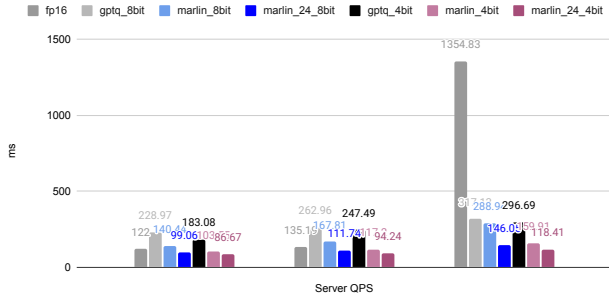


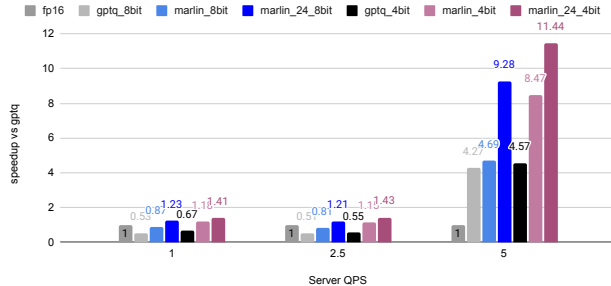
vLLM Server: TTFT Yi-34B-Chat: GPTQ/Marlin/Marlin24 4/8 bit

A100 GPU, 256 prompt, 128 new tokens



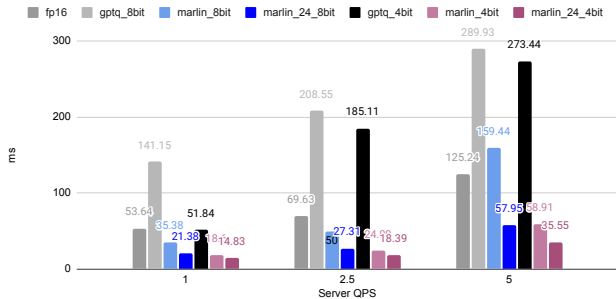
vLLM Server: TTFT speedup for Yi-34B-Chat: GPTQ/Marlin/Marlin24 4/8 bit

A100 GPU, 256 prompt, 128 new tokens



vLLM Server: TPOT Yi-34B-Chat: GPTQ/Marlin/Marlin24 4/8 bit

A100 GPU, 256 prompt, 128 new tokens



vLLM Server: TPOT speedup for Yi-34B-Chat: GPTQ/Marlin/Marlin24 4/8 bit

A100 GPU, 256 prompt, 128 new tokens

