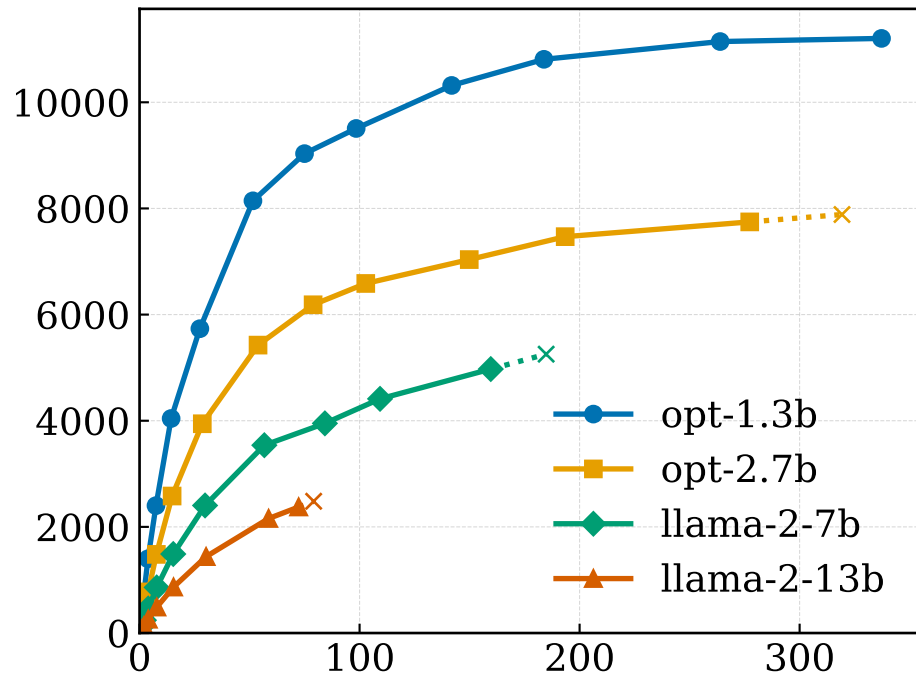
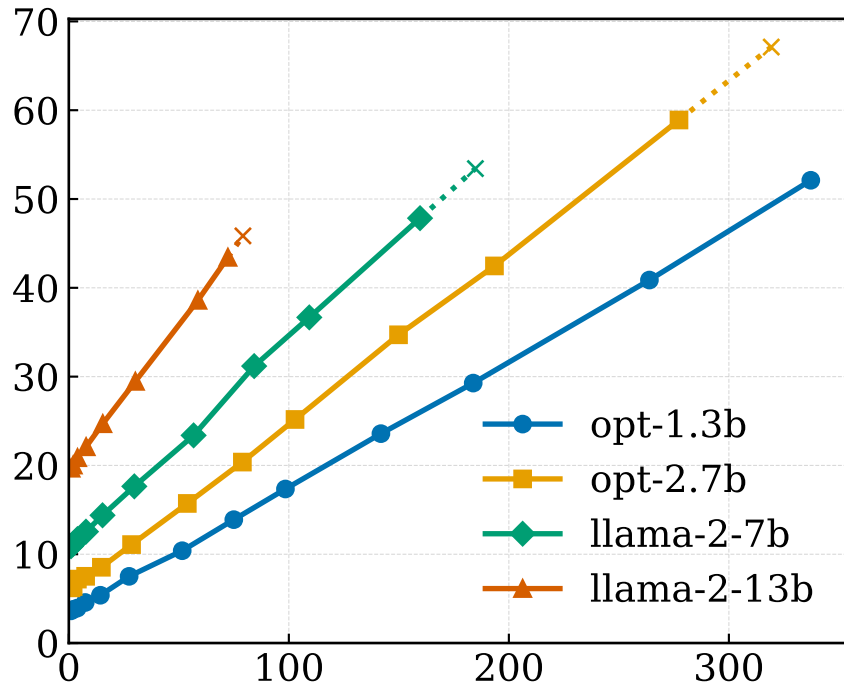


Throughput (tokens/s)



Average Batch Size (reqs)

Latency (ms)



Average Batch Size (reqs)