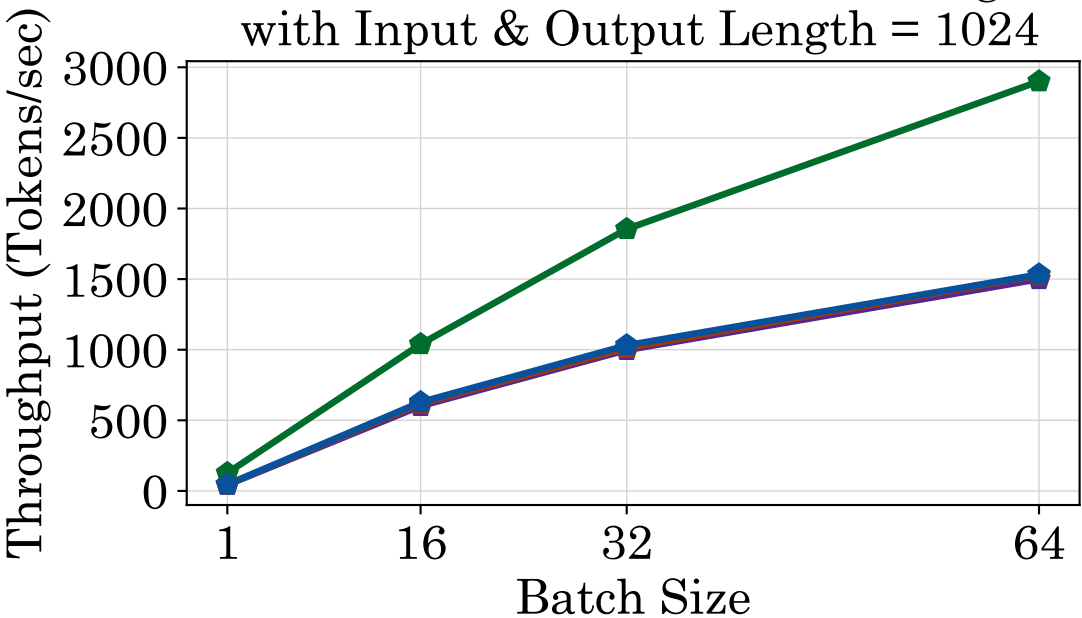


70B Models on Four MI250 GPUs using vLLM
with Input & Output Length = 1024



Model

Qwen2-72B

LLaMA-3-70B

Mixtral-8x7B

LLaMA-2-70B