

vLLM: Batch Size vs Input/Output Length of LLaMA-3-8B on a Single A100 GPU

