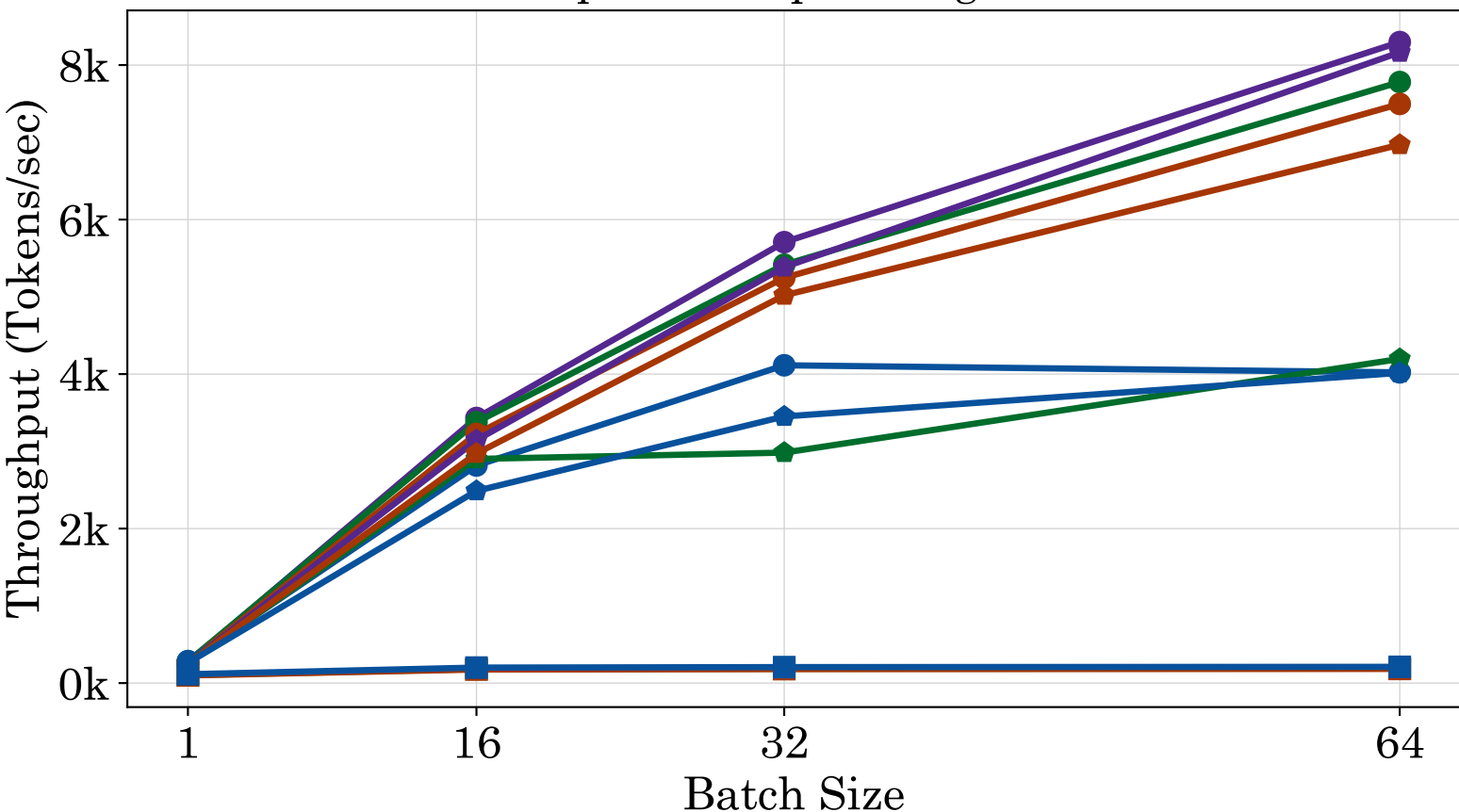


H100: Framework Comparison of 7B Models on One GPU with Input & Output Length = 1024



Framework & Model

- TRT-LLM Qwen2-7B
- TRT-LLM Mistral-7B
- TRT-LLM LLaMA-3-8B
- TRT-LLM LLaMA-2-7B
- vLLM Qwen2-7B
- vLLM Mistral-7B
- vLLM LLaMA-3-8B
- vLLM LLaMA-2-7B
- llama.cpp Mistral-7B
- llama.cpp LLaMA-3-8B
- llama.cpp LLaMA-2-7B