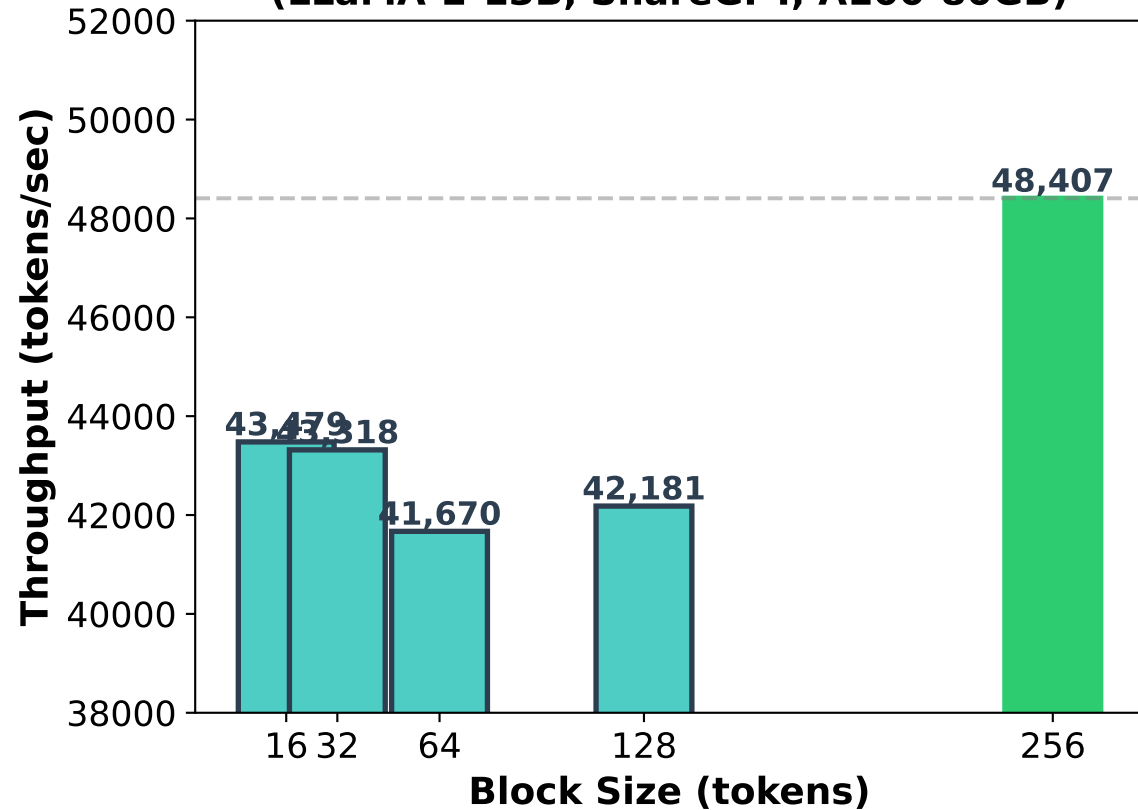


Throughput by Block Size
(LLaMA-2-13B, ShareGPT, A100-80GB)



Memory Efficiency by Block Size
(LLaMA-2-13B, ShareGPT, A100-80GB)

