

DeepSeek-R1-Distill-Llama-70B
TP=2, PP=4, Input=2048, Output=512
Instance=2x g6e.12xlarge, GPU=L40S
Elapsed=90.5s, GPU Time=90.2s

