

# Efficient LLMs

Yatin Nandwani  
Research Scientist, IBM Research



Semester 1,  
2025-2026

Large Language Models: Introduction and Recent Advances

ELL881 · AIL821

- How to train big models on big data?
- What's in the GPU memory during training?
  - i. Model weights ; ii. Param gradients iii. Optim states iv. Activations

• What is the size of params / grads / optim states?

• What is the size of activations?

• How to reduce activation memory?

- Activation re-computation aka Gradient checkpointing

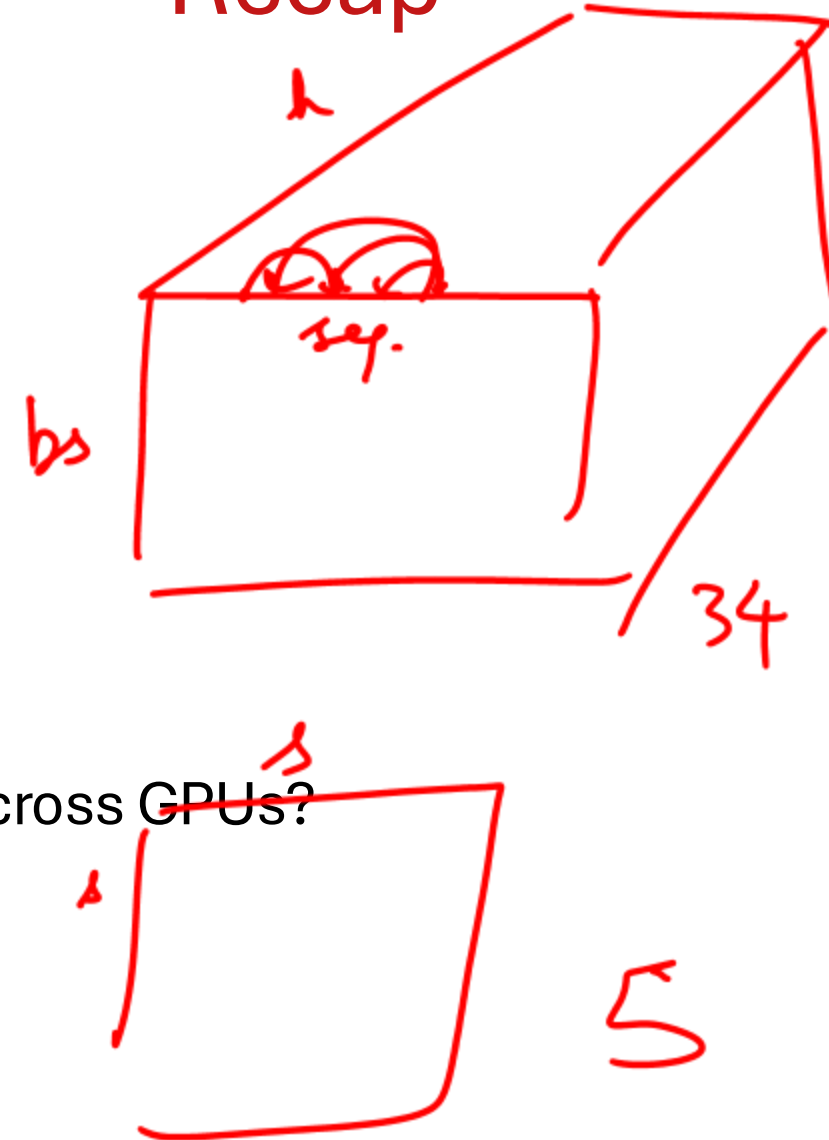
• How to increase batch size?

- Gradient accumulation (run fwd / bwd  $k$  times before `optim.step()`)

• Can we parallelize grad. Accumulation?

• Can we shard the optim. states, grads, and model params across GPUs?

## Recap



**Command:** `torchrun --nproc_per_node 2 train.py`

```
from torch.distributed.fsdp import fully_shard, FSDPModule
model = Transformer()
for layer in model.layers:
    fully_shard(layer)
fully_shard(model)
```

```
assert isinstance(model, Transformer)
assert isinstance(model, FSDPModule)
print(model)
# FSDPTransformer(
#   (tok_embeddings): Embedding(...)
#   ...
#   (layers): 3 x FSDPTransformerBlock(...)
#   (output): Linear(...)
# )
```

```
for _ in range(epochs):
    x = torch.randint(0, vocab_size, (batch_size, seq_len), device=device)
    loss = model(x).sum()
    loss.backward()
    optim.step()
    optim.zero_grad()
```

```
from torch.utils.data import DataLoader
```

```
train_loader = DataLoader(  
    dataset, ✓  
    batch_size=32,  
)
```

for batch in train\_loader:

\_\_\_\_\_

```
from torch.utils.data import DataLoader

train_loader = DataLoader(
    dataset,
    batch_size=32,
)
```

```
1  from torch.utils.data import DataLoader, DistributedSampler
2  import torch.distributed as dist
3
4  dist.init_process_group("nccl")
5  rank = dist.get_rank() ✓
6  world_size = dist.get_world_size()
7
8  train_sampler = DistributedSampler(
9      dataset, num_replicas=world_size, rank=rank,
10     shuffle=True, # shuffle at each epoch
11 )
12 train_loader = DataLoader(
13     dataset,
14     batch_size=local_batch_size, # per-GPU batch size
15     sampler=train_sampler,
16 )
```

huggingface.co/docs/trl/en/sft\_trainer

personal | IBM | A Butterfly Valley | ... | AIoT-MLSys-Lab/E... | Why is everyone u... | The Ultra-Scale Pl... | (46) How LLMs us...

• TRL ▾

Search documentation 🔍

V0.22.1 ▾ EN ▾ 15,313

- DDPO
- DPO
- Online DPO
- GKD
- GRPO
- KTO
- Nash-MD
- ORPO
- PPO
- PRM
- Reward
- RLOO
- SFT**
- Iterative SFT
- XPO

Model Classes

Model Utilities

Best of N Sampling

Judges

## SFT Trainer

All models SFT smol course Chapter 1

### Overview

TRL supports the Supervised Fine-Tuning (SFT) Trainer for training language models.

This post-training method was contributed by [Younes Belkada](#).

### Quick start

This example demonstrates how to train a language model using the [SFTTrainer](#) from TRL. We train a [Qwen 3 0.6B](#) model on the [Capybara dataset](#), a compact, diverse multi-turn dataset to benchmark reasoning and generalization.

```
from trl import SFTConfig, SFTTrainer
from datasets import load_dataset

trainer = SFTTrainer(
    model="Qwen/Qwen3-0.6B",
    train_dataset=load_dataset("trl-lib/Capybara", split="train"),
)
trainer.train()
```

trl

# SFTTrainer

**torchrun** is low-level  
PyTorch-native

**accelerate** is high-level  
and automates much of the  
distributed setup.

[https://huggingface.co/docs/trl/en/sft\\_trainer](https://huggingface.co/docs/trl/en/sft_trainer)

```

accelerate launch --config_file <path/to/acc/config> trl/scripts/sft.py \
  --model_name_or_path Qwen/Qwen2-0.5B \
  --dataset_name trl-lib/Capybara \
  --learning_rate 2.0e-5 \
  --num_train_epochs 1 \
  --per_device_train_batch_size 2 \
  --gradient_accumulation_steps 8 \
  --gradient_checkpointing \
  --eos_token '<|im_end|>' \
  --eval_strategy steps \
  --eval_steps 100 \
  --output_dir Qwen2-0.5B-SFT \

```

SFT Logis

$$y = Wx$$

$$\Delta W = \left( \frac{\partial \mathcal{L}}{\partial y} \right) x^T$$

# SFTTrainer

**torchrun** is low-level  
PyTorch-native

**accelerate** is high-level  
and automates much of the  
distributed setup.

[https://huggingface.co/docs/trl/en/sft\\_trainer](https://huggingface.co/docs/trl/en/sft_trainer)

[https://huggingface.co/docs/trl/main/en/distributing\\_training](https://huggingface.co/docs/trl/main/en/distributing_training)

<https://github.com/huggingface/trl/blob/main/trl/scripts/sft.py>

# Recap

- How to train big models on big data?
- What's in the GPU memory during training?
  - i. Model weights ; ii. Param gradients iii. Optim states iv. Activations
- What is the size of params / grads / optim states?
- What is the size of activations?
- How to reduce activation memory?
  - Activation re-computation *aka* Gradient checkpointing
- How to increase batch size?
  - Gradient accumulation ( run fwd / bwd  $k$  times before `optim.step()`)
- Can we parallelize grad. Accumulation?
- Can we shard the optim. states, grads, and model params across GPUs?



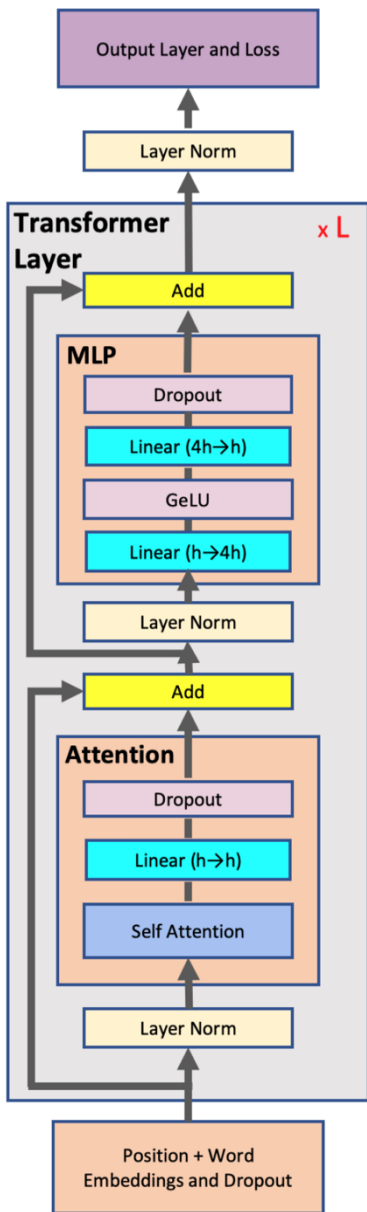
- How to train big models on big data?
- What's in the GPU memory during training?
  - i. Model weights ; ii. Param gradients iii. Optim states iv. Activations
- What is the size of params / grads / optim states?
- What is the size of activations?
- How to reduce activation memory?
  - Activation re-computation *aka* Gradient checkpointing
- How to increase batch size?
  - Gradient accumulation ( run fwd / bwd  $k$  times before `optim.step()`)
- Can we parallelize grad. Accumulation?
- Can we shard the optim. states, grads, and model params across GPUs?
- Still not sufficient for big models (e.g 70B Llama). Can we split one input sequence across GPUs?

## Recap



# Let us revisit activation memory





<b>Total</b>	$34 * seq * bs * h + 5 * n_{heads} * seq^2 * bs$
<b>MLP Block</b>	$19 * seq * bs * h$
D/o mask	$1 * seq * bs * h$
Linear ( $4h \rightarrow h$ )	$8 * seq * bs * h$
GELU	$8 * seq * bs * h$
Linear ( $h \rightarrow 4h$ )	$2 * seq * bs * h$
<b>Layer Norm</b>	$2 * seq * bs * h$
<b>Attention Block</b>	$11 * seq * bs * h + 5 * n_{heads} * seq^2 * bs$
<b>Layer Norm</b>	$2 * seq * bs * h$

# Memory for Activations

$$m_{act} = L * \left( \begin{array}{c} 34 * seq * bs * h \\ + \\ 5 * n_{heads} * seq^2 * bs \end{array} \right)$$

- Scales Linearly with batch size
- Quadratically with the sequence length

[Korthikanti etal. 2022, Reducing Activation Recomputation in Large Transformer Models](#)



# MLP Block

Linear ( $4h \rightarrow h$ )	$8 * seq * bs * h$
GELU	$8 * seq * bs * h$
Linear ( $h \rightarrow 4h$ )	$2 * seq * bs * h$

- neurons are independent of
- Can we compute them independently on different GPUs?
- How will it split the weight matrix?

Inputs

$x_1$

$x_h$

$W_{11}$

$u_1$

$u_{4h}$

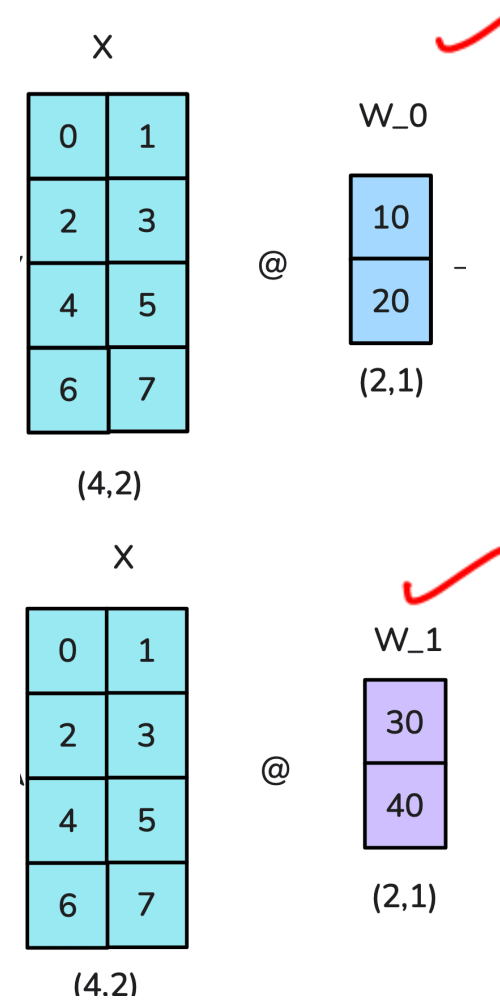
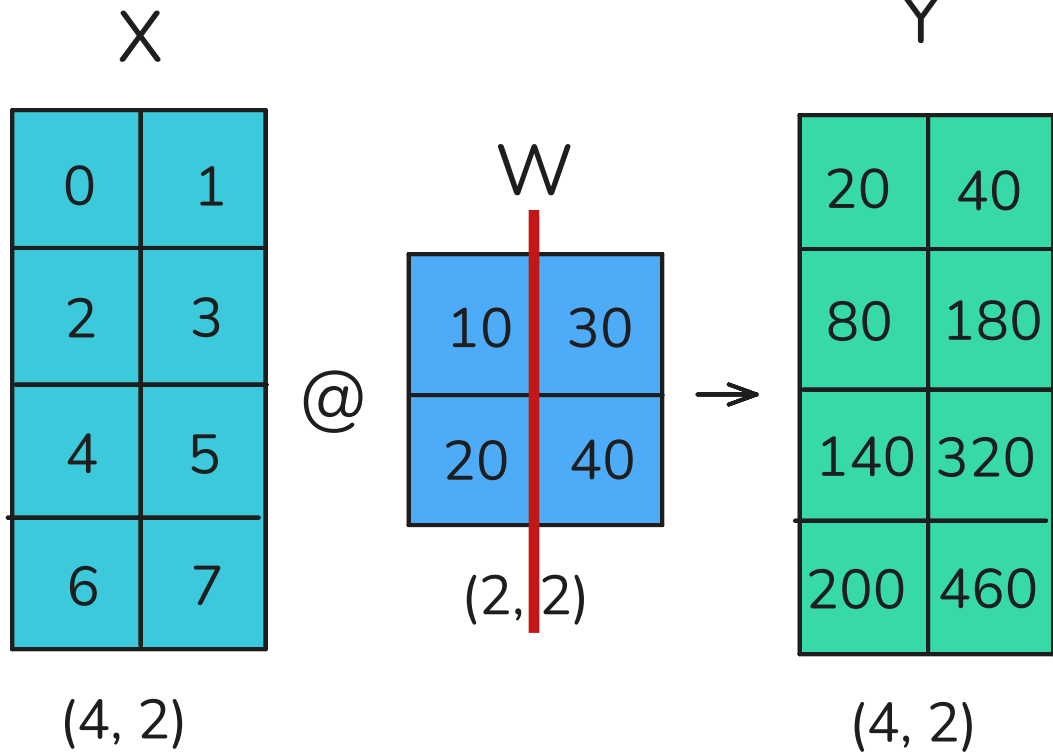
$o_1$

$o_h$

$$1. \quad A \cdot B = A \cdot [B_1 \quad B_2 \quad \dots] = [AB_1 \quad AB_2 \quad \dots]$$

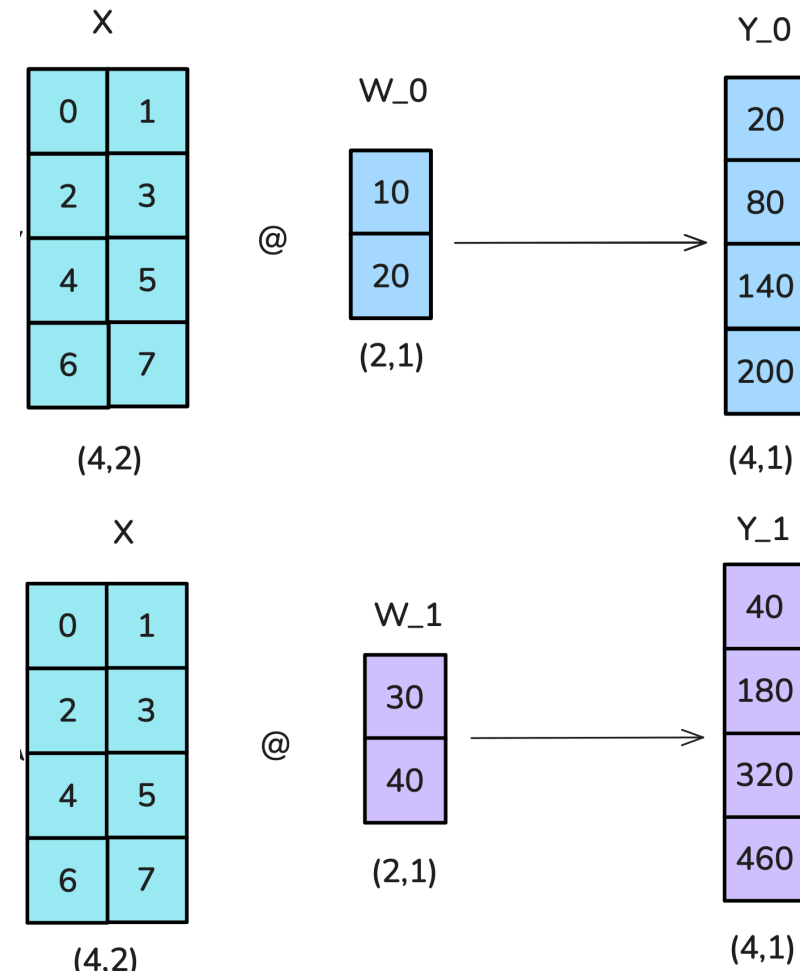
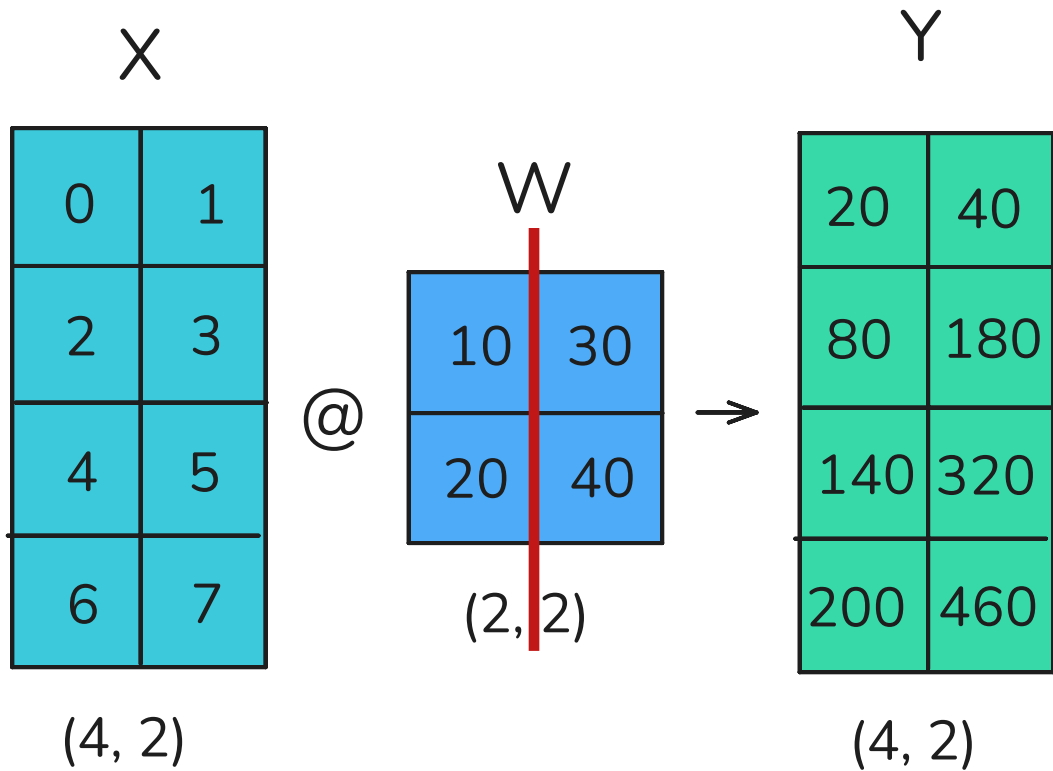


# Column Linear



$$1. \quad A \cdot B = A \cdot [B_1 \quad B_2 \quad \dots] = [AB_1 \quad AB_2 \quad \dots]$$

# Column Linear



$$1. \quad A \cdot B = A \cdot [B_1 \quad B_2 \quad \dots] = [AB_1 \quad AB_2 \quad \dots]$$

# MLP Block

Linear ( $4h \rightarrow h$ )	$8 * seq * bs * h$
GELU	$8 * seq * bs * h$
Linear ( $h \rightarrow 4h$ )	$2 * seq * bs * h$

- Focus on the 2<sup>nd</sup> Linear Layer and output neurons
- Both   and   neurons contribute to all the output neurons

- Can we split the computation on two GPUs?

Inputs

$x_1$

$\vdots$

$x_h$

$W_{11}$

$u_1$

$u_{4h}$

$o_1$

$\vdots$

$o_h$

$$1. A \cdot B = A \cdot [B_1 \ B_2 \ \dots] = [AB_1 \ AB_2 \ \dots]^T$$

$$2. A \cdot B = [A_1 \ A_2 \ \dots] \begin{bmatrix} B_1 \\ B_2 \\ \vdots \end{bmatrix} = \sum_{i=1}^n A_i B_i$$

Handwritten notes for equation 2:

- $u^T m = w_1 m_{1,1-h} + w_2 m_{1,h+1-2h} + w_3 m_{1,2h+1-3h} + w_4 m_{1,3h+1-4h}$
- $GPU1$  (under  $A_1$ )
- $GPU2$  (under  $A_2$ )
- $GPU4$  (under  $A_4$ )



$$\underbrace{\begin{bmatrix} w_1 & w_2 \end{bmatrix}}_{1 \times 20} \underbrace{M}_{20 \times 30} = \underbrace{w}_{1 \times 20} \underbrace{\begin{bmatrix} M_1 \\ M_2 \\ M_3 \\ M_4 \end{bmatrix}}_{20 \times 30} = \underbrace{w_1 M_1}_{1 \times 30}$$

$$\begin{bmatrix} w_1 & w_2 \end{bmatrix} \begin{bmatrix} 10 & M_1 \\ 10 & M_2 \end{bmatrix}_{10 \times 30} = \underbrace{w^T m^1}_{1 \times 30} = \begin{bmatrix} w_1^T & w_2^T \end{bmatrix} \begin{bmatrix} m_1^1 \\ m_2^1 \end{bmatrix}_{10}$$

X

0	1
2	3
4	5
6	7

(4, 2)

@

W

10	30
20	40

(2, 2)




(4, 2)

@

W

10	30
20	40

(2, 2)



Y2


(4, 2)

$$o_1 = \underline{w_1^T m_1^1} + \underline{w_2^T m_2^1}$$

$$o_2 = \underline{w_1^T m_1^2} + \underline{w_2^T m_2^2}$$

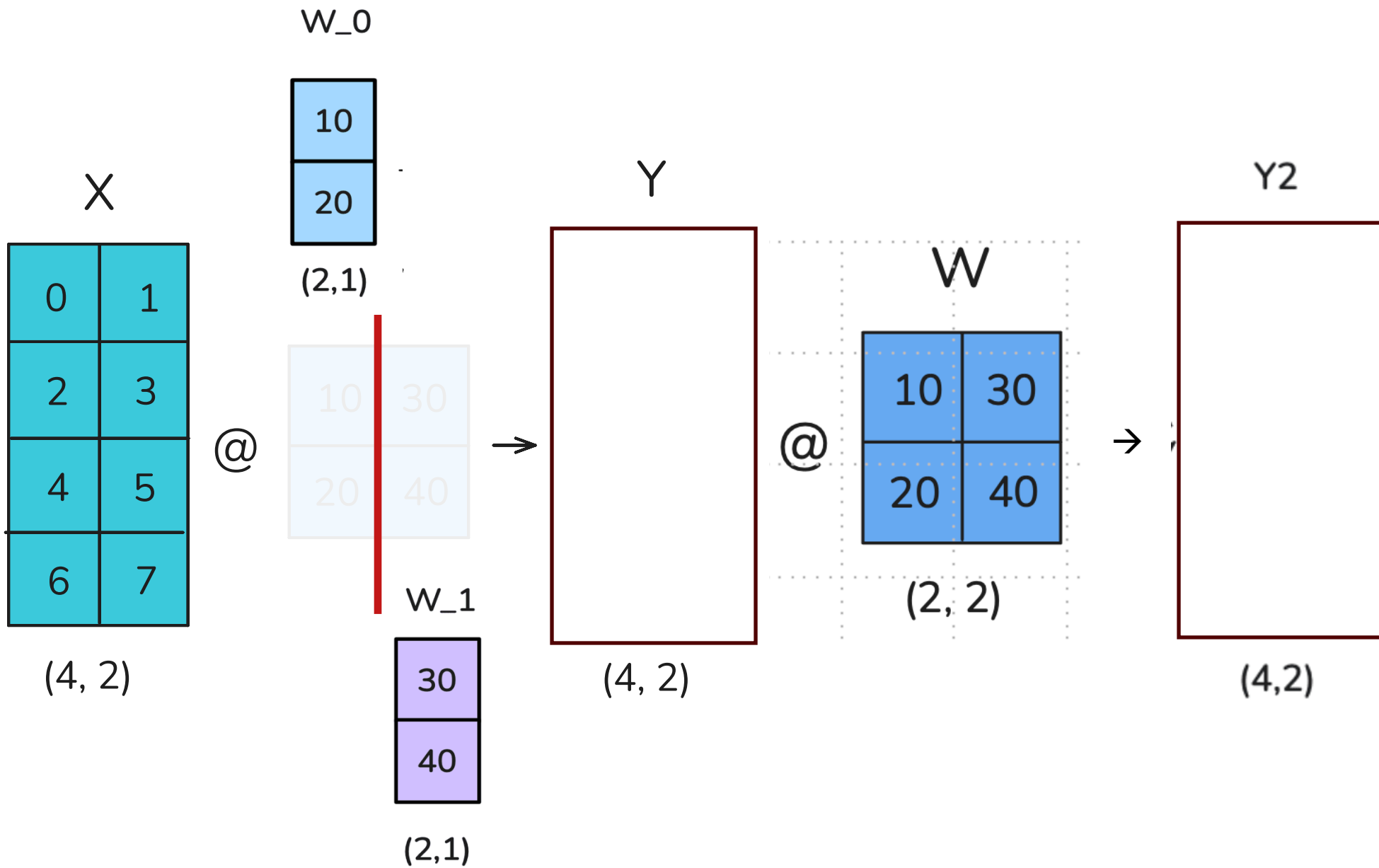
$$\left( \begin{array}{cc} w_1^T M_1 & 1 \times 10 \quad 10 \times 30 \\ + \\ w_1^T M_2 & 1 \times 10 \quad 10 \times 30 \end{array} \right) \begin{bmatrix} \text{ } \\ \text{ } \end{bmatrix}_{30}$$

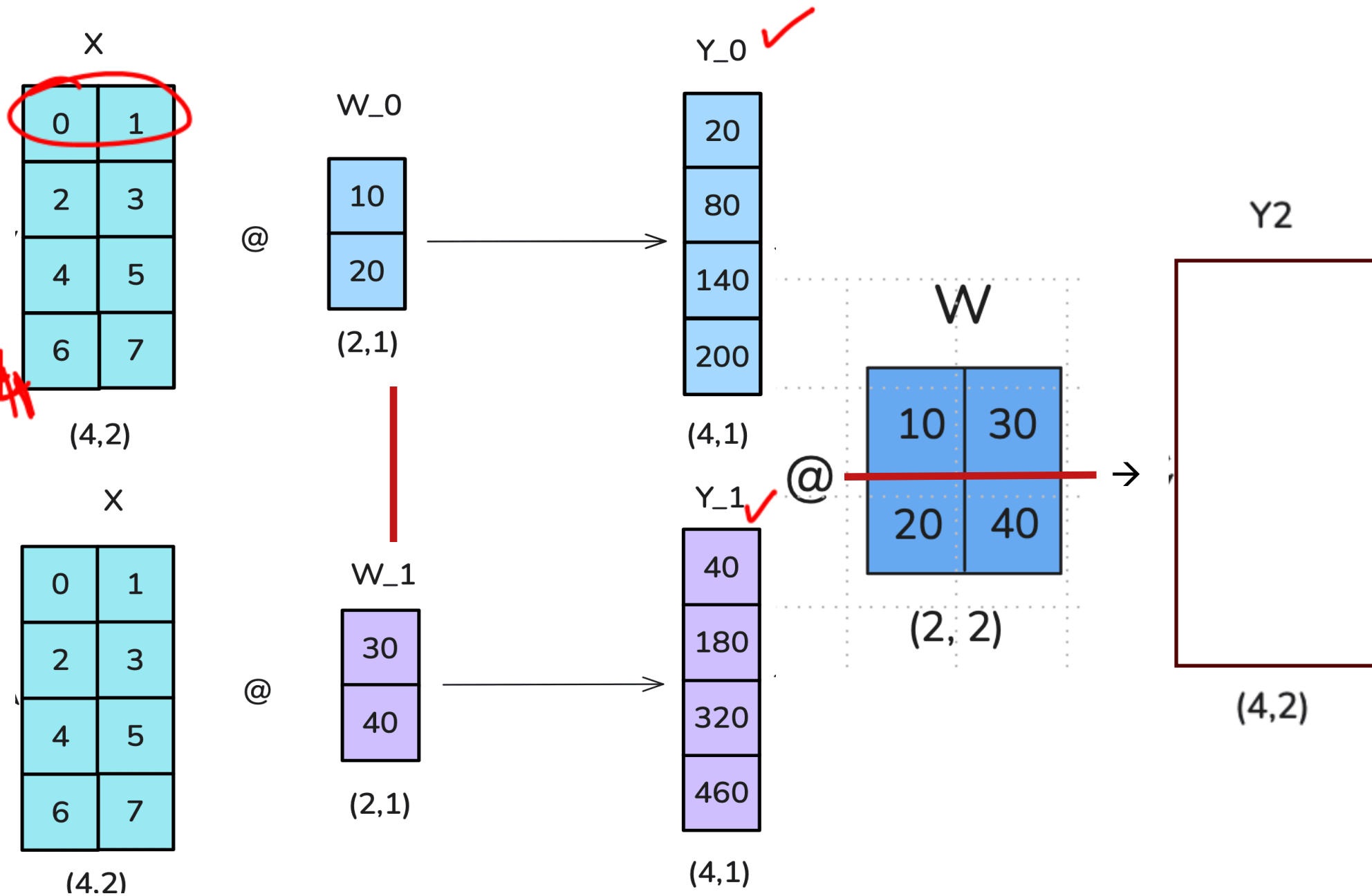
1.  $A \cdot B = A \cdot [B_1 \ B_2 \ \dots] = [AB_1 \ AB_2 \ \dots]$

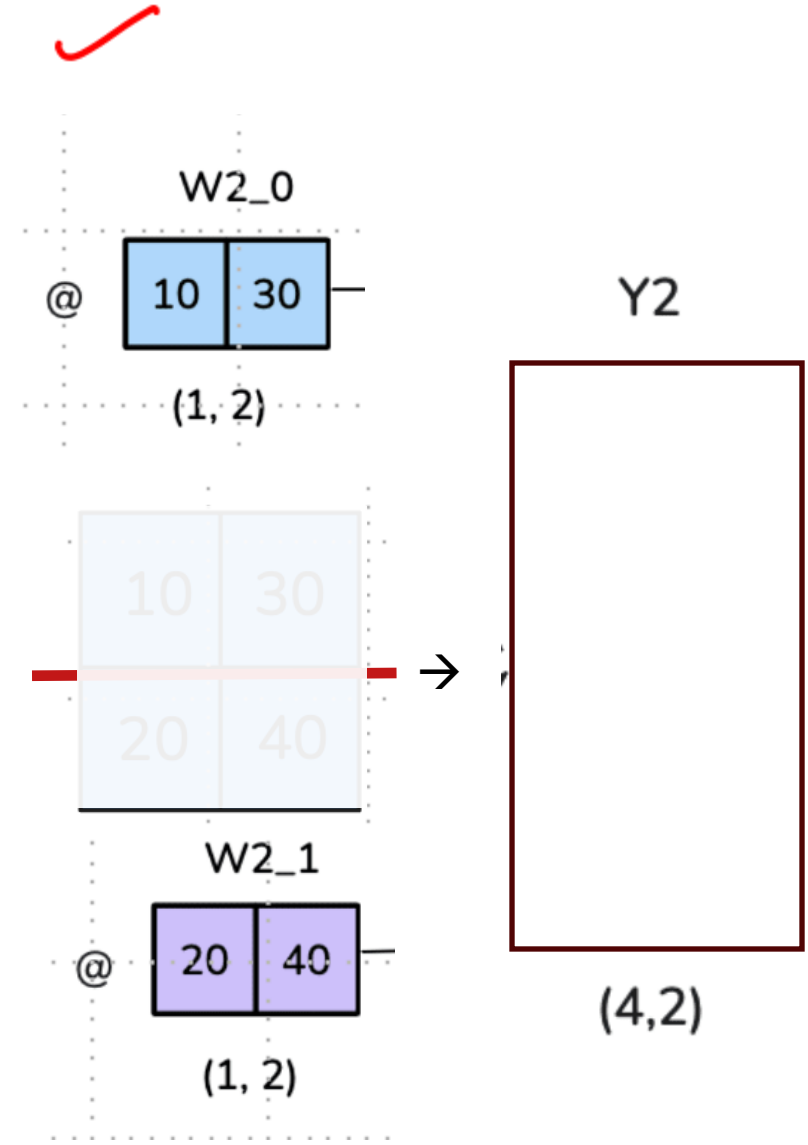
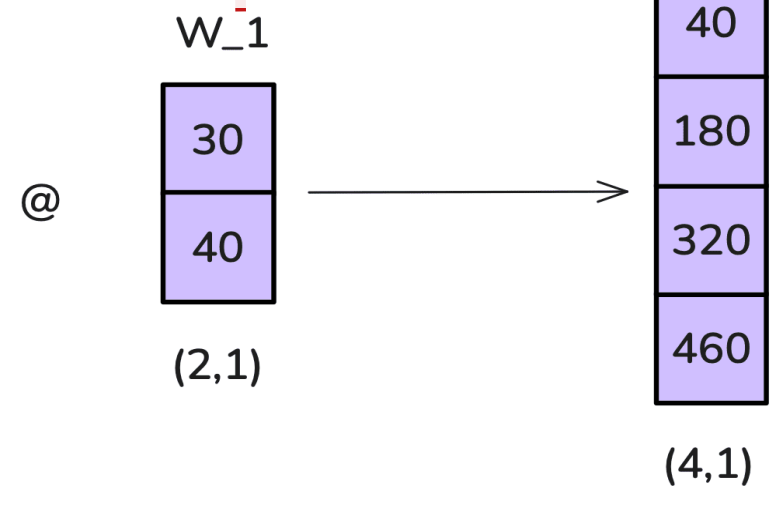
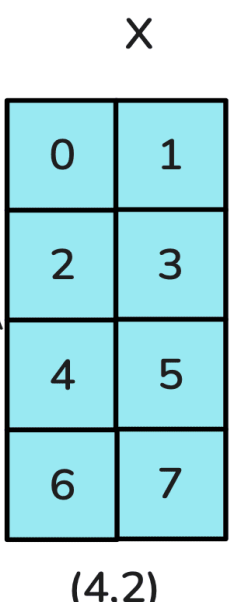
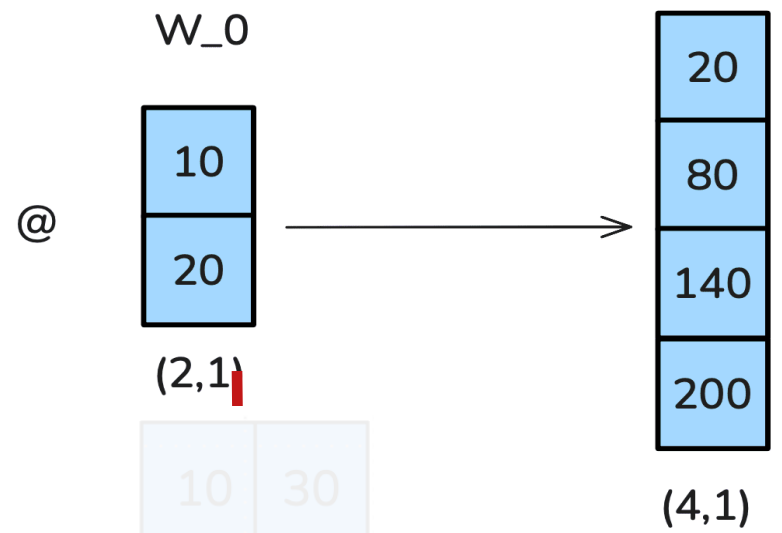
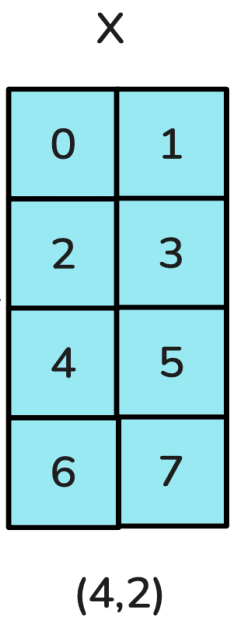
2.  $A \cdot B = [A_1 \ A_2 \ \dots] \begin{bmatrix} B_1 \\ B_2 \\ \vdots \end{bmatrix} = \sum_{i=1}^n A_i B_i$

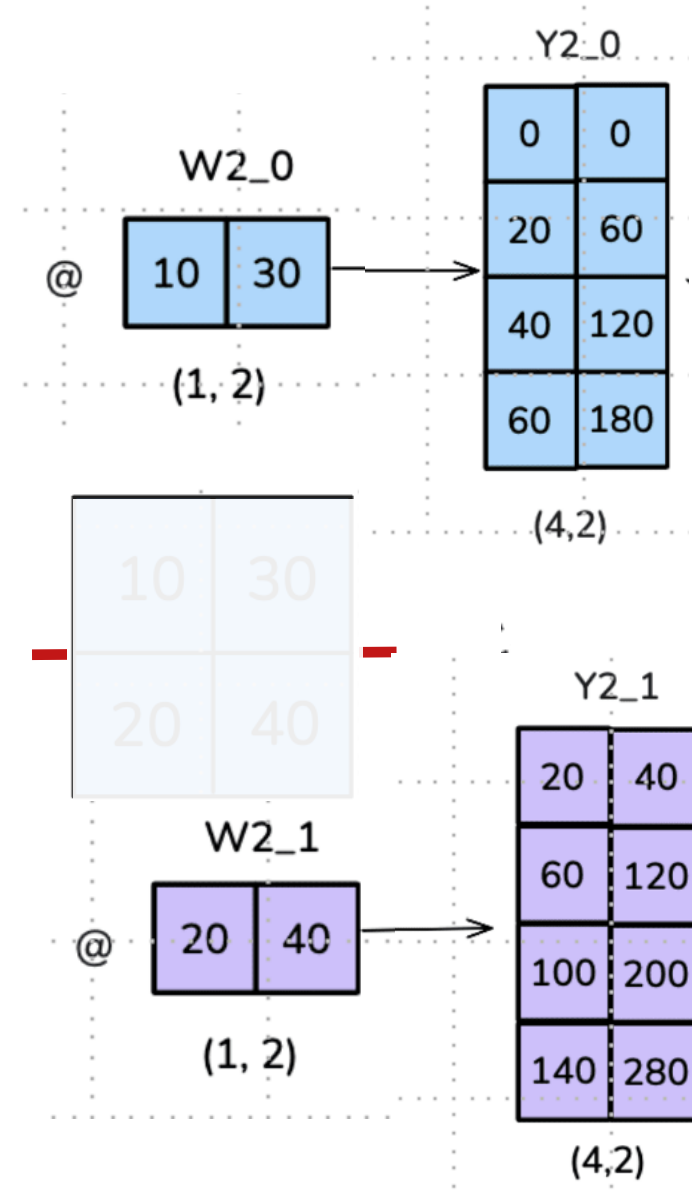
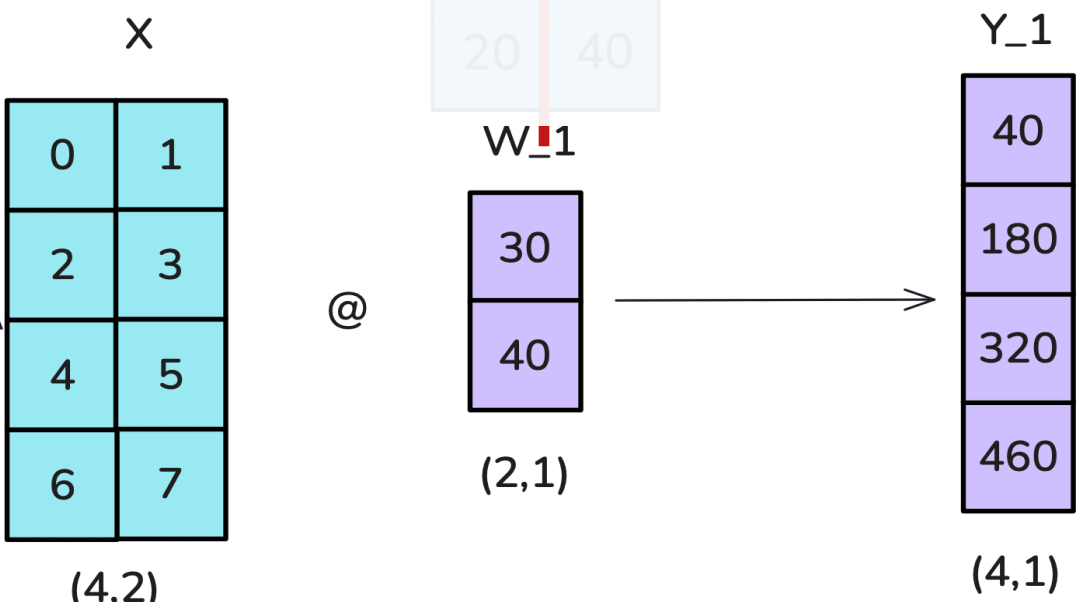
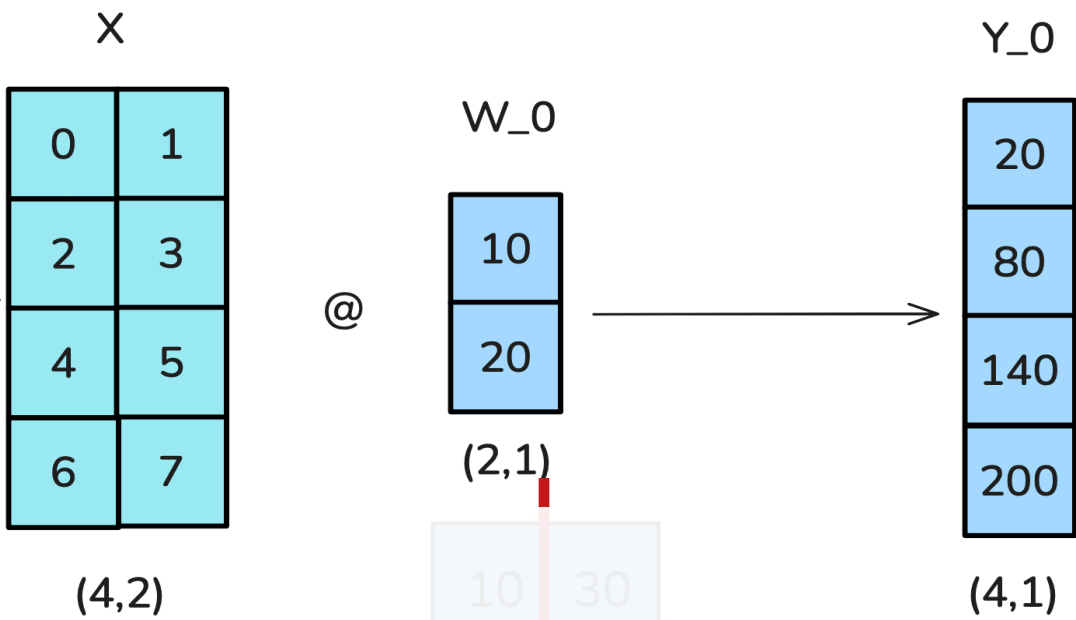


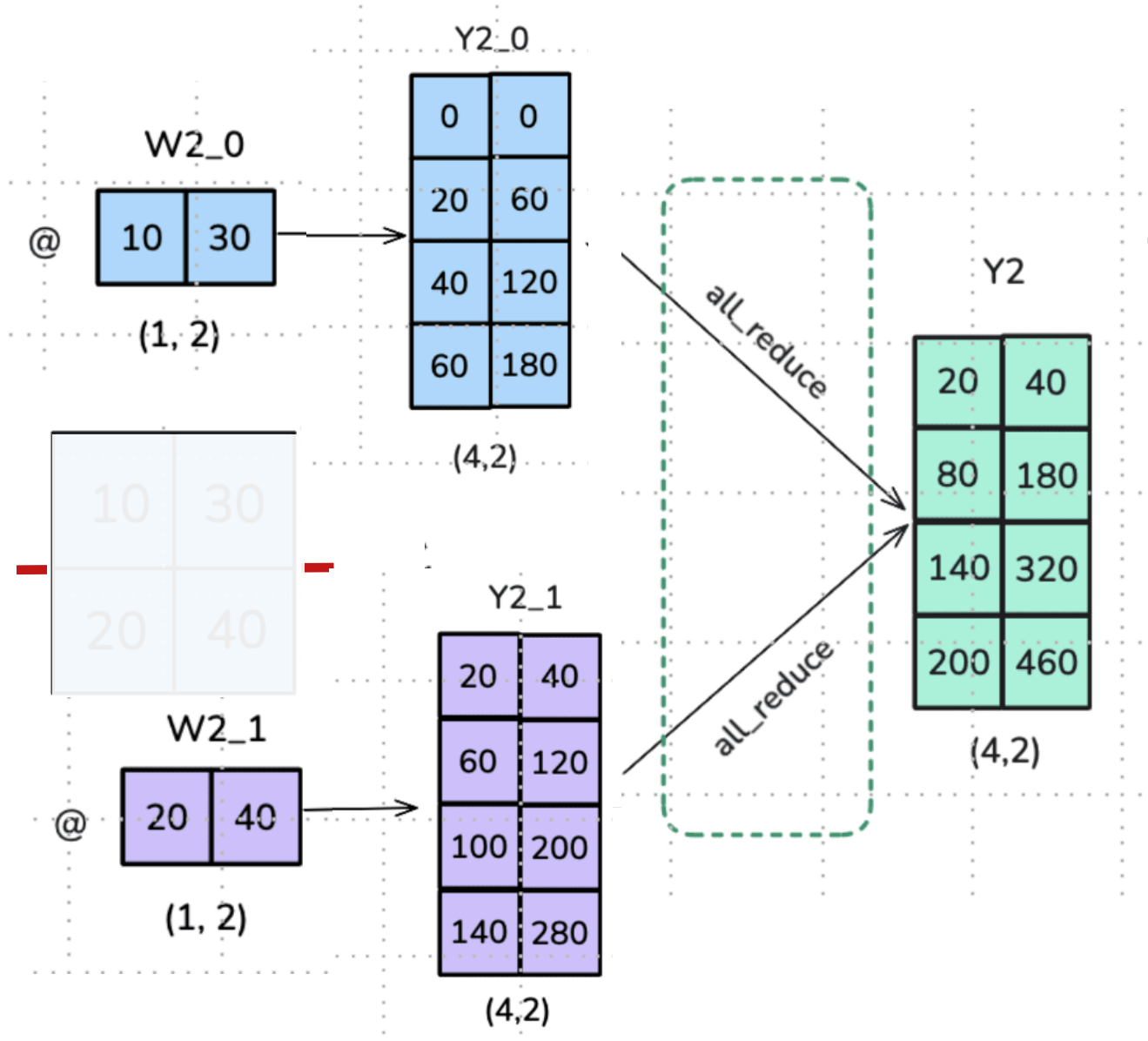
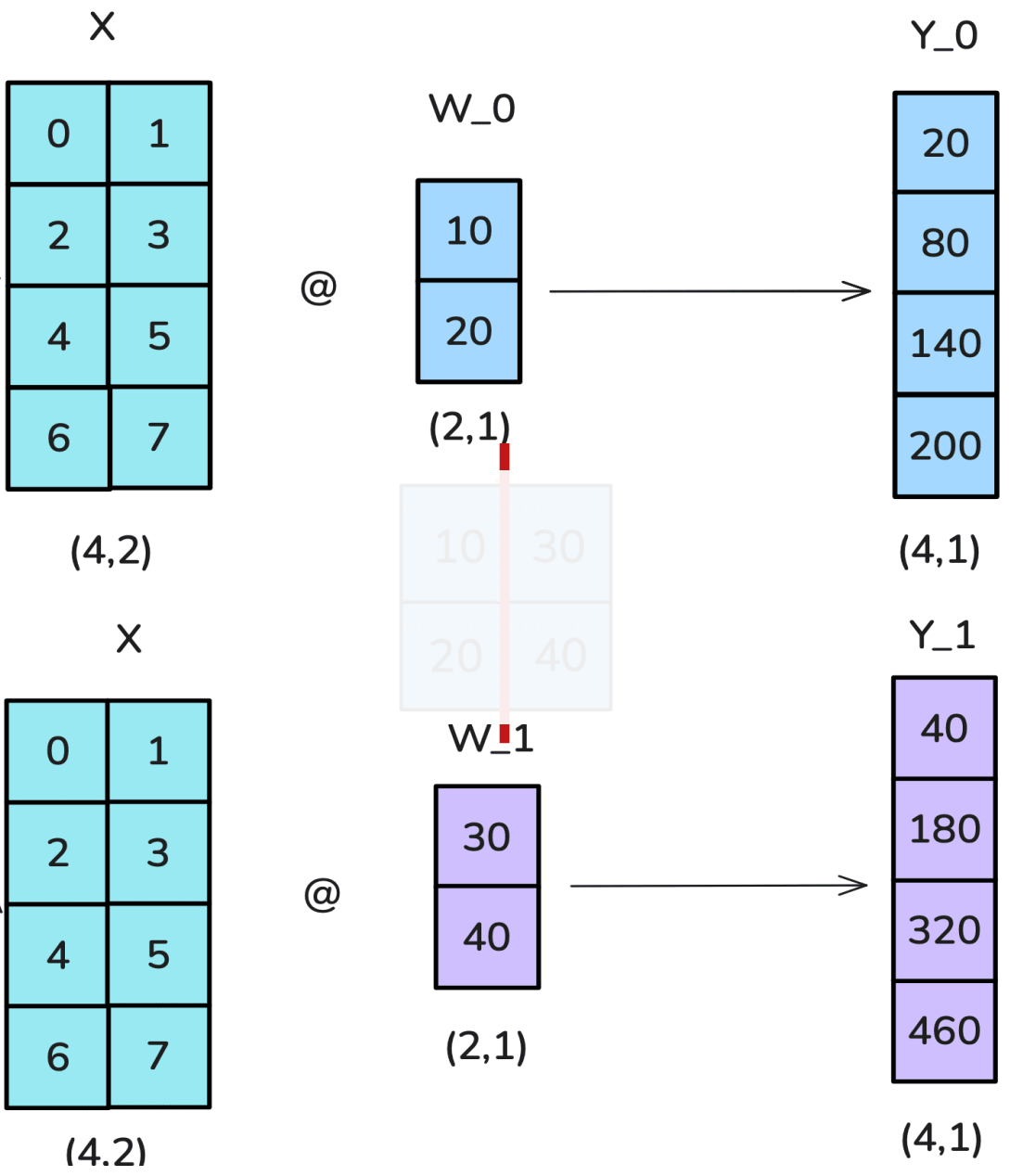




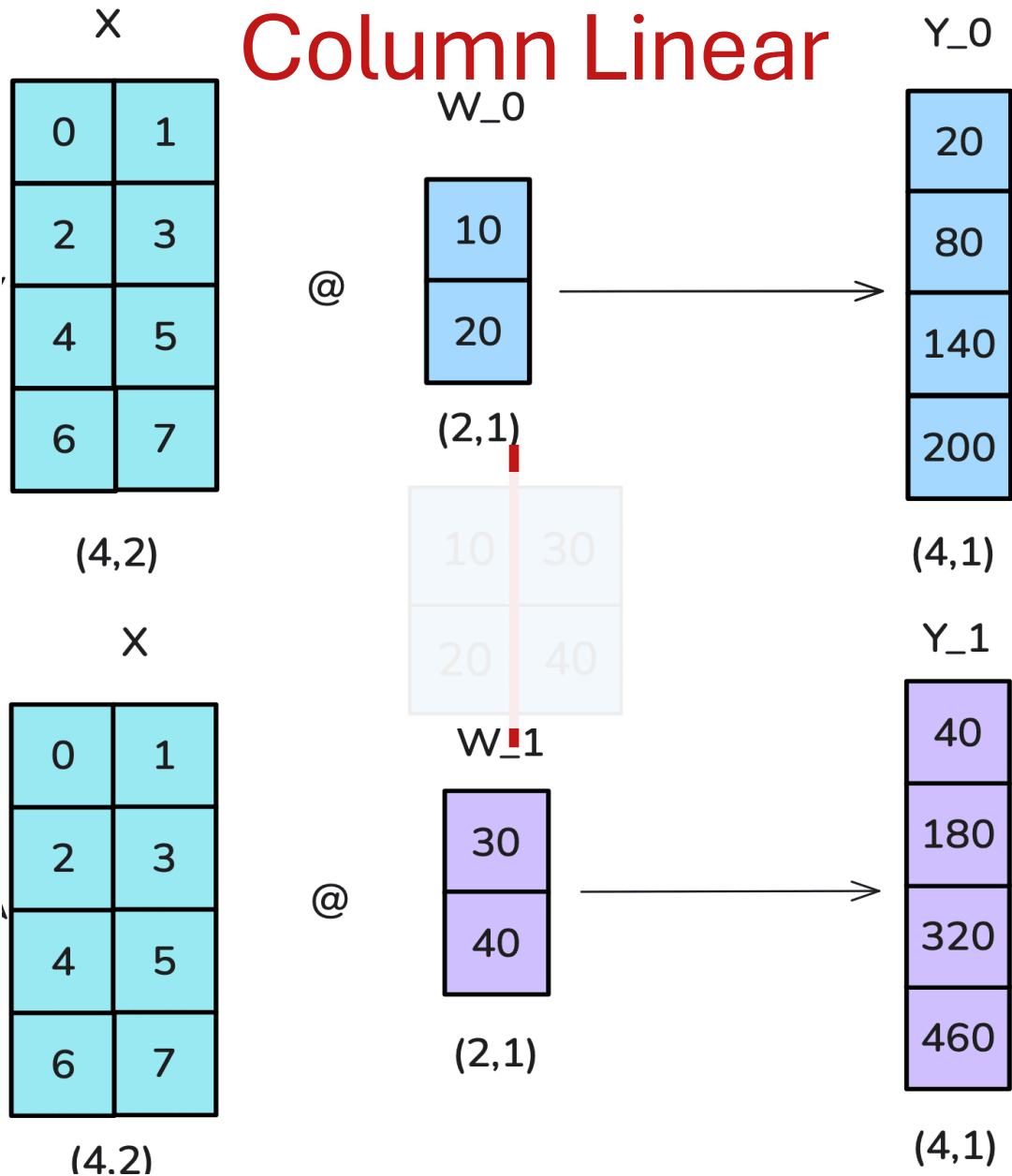




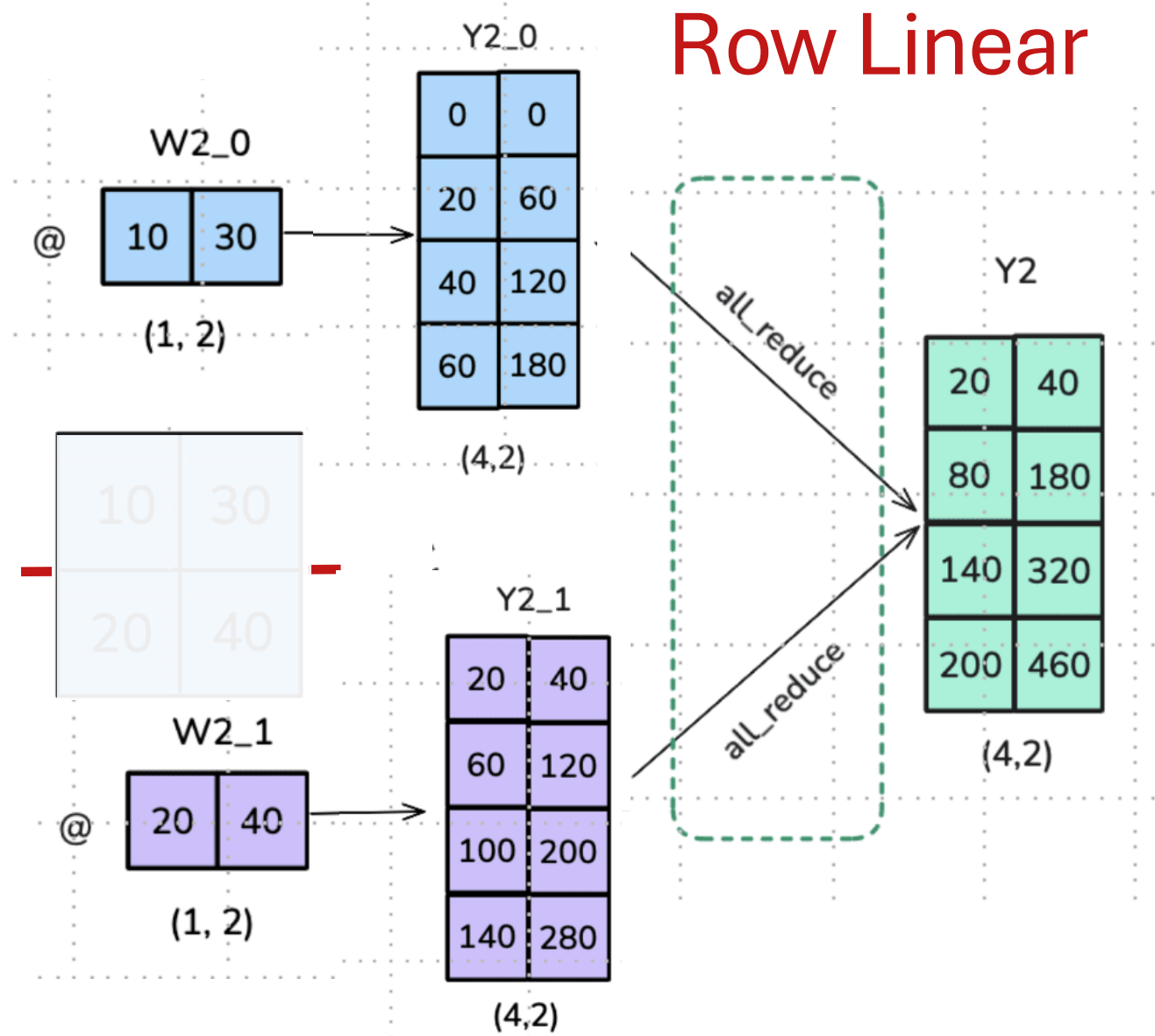




# Column Linear



# Row Linear



# Tensor Parallelism in MLP Block

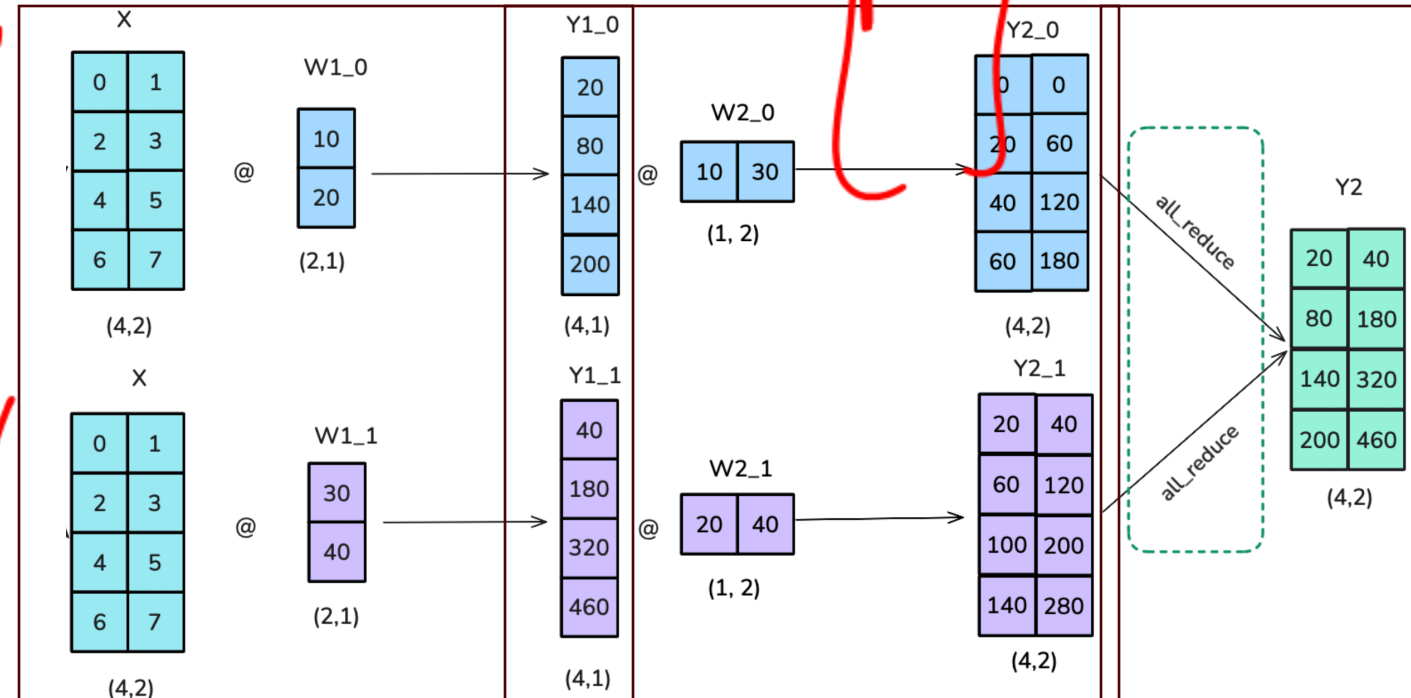
- Use **column-linear** to split the 1<sup>st</sup> layer:

- compute different neurons on different GPUs
- Each GPU compute  $\frac{4h}{TP}$  activations

- Use **row-linear** to split the 2<sup>nd</sup> layer:

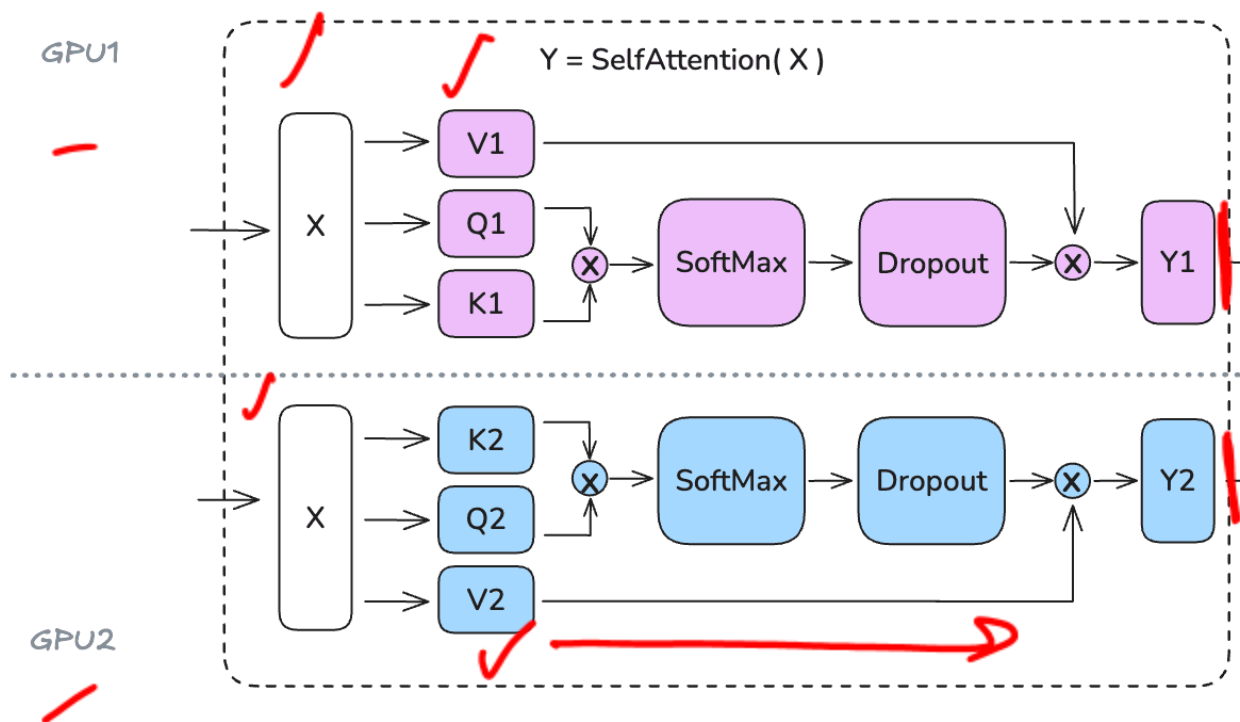
- Each GPU acts on different neurons and computes partial output
- No need to communicate the intermediate activations across GPUs → reduction in activation memory!

- Use **all\_reduce** to communicate the partial outputs



# Tensor Parallelism in MHA Block

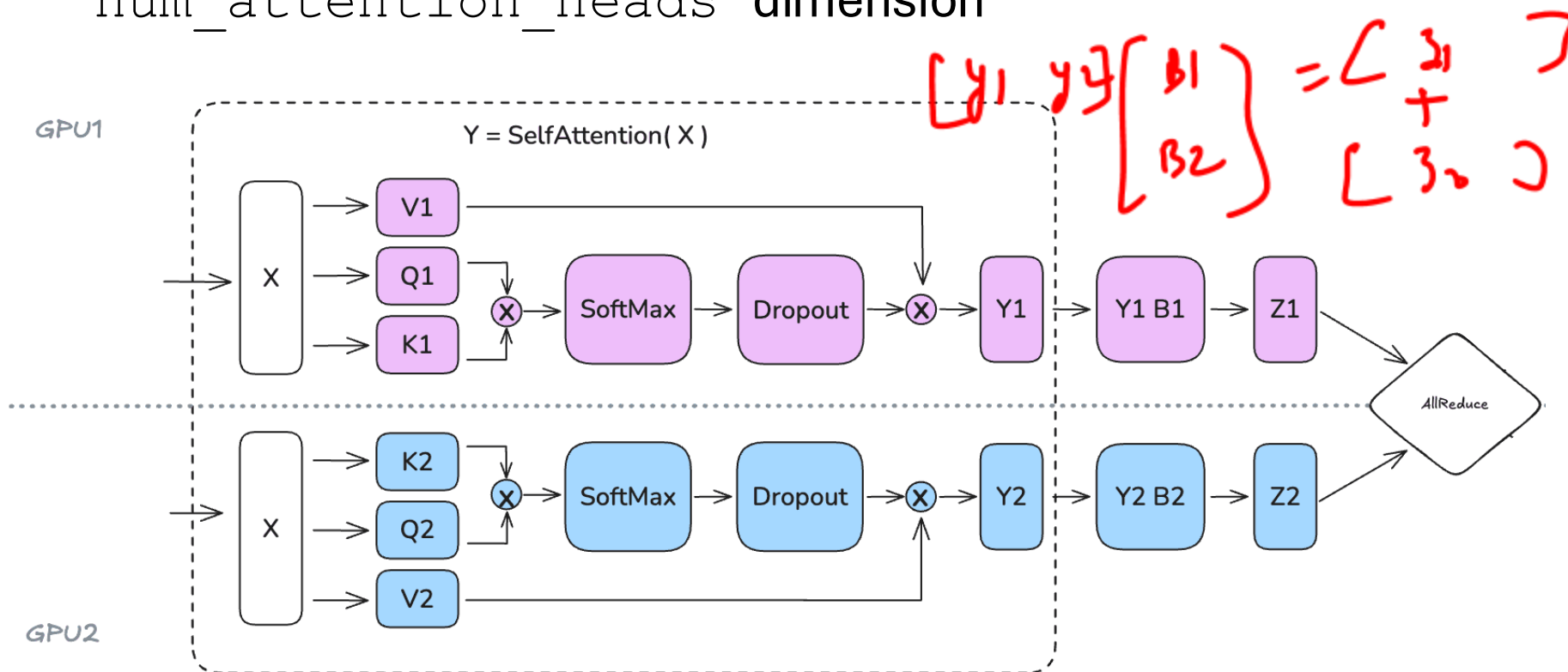
- Parallelize different heads on different GPUs – i.e. along `num_attention_heads` dimension
- Same as splitting  $K, Q, V$  matrices in column-parallel

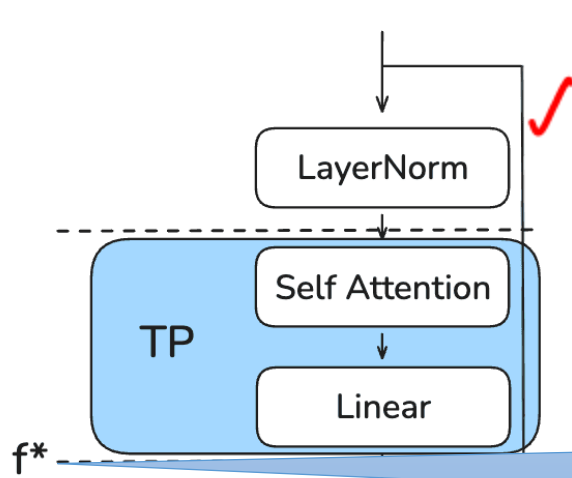




# Tensor Parallelism in Attention Block

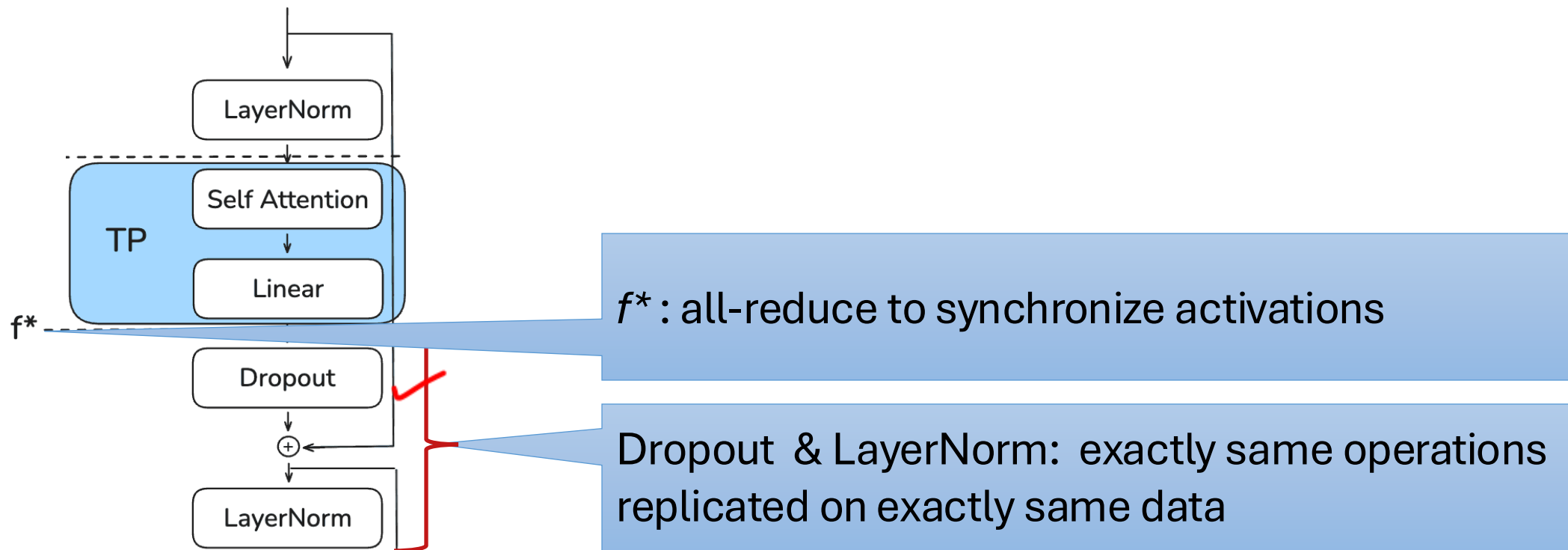
- Parallelize different heads on different GPUs – i.e. along `num_attention_heads` dimension
- Same as splitting  $K, Q, V$  matrices in column-parallel

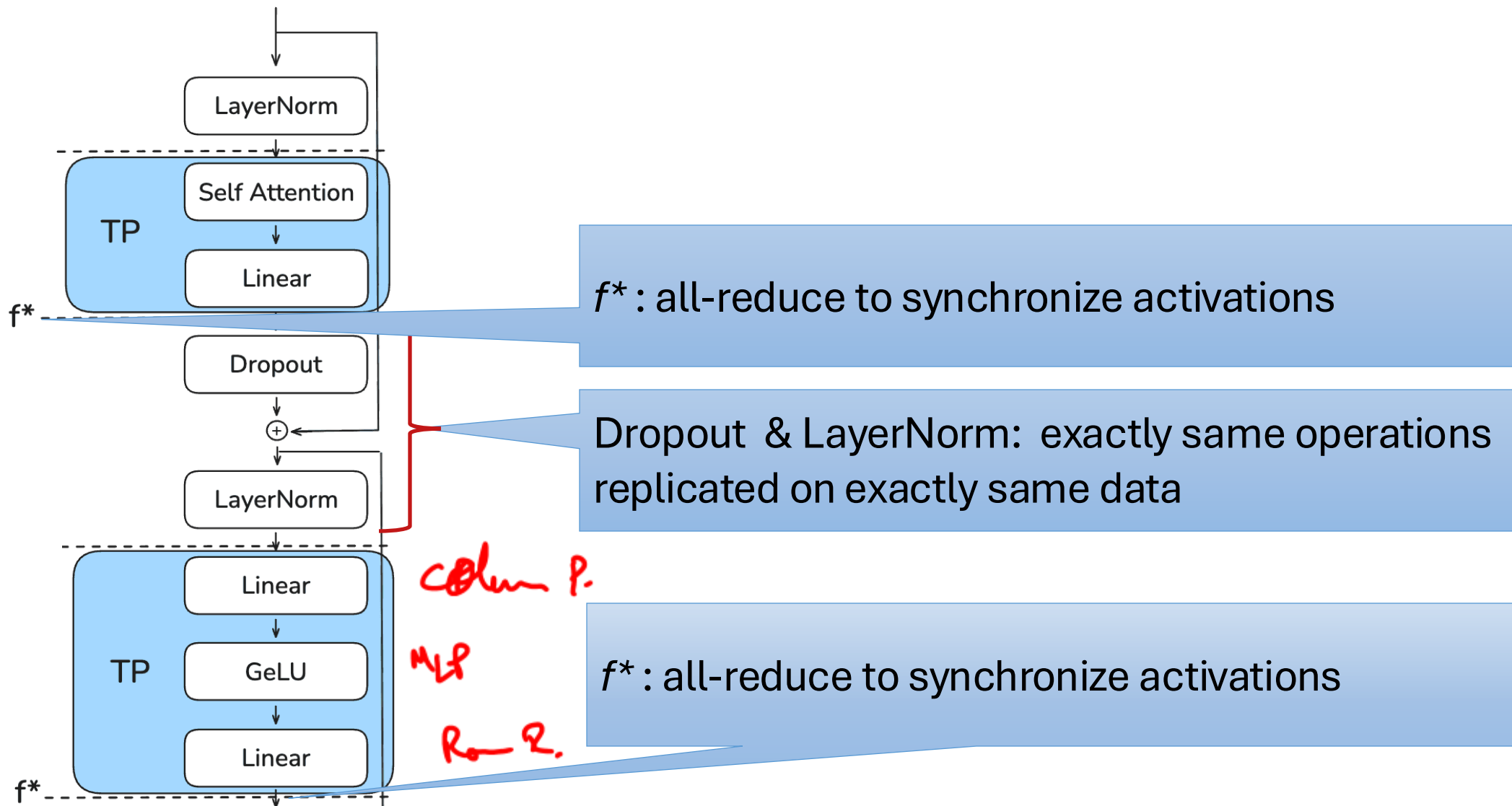


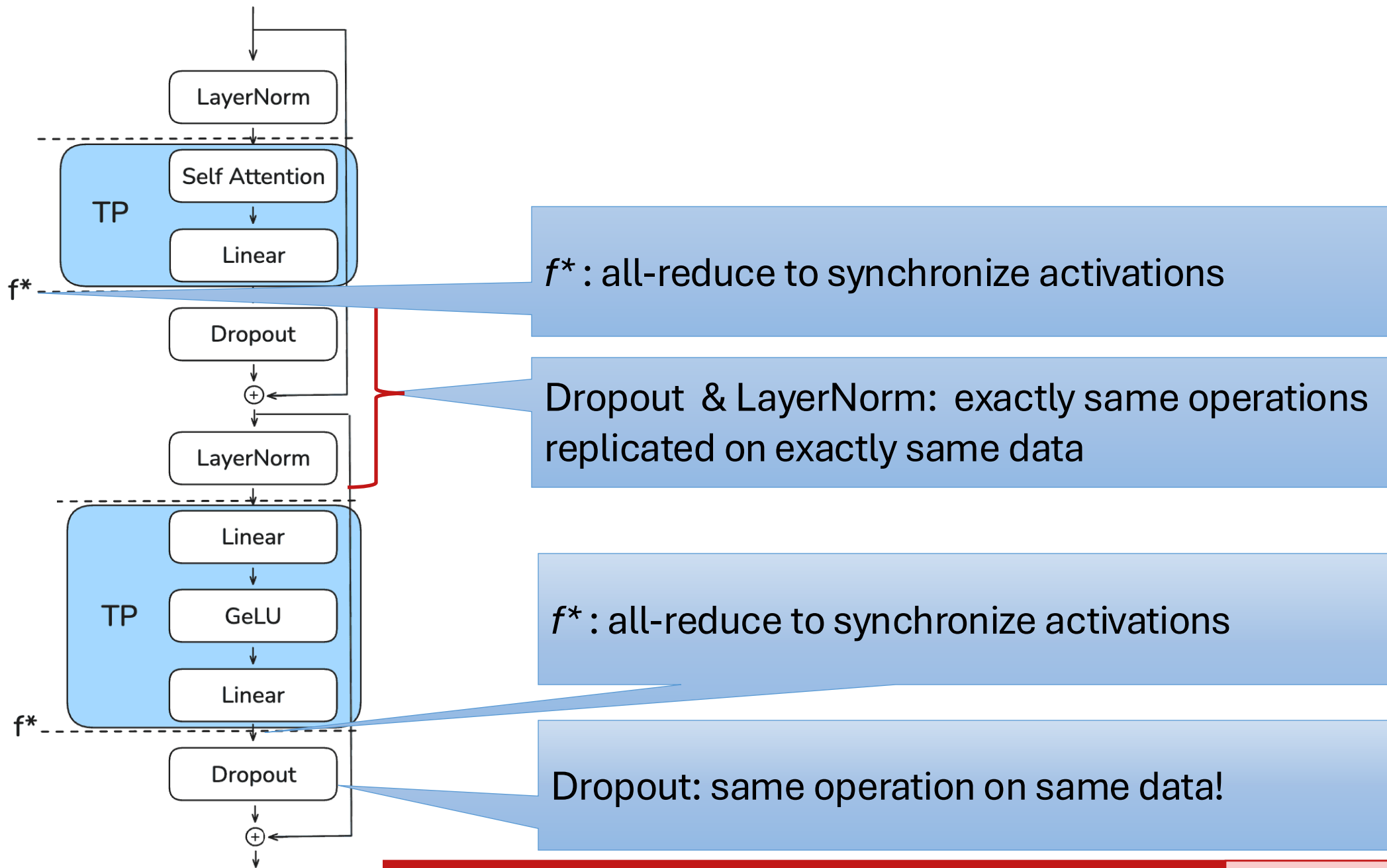


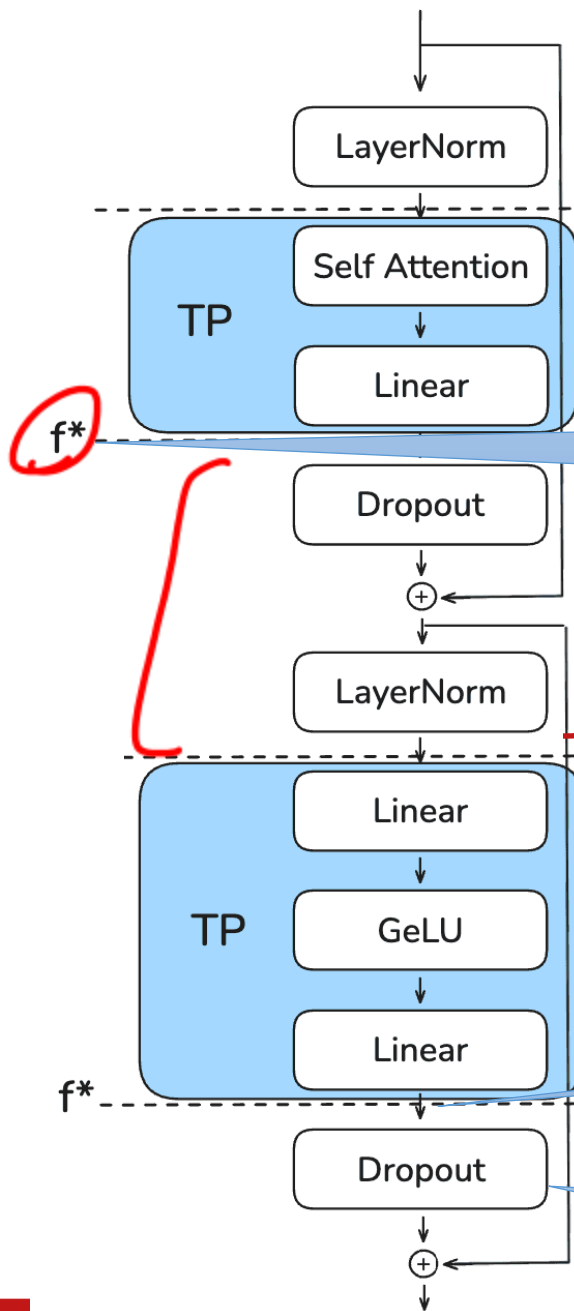
$f^*$  : all-reduce to synchronize activations

1. Synchronization not overlapping with computation
2. “Exposed communication” overhead is necessary to combine partial results across tensor-parallel ranks before the final LayerNorm can be applied.









Dropout and LayerNorm – doing same operation on same data on all TP GPUs!

$f^*$  : all-reduce to synchronize activations

Dropout & LayerNorm: exactly same operations replicated on exactly same data

$f^*$  : all-reduce to synchronize activations

Dropout: same operation on same data!

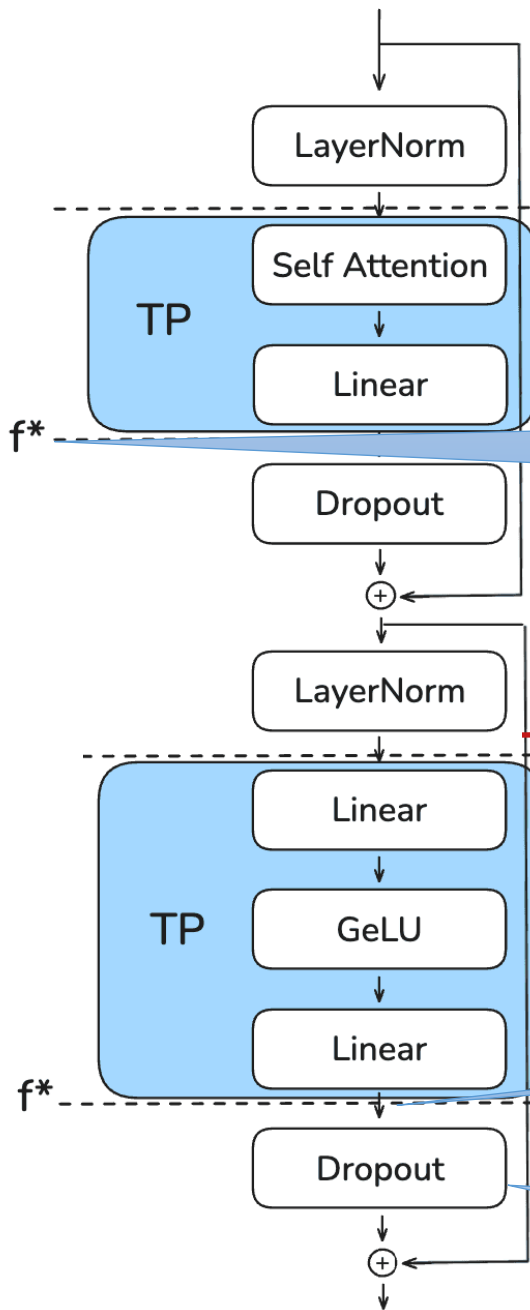
Can we parallelize dropout and LayerNorm?

$f^*$  : all-reduce to synchronize activations

Dropout & LayerNorm: exactly same operations replicated on exactly same data

$f^*$  : all-reduce to synchronize activations

Dropout: same operation on same data!



# Sequence Parallel – parallelizing dropout & LayerNorm

<b>Total</b>	$34 * seq * bs * h + 5 * n_{heads} * seq^2 * bs$
<b>MLP Block</b>	$19 * seq * bs * h$
D/o mask	$1 * seq * bs * h$
Linear ( $4h \rightarrow h$ )	$8 * seq * bs * h$
GELU	$8 * seq * bs * h$
Linear ( $h \rightarrow 4h$ )	$2 * seq * bs * h$
<b>Layer Norm</b>	$2 * seq * bs * h$

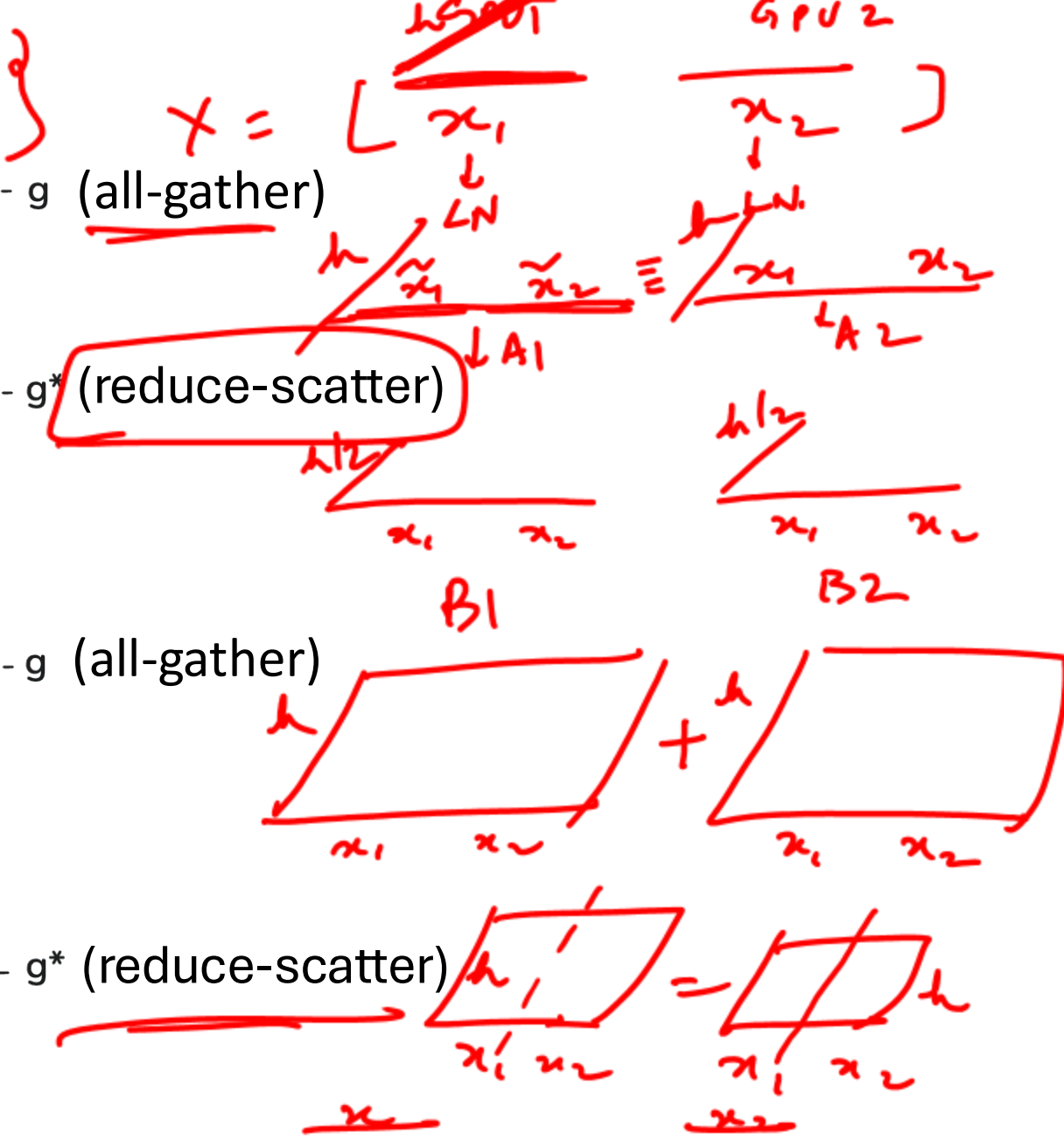
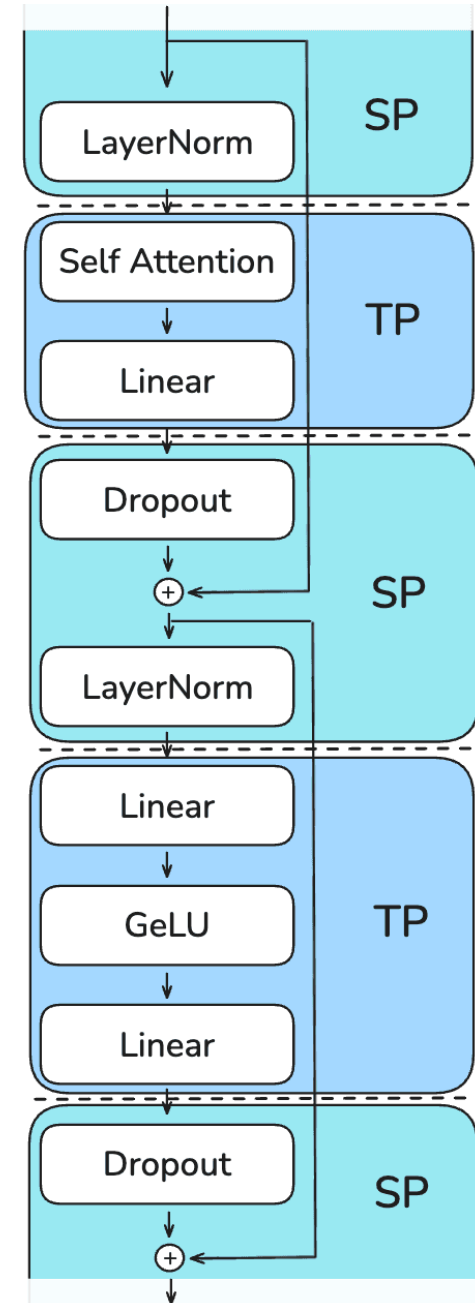
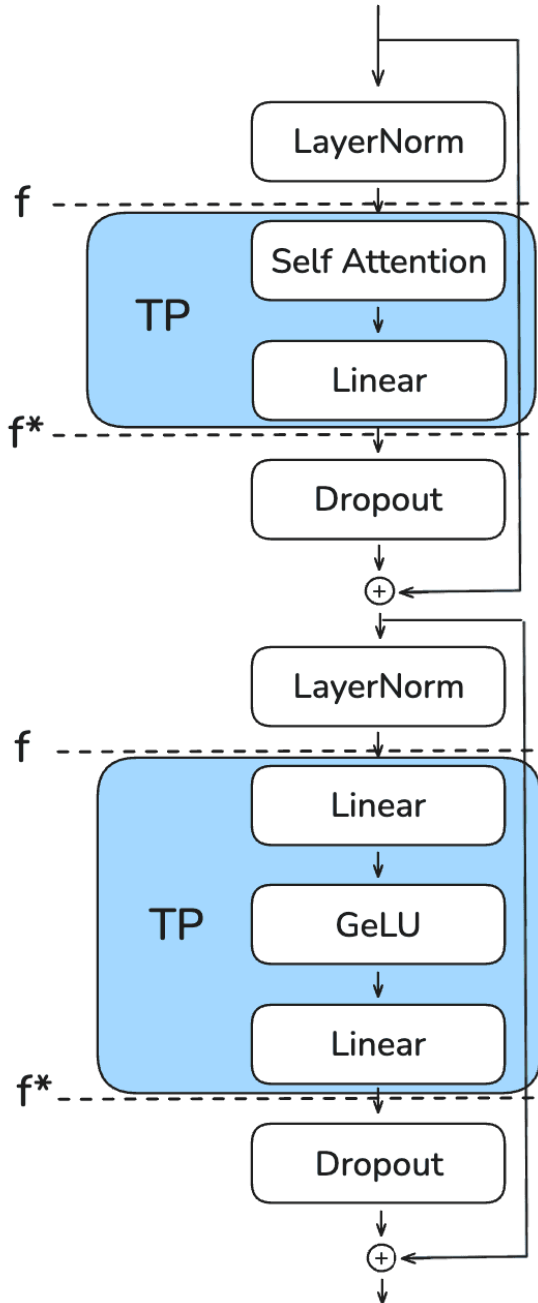
<b>Attention Block</b>	$11 * seq * bs * h + 5 * n_{heads} * seq^2 * bs$
------------------------	--

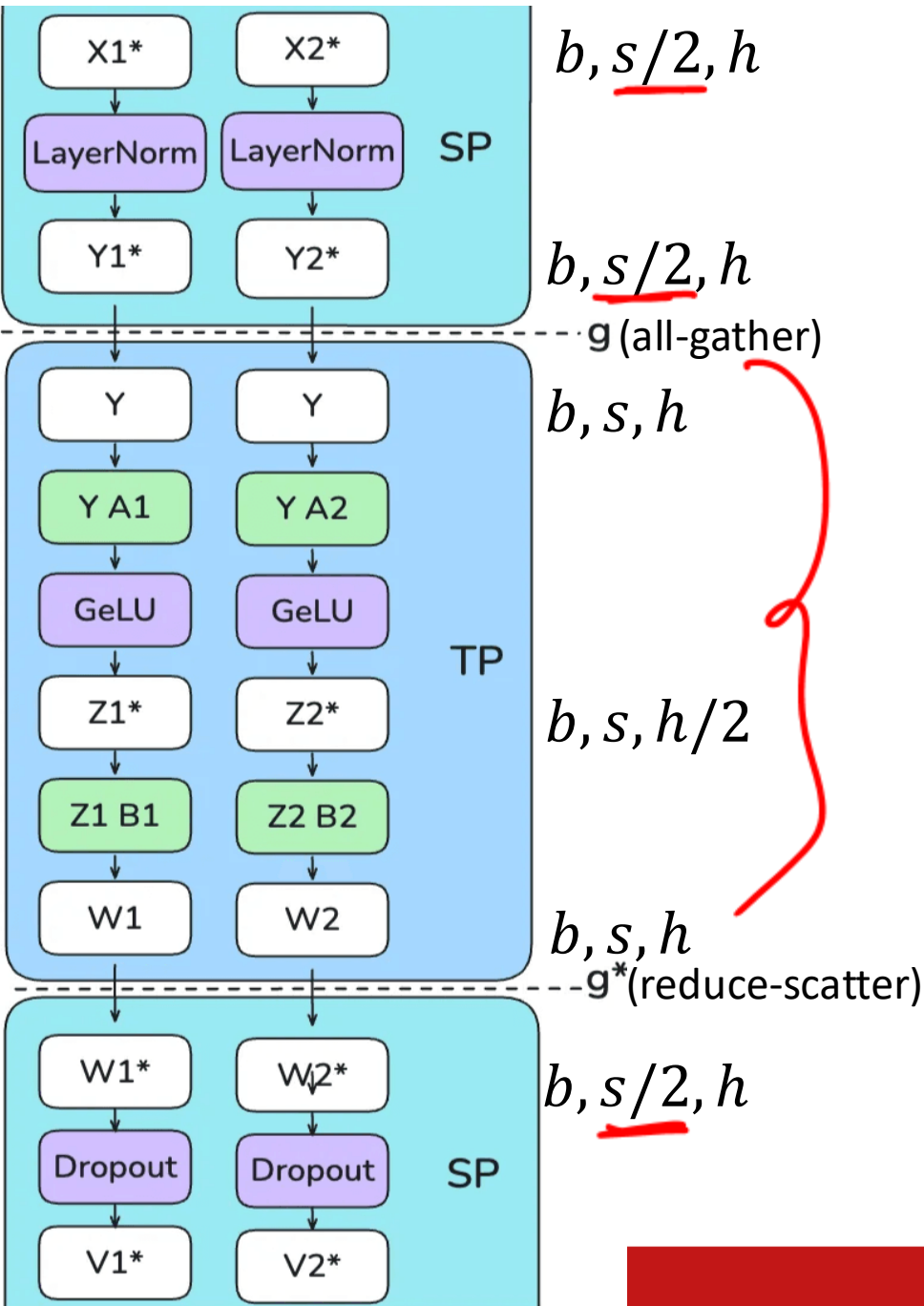
<b>Layer Norm</b>	$2 * seq * bs * h$
-------------------	--------------------

- In DP, we parallelize along the “batch dim” ( $bs$ )
- In TP, we parallelize along the “hidden dim” ( $h$ )
- In SP, we parallelize along the input sequence dimension ( $seq$ )









### Initial LayerNorm layer (SP region)

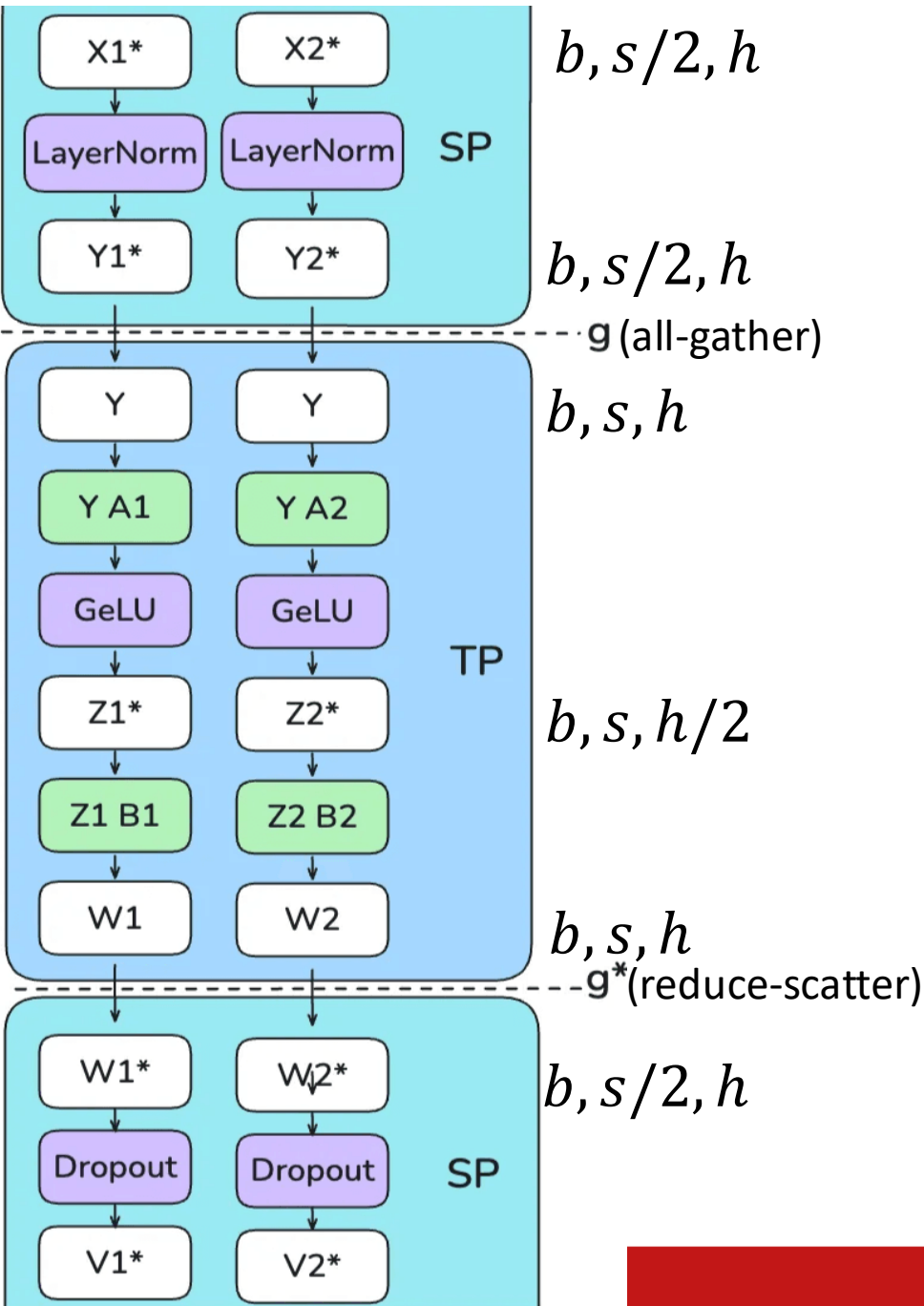
- Input tensors  $X1^*$  and  $X2^*$  ( $b, s/2, h$ ) enter, already split across the sequence dimension.
- Each GPU computes LayerNorm independently on its sequence chunk, giving  $Y1^*$  and  $Y2^*$ .

### First transition (SP $\rightarrow$ TP)

- $g$  operation (all-gather) combines  $Y1$  and  $Y2$  back to full sequence length.
- Restores  $Y$  ( $b, s, h$ ) since column-linear layers need the full hidden dimension  $h$ .

### First linear layer (TP region)

- $A1$  and  $A2$  are column-linear layers, so they split  $Y$  along the hidden dimension.
- GELU is applied independently on each GPU.
- $Z1^*$  and  $Z2^*$  are ( $b, s, h/2$ ).



### Initial LayerNorm layer (SP region)

- Input tensors  $X1^*$  and  $X2^*$  ( $b, s/2, h$ ) enter, already split across the sequence dimension.
- Each GPU computes LayerNorm independently on its sequence chunk, giving  $Y1^*$  and  $Y2^*$ .

### First transition (SP $\rightarrow$ TP)

- $g$  operation (all-gather) combines  $Y1$  and  $Y2$  back to full sequence length.
- Restores  $Y$  ( $b, s, h$ ) since column-linear layers need the full hidden dimension  $h$ .

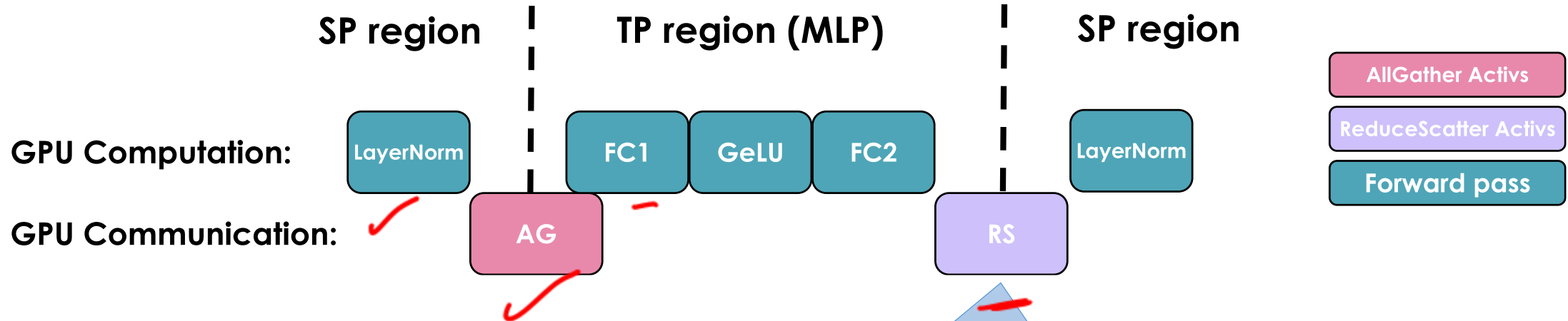
### Second linear layer (TP region)

- $B1$  and  $B2$  are row-linear layers, so they restore the hidden dimension.
- $W1$  and  $W2$  are ( $b, s, h$ ) that need to be summed together.

### Final transition (TP $\rightarrow$ SP)

- $g^*$  operation (reduce-scatter) reduces for previous row-linear correctness while scattering along the sequence dimension.
- $W1^*$  and  $W2^*$  are ( $b, s/2, h$ )

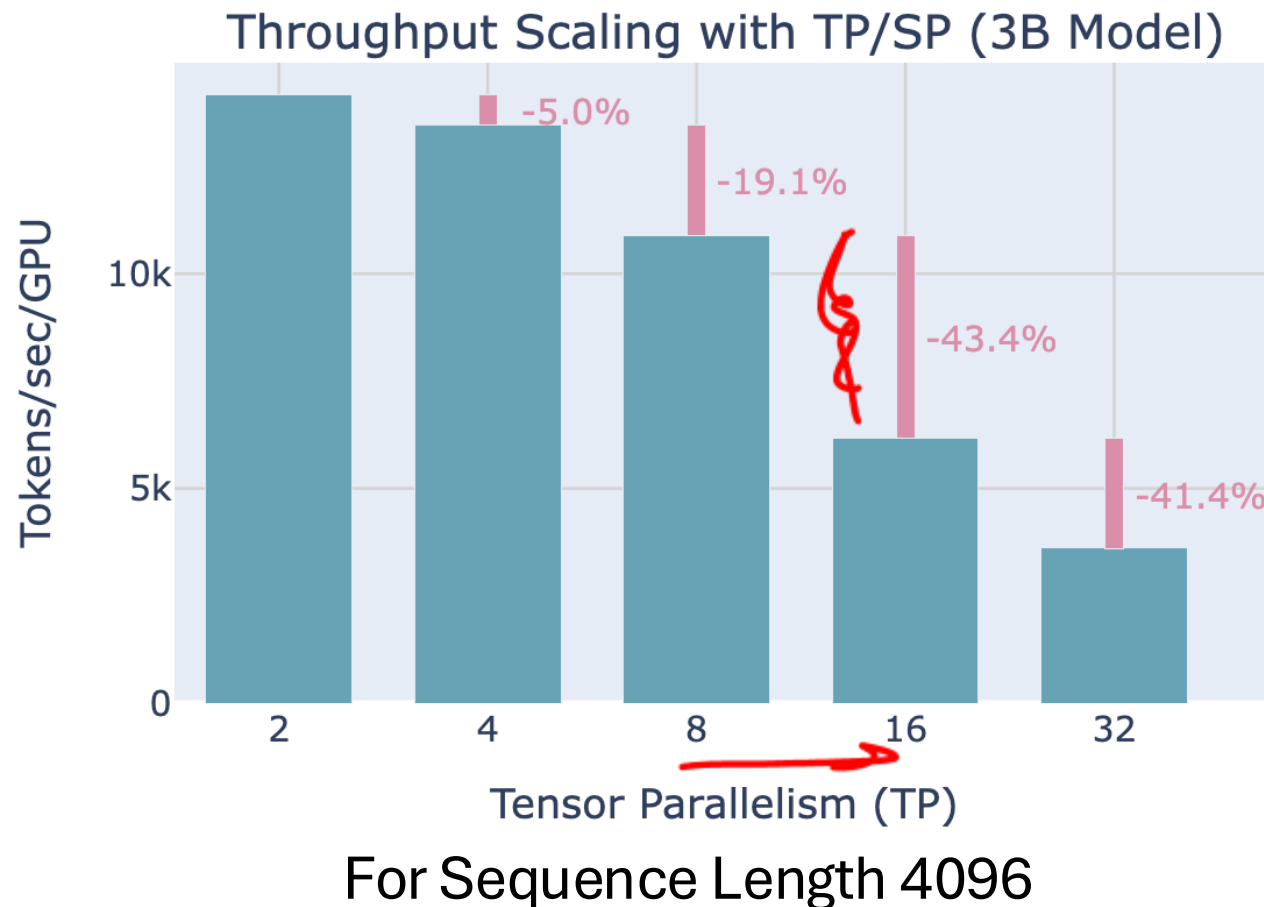
# Computation-communication timeline for MLP Layer



1. Synchronization not overlapping with computation
2. “Exposed communication” overhead is necessary to combine partial results across tensor-parallel ranks before the LayerNorm can be applied.

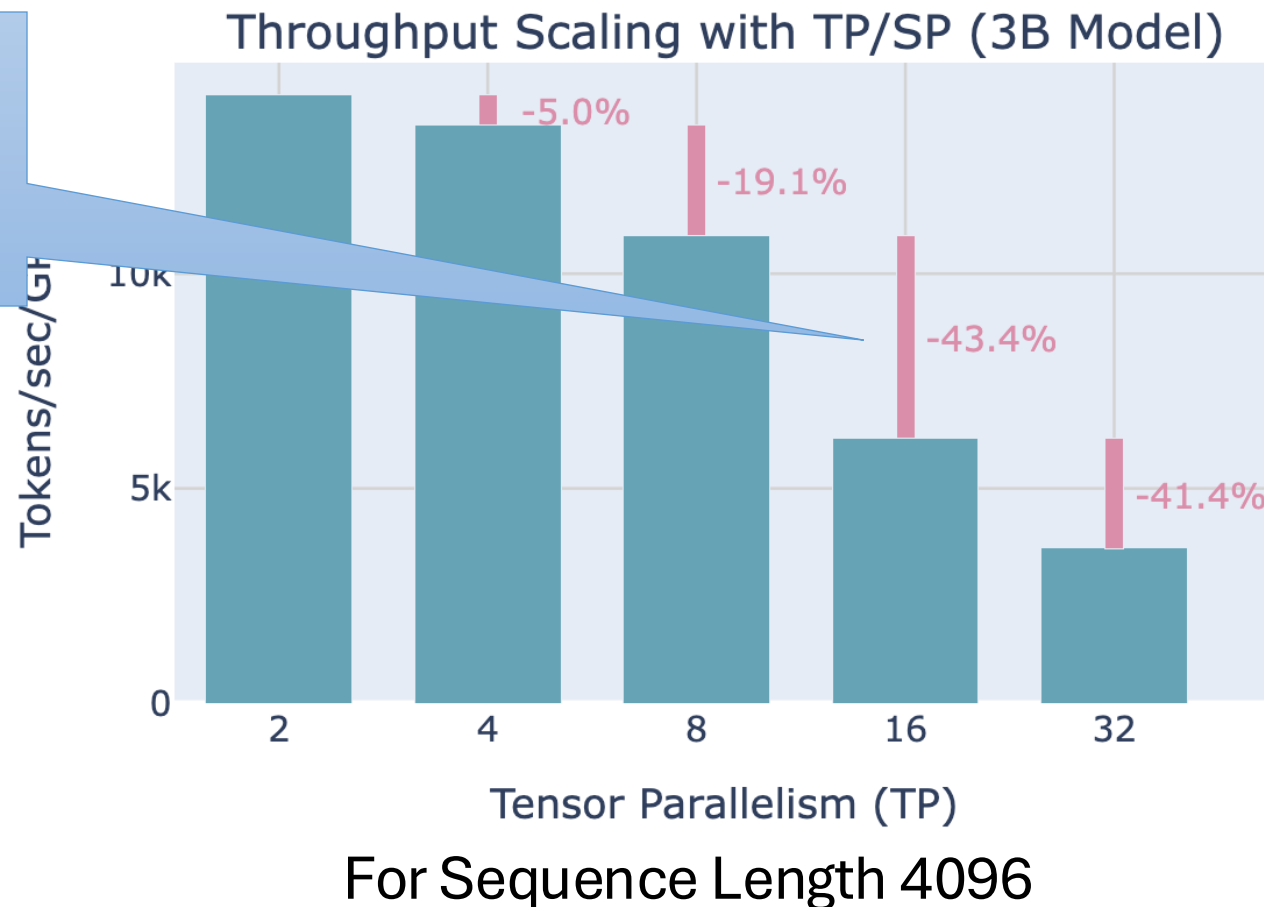
# Tensor+Sequence Parallelism - Tradeoffs

- Intermediate activations sharded across GPUs.
- **TP**: Reduces activation memory for matrix multiplication
- **SP**: Reduces activation memory for LayerNorm & dropout
- Need to gather full activations for LayerNorm
- Introduces significant communication overhead
- Introduces “*exposed communication*”

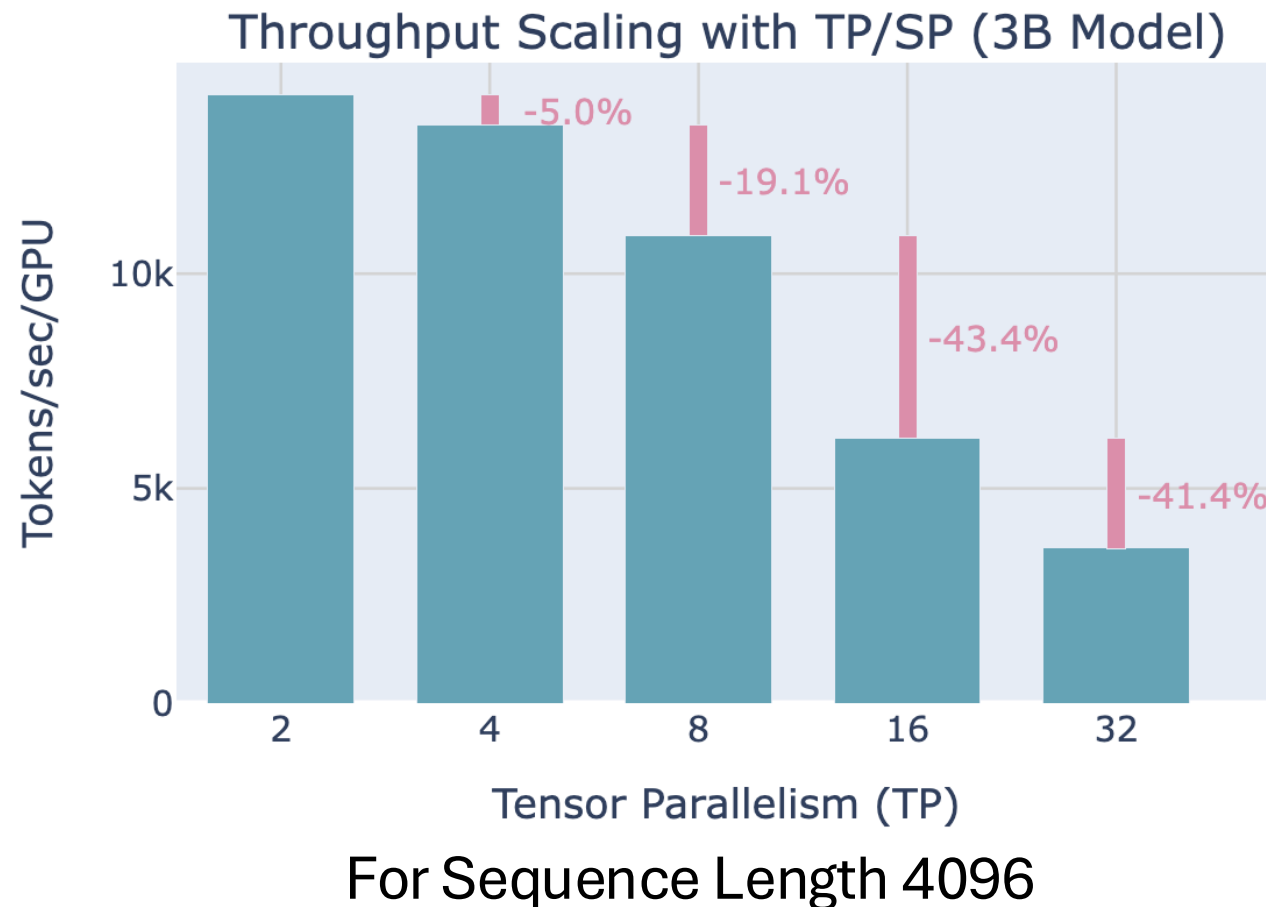
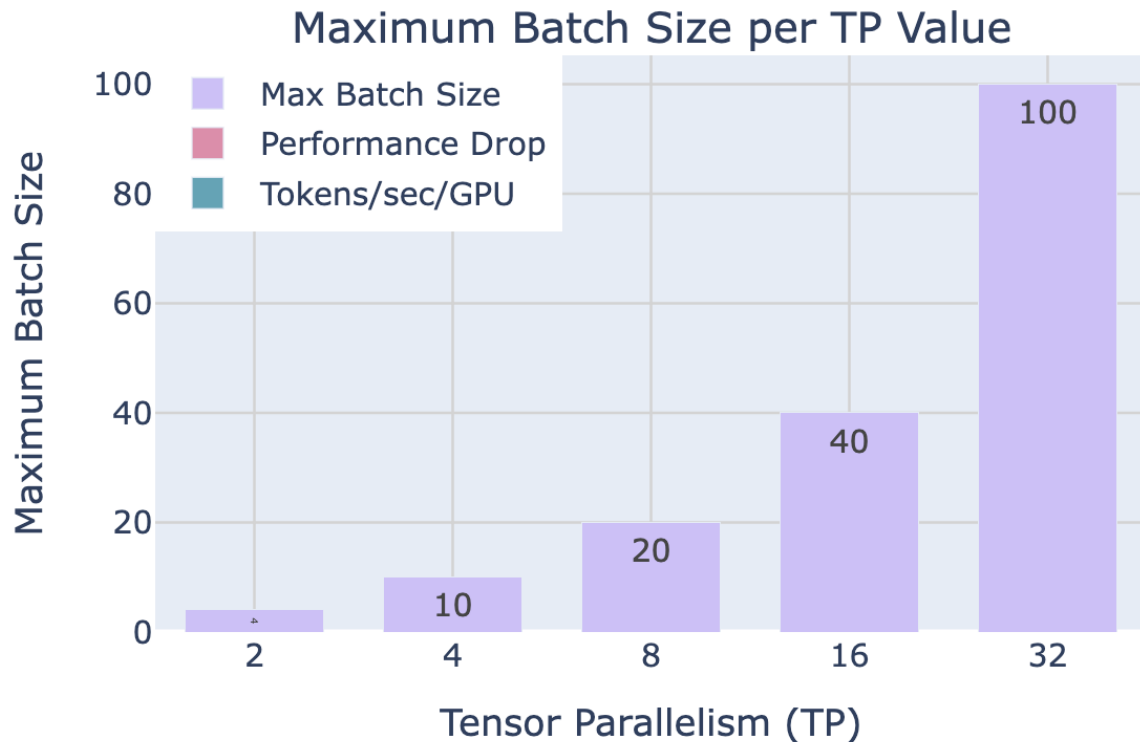


# Tensor+Sequence Parallelism - Tradeoffs

- TP leverages fast NVLink interconnects within a node.
- Slower network connections across nodes results in huge throughput drop.



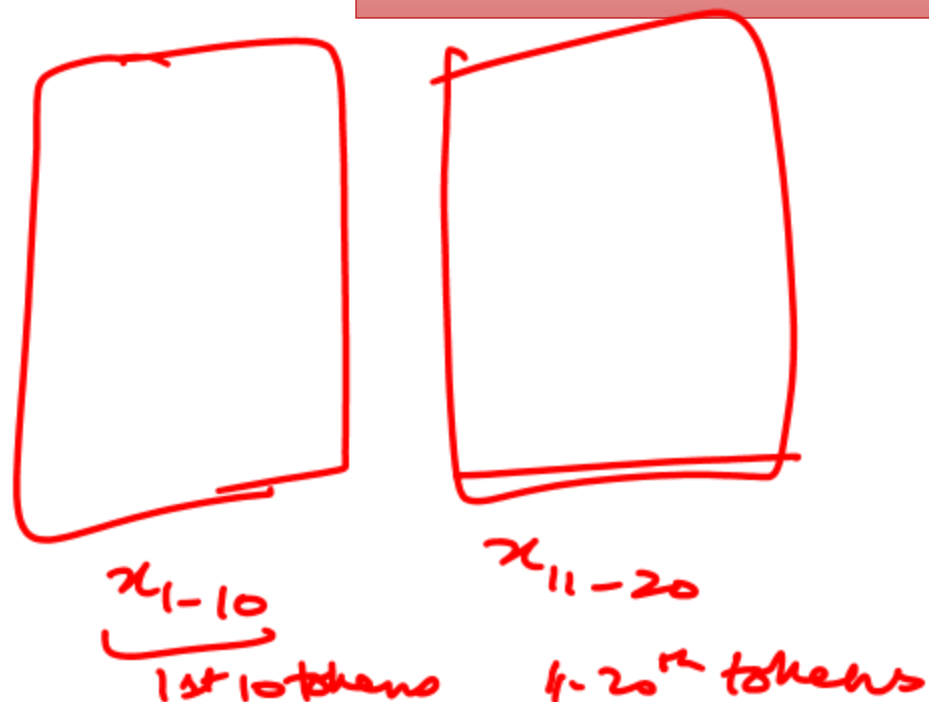
# Tensor+Sequence Parallelism - Tradeoffs



# Tensor+Sequence Parallelism - Limitations

1. If we scale the sequence length the activation memory will still blow up in the TP region

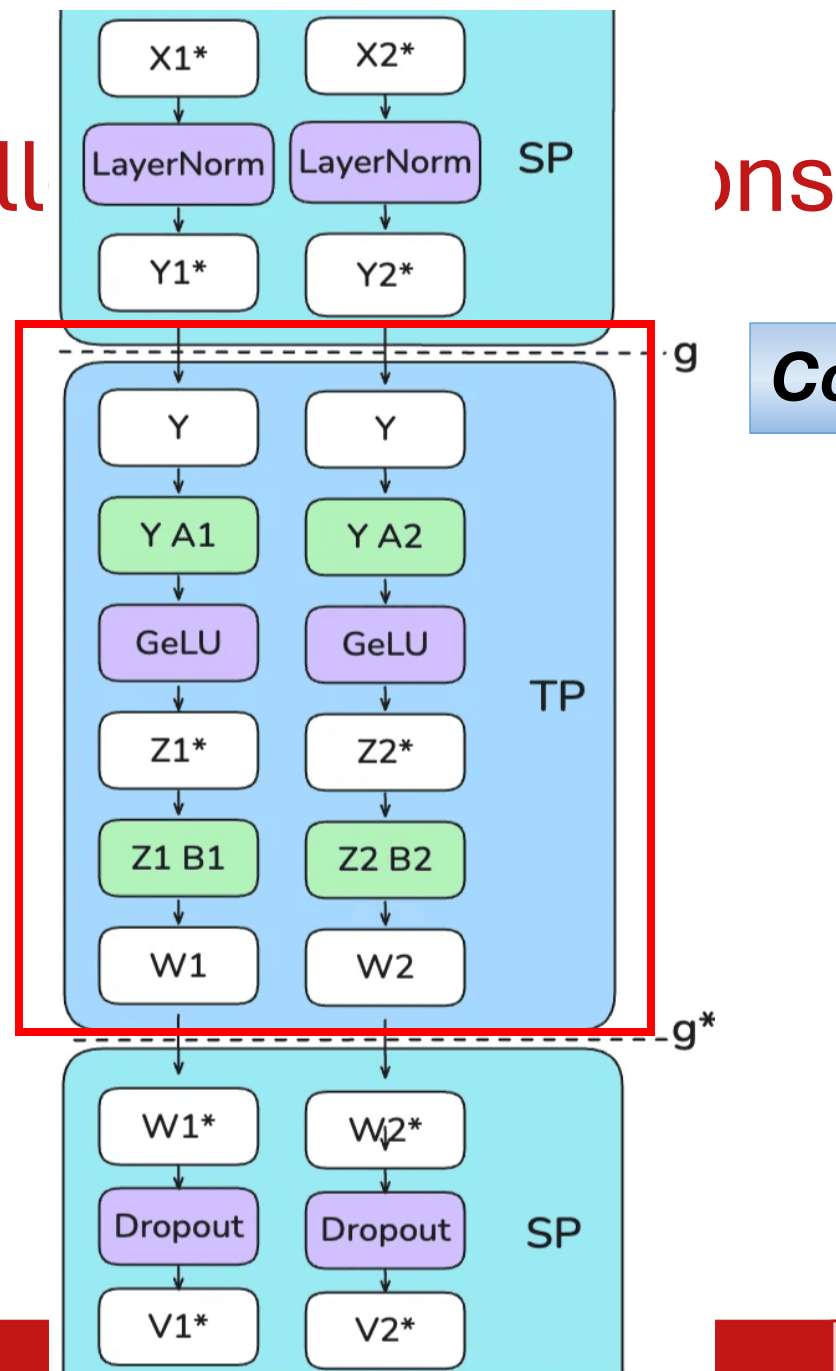
- In SP, we split one sequence into chunks and process each chunk in parallel.
- Can we do the same for MLP Layer?
- How about Attention Layer?





# Tensor+Sequence Parallel

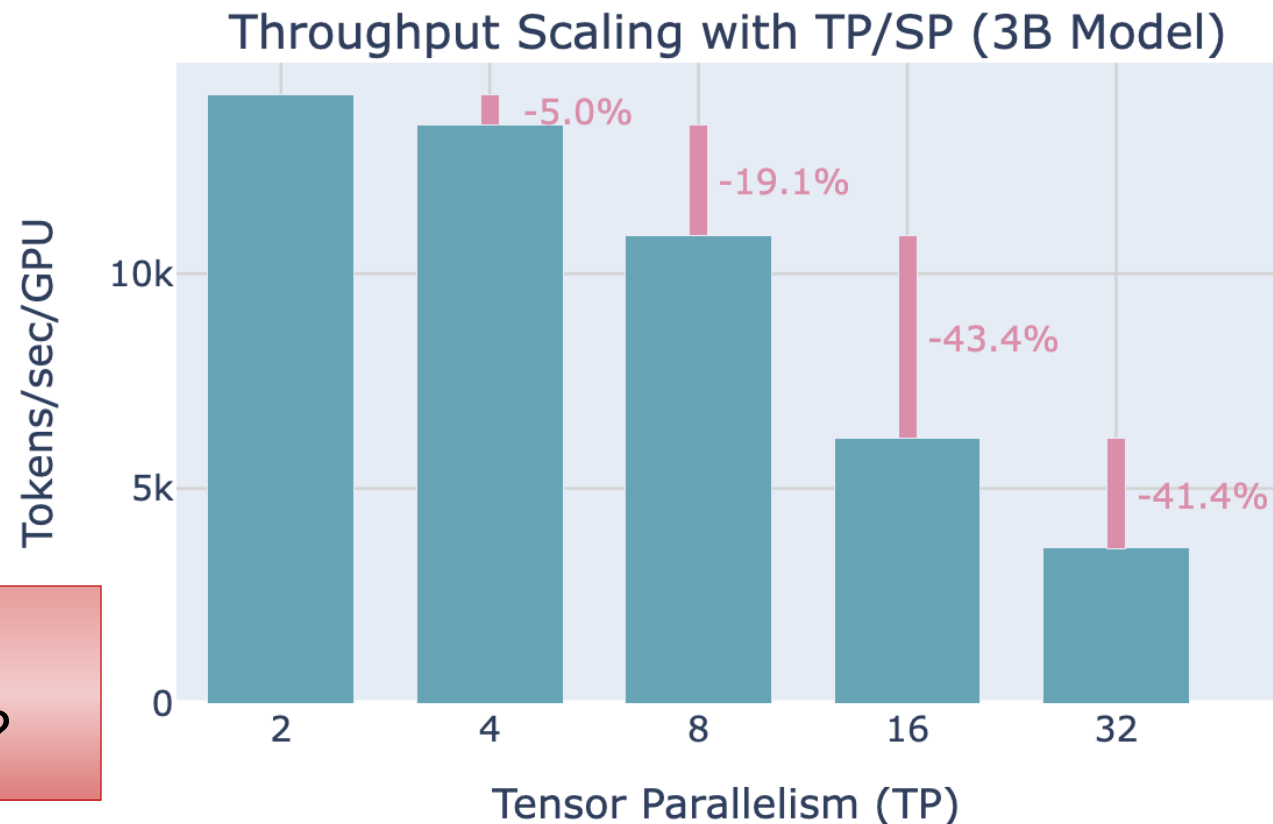
1. If we scale the sequence length the activation memory will still blow up in the TP region



**Context Parallelism**

# Tensor+Sequence Parallelism - Limitations

1. If we scale the sequence length the activation memory will still blow up in the TP region
2. If the model is too big to fit with TP=8 we will see a massive slowdown due to the inter-node connectivity.
  - In TP, we split a “Tensor” across GPUs.
  - How about splitting layers across GPUs?



**Pipeline Parallelism**

# Tensor+Sequence Parallelism - Limitations

1. If we scale the sequence length the activation memory will still blow up in the TP region
2. If the model is too big to fit with TP=8 we will see a massive slowdown due to the inter-node connectivity.

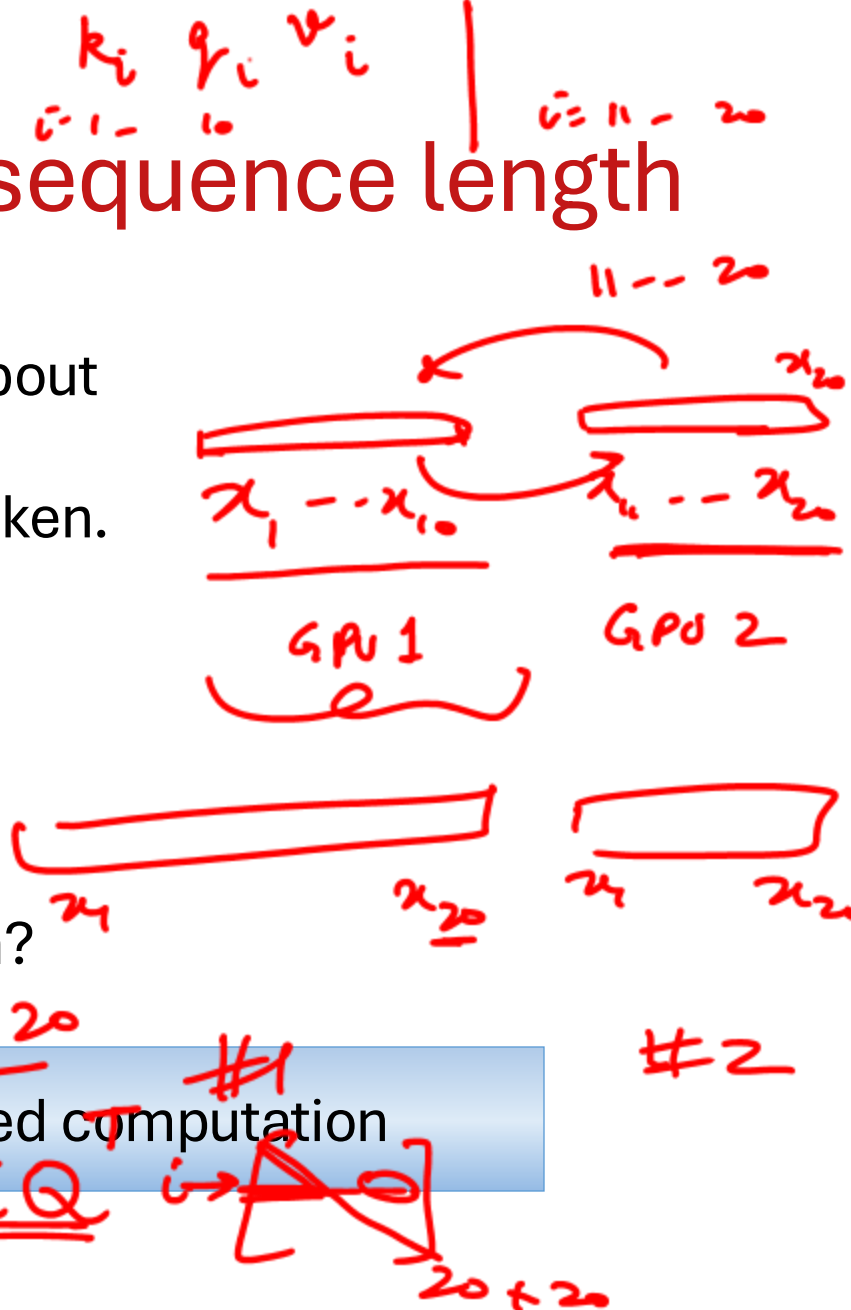
***Context Parallelism***

***Pipeline Parallelism***



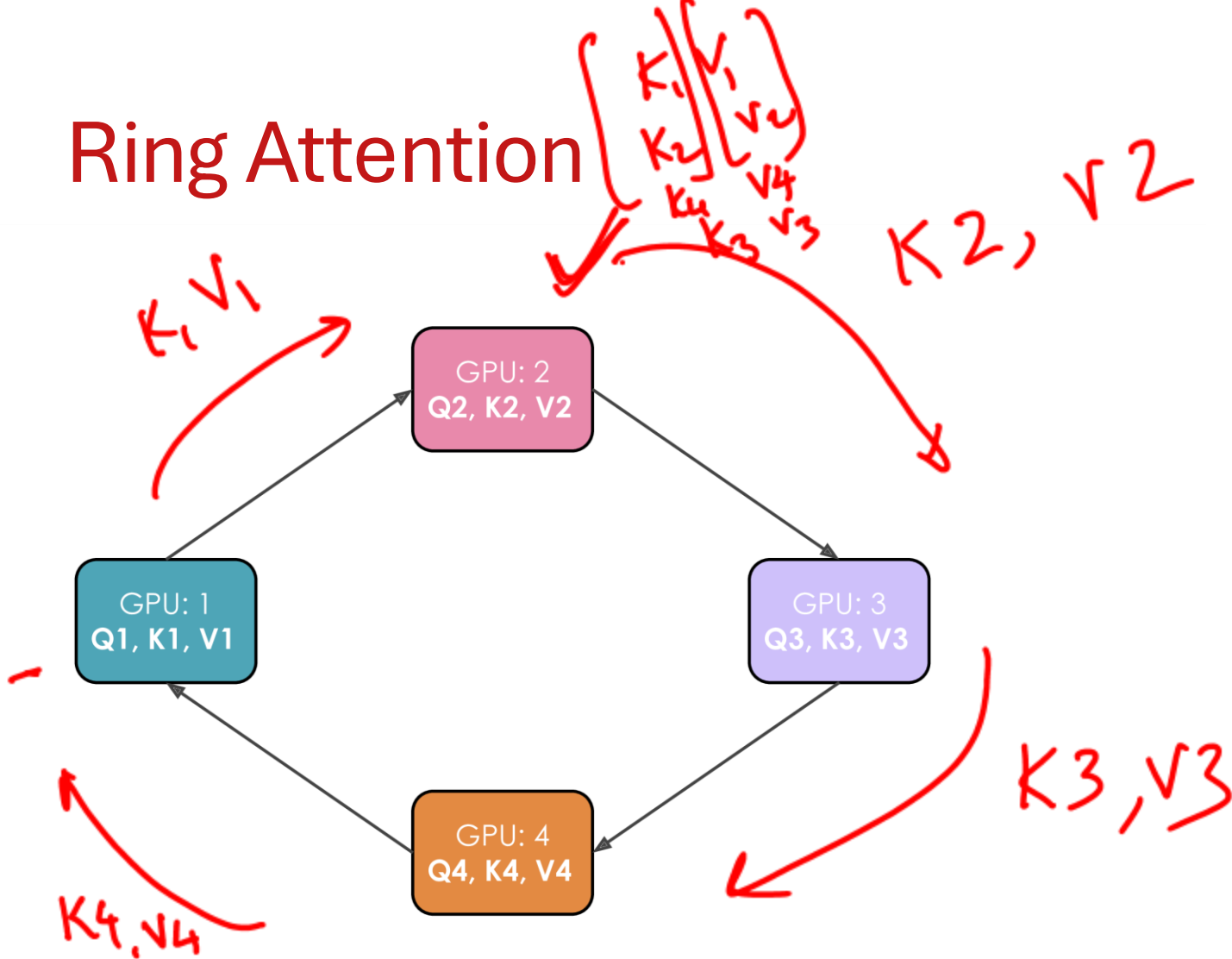
# Context Parallelism – partition along sequence length

- For MLP layers, its exactly same as SP for LayerNorm & Dropout
- In Attention Layer, each token has to attend on every other token.
- But tokens in a different chunk are on a different GPU!
- Full communication b/w GPUs?
- Can we somehow overlap computation with communication?



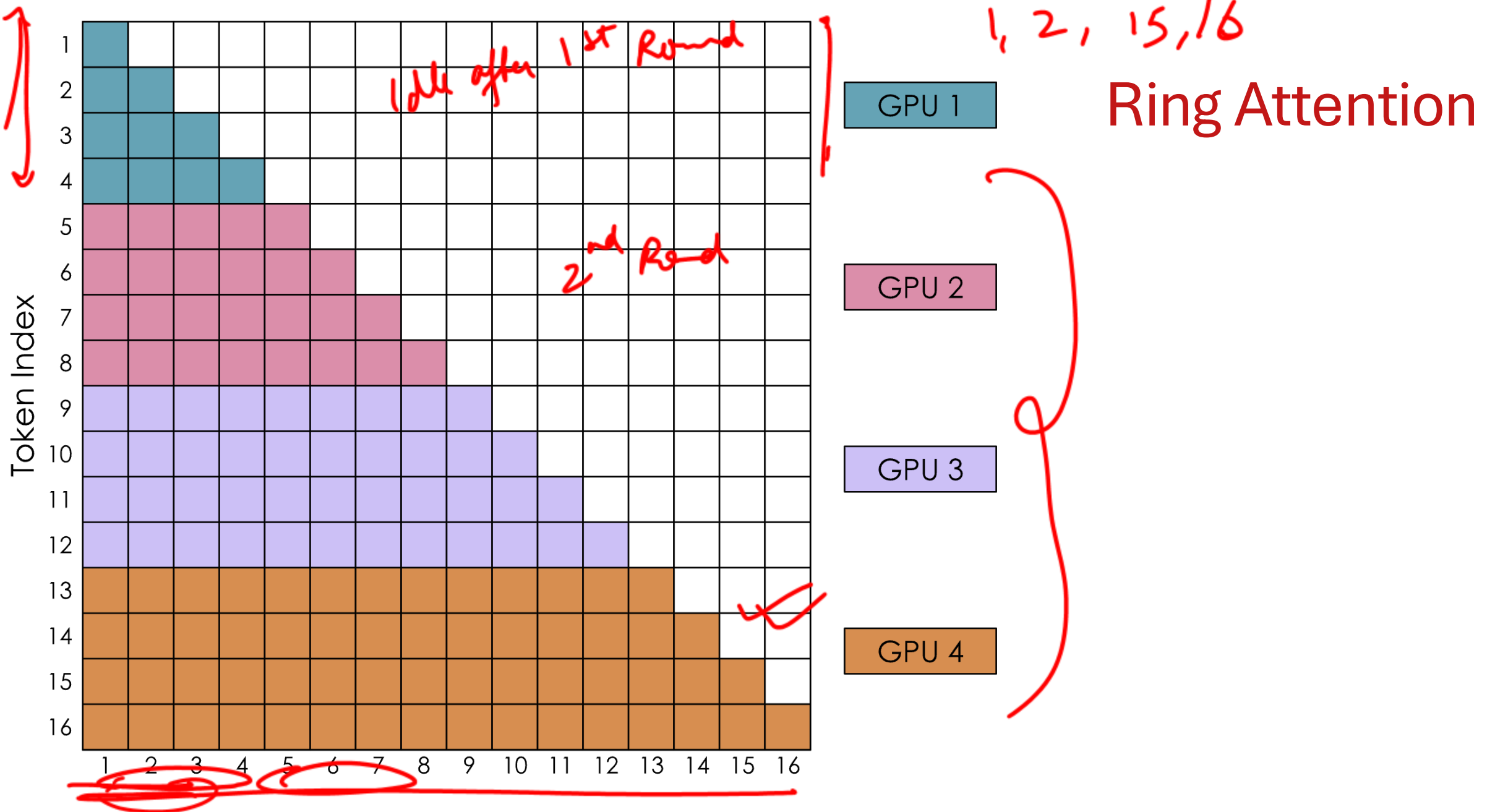
Ring Attention - online softmax computation + overlapped computation

# Ring Attention

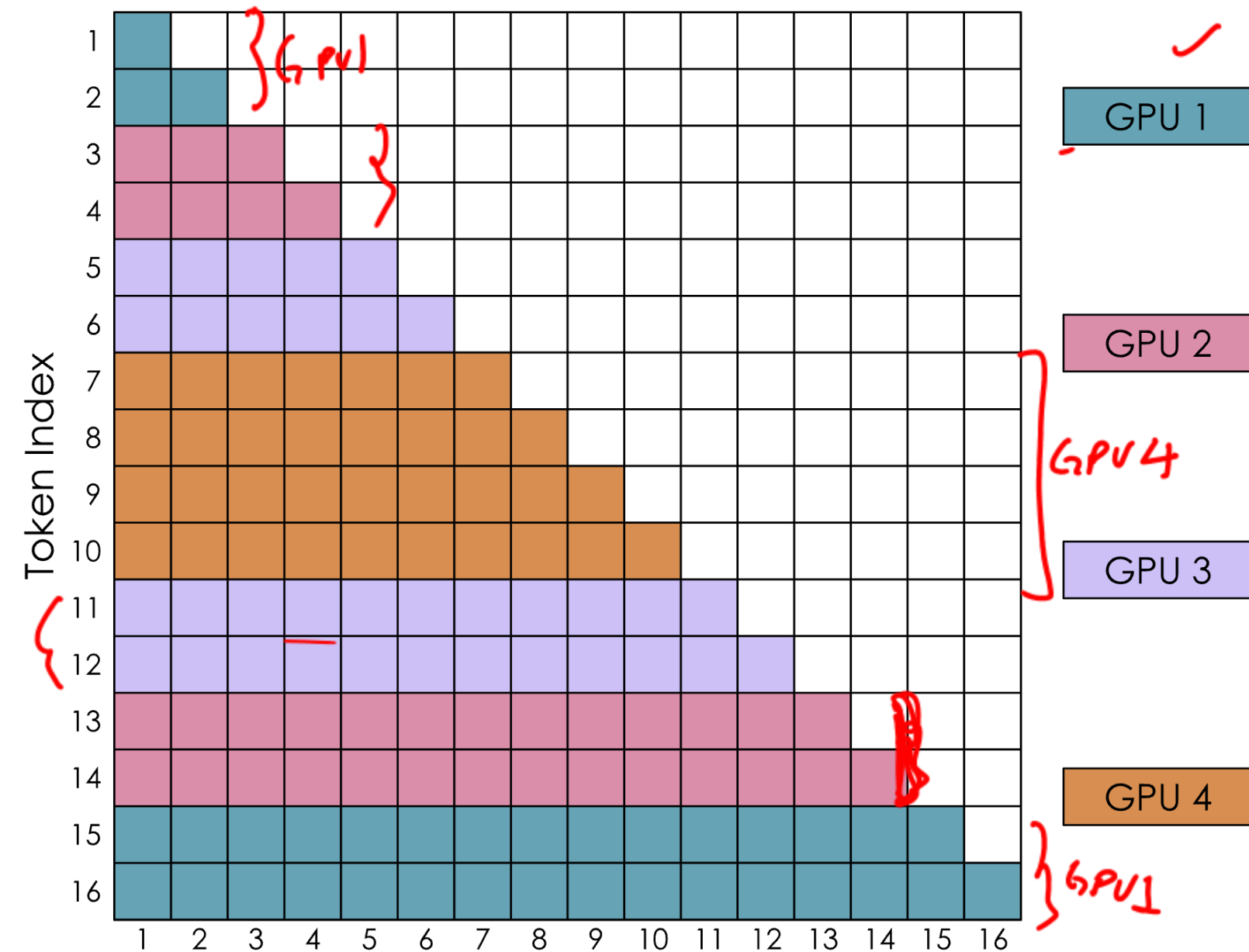


Ring Attention with Blockwise Transformers for Near-Infinite Context H. Liu, M. Zaharia, P. Abbeel. 2023 [\[PDF\]](#)





# Ring Attention



# Computation-communication timeline



## All-to-all (ring) implementation:

- GPUs exchange K/V pairs in a ring-like pattern, one chunk at a time.
- More memory-efficient, as each GPU only needs to store one additional chunk temporarily.
- Communication is spread out and overlapped with computation, though with some additional base latency overhead from multiple communication steps.



# Comparing with Naïve *all-gather* Implementation

## Attention

### All-gather implementation

GPU Computation:

$\text{Attn}(Q_i, K_i, V_i)$

$\text{Attn}(Q_i, K_{i+1}, V_{i+1})$

$\text{Attn}(Q_i, K_{i+2}, V_{i+2})$

AllGather Activs

GPU Communication:

$\text{AG}(K, V)$

Forward pass

## Attention

### All-to-all (ring) implementation

GPU Computation:

$\text{Attn}(Q_i, K_i, V_i)$

$\text{Attn}(Q_i, K_{i+1}, V_{i+1})$

$\text{Attn}(Q_i, K_{i+2}, V_{i+2})$

P2P Activs

GPU Communication:

Fetch  $K_{i+1}, V_{i+1}$

Fetch  $K_{i+2}, V_{i+2}$

Fetch  $K_{i+3}, V_{i+3}$

Forward pass



# Tensor+Sequence Parallelism - Limitations

1. If we scale the sequence length the activation memory will still blow up in the TP region

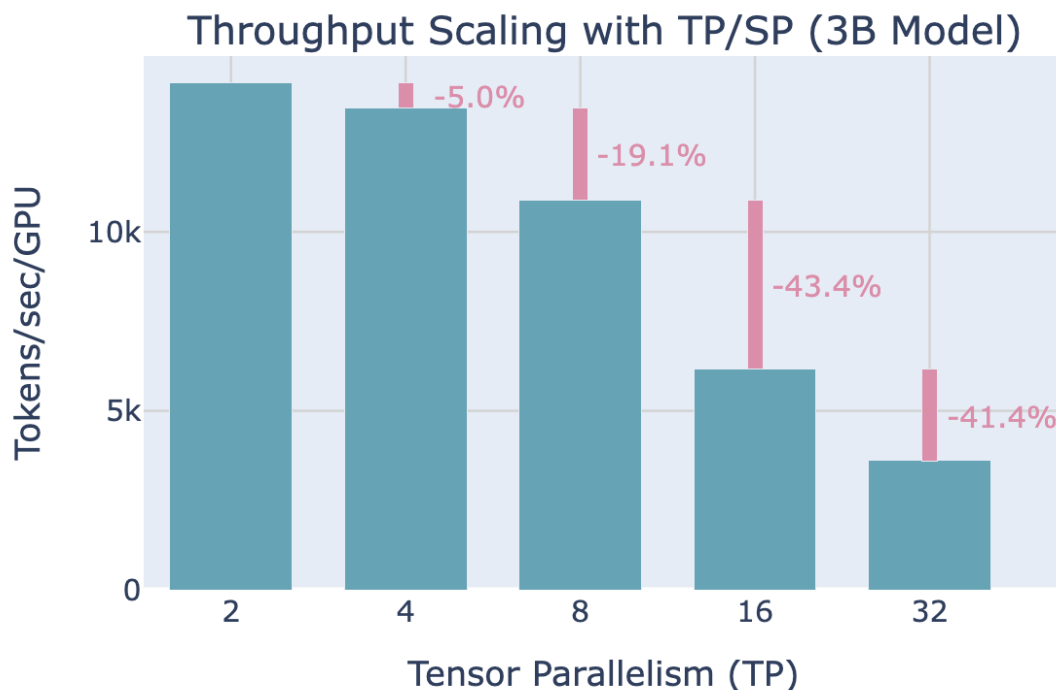
## ***Context Parallelism***

- TP: Split a model across one node to tame large models
  - CP: Tame the activation explosion with long sequences.
- 
- TP: doesn't scale well across nodes



# Tensor+Sequence Parallelism - Limitations

1. If we scale the sequence length the activation memory will still blow up in the TP region



## ***Context Parallelism***

- TP: Split a model across one node to tame large models
- CP: Tame the activation explosion with long sequences.

- TP: doesn't scale well across nodes

# Tensor+Sequence Parallelism - Limitations

1. If we scale the sequence length the activation memory will still blow up in the TP region

## ***Context Parallelism***

- TP: Split a model across one node to tame large models
  - CP: Tame the activation explosion with long sequences.
- 
- TP: doesn't scale well across nodes
  - How about splitting layers across GPUs?



# Tensor+Sequence Parallelism - Limitations

1. If we scale the sequence length the activation memory will still blow up in the TP region
2. If the model is too big to fit with TP=8 we will see a massive slowdown due to the inter-node connectivity.

## ***Pipeline Parallelism***

## ***Context Parallelism***

- TP: Split a model across one node to tame large models
- CP: Tame the activation explosion with long sequences.

- TP: doesn't scale well across nodes
- How about splitting layers across GPUs?

