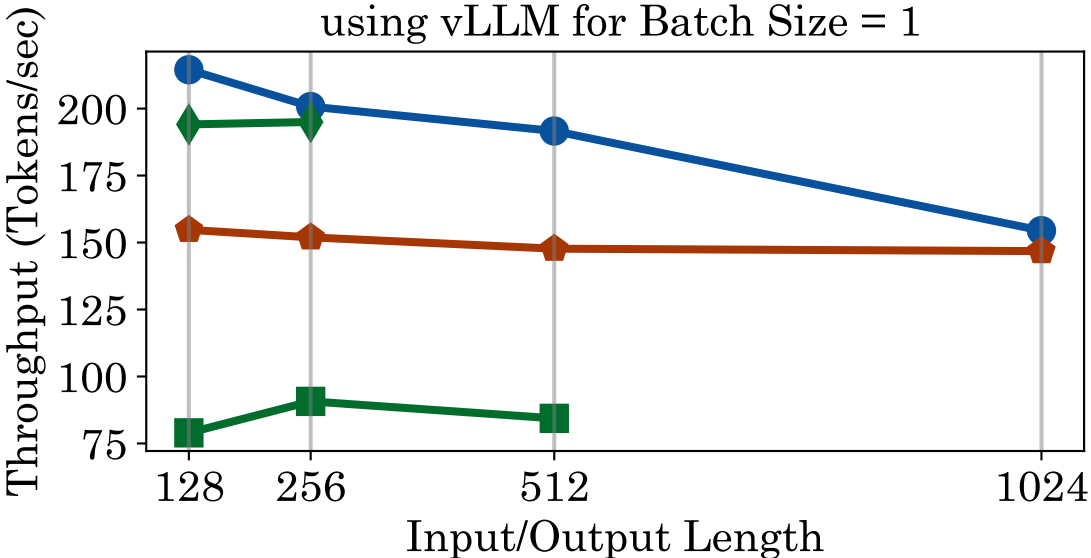


Speculative Decoding (SD) on One A100 GPU using vLLM for Batch Size = 1



Models & SD

- LLaMA-2-7B w SD
- ◆ Mixtral-8x7B w SD
- ◆ LLaMA-2-7B w/o SD
- Mixtral-8x7B w/o SD