

vLLM: 7B Models on GPUs with Batch Size = 32
and Input/Output Length = 2048

