



t: throughput (tokens/s)

SmallRequest

MediumRequest

SmallRequest

