

KV Cache Block Size vs Batch Size for
Input/Output Length = 1024
on LLaMA-3-8B on a Single A100 GPU

