**LLM Inference Performance Comparison (LLaMA-2-13B)**

LLMIR outperforms vLLM by 22.4%