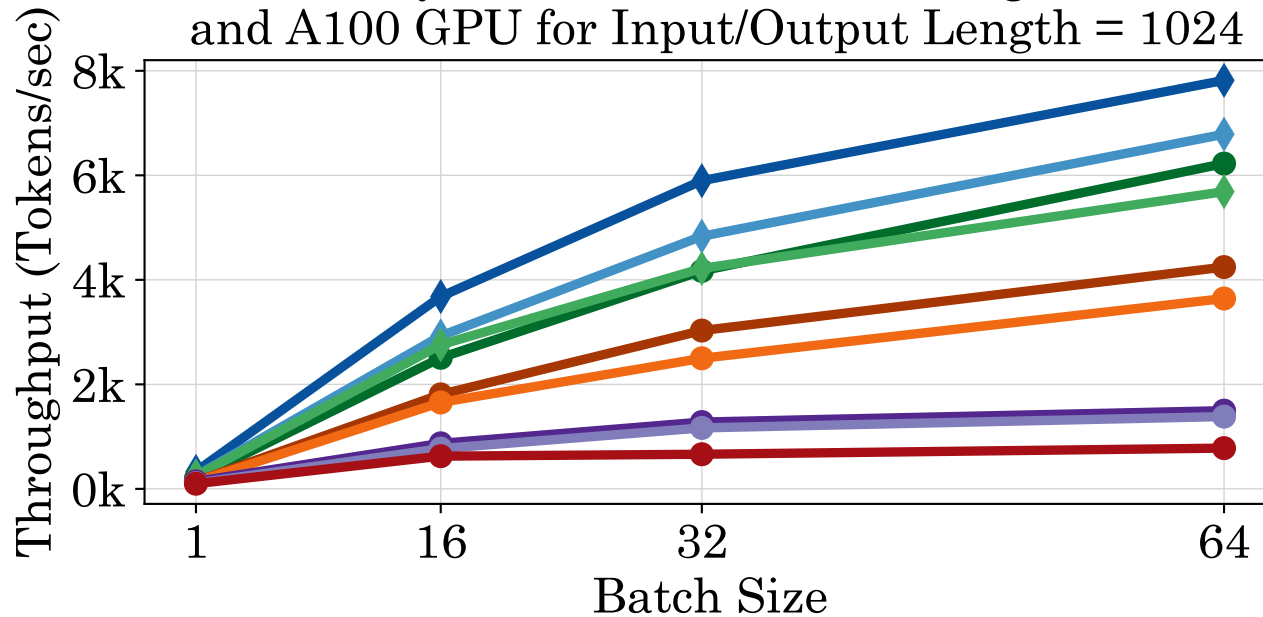


LLaMA-3-8B: Quantization Benchmarking on One H100
and A100 GPU for Input/Output Length = 1024



Hardware, Framework & {Weight Precision, KV Cache Precision}

- | | | |
|-----------------------------|-----------------------------|----------------------------|
| ◆ H100 vLLM {fp8, fp8} | ◆ H100 vLLM {fp16, fp8} | ● A100 TRT-LLM {int8, fp8} |
| ◆ H100 vLLM {fp16, fp16} | ● A100 TRT-LLM {fp16, int8} | ● A100 TRT-LLM {fp16, fp8} |
| ● A100 TRT-LLM {int8, int8} | ● A100 vLLM {fp16, fp16} | ● A100 vLLM {fp16, fp8} |