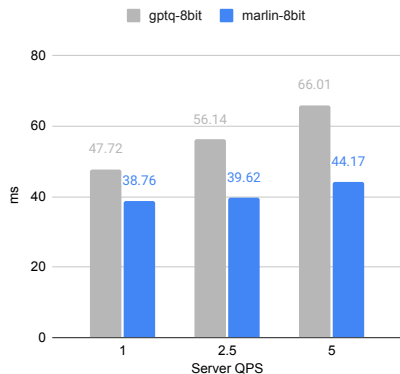


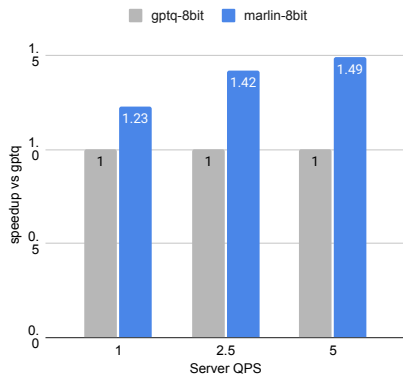
vLLM Server: TTFT, Llama-3-8B 8-bit GPTQ

A100 GPU, 256 prompt, 128 new tokens



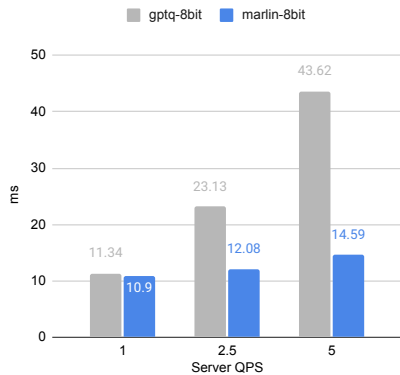
vLLM Server: TTFT speedup, Llama-3-8B 8-bit GPTQ

A100 GPU, 256 prompt, 128 new tokens



vLLM Server: TPOT, Llama-3-8B 8-bit GPTQ

A100 GPU, 256 prompt, 128 new tokens



vLLM Server: TPOT speedup, Llama-3-8B 8-bit GPTQ

A100 GPU, 256 prompt, 128 new tokens

