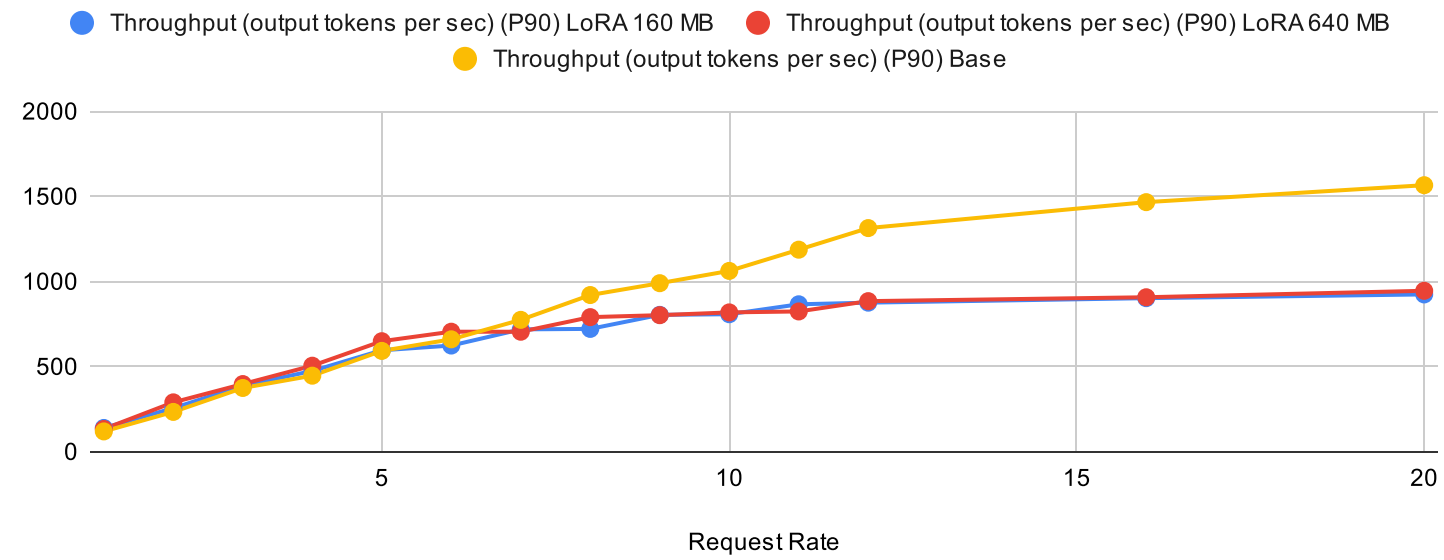


Throughput (output tokens per sec) (P90) LoRA 160 MB, Throughput (output tokens per sec) (P90) LoRA 640 MB and Throughput (output tokens per sec) (P90) Base



Latency (sec) (P90) LoRA 160 MB, Latency (sec) (P90) LoRA 640 MB and Latency (sec) (P90) Base

