

Hardware: a100 Inference: vLLM

Throughput (Tokens per seconds)

