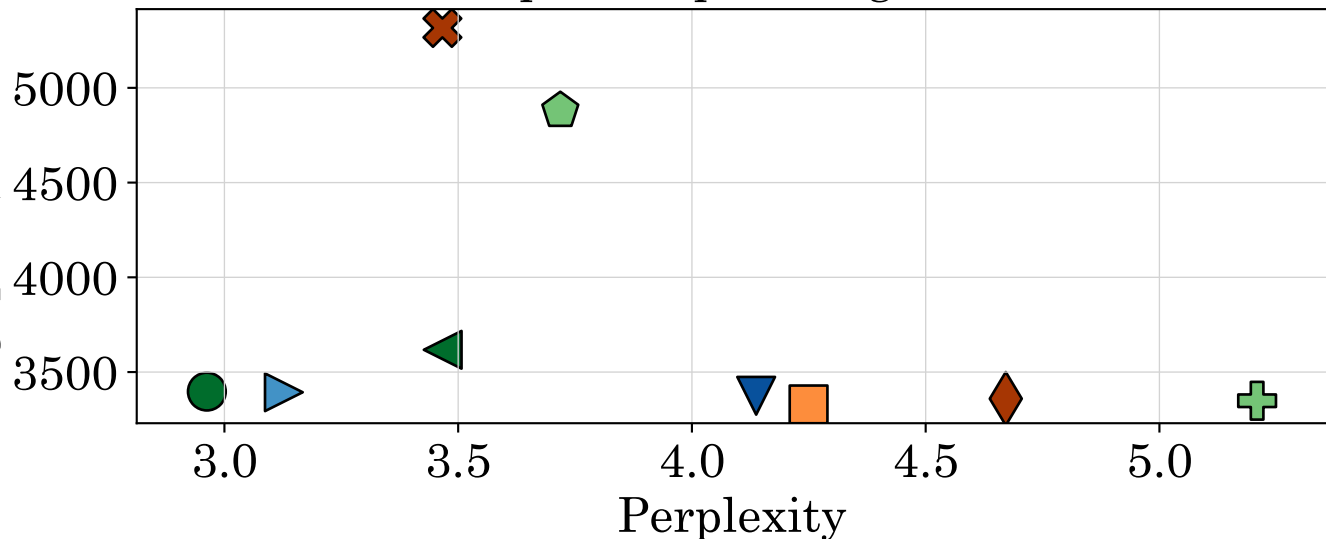


Perplexity vs Throughput Comparison of ~7B Models
using vLLM on One H100 GPU for Batch Size = 32
and Input/Output Length = 1024

Throughput (Tokens/sec)



Hardare & Model

