



● Baseline

▲ NF4

■ Int8

Techniques (Averaged Across All Models)

✖ SnapKV

+ PyramidKV

★ StreamingLLM

▼ DuoAttn