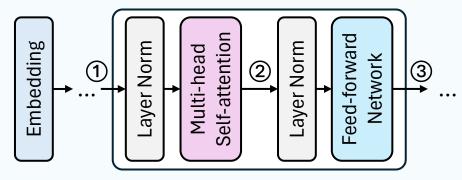
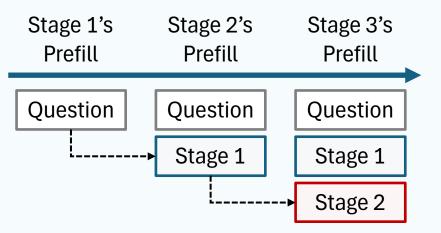
Offline Configuration

■ Step 1: Estimate Layer Importance



MHSA importance: cos(①, ②)
FFN importance: cos(②,③)

Accumulate Importance across Stages:



Set Target Accuracy (Acc).

From Slowest to Fastest Stage, Repeat:

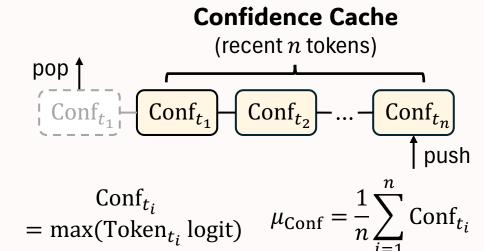
Within Target Acc, Find Best Latency.

# Skip	Acc (%)	Latency
0	51.7 ✓	0.3192 X
1	50.6	0.3157
2	50.5	0.3012
3	51.3 ✓	0.3029 X
4	51.4 🗸	0.2986 🗸
5	49.5	0.2885
6	49.4	0.3103
7	48.6	0.3033
8	49.4	0.3202
•		

2 Set **# Skip** for this Stage.

Online Adjustment

% Step 3: Generation Early Exit



Terminate Generation if Confidence is Low:

