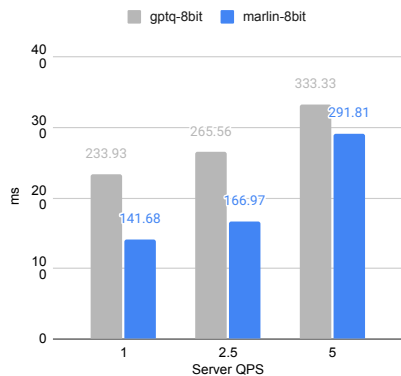
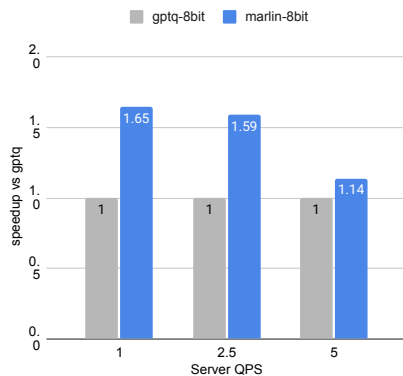


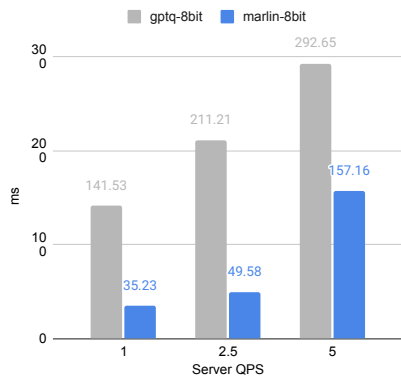
vLLM Server: TTFT, Yi-34B-Chat 8-bit GPTQ
A100 GPU, 256 prompt, 128 new tokens



vLLM Server: TTFT speedup, Yi-34B-Chat 8-bit GPTQ
A100 GPU, 256 prompt, 128 new tokens



vLLM Server: TPOT, Yi-34B-Chat 8-bit GPTQ
A100 GPU, 256 prompt, 128 new tokens



vLLM Server: TPOT speedup, Yi-34B-Chat 8-bit GPTQ
A100 GPU, 256 prompt, 128 new tokens

