**Inference Benchmark: Framework Comparison**

**Average Latency (lower is better)**

Transformers: 1.501s
Unsloth: 0.918s
vLLM: 0.246s
Ollama: 0.440s

**Throughput (higher is better)**

Transformers: 95
Unsloth: 136
vLLM: 346
Ollama: 244

**Peak GPU Memory (lower is better)**

GPU Limit (12GB)
Transformers: 1.00
Unsloth: 1.14
vLLM: 10.90
Ollama: 1.99