

SoftMoE

4 experts

1 expert scaled x4

1 expert

ExpertChoice

1 expert scaled x4

Baseline

Penultimate layer scaled x4

