

Hardware: tesla Inference: vLLM

Throughput (Tokens per seconds)

