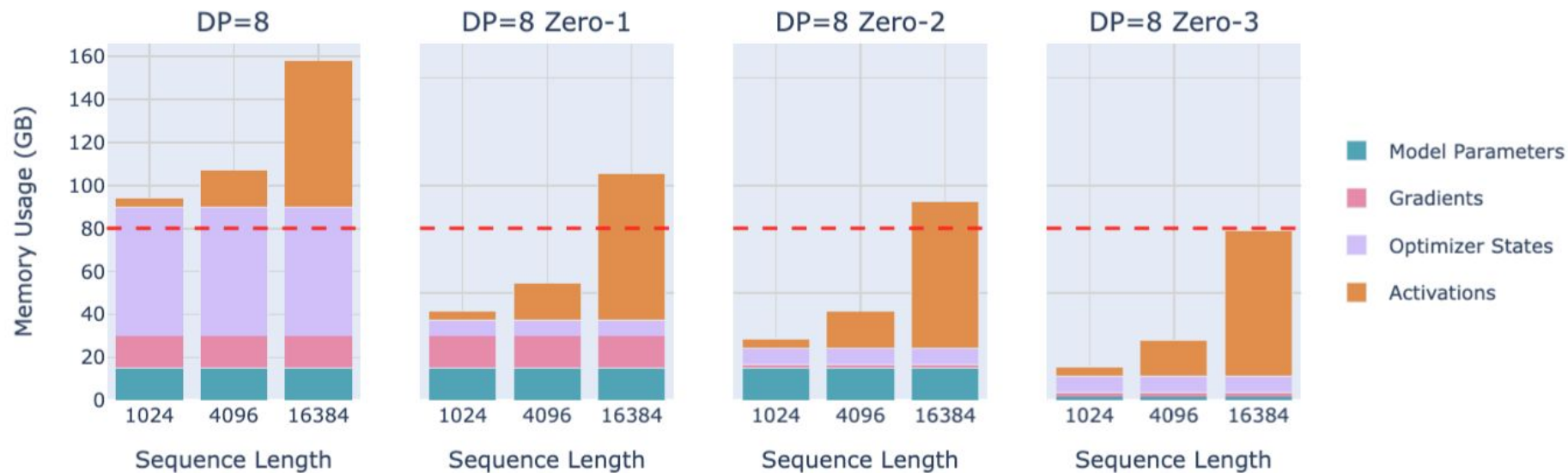


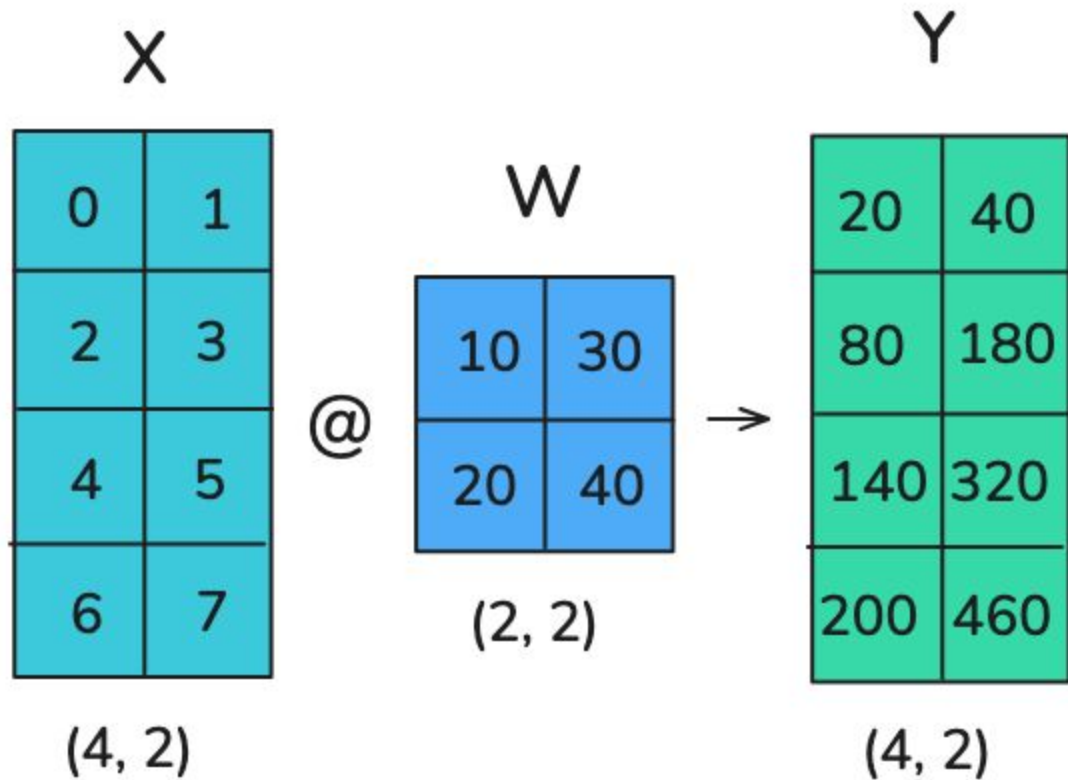
# Tensor Parallelism

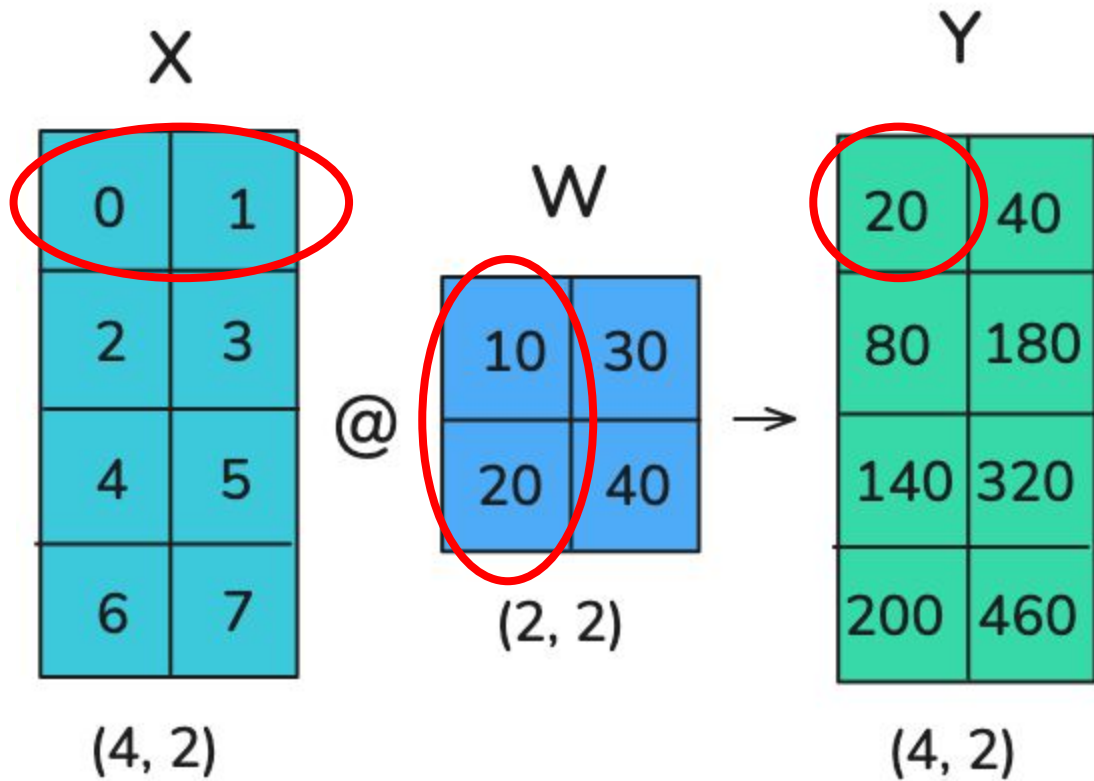
made with  for “Little ML book club”

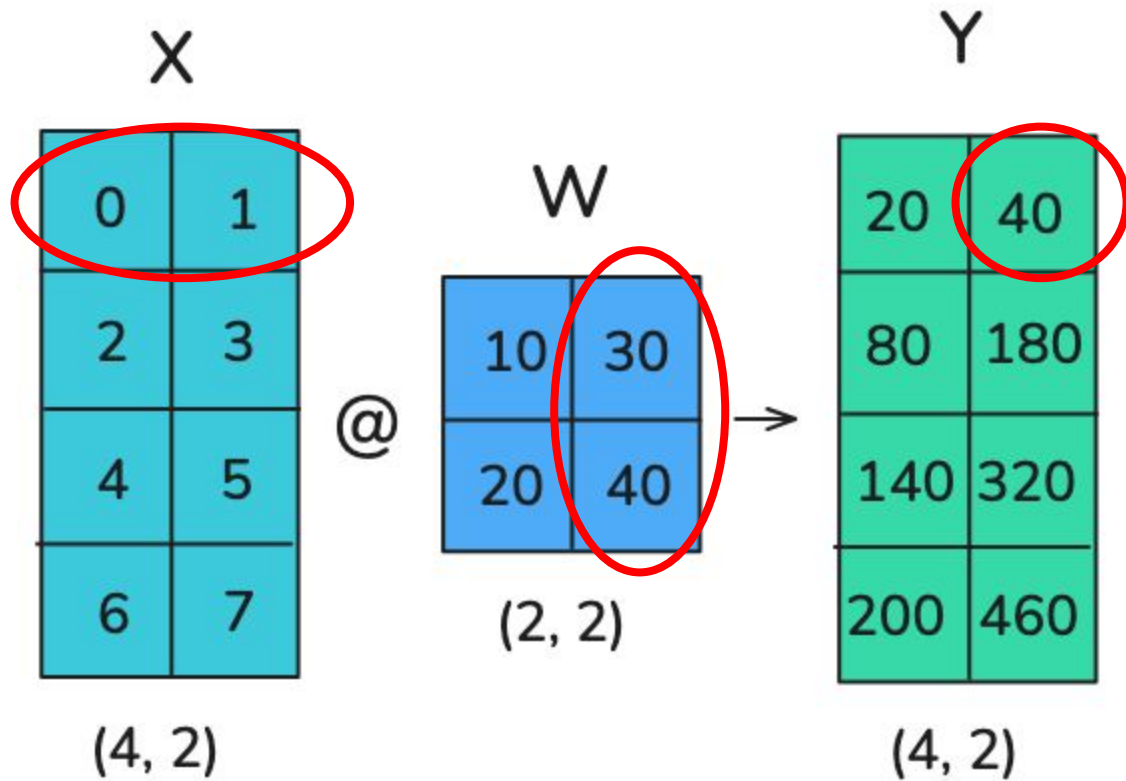
## Memory Usage for 8B Model



mom,  
but I want to play  
with veeeeery long sequences  
and  
veeeeeery large models







X

0	1
2	3
4	5
6	7

(4, 2)

@

W

10	30
20	40

(2, 2)



Y

20	40
80	180
140	320
200	460

(4, 2)

X

0	1
2	3
4	5
6	7

(4, 2)

@

W

10	30
20	40

(2, 2)

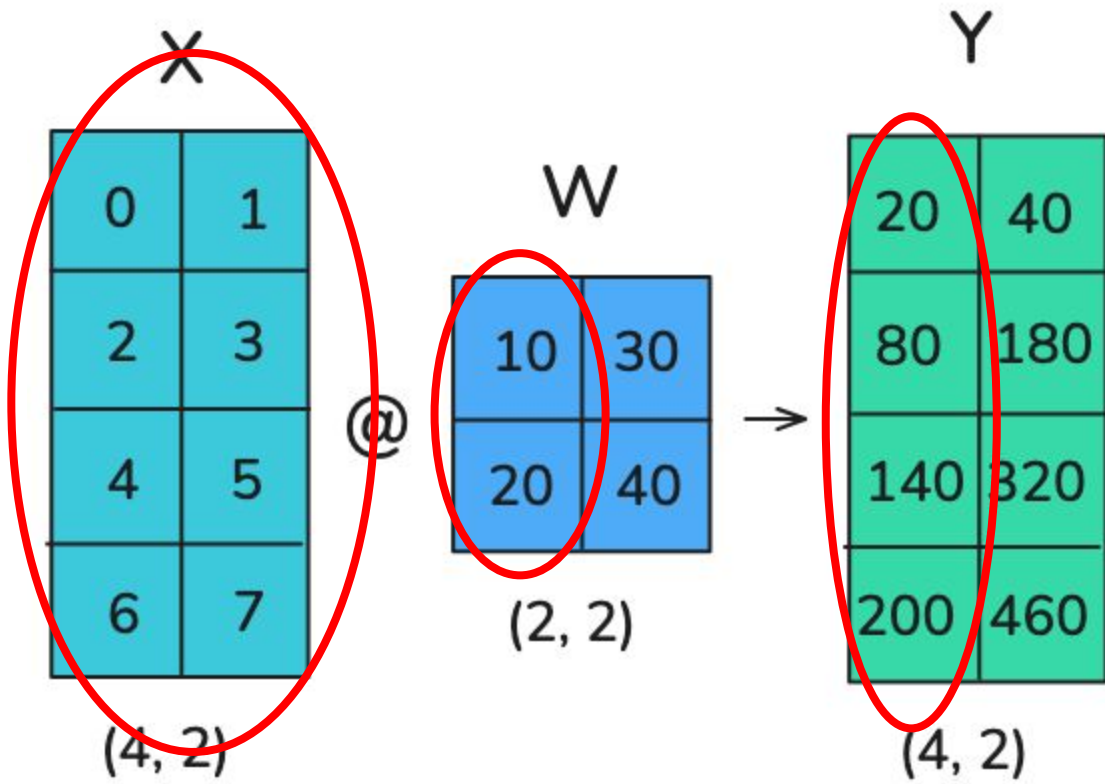


Y

20	40
80	180
140	320
200	460

(4, 2)





X

0	1
2	3
4	5
6	7

(4, 2)

@

W

10	30
20	40

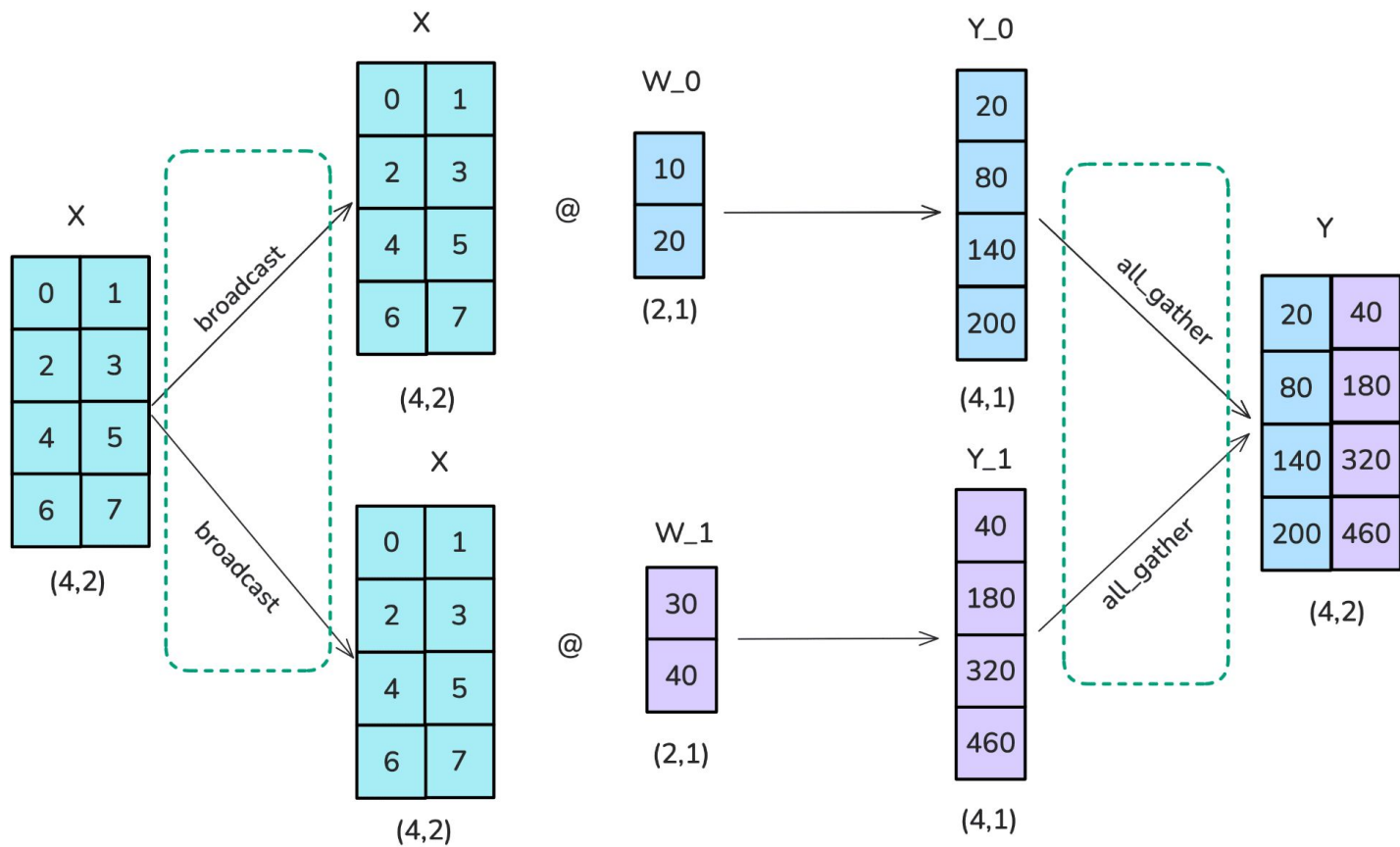
(2, 2)



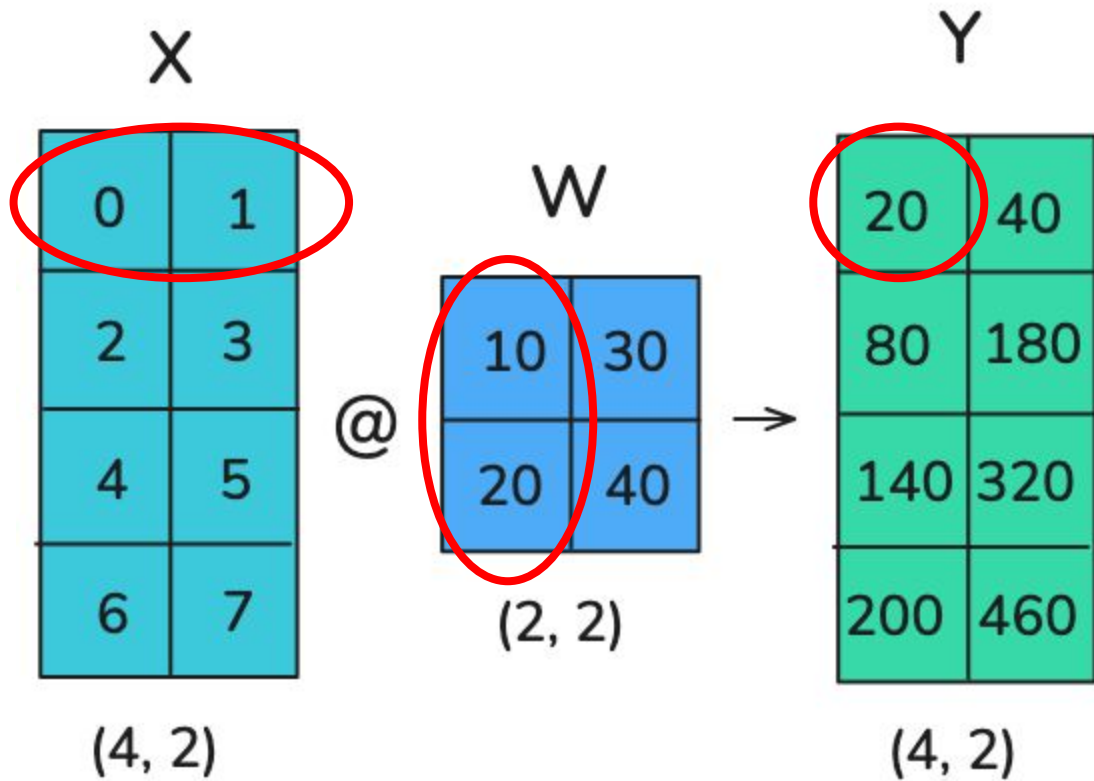
Y

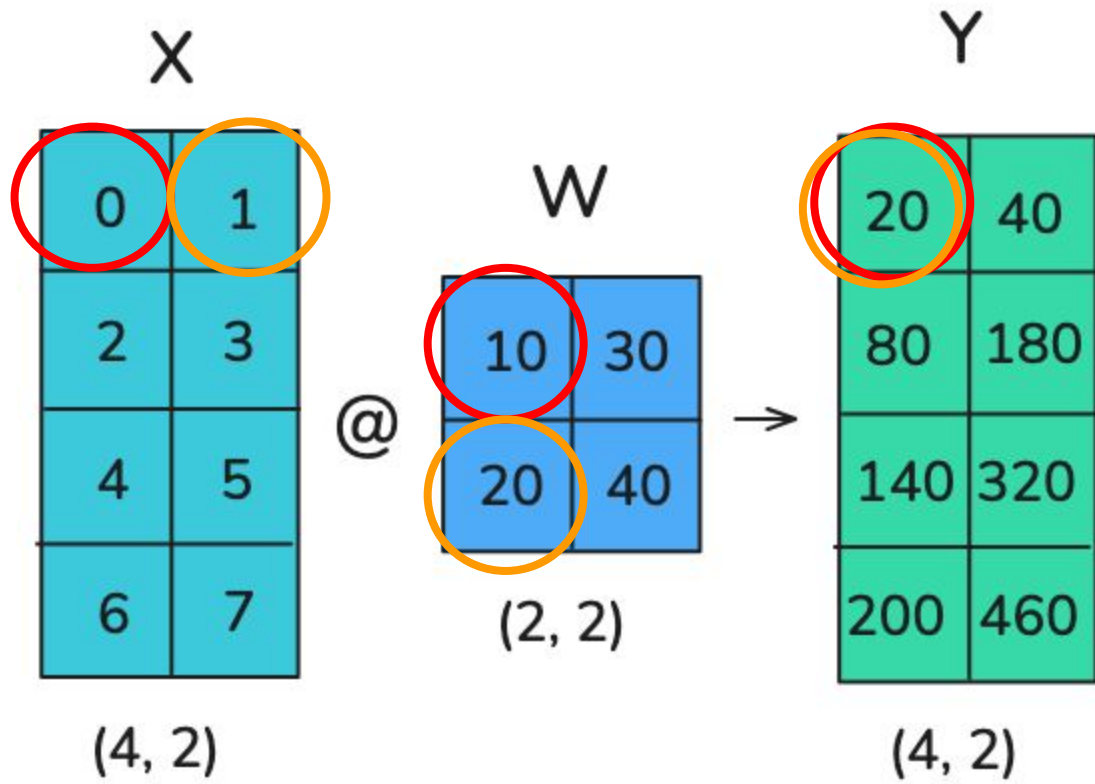
20	40
80	180
140	320
200	460

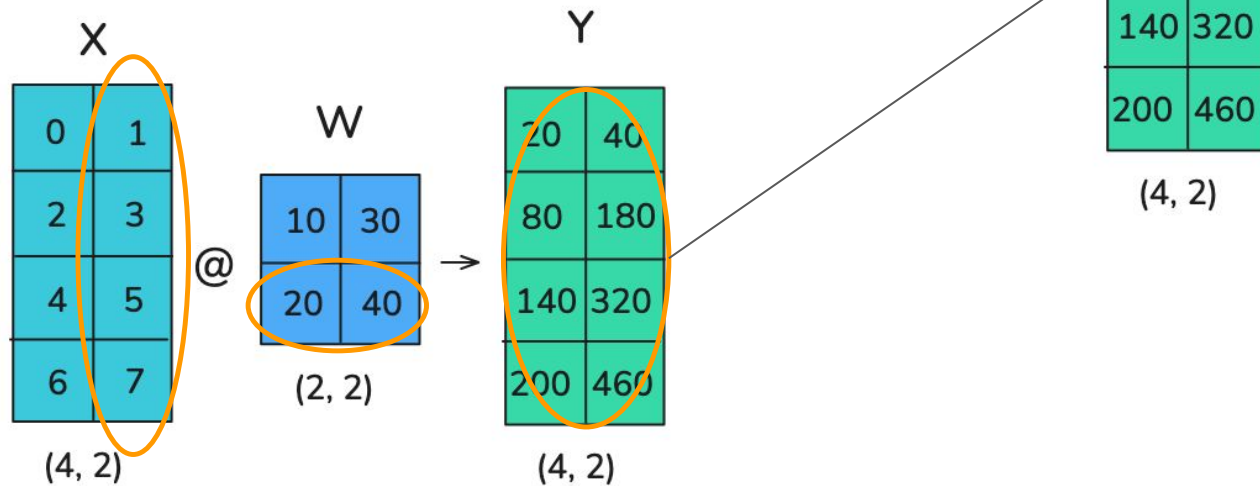
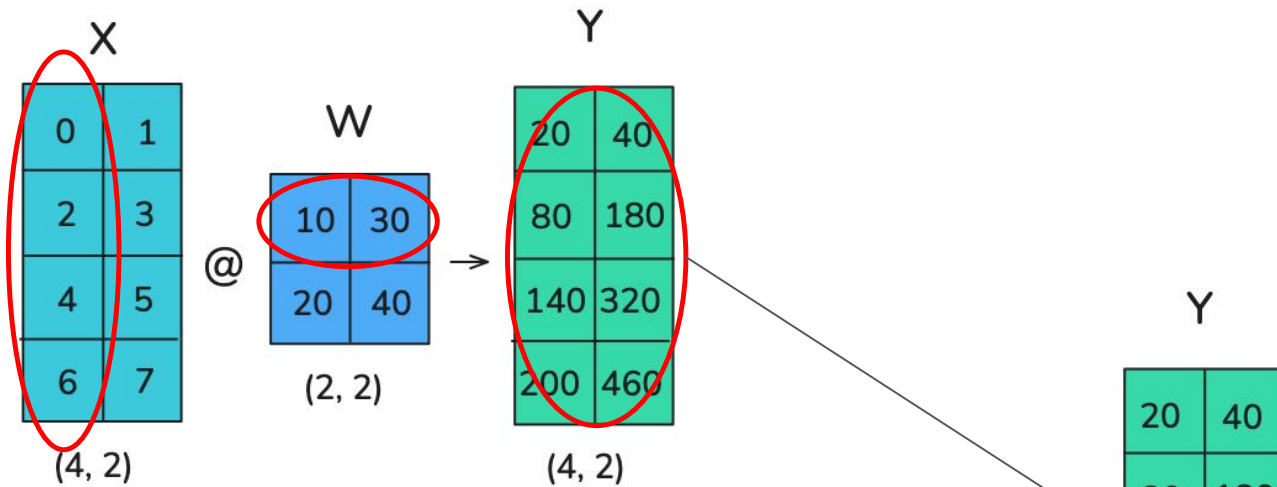
(4, 2)

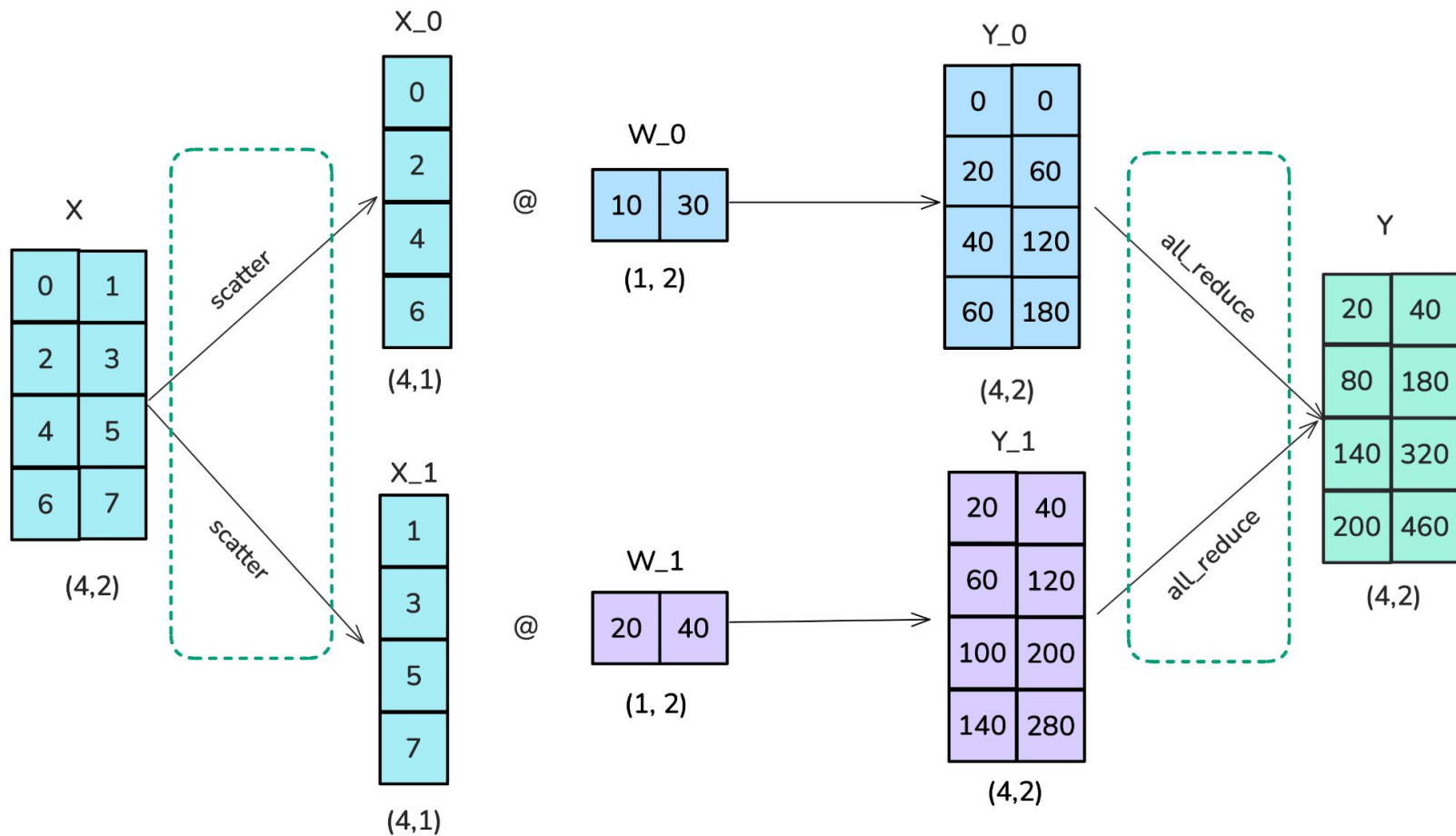


Column linear

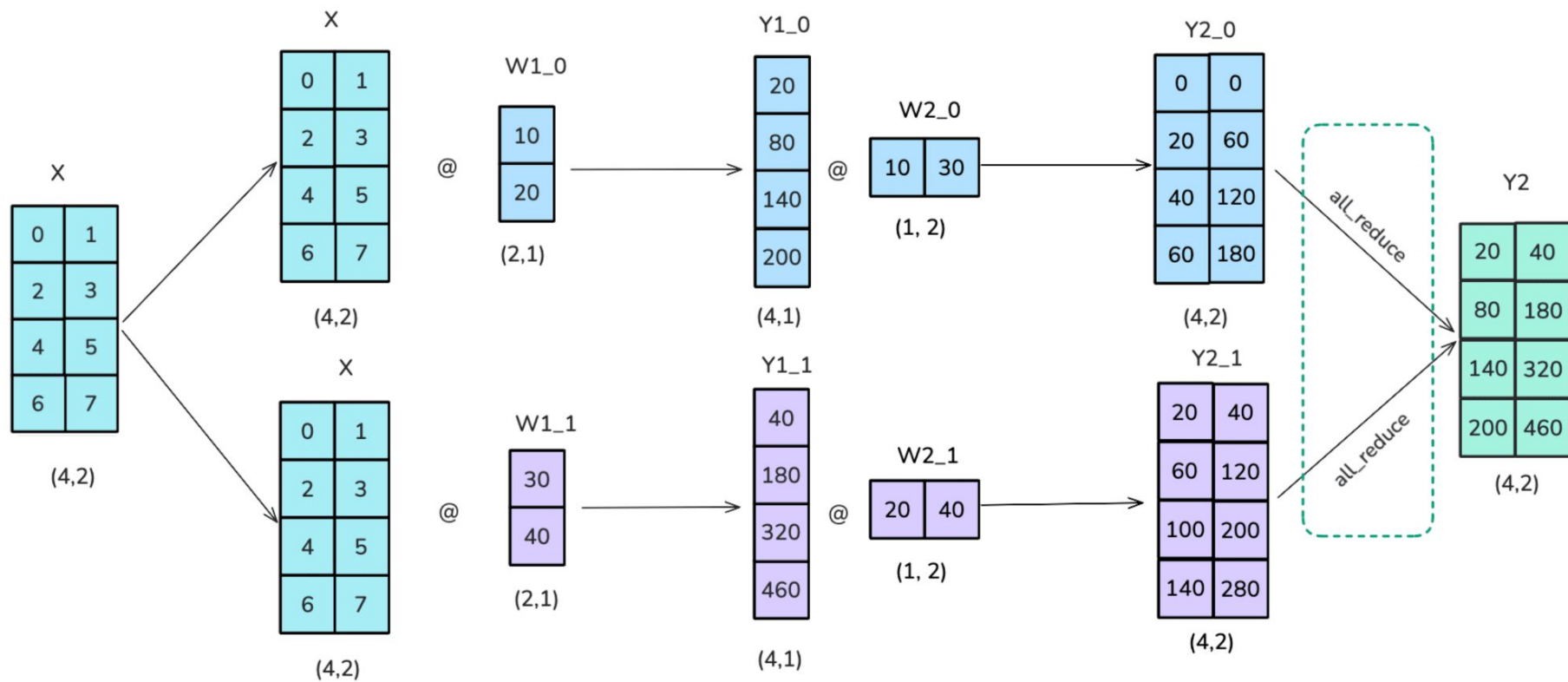








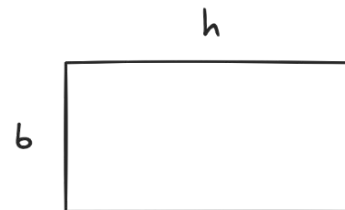
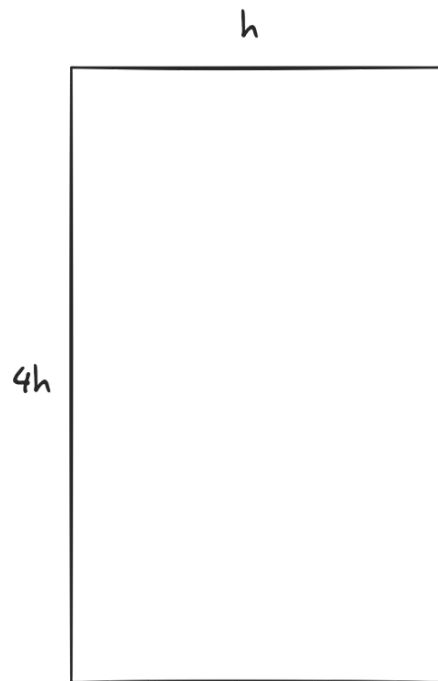
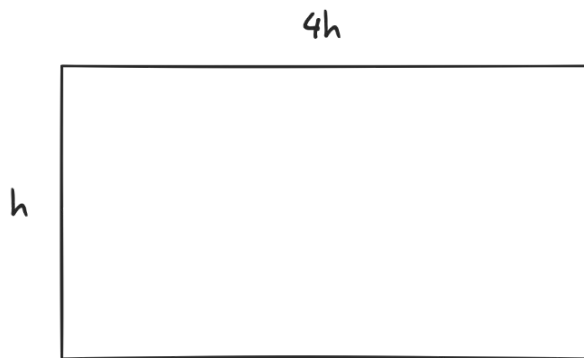
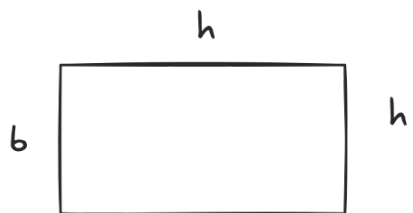
Row linear

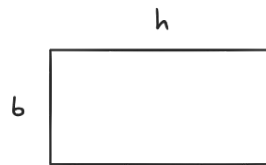
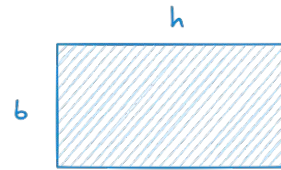
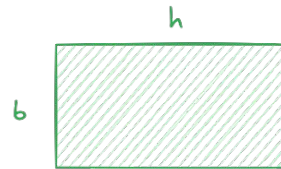
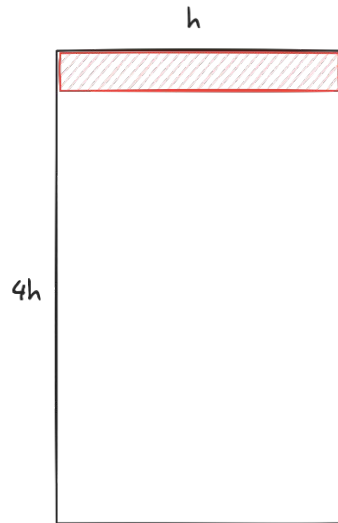
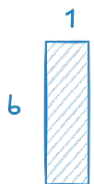
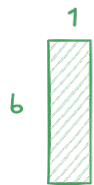
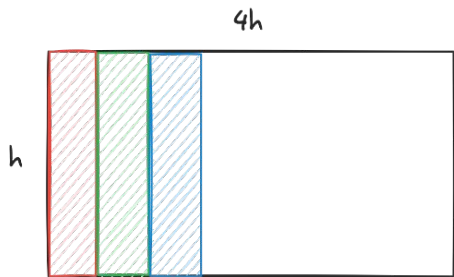
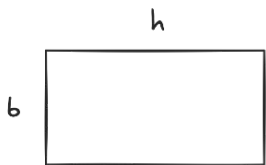


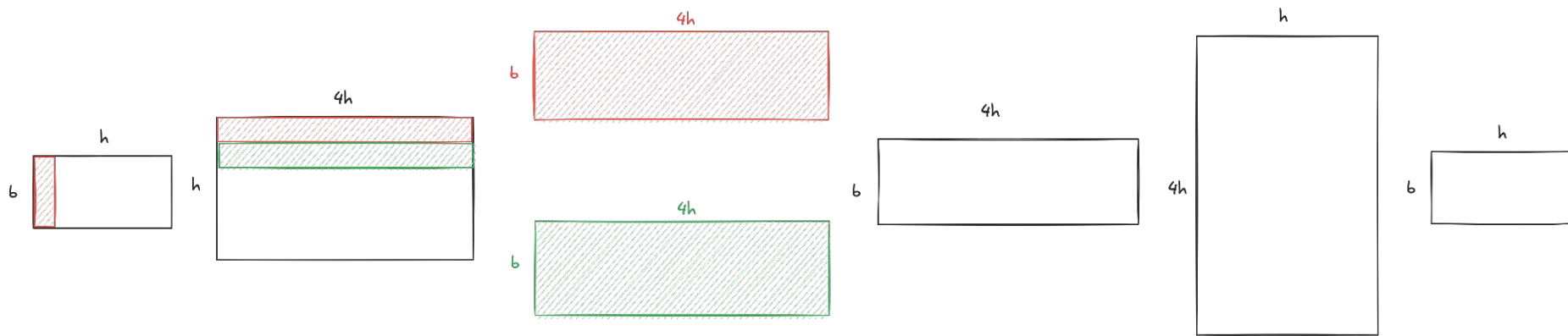
Tensor parallelism with column linear + row Linear

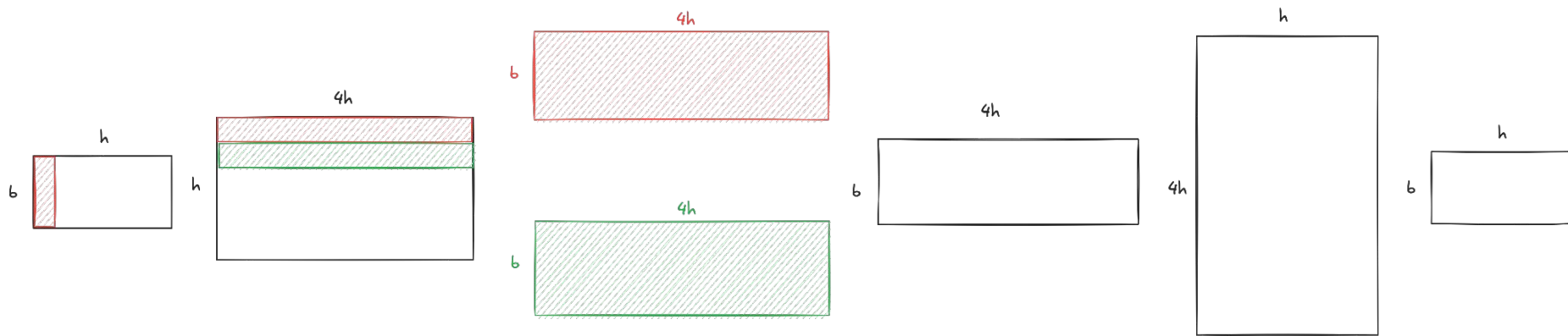


why combine 2 operations?

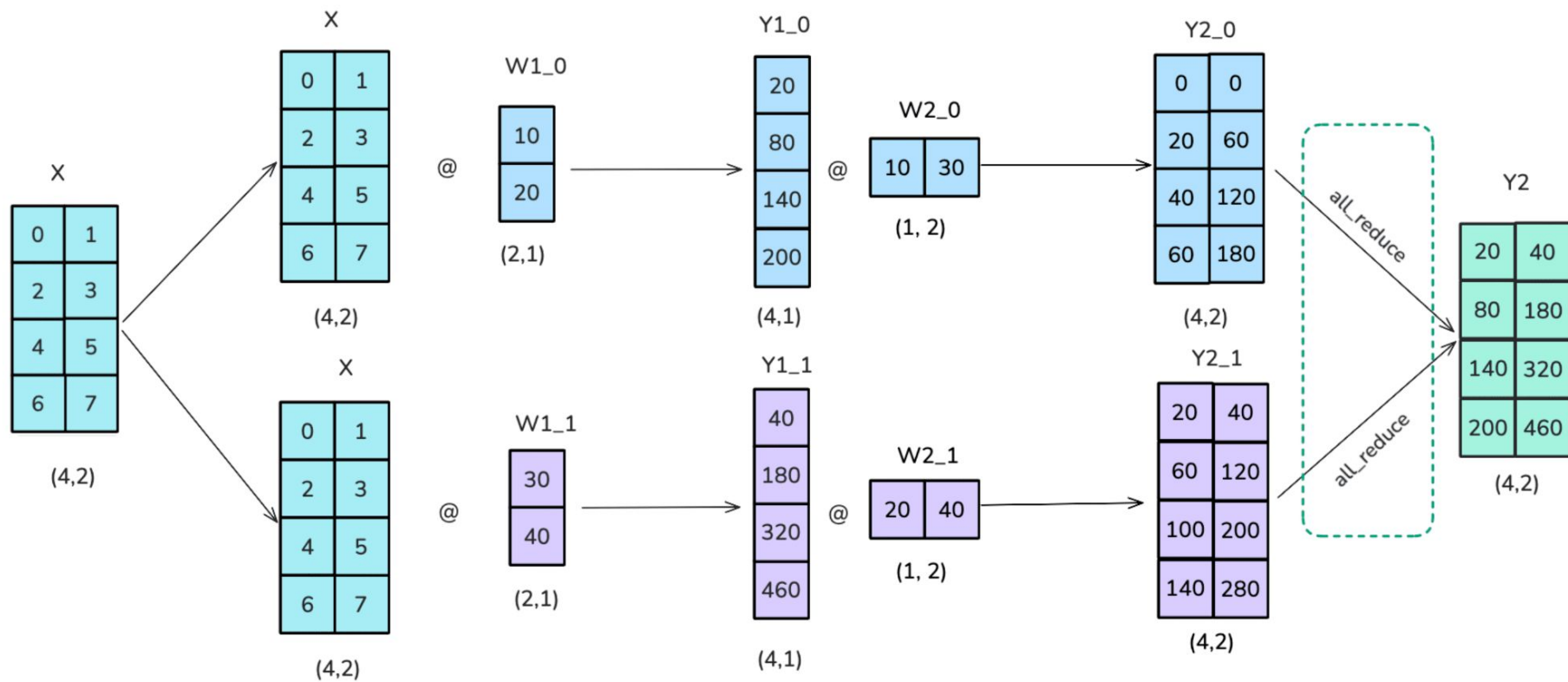








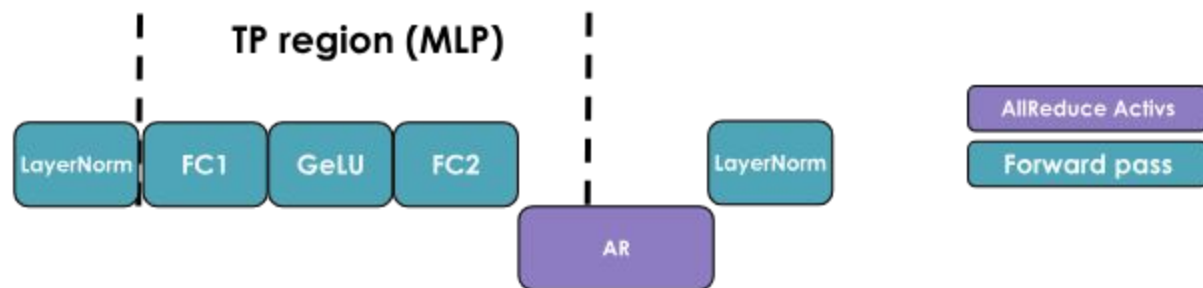
$$b * h \ll b * 4h + \text{all\_reduce}$$



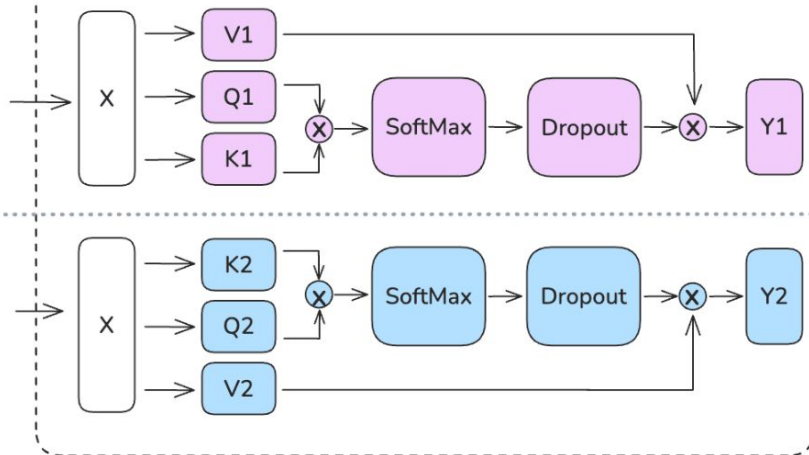
Tensor parallelism with column linear + row Linear

GPU Computation:

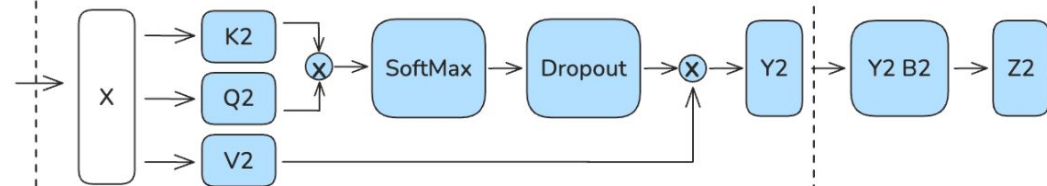
GPU Communication:



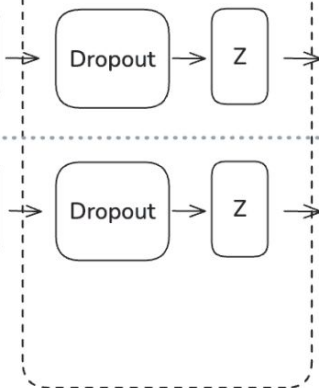
GPU1

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

GPU2

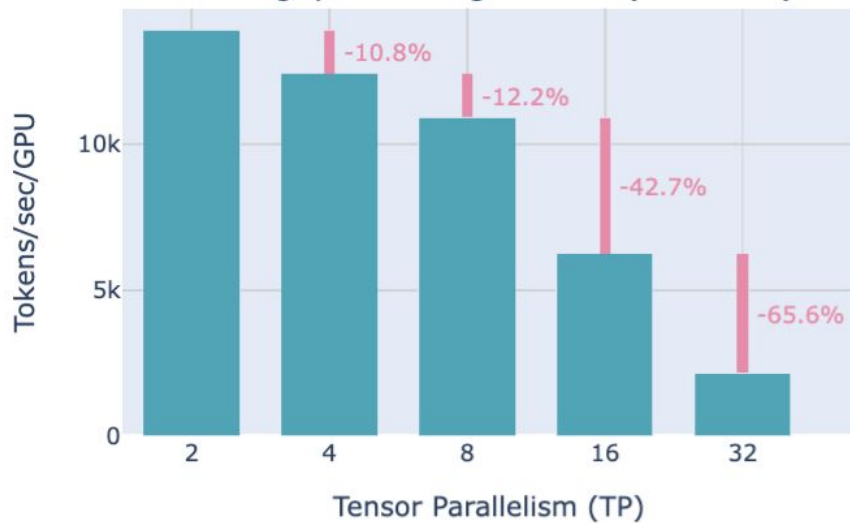


AllReduce

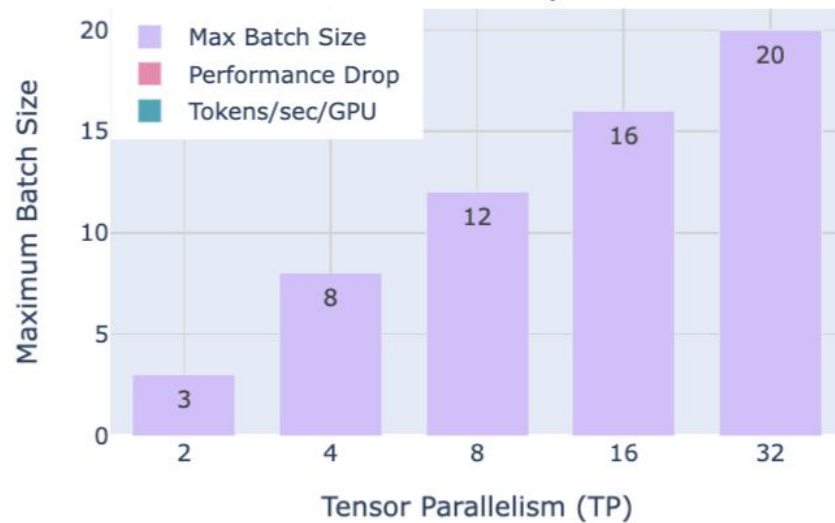
 $Z = \text{Dropout}(Y B)$ 



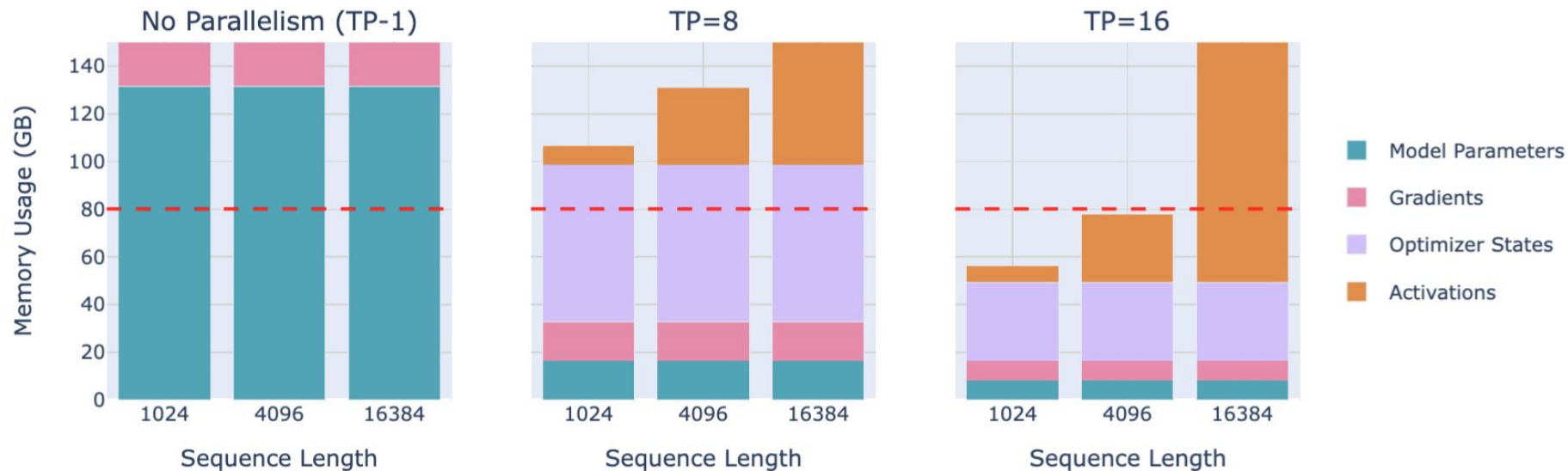
### Throughput Scaling with TP (3B Model)



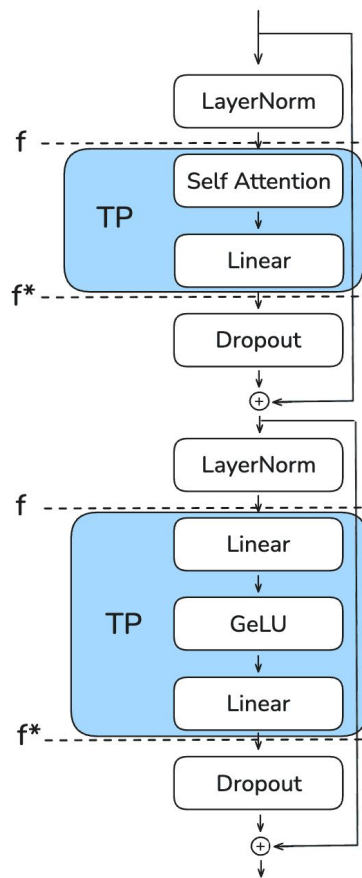
### Maximum Batch Size per TP Value



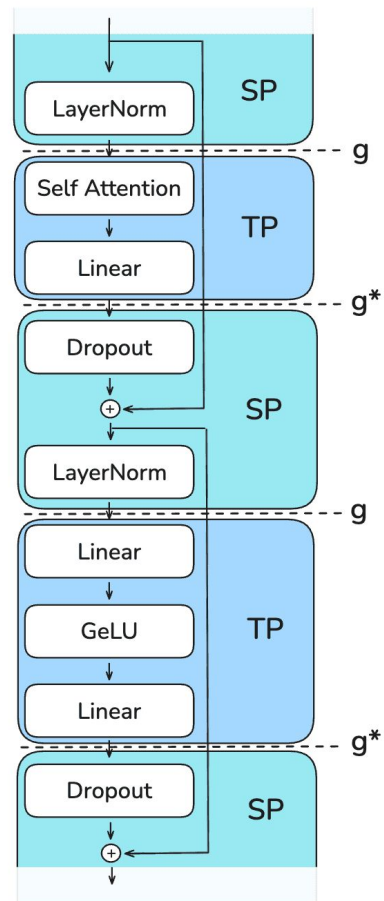
## Memory Usage for 70B Model

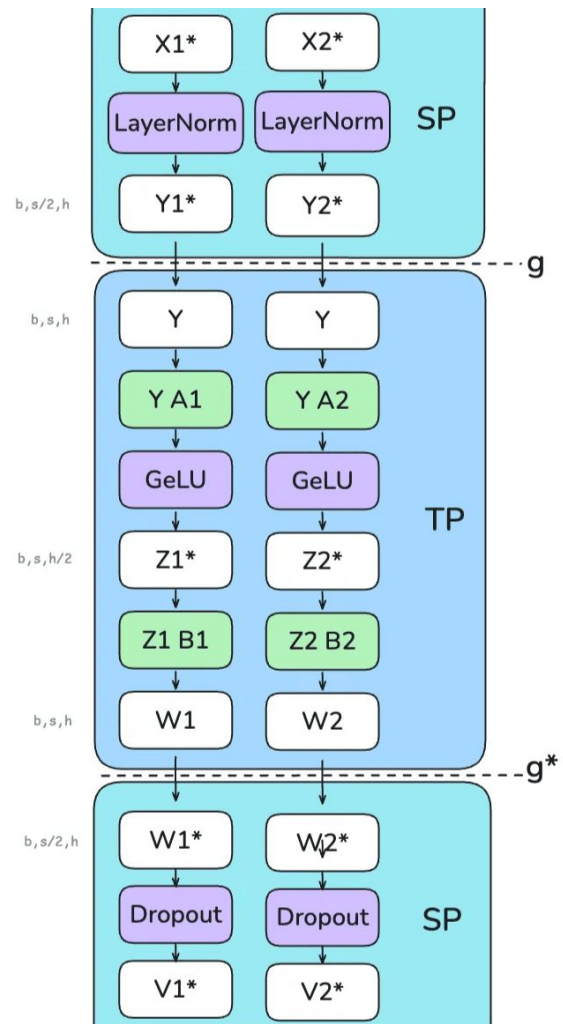


Tensor Parallel



Tensor + Sequence Parallel





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





Still, there are two limits to TP+SP: if we scale the sequence length the activation memory will still blow up in the TP region, and if the model is too big to fit with TP=8 we will see a massive slowdown due to the inter-node connectivity.



see you next time