DeepSeek-R1-Distill-Llama-70B
TP=2, PP=4, Input=8192, Output=2048
Instance=2x g6.12xlarge, GPU=L4
Elapsed=387.4s, GPU Time=387.0s