

Perplexity vs Throughput Comparison of ~7B Models
using vLLM on One A100 GPU for Batch Size = 32
and Input/Output Length = 1024

