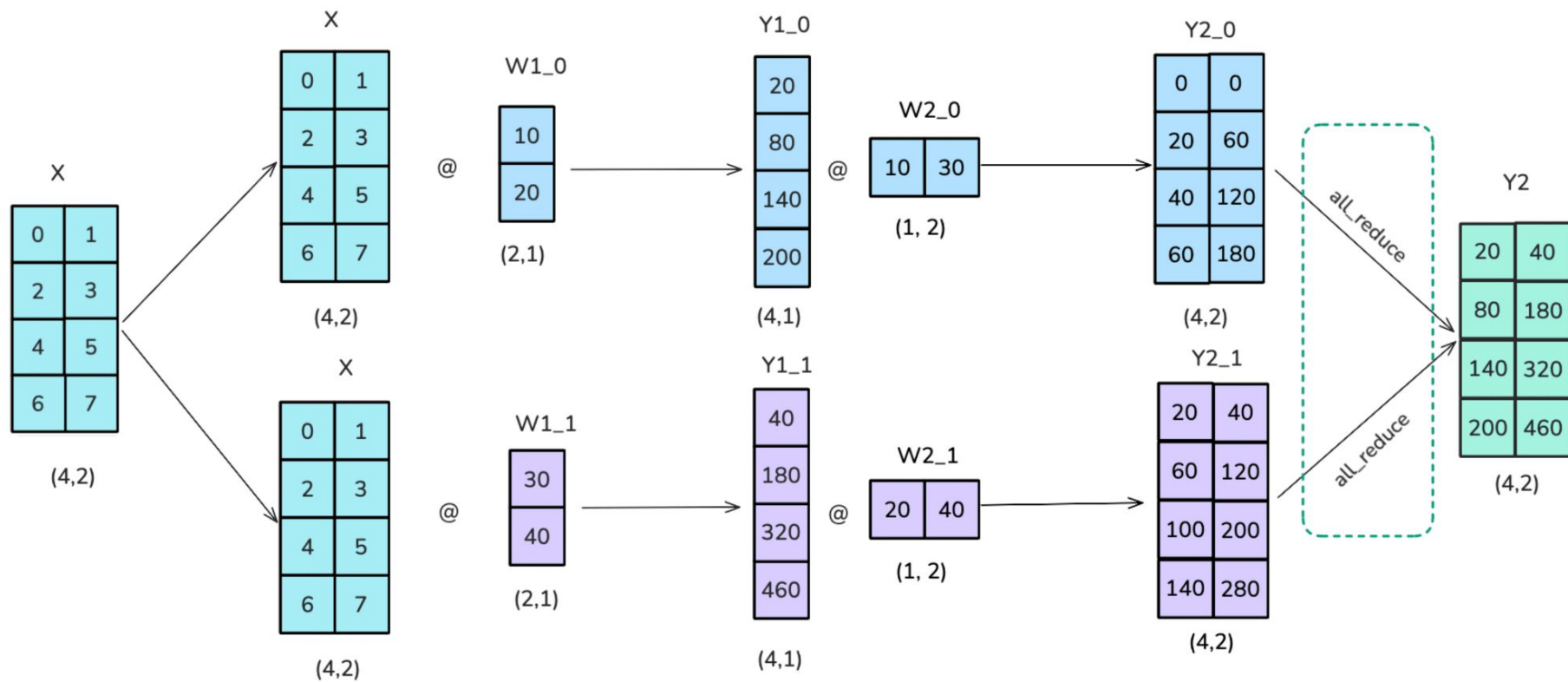


# Sequence Parallelism

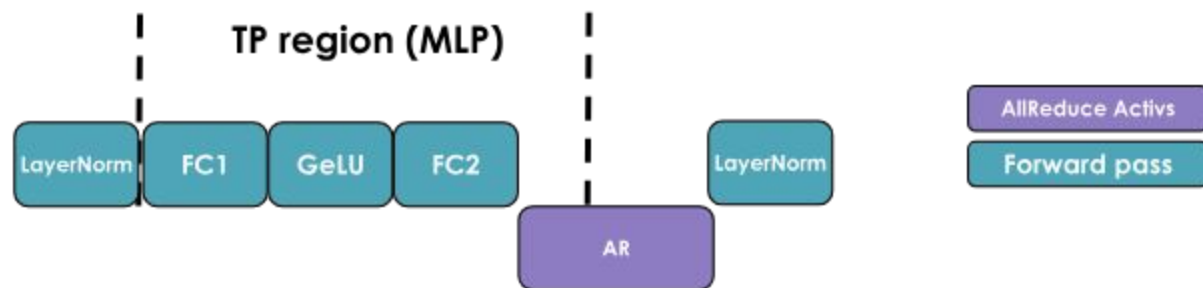
made with  for “Little ML book club”



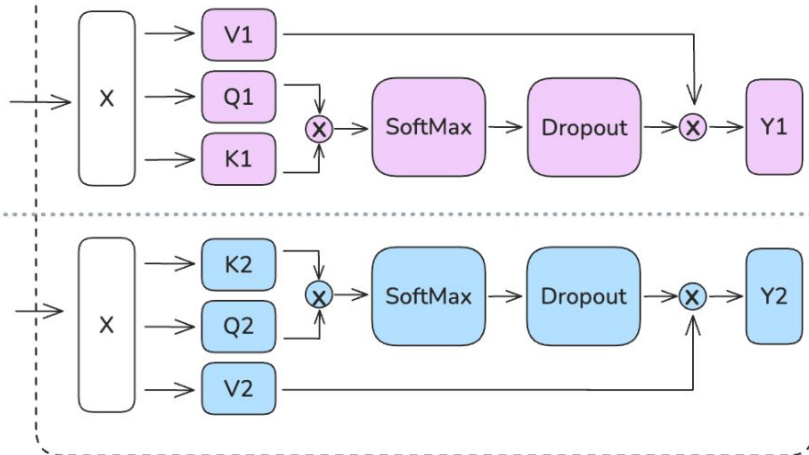
Tensor parallelism with column linear + row Linear

GPU Computation:

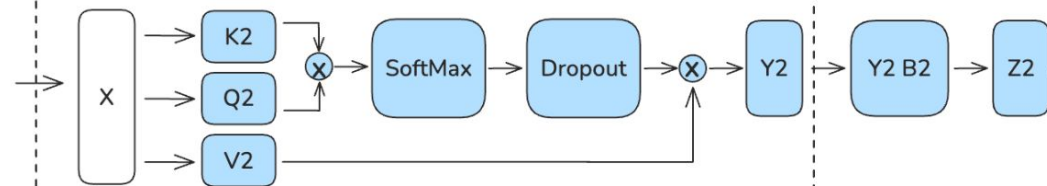
GPU Communication:



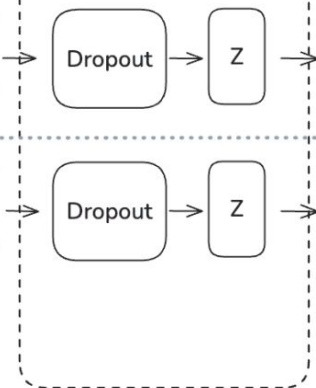
GPU1

 $Y = \text{SelfAttention}(X)$ 

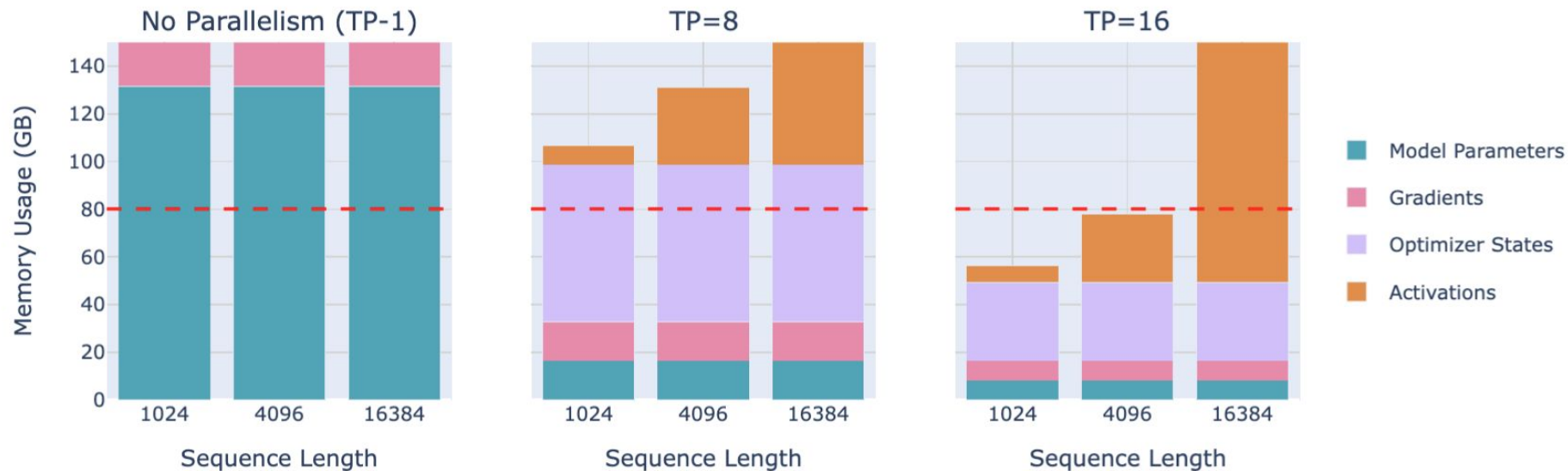
GPU2



AllReduce

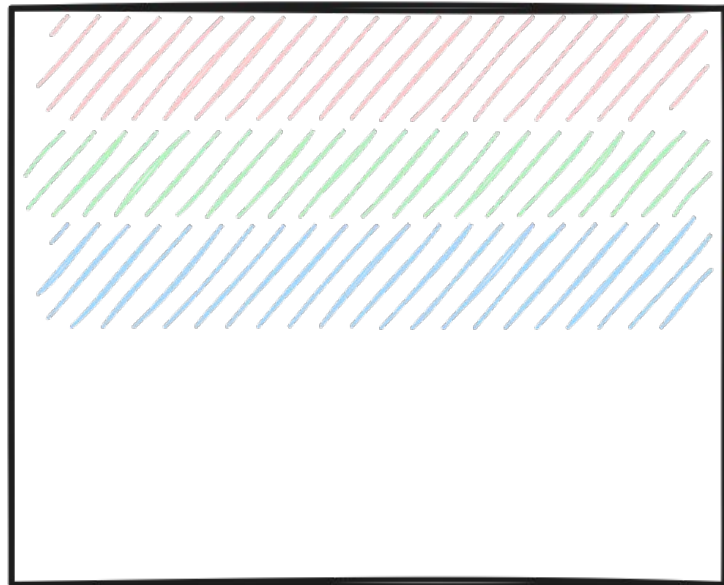
 $Z = \text{Dropout}(Y B)$ 

## Memory Usage for 70B Model

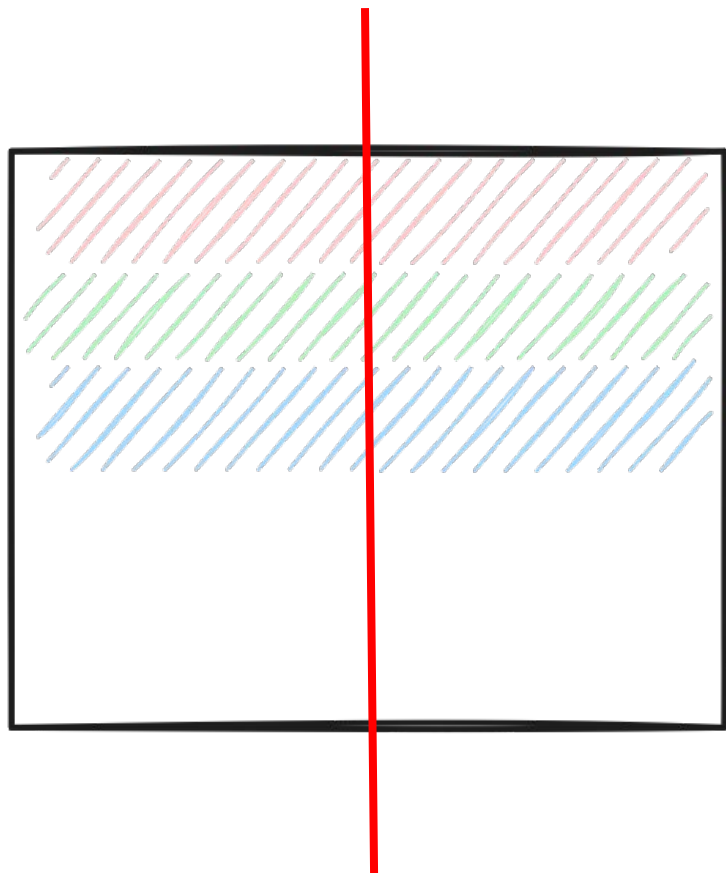


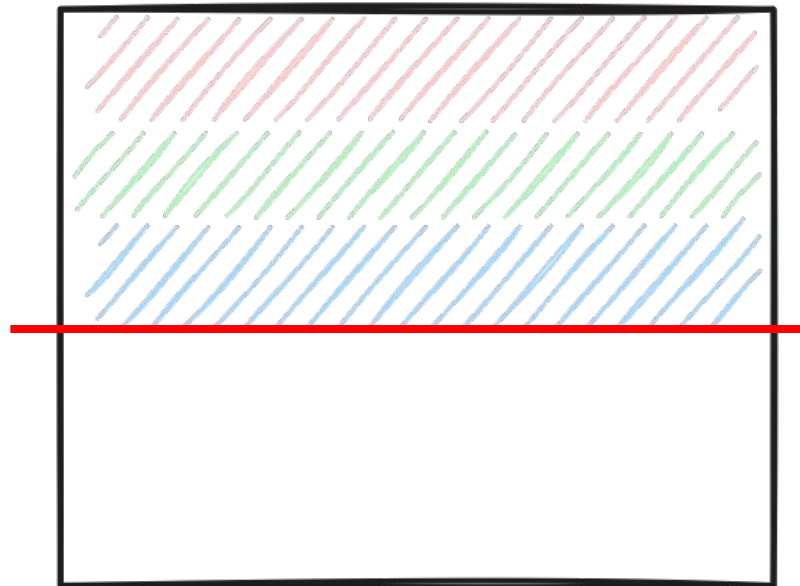
- layer norm
- dropout

- layer norm
- dropout ???



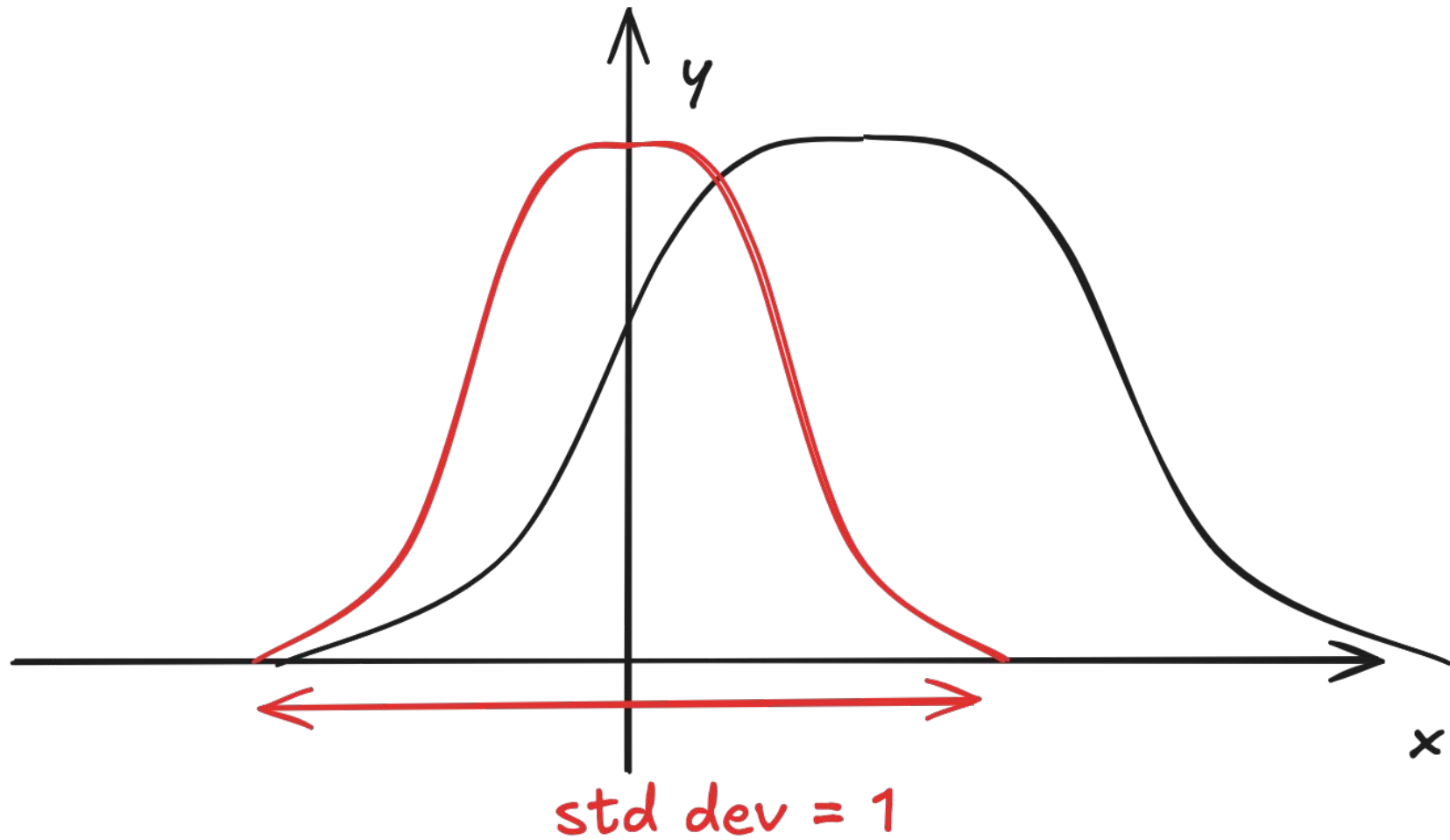






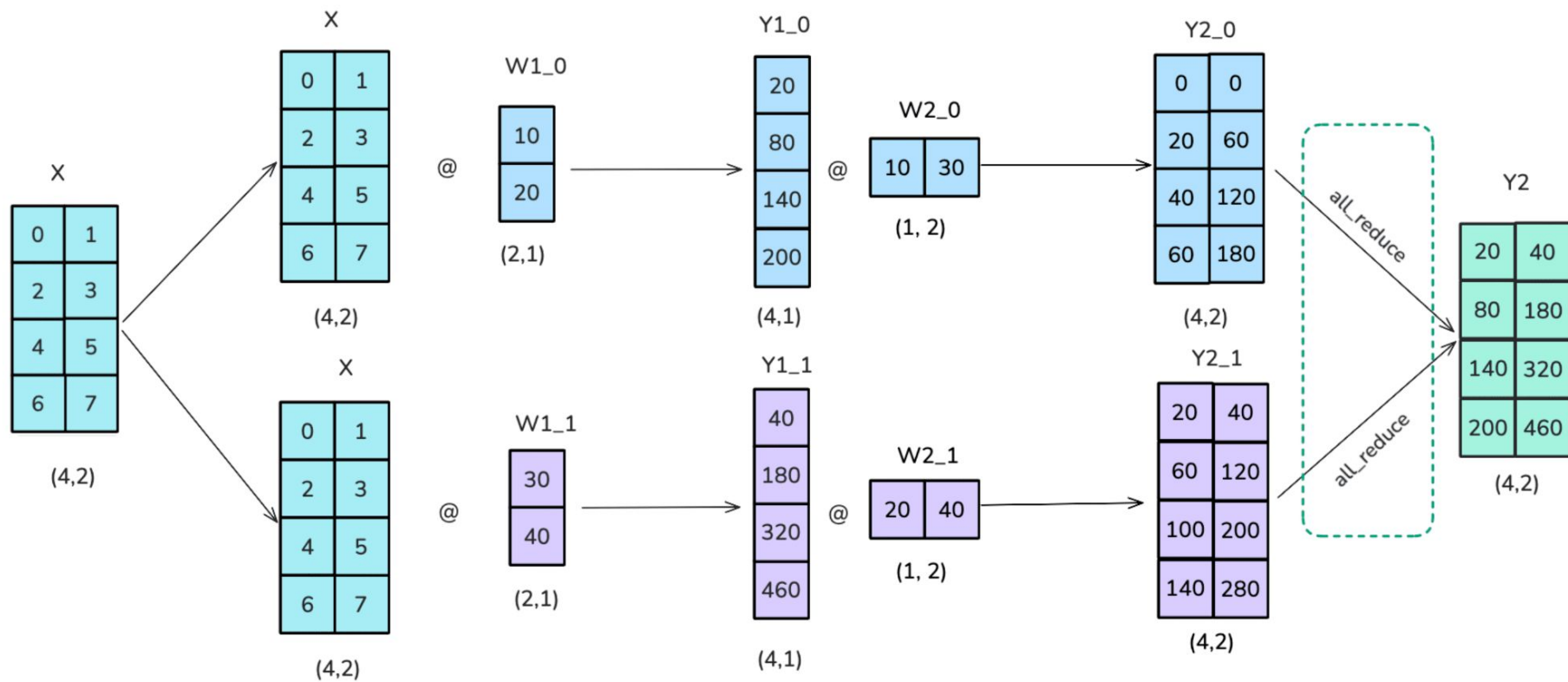
$$\text{LayerNorm}(x) = \gamma \cdot \frac{x - \mu}{\sqrt{\sigma^2 + \epsilon}} + \beta$$

where  $\mu = \text{mean}(x)$  and  $\sigma^2 = \text{var}(x)$  are computed across hidden dimension  $h$ .

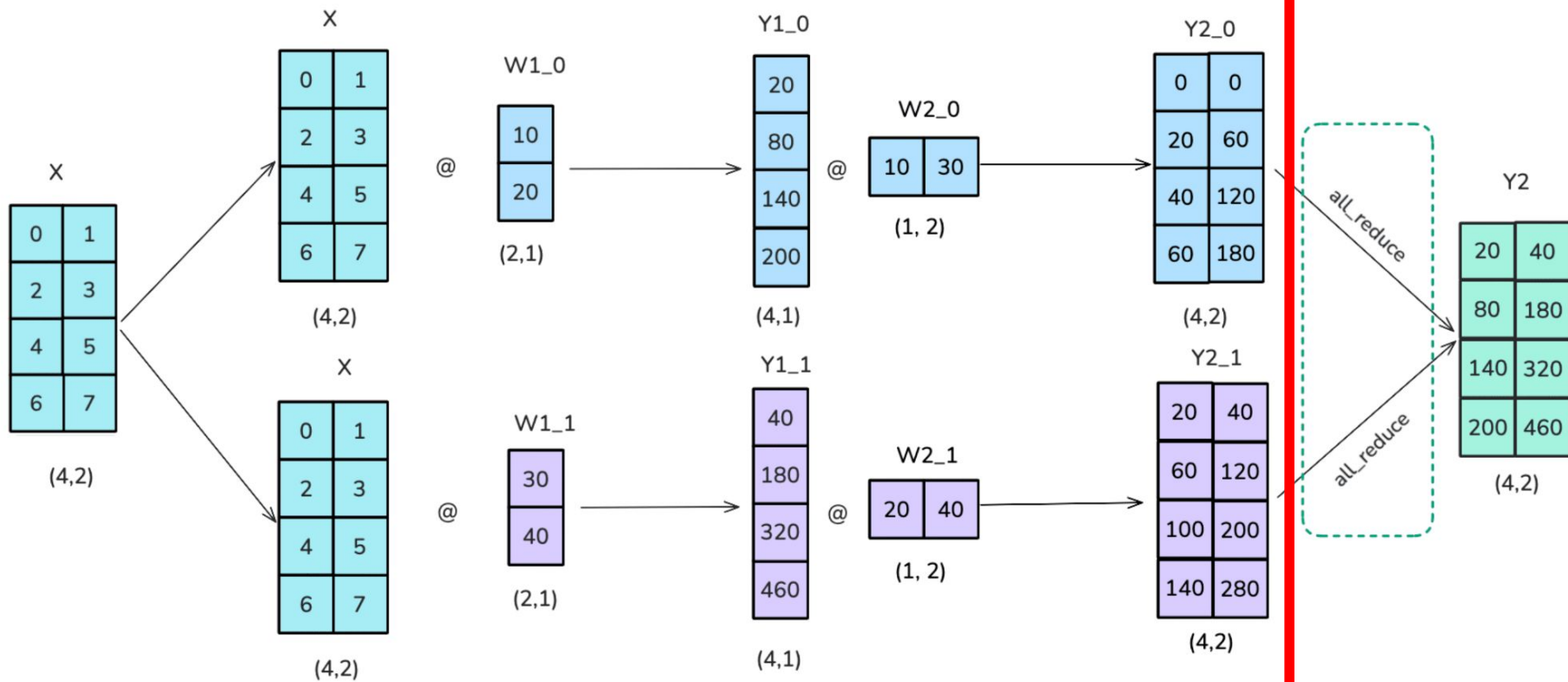




```
(batch, seq_len, d_model)
```



Tensor parallelism with column linear + row Linear



Tensor parallelism with column linear + row Linear



```
Full tensor: (batch, seq_len, d_model)
```

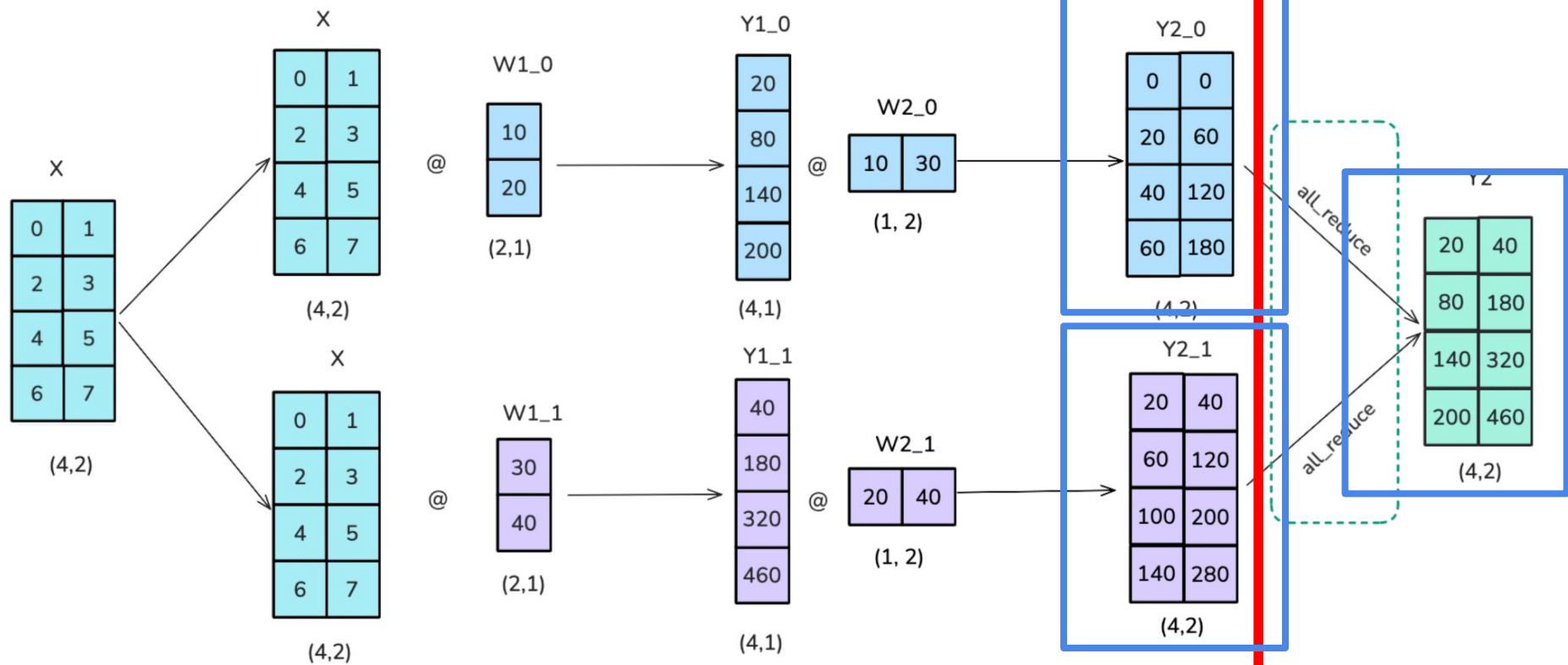
↑

SHARDED across GPUs

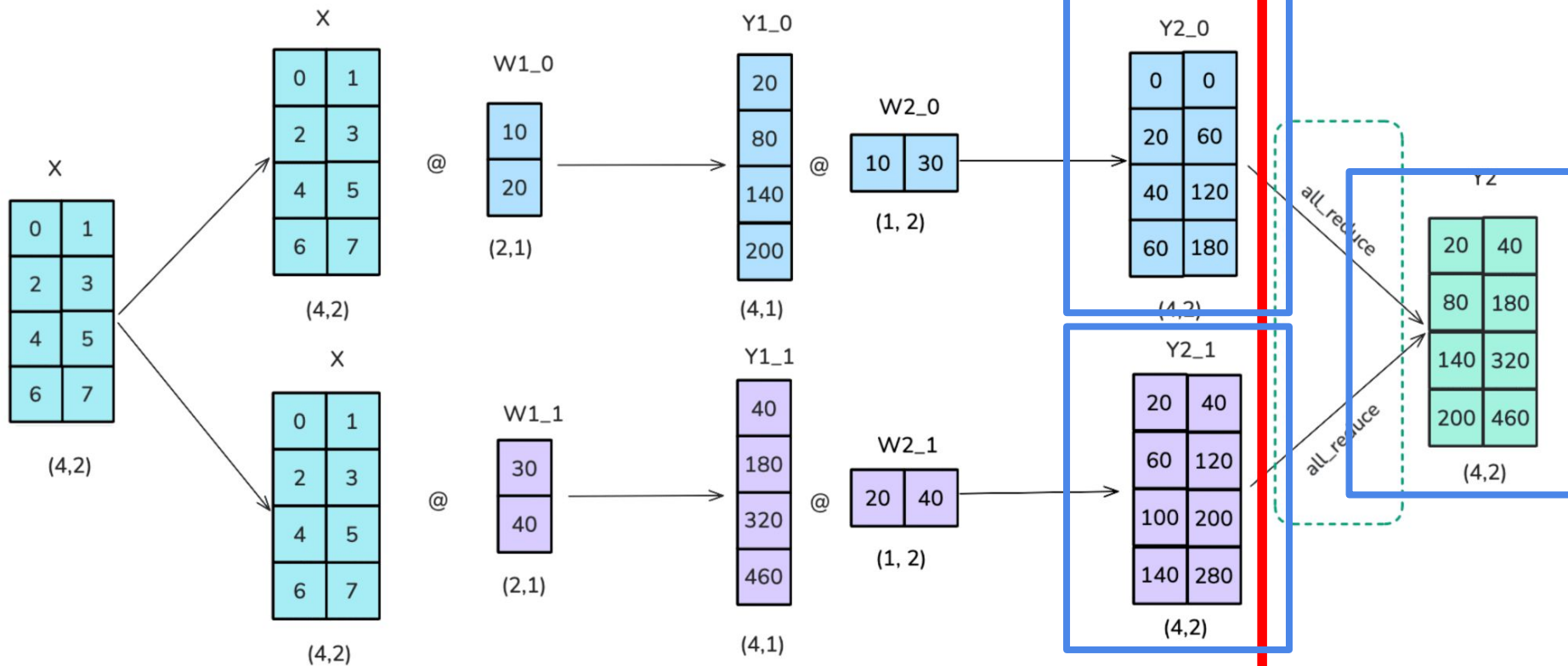
```
GPU 0: x[:, :, 0:D//2]      # all tokens, first half of features  
GPU 1: x[:, :, D//2:D]     # all tokens, second half of features
```



hmmm....



Tensor parallelism with column linear + row Linear



Tensor parallelism with column linear + row Linear

size of  $Y2_0$  == size of  $YZ$

	Standard TP (The Diagram)	Sequence Parallelism
Operation	All-Reduce	Reduce-Scatter
GPU 0 Output	$\begin{bmatrix} 20 & 40 \\ 80 & 180 \\ 140 & 320 \\ 200 & 460 \end{bmatrix}$	$\begin{bmatrix} 20 & 40 \\ 80 & 180 \end{bmatrix}$
GPU 1 Output	$\begin{bmatrix} 20 & 40 \\ 80 & 180 \\ 140 & 320 \\ 200 & 460 \end{bmatrix}$	$\begin{bmatrix} 140 & 320 \\ 200 & 460 \end{bmatrix}$
Memory Used	100% (Full Matrix)	50% (Sharded Matrix)



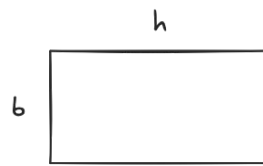
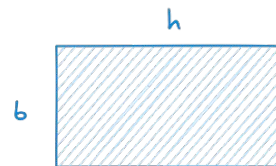
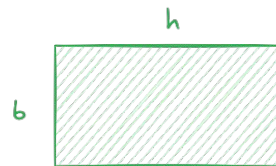
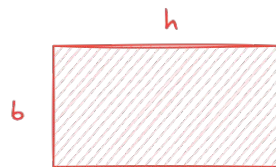
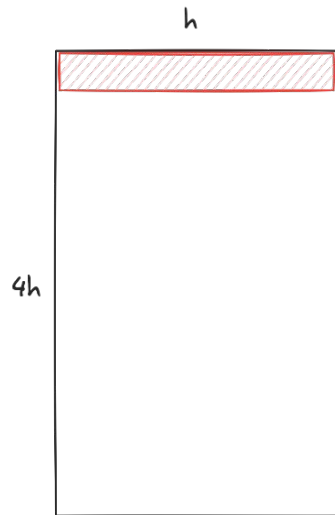
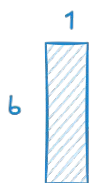
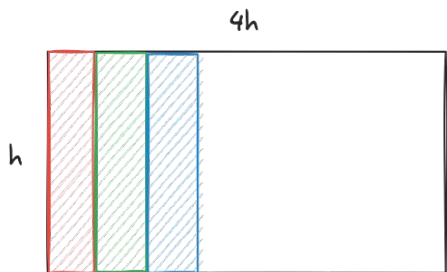
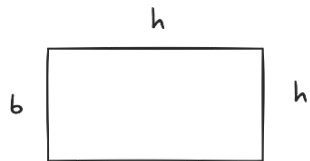
Full tensor: (batch, seq\_len, d\_model)

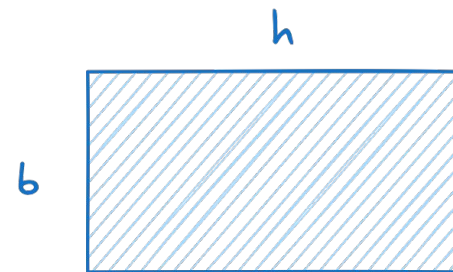
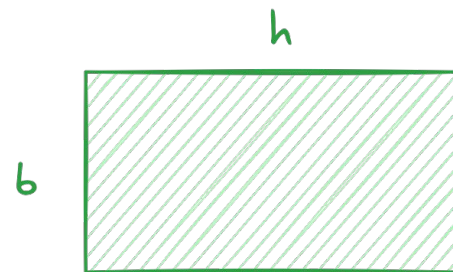
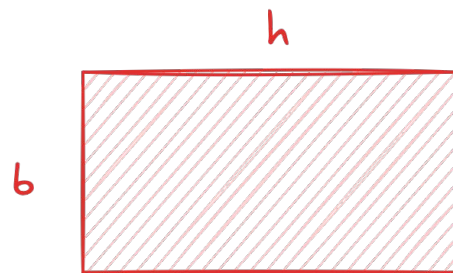
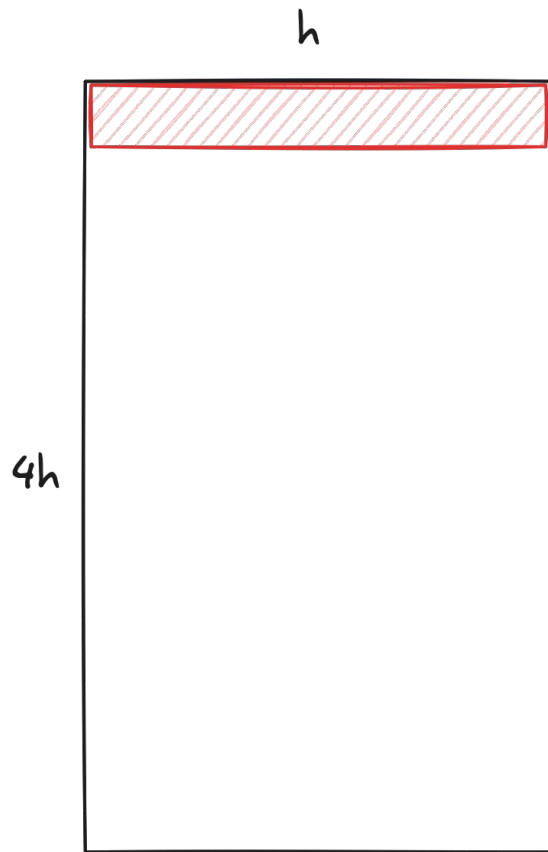
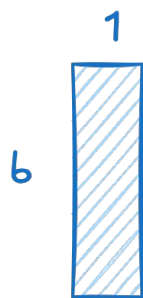
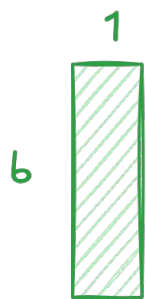
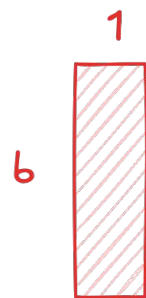
↑

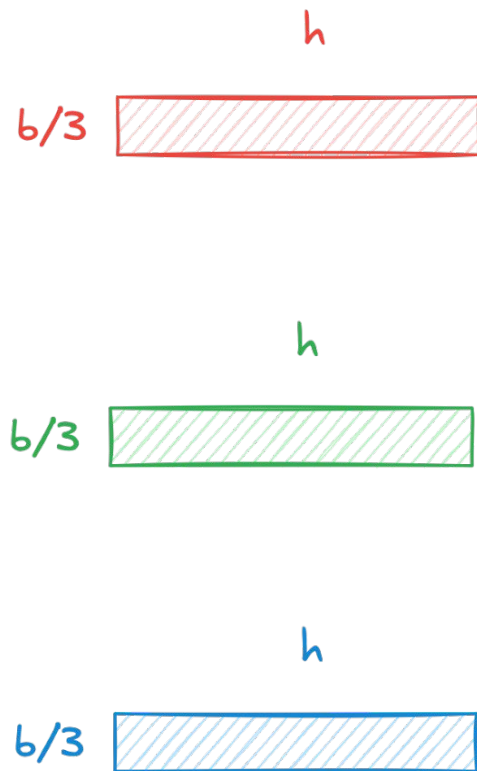
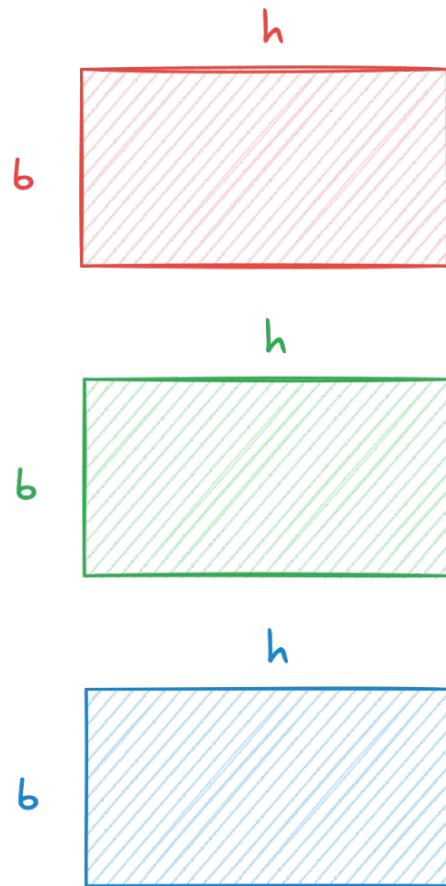
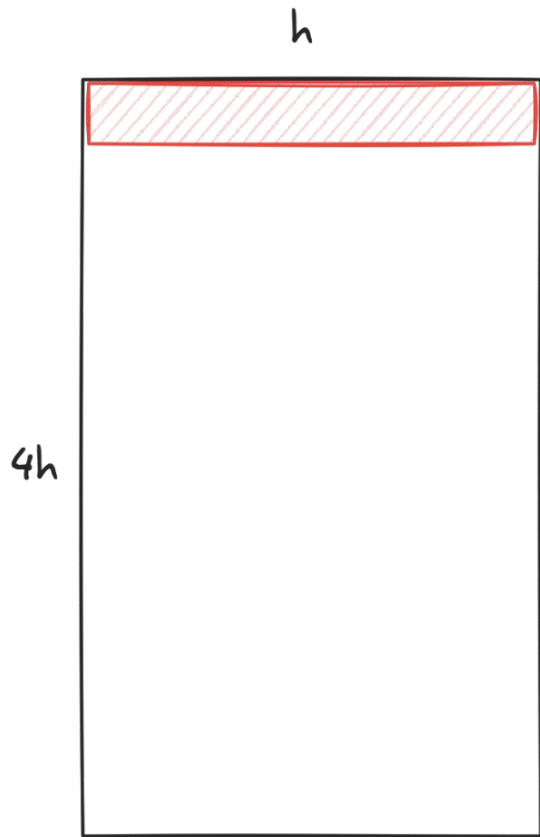
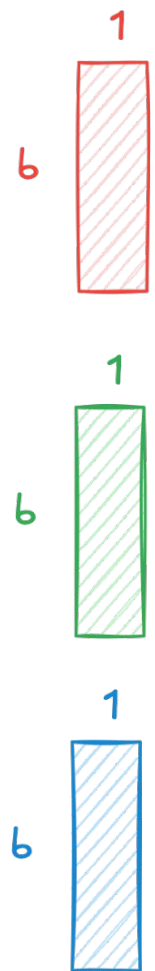
SHARDED across GPUs

GPU 0: x[:, :, 0:2D]           # all tokens, partials of all features

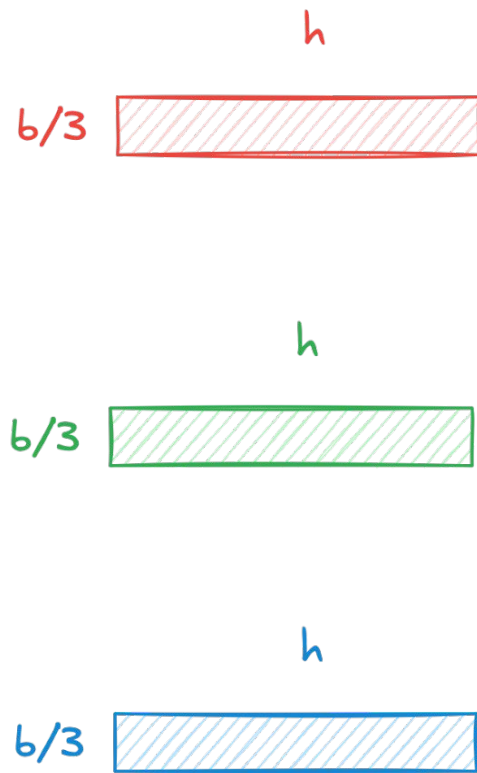
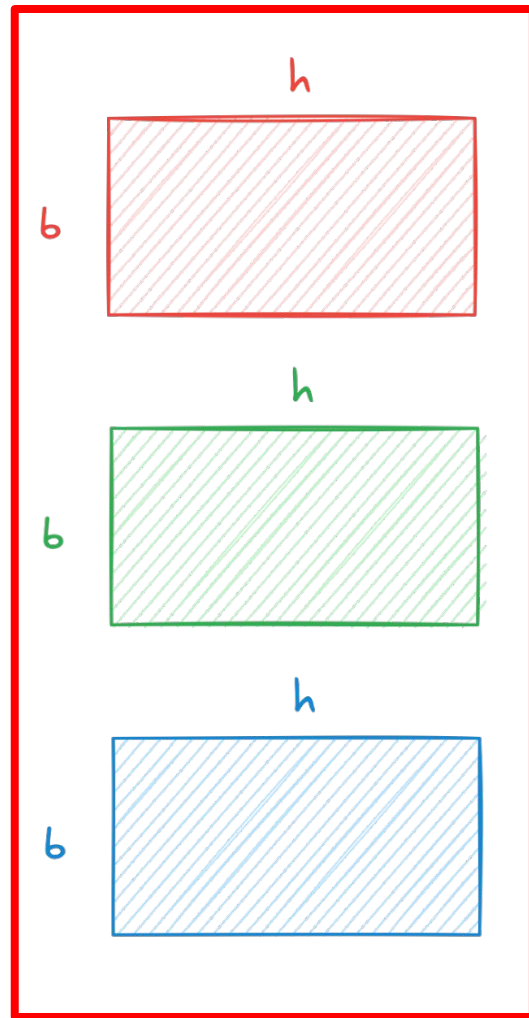
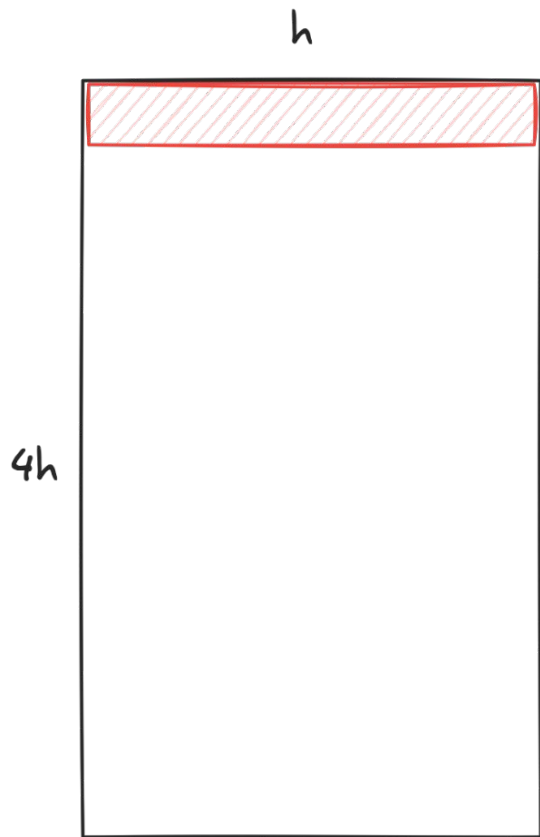
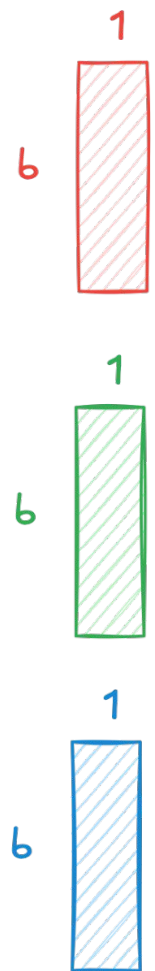
GPU 1: x[:, :, 0:2D]           # all tokens, partials of all features

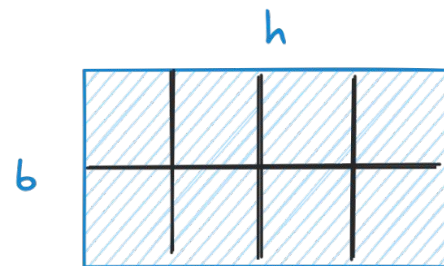
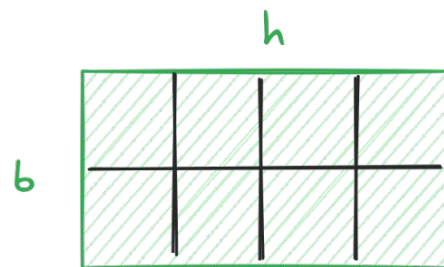
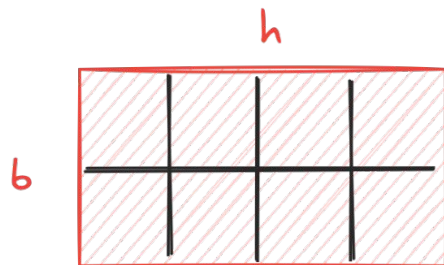
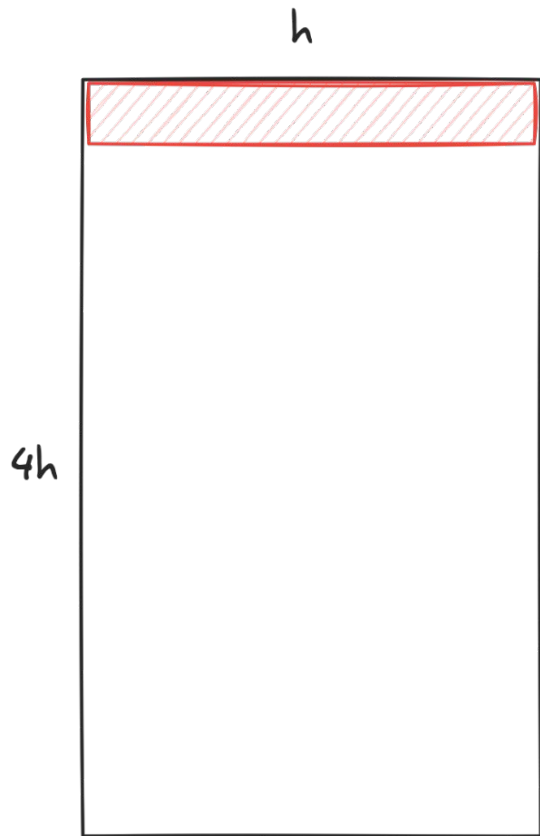
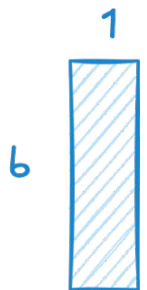
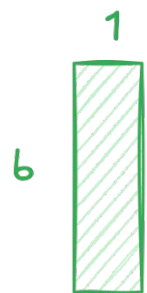
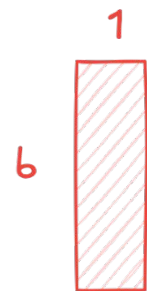










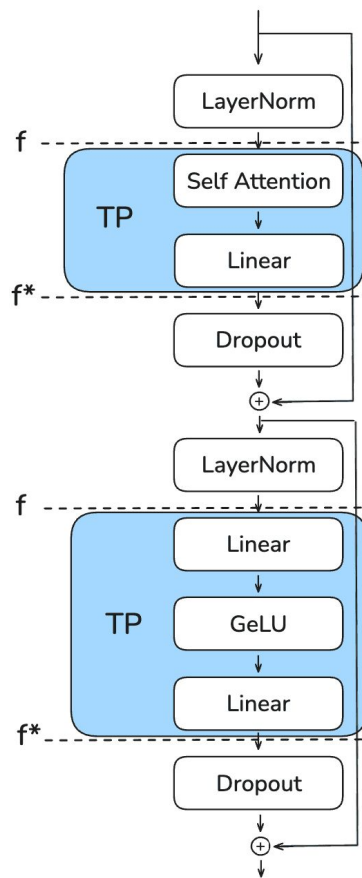


dropout ???

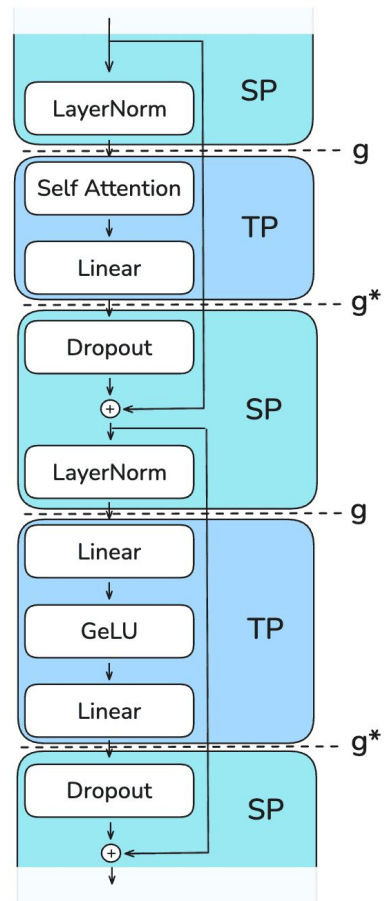
$$y = x \cdot \text{RandomMask}$$

$$y = x \cdot \text{RandomMask}$$
The word "RandomMask" is displayed in a black serif font. It is overlaid with four adjacent rectangular boxes of different colors: blue, red, green, and magenta. The blue box covers the first two letters "Ra", the red box covers the next two "ndom", the green box covers the next two "Ma", and the magenta box covers the final two "sk".

Tensor Parallel



Tensor + Sequence Parallel



### 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)

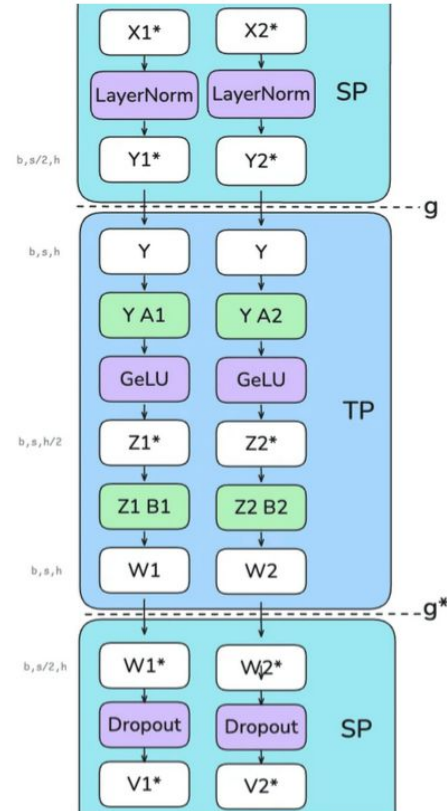
- $A1$  and  $A2$  are column-linear layers, so they split  $Y$  along the hidden dimension.
- GELU is applied independently on each GPU.
- $Z1^*$  and  $Z2^*$  are ( $b, s, h/2$ ).

### 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)

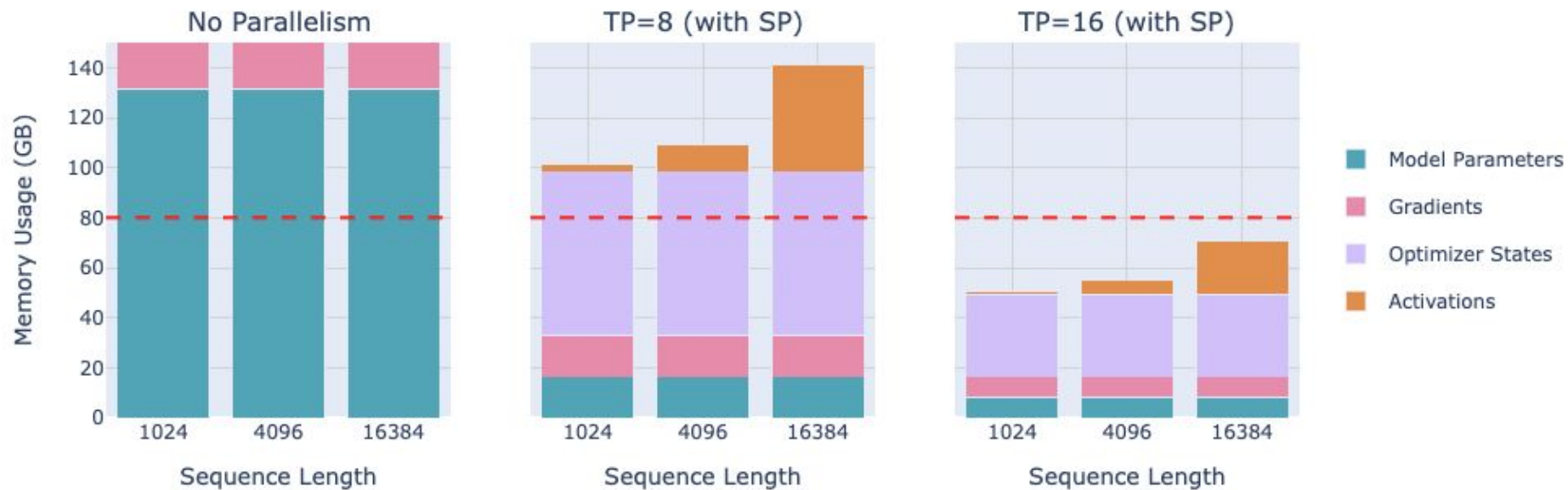
- $g^*$  operation (reduce-scatter) reduces for previous row-linear correctness while scattering along the sequence dimension.
- $W1^*$  and  $W2^*$  are ( $b, s/2, h$ ).



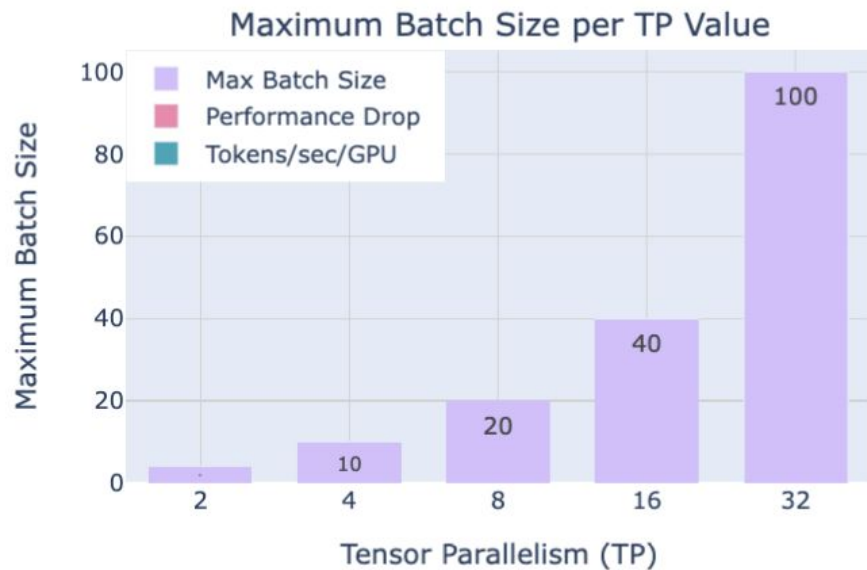
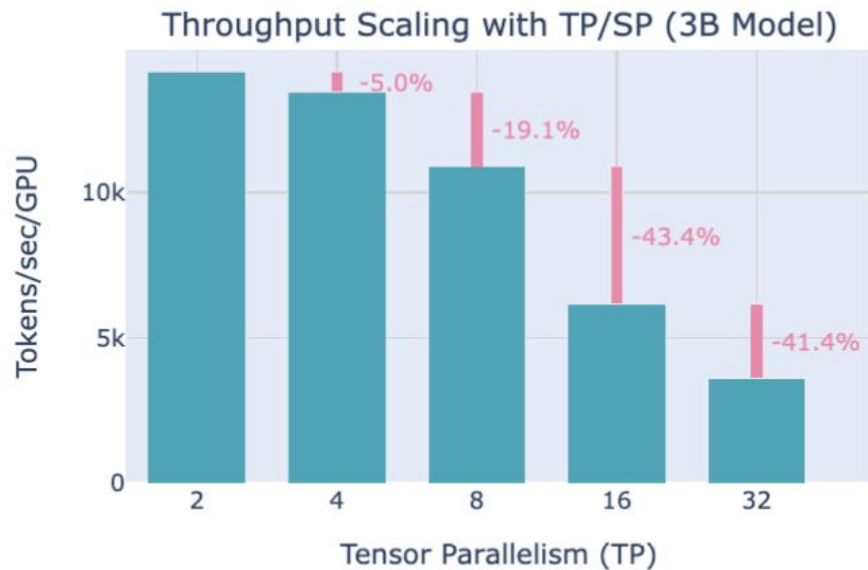
Region	TP only	TP with SP
Enter TP (column-linear)	$h$ : sharded (weight_out is sharded) $s$ : full	$h$ : sharded (weight_out is sharded) $s$ : <b>all-gather</b> to full
TP region	$h$ : sharded $s$ : full	$h$ : sharded $s$ : full
Exit TP (row-linear)	$h$ : full (weight_out is full + <b>all-reduce</b> for correctness) $s$ : full	$h$ : full (weight_out is full + <b>reduce-scatter</b> for correctness) $s$ : <b>reduce-scatter</b> to sharded
SP region	$h$ : full $s$ : full	$h$ : full $s$ : sharded



## Memory Usage for 70B Model







limits

- veeery large sequences blow up in TP
- inter-node connectivity tax for big models

see you next time