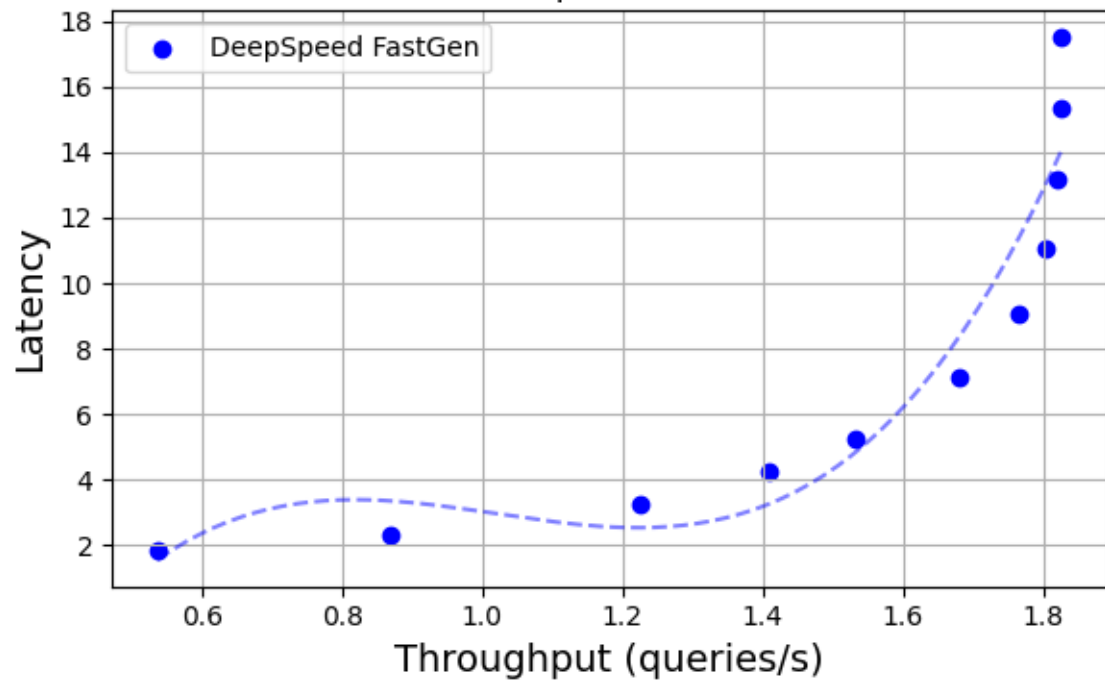


Model Llama 2 7B, Prompt: 2600, Generation: 60, TP: 1



Effective throughput (SLA prompt: 512 tokens/s, generation: 4 tokens/s)
Llama 2 7B Prompt: 2600, Generation: 60, TP: 1

