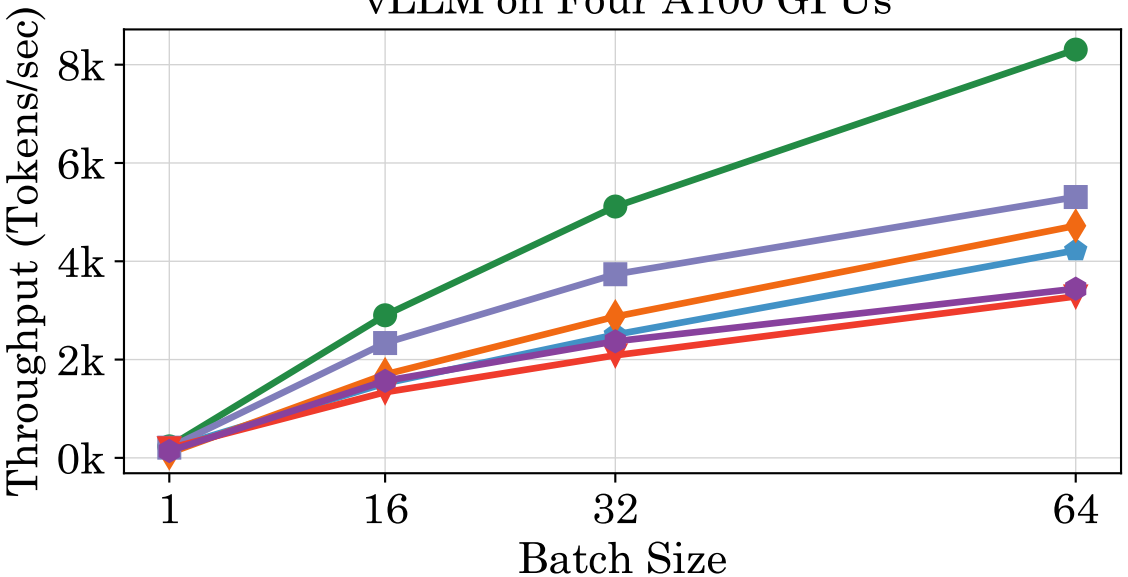


# Mixtral-8x7B: TRT-LLM vs DS-MII vs vLLM on Four A100 GPUs



Input/Output Length & Framework

- 128 TRT-LLM
- ◆ 128 DS-MII
- ▼ 2048 vLLM
- ◆ 128 vLLM
- 2048 TRT-LLM
- ◆ 2048 DS-MII