

Pipeline Parallelism

made with ❤️ for “Little ML book club”

Attention

GPU Computation:

Attn(Q_i, K_i, V_i)

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

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

AllGather Activs

GPU Communication:

AG(K,V)

Forward pass

Attention

GPU Computation:

Attn(Q_i, K_i, V_i)

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

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

P2P Activs

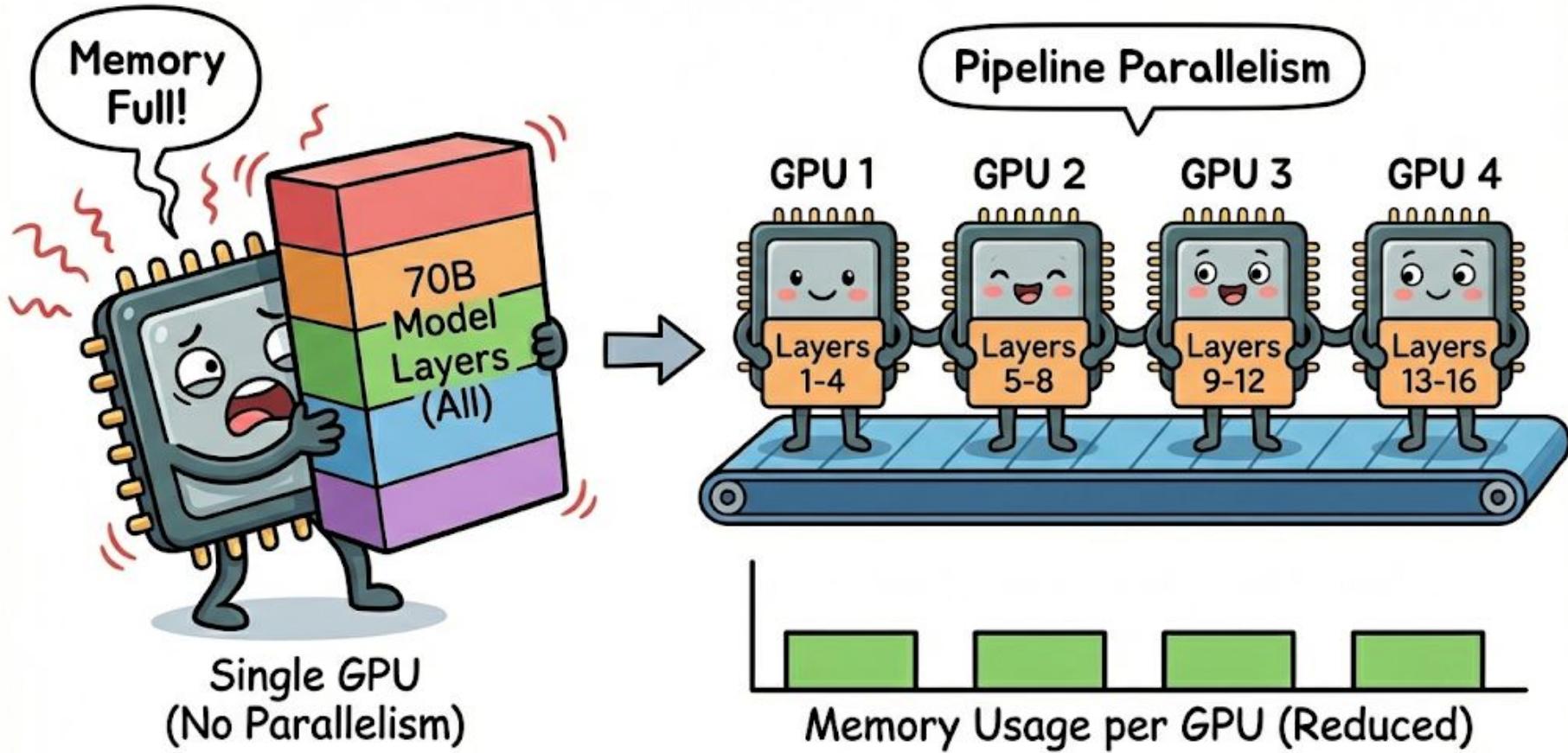
GPU Communication:

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

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

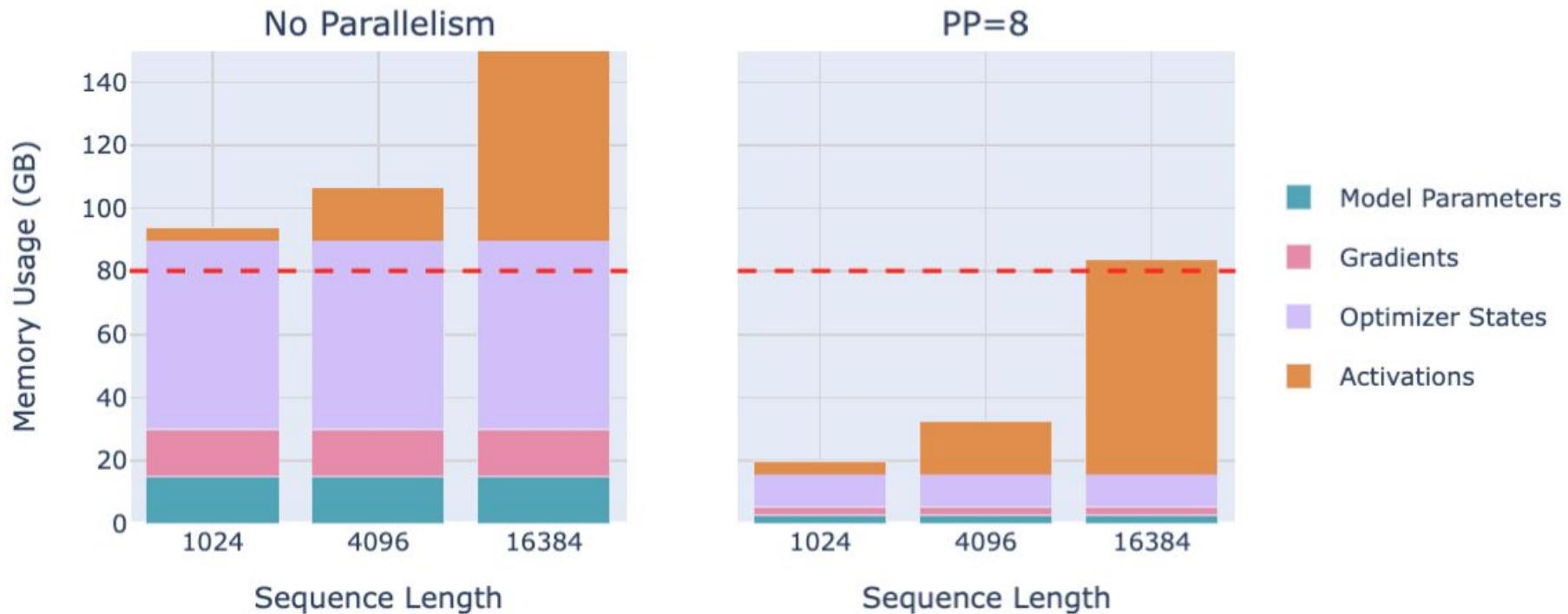
Fetch K_{i+3}, V_{i+3}

Forward pass



Efficiently splits model layers across multiple GPUs

Memory Usage for 8B Model

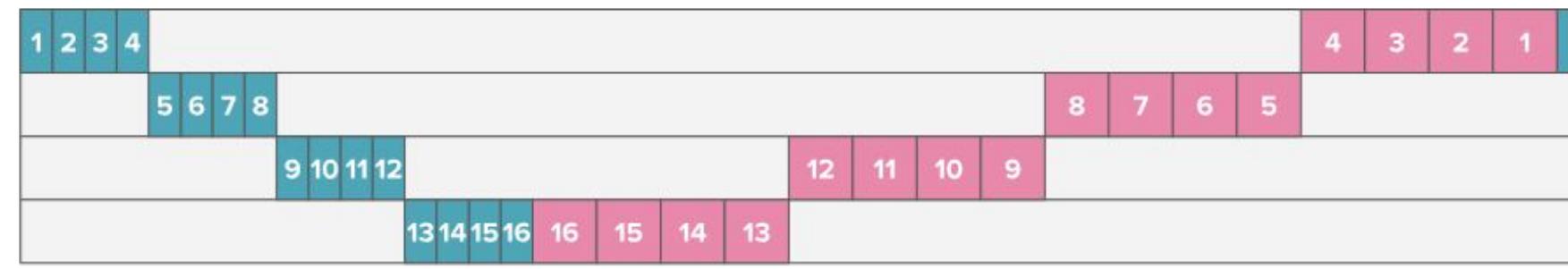




Note

This is because each GPU needs to perform PP forward passes before starting the first backward pass. Since each GPU handles $1/PP$ of the layers but needs to process PP micro-batches before the first backward, it ends up storing $PP \times (activs/PP) \approx activs$, which means the activation memory requirement remains roughly the same as without pipeline parallelism.

GPU



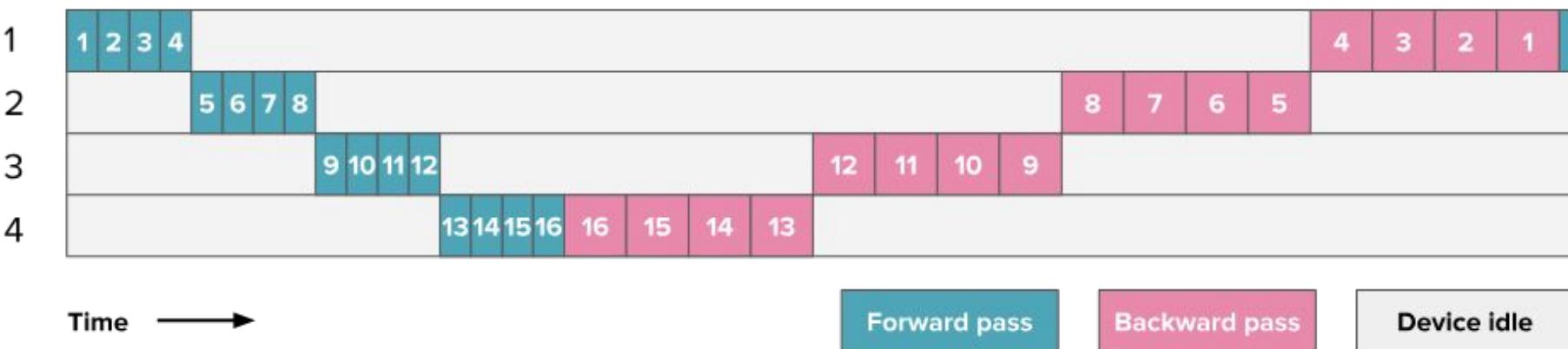
Time →

Forward pass

Backward pass

Device idle

GPU



GPU activation memory = 1/4 of total activation memory

GPU



Time →

Forward pass

Backward pass

Device idle

GPU



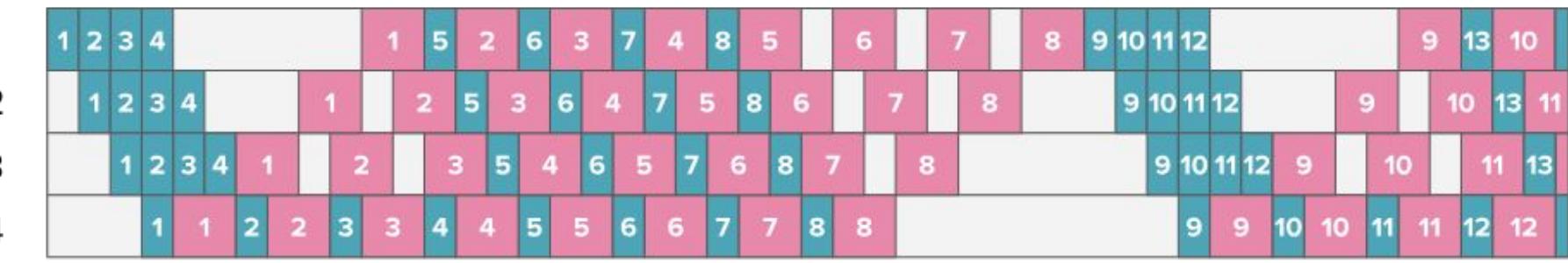
Forward pass

Backward pass

Device idle

GPU activation memory = 8/4 of total activation memory

GPU



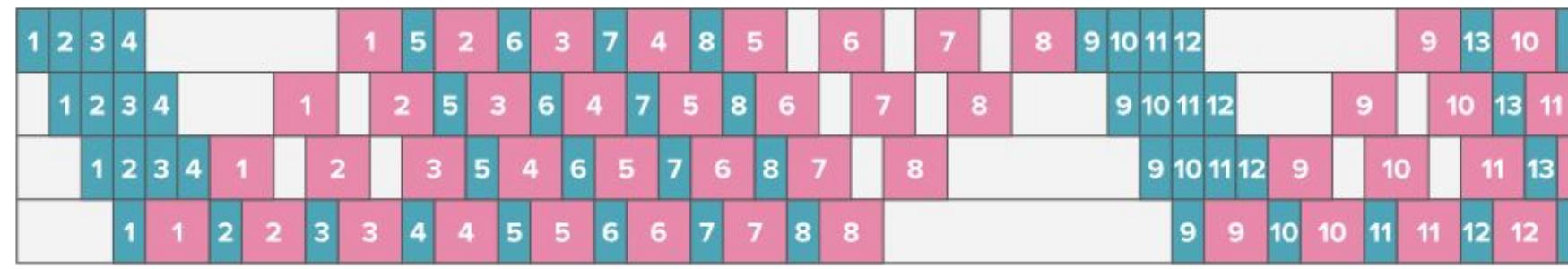
Time →

Forward pass

Backward pass

Device idle

GPU



Time →

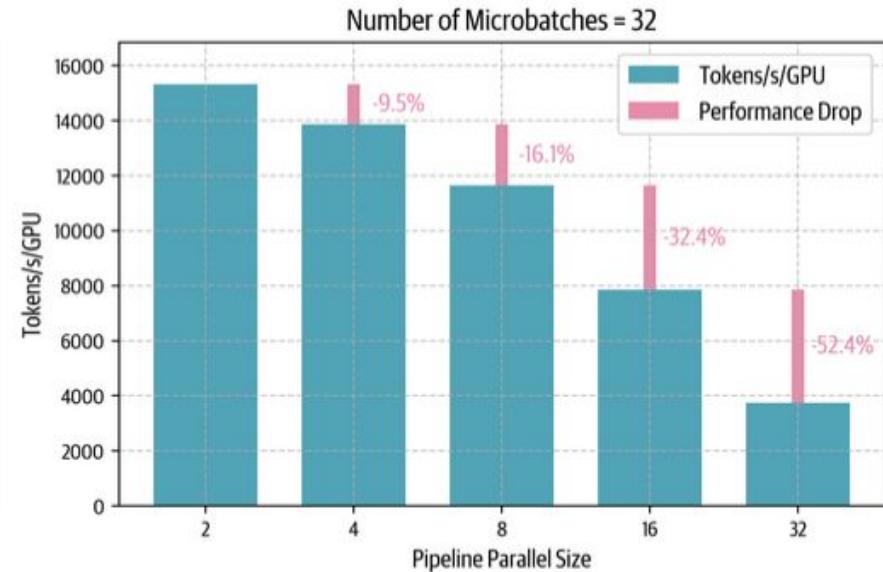
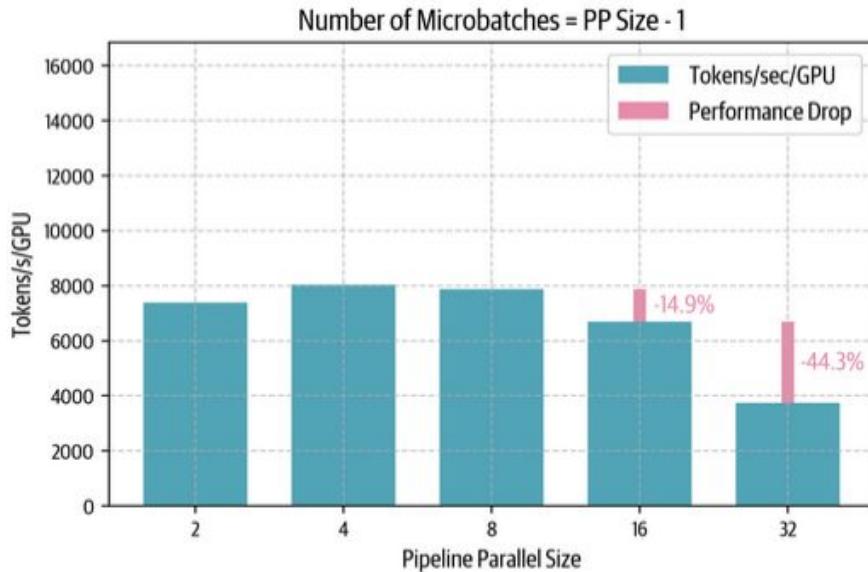
Forward pass

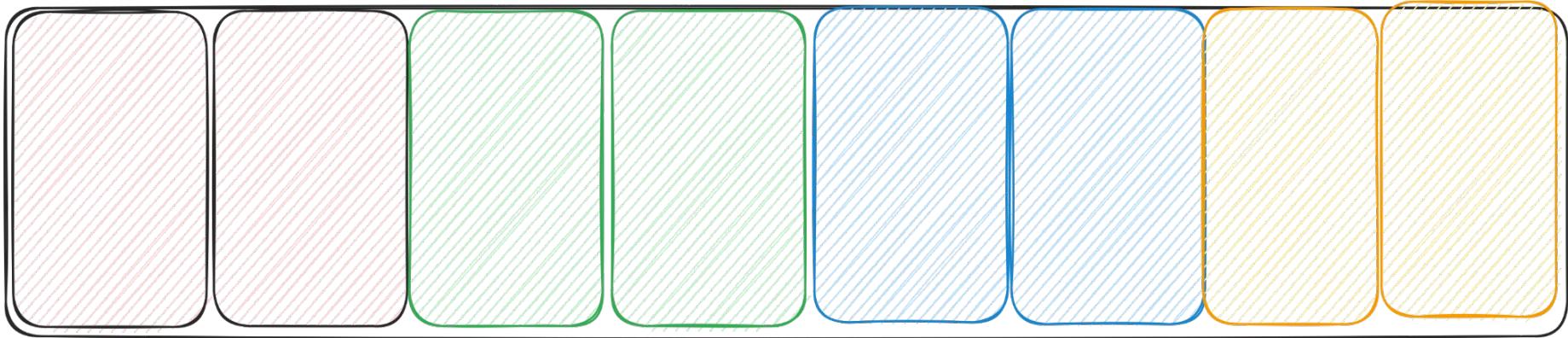
Backward pass

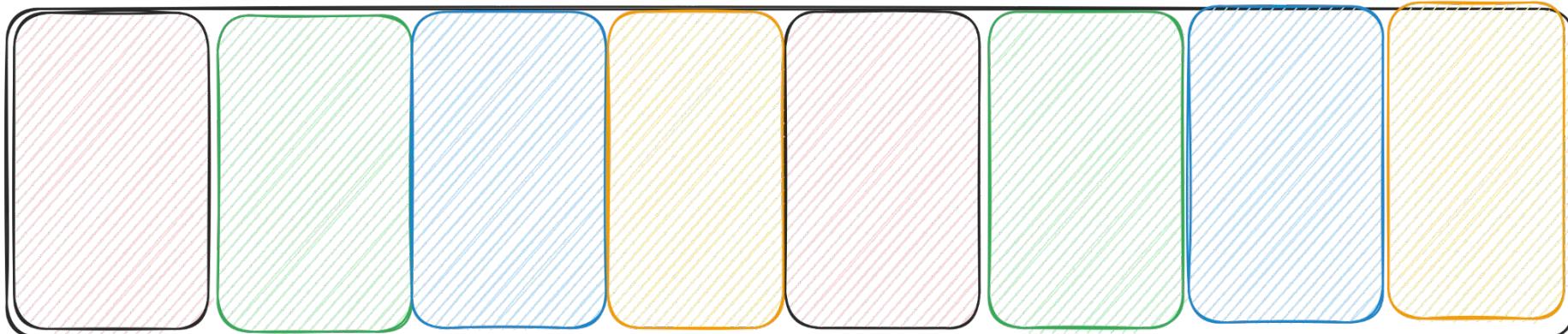
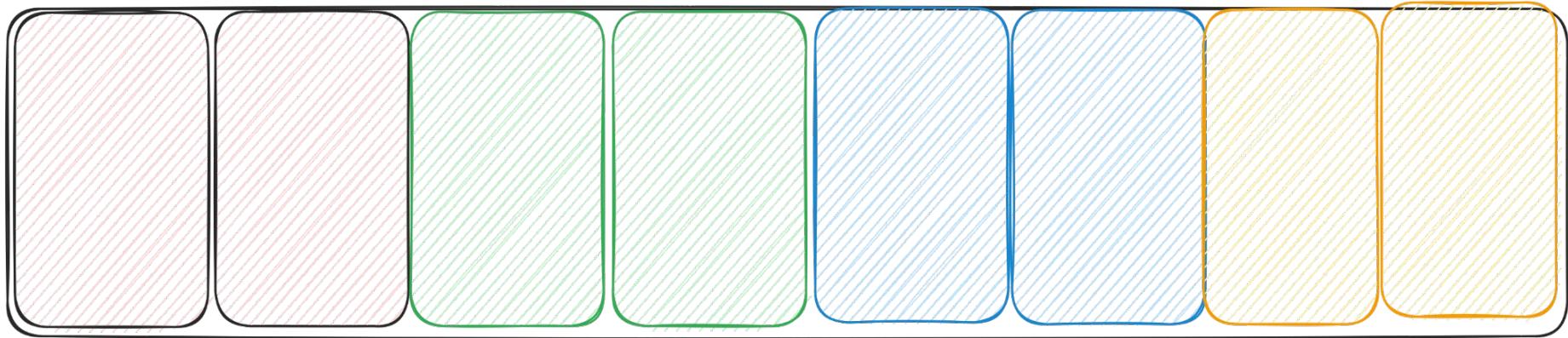
Device idle

GPU activation memory = 4/4 of total activation memory

Throughput Scaling with Pipeline Parallelism (1F1B schedule)







GPU



GPU



GPU activation memory = 4/4 of total activation memory

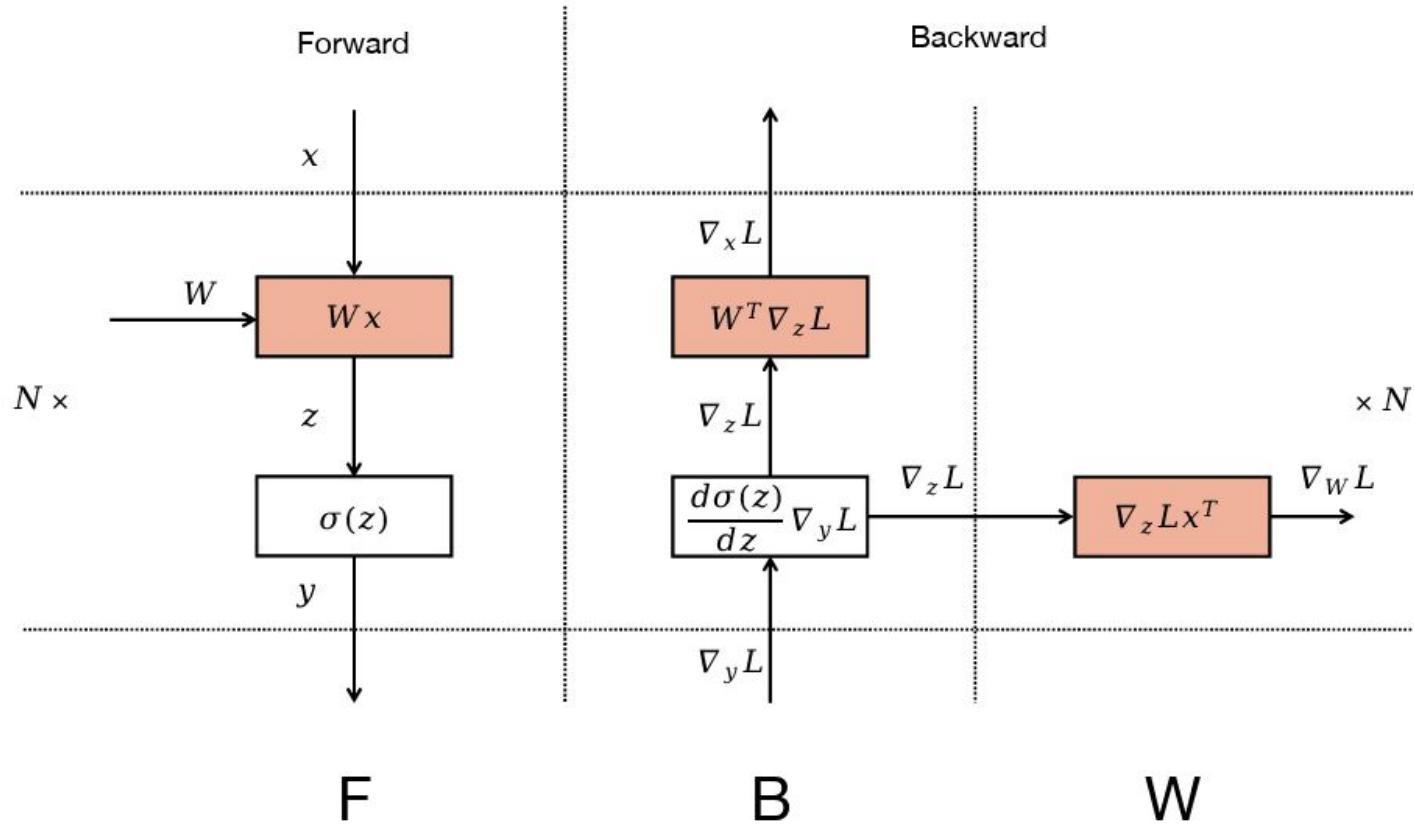


Figure 1: Computation Graph for MLP.

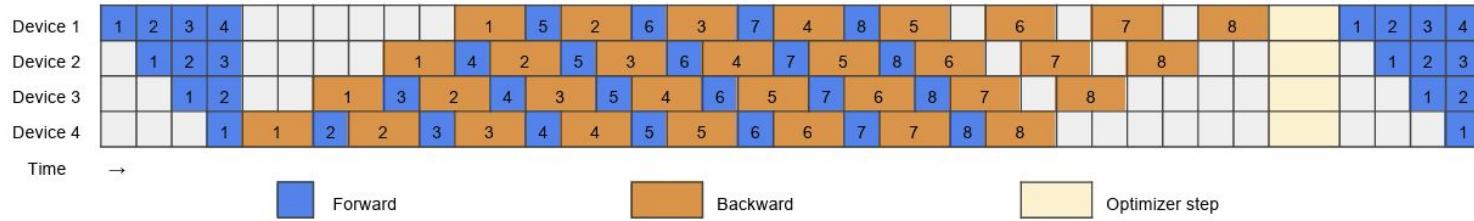


Figure 2: 1F1B pipeline schedule.

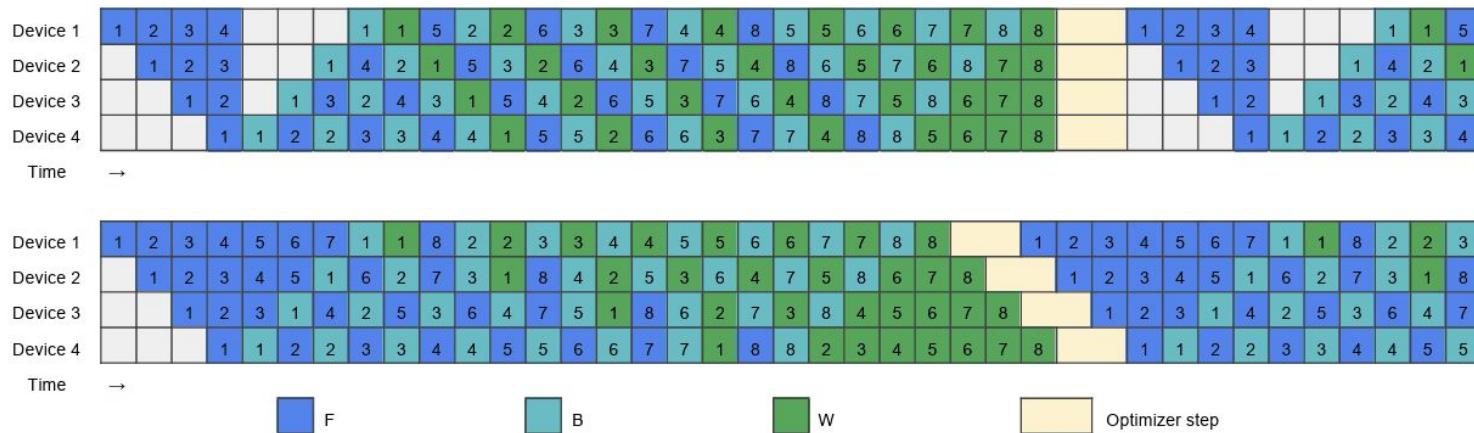
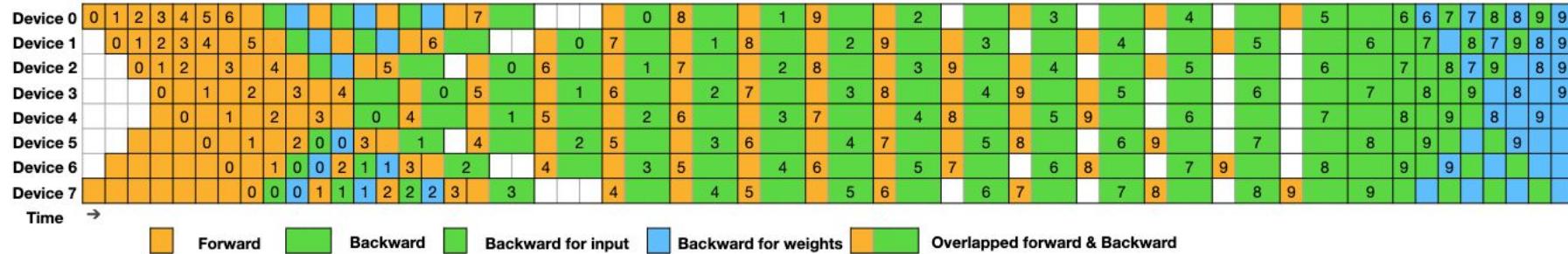


Figure 3: Handcrafted pipeline schedules, top: ZB-H1; bottom: ZB-H2



Method	Bubble	Parameter	Activation
1F1B	$(PP - 1)(F + B)$	$1\times$	PP
ZB1P	$(PP - 1)(F + B - 2W)$	$1\times$	PP
DualPipe (Ours)	$(\frac{PP}{2} - 1)(F\&B + B - 3W)$	$2\times$	$PP + 1$

see you next time