

Throughput (tokens/s)

opt-1.3b llama-2-7b
opt-2.7b llama-2-13b

