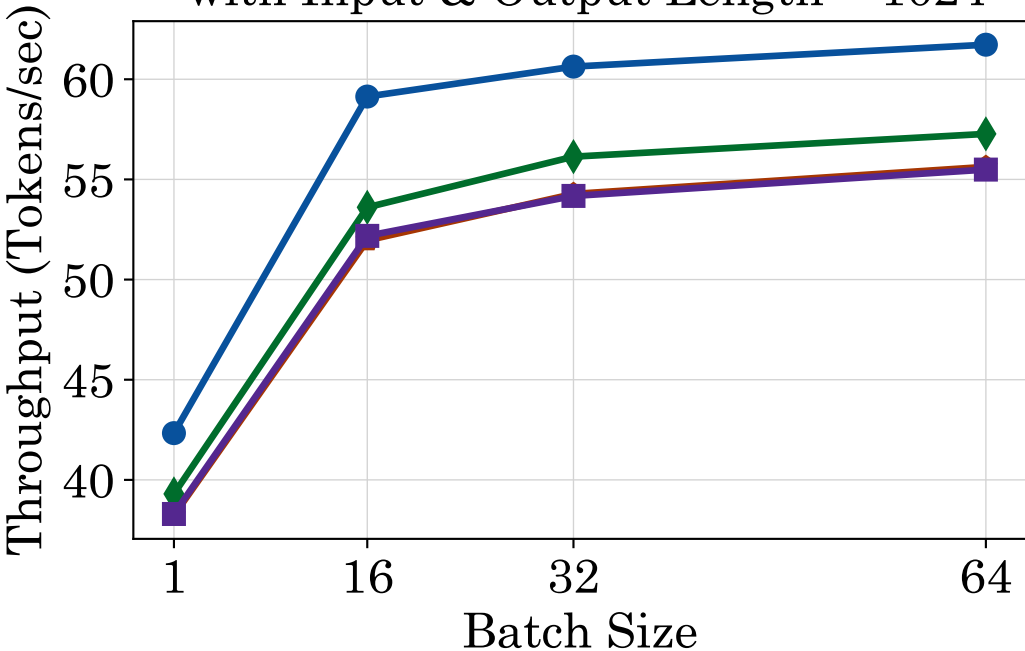


7B Models on One MI250 GPU using llama.cpp
with Input & Output Length = 1024



- Model
- LLaMA-2-7B
 - ◆ Mistral-7B
 - ◆ LLaMA-3-8B
 - Qwen2-7B