**LLM Inference Throughput Comparison**
**NVIDIA A10 GPU (24GB)**

Throughput (tokens/second)

Legend:
- Llama-3-8B
- Mistral-7B

NIM Baseline

| Inference Framework | Llama-3-8B | Mistral-7B |
|---|---|---|
| NVIDIA NIM (TensorRT-LLM) | 31.7 | 33.6 |
| SGLang (RadixAttention) | 28.6 | 30.2 |
| HuggingFace TGI (FlashAttention) | 28.1 | 29.6 |
| vLLM (PagedAttention) | 28.0 | 29.0 |