

Normalized Speedup

5
4
3
2
1
0

MSCOCO
34B

Baseline w/o SDPA
gives Out of Memory

Vizwiz
34B

Baseline w/o SDPA
gives Out of Memory

Coco_Image
34B

HumanEval
34B

S2ST

S2TT

T2ST

T2TT

Workloads

