

Throughput (tokens/s)



KV Cache Maximum Usage (%)