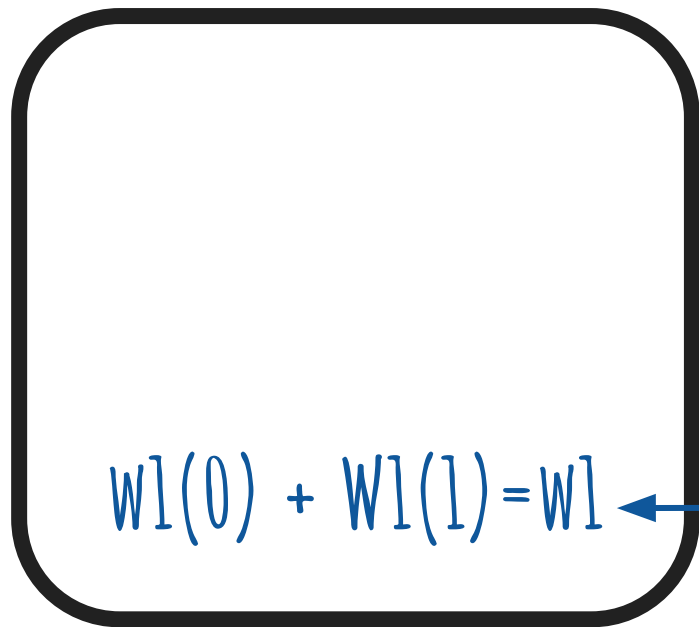
↑
BATCH #2



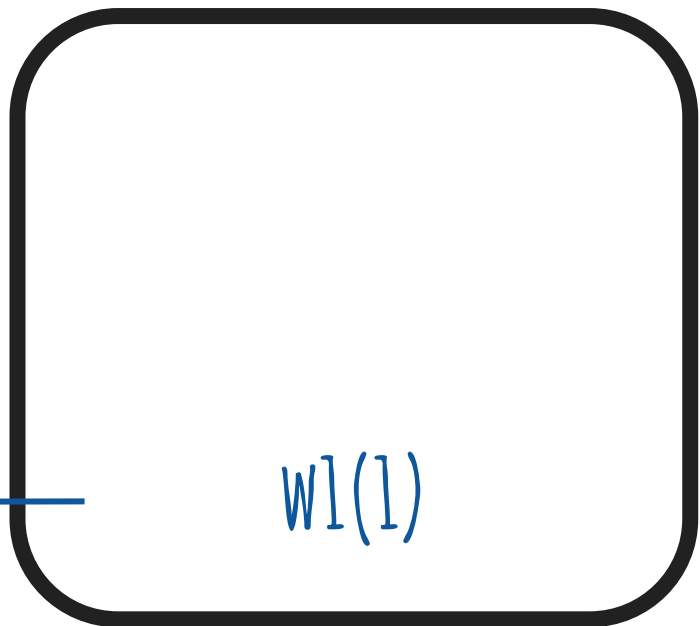
FORWARD PASS: FSDP

GPU 0 (SHARD 0)

GPU 1 (SHARD 1)



↑ BATCH #1



↑ BATCH #2

T

FORWARD PASS: FSDP

GPU 0 (SHARD 0)

GPU 1 (SHARD 1)

$$W2(0) + W2(1) = W2$$

$$W2(1)$$

↑ BATCH #1

↑ BATCH #2

T

FORWARD PASS: FSDP

GPU 0 (SHARD 0)

GPU 1 (SHARD 1)

$$w_3(0) + w_3(1) = w_3$$

$$w_3(1)$$

↑
BATCH #1

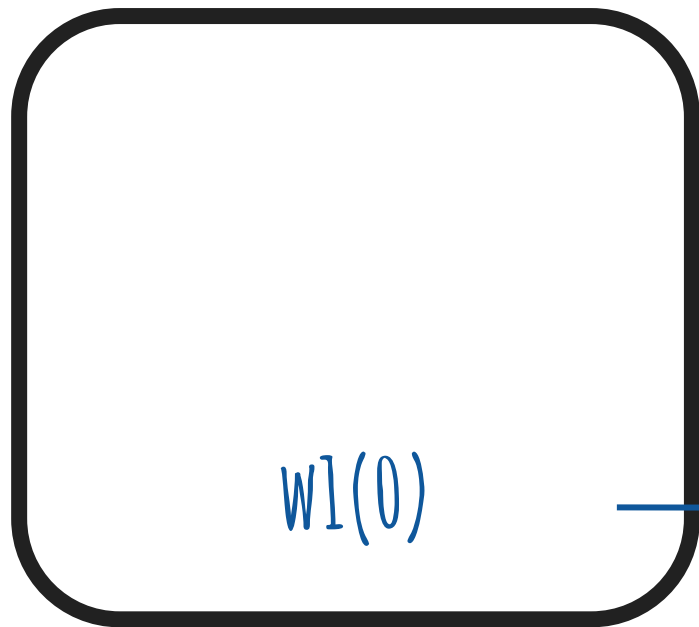
↑
BATCH #2

T

FORWARD PASS: FSDP

GPU 0 (SHARD 0)

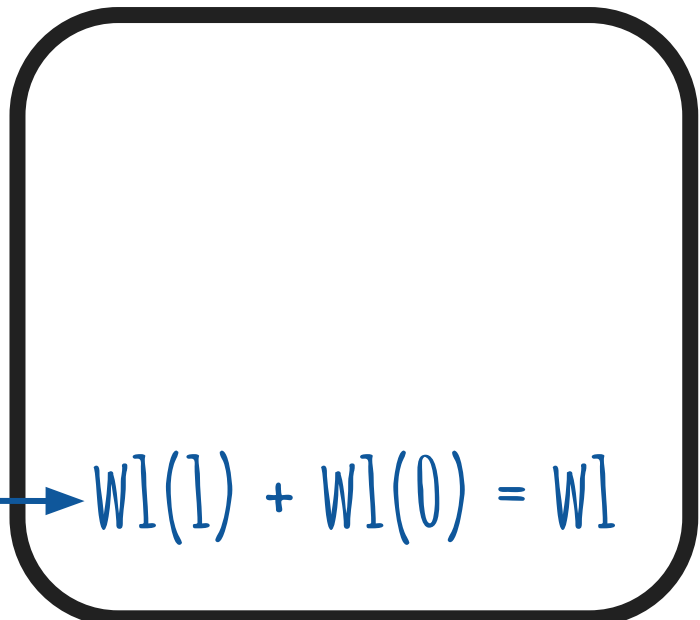
GPU 1 (SHARD 1)



$wl(0)$



BATCH #1



$wl(1) + wl(0) = wl$



BATCH #2

T

BACKWARD PROPAGATION: FSDP

GPU 0 (SHARD 0)

GPU 1 (SHARD 1)

$w_3(0)$, $G_3(0)$, $o_3(0)$

$w_2(0)$, $G_2(0)$, $o_2(0)$

$w_1(0)$, $G_1(0)$, $o_1(0)$

↑
BATCH #1



$w_3(1)$, $G_3(1)$, $o_3(1)$

$w_2(1)$, $G_2(1)$, $o_2(1)$

$w_1(1)$, $G_1(1)$, $o_1(1)$

↑
BATCH #2

J

BACKWARD PROPAGATION: FSDP

GPU 0 (SHARD 0)

GPU 1 (SHARD 1)

$w_3(0)$, $G_3(0)$, $o_3(0)$

$w_2(0)$, $G_2(0)$, $o_2(0)$

$w_1(0)$, $G_1(0)$, $o_1(0)$

$w_3(1)$, $G_3(1)$, $o_3(1)$

$w_2(1)$, $G_2(1)$, $o_2(1)$

$w_1(1)$, $G_1(1)$, $o_1(1)$

↑ BATCH #1

↑ BATCH #2

T

BACKWARD PROPAGATION: FSDP

GPU 0 (SHARD 0)

$w_3(0)$, $G_3(0)$, $o_3(0)$

$w_2(0)$, $G_2(0)$, $o_2(0)$

$w_1(0)$, $G_1(0)$, $o_1(0)$

↑
BATCH #1

GPU 1 (SHARD 1)

$w_3(1)$, $G_3(1)$, $o_3(1)$

$w_2(1)$, $G_2(1)$, $o_2(1)$

$w_1(1)$, $G_1(1)$, $o_1(1)$

↑
BATCH #2



T

COMPARING APPROACHES

MODEL PARALLEL (MP)

LOW COMMUNICATION OVERHEAD

LOW UTILIZATION FOR MULTI-GPU

DISTRIBUTED DATA PARALLEL (DDP)

MODEST COMMUNICATION OVERHEAD
(GRADIENTS)

FULLY SHARDED DATA PARALLEL (FSDP)

HIGH COMMUNICATION OVERHEAD
(WEIGHTS + GRADIENTS)



GPU SETUP

1. YOU CAN ALWAYS DO *MODEL PARALLEL (MP)*

2. DOES MODEL FIT ON A SINGLE GPU?

YES (SINGLE GPU)

MORE SPEED WITH

DISTRIBUTED DATA PARALLEL (DDP)

NO (NEED MULTIPLE GPUS)

MORE SPEED WITH

FULLY SHARDED DATA PARALLEL (FSDP)



A NOTE ON DEEPSPEED (MICROSOFT)

DEEPSPEED STAGE 1 \leftrightarrow PYTORCH FSDP

SHARD 01

DEEPSPEED STAGE 2 \leftrightarrow PYTORCH FSDP

SHARD G1, 01

DEEPSPEED STAGE 3 \leftrightarrow PYTORCH FSDP

SHARD W1, G1, 01



ACCELERATE

A WRAPPER AROUND LIBRARIES:

- SUPPORTS PYTORCH DDP + FSDP
 - SUPPORTS DEEPSPEED

CONFIGURE WITH 'ACCELERATE CONFIG'



RESULTS - TINYLLAMA 1.1B W/ LORA

4X A6000 ADA (48GB)

	GPU UTILIZATION	TRAINING TIME
MODEL PARALLEL	24%	3M53S
DDP	90%	59S

1024 ROWS, 1024 CONTEXT, BATCH 16, GA 2

T

RESULTS - CODELLAMA 34B W/ LORA

4X A6000 ADA (48GB)

	GPU UTILIZATION	TRAINING TIME
MODEL PARALLEL	~20%	47 MINS
FSDP	~100%	30 MINS

1024 ROWS, 1024 CONTEXT, BATCH 4, GA 4



K V CACHE

	key("The") · query ("quick")	The
The	quick	



K V CACHE

	key("The") · query ("quick")	key("The") · query ("brown")	The
		key("quick") · query ("brown")	quick
The	quick	brown	

