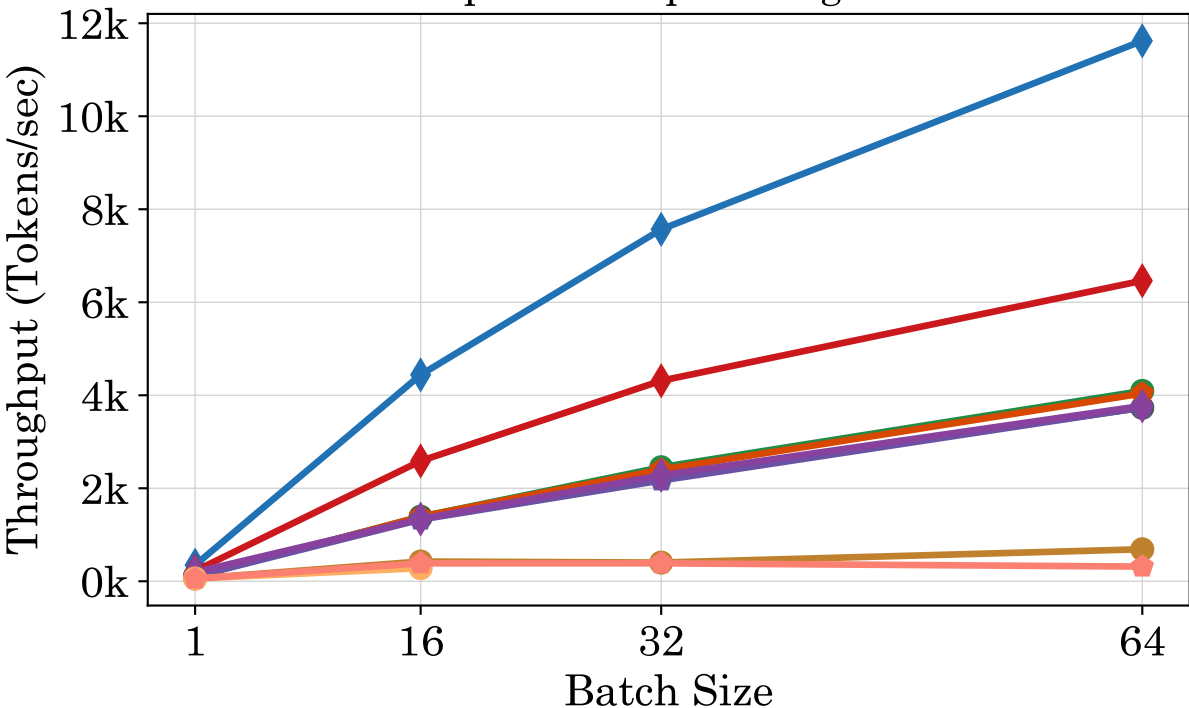


70B Models on four A100 and H100 GPUs
with Input & Output Length = 1024



Hardware, Framework & Model

- | | |
|-----------------------------|-----------------------------|
| ◆ H100 TRT-LLM Mixtral-8x7B | ◆ A100 TRT-LLM Mixtral-8x7B |
| ● H100 TRT-LLM LLaMA-2-70B | ◆ A100 vLLM Mixtral-8x7B |
| ● H100 vLLM LLaMA-2-70B | ● A100 TRT-LLM LLaMA-2-70B |
| ◆ H100 TRT-LLM LLaMA-3-70B | ● A100 vLLM LLaMA-2-70B |
| ◆ H100 vLLM LLaMA-3-70B | ◆ A100 TRT-LLM LLaMA-3-70B |