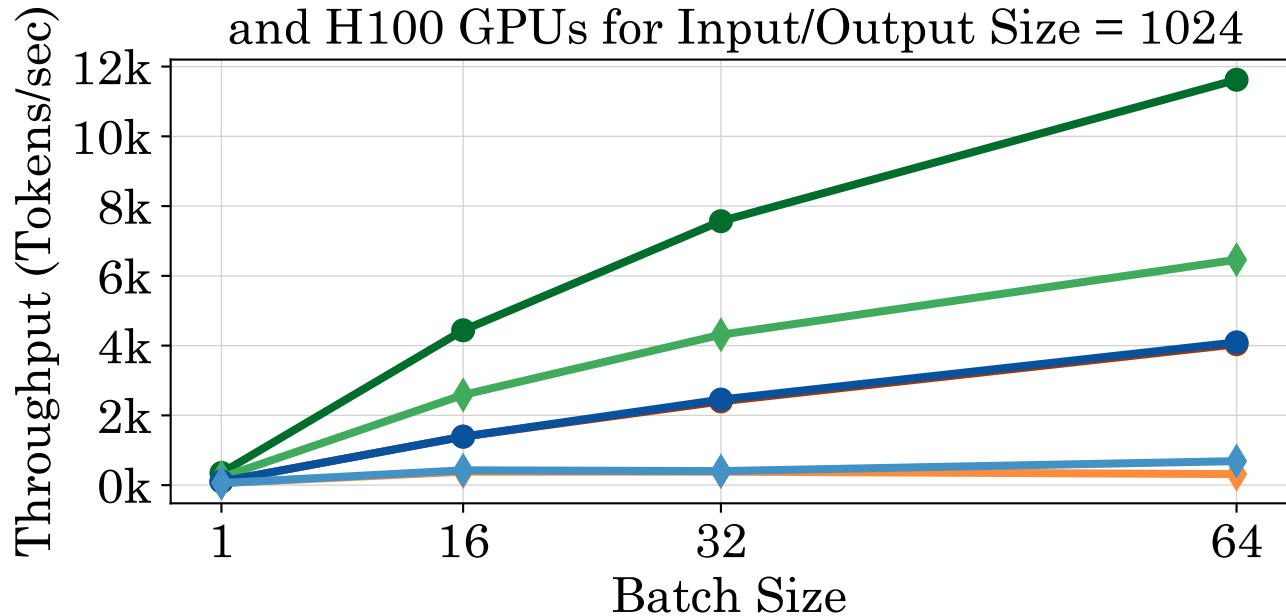


TensorRT-LLM: MoE and 70B Models on Four A100
and H100 GPUs for Input/Output Size = 1024



Hardware & Model

- H100 Mixtral-8x7B
- ◆ A100 Mixtral-8x7B
- H100 LLaMA-2-70B
- H100 LLaMA-3-70B
- ◆ A100 LLaMA-2-70B
- ◆ A100 LLaMA-3-70B