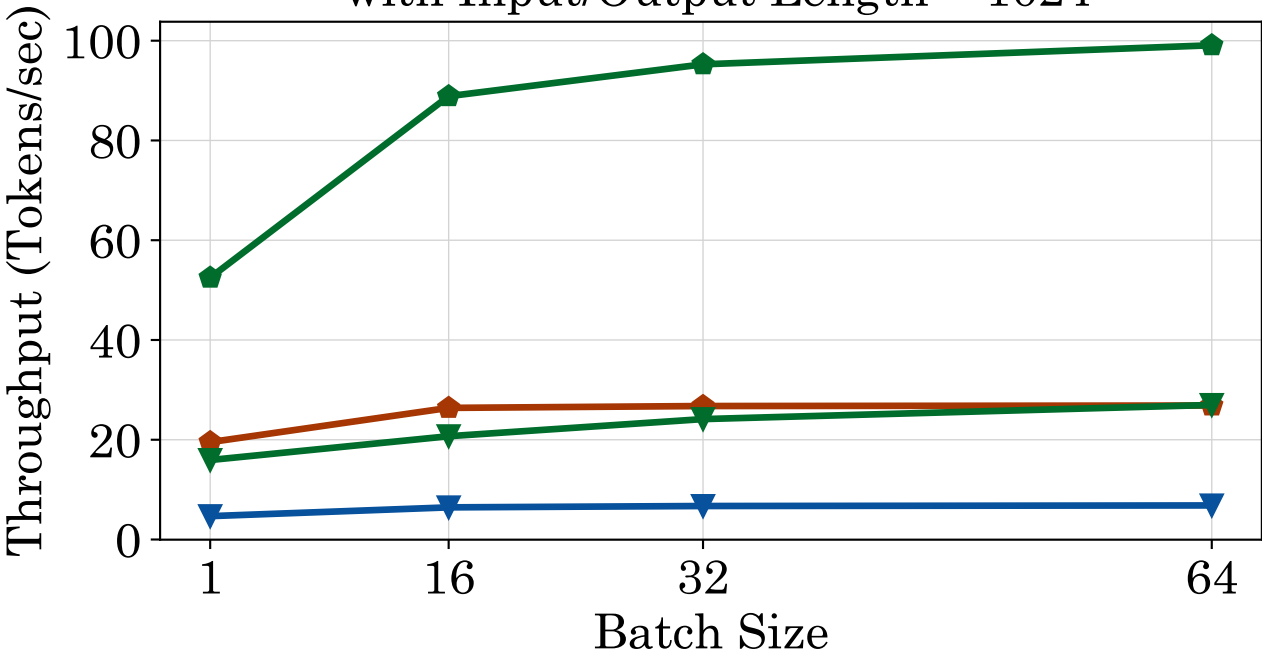


# llama.cpp: 70B Models on Four GPUs with Input/Output Length = 1024



Hardware & Model

- ◆ H100 Mixtral-8x7B
- ▼ MI250 Mixtral-8x7B
- ◆ H100 LLaMA-3-70B
- ▼ MI250 LLaMA-2-70B