

# Ptype tp Model meta-llama/Llama-2-13b-chat-hf

P99 Token Latency (s/token)

