DeepSeek-R1-Distill-Llama-70B
TP=4, PP=2, Input=2048, Output=512
Instance=2x g5.12xlarge, GPU=A10G
Elapsed=426.4s, GPU Time=426.3s