

Inference throughput: Baseline vs INT4

