一、 SGlang 框架概述

SGlang 是一个高性能推理引擎,专为混合专家(MoE)语言模型设计,如 DeepSeek-R1。它支持多节点张量并行计算,能够在多台机器上协同工作,从而 满足大规模模型的部署需求。此外,SGlang 还支持 FP8(W8A8)和 KV 缓存 优化,并通过 Torch Compile 技术进一步提升推理效率。

工具名称	性能表现	易用性	适用场 景	硬件需求	模型支持	部署方式	系统支持
SGLang	零开销批处 理提升1.1倍 吞吐量,缓 存感知负载 均衡提升1.9 倍,结构化 输出提速10 倍	需一定技术基础, 但提供完整API和示例	企推务并 景要化的级服高场需构出用	推荐 A100/H100, 支持多GPU部 署	全面支持主 流大模型, 特别优化 DeepSeek 等模型	Docker、 Python包	Linux
Ollama	继承 Ilama.cpp 的 高效推理能 力,提供便 捷的模型管 理和运行机 制	小白友好,提供图形界程序 一和命令支持RESTAPI	个发意证生学日答意等轻应景人者验、辅习常、写个量用开创 学助、问创作人级场	与 llama.cpp相同,但提供更简便的资源管理	模型库丰 富,涵盖 1700多款, 支持键下载 安装	独立应用 程序、 Docker、 REST API	Windows、 macOS、 Linux
VLLM	高效性能	较为复杂	研究开 发与商 业	CPU/GPU	广泛的模型 支持	本地部 署、容器 化	Linux, macOS, Windows
LLaMA.cpp	极高性能	相对复杂	大规模 商业应 用与研究	高端GPU	专属LLaMA 模型	本地部署、分布式部署	Linux, Windows

二、SGLang 部署 Qwen3

(一) 服务器硬件配置

- 服务器 1: CPU: 英特尔至强 Max 9468*2、GPU: HGX H20(96GB)*8、内存: 64GB*32、存储: 3.84T Nvme*4
- 服务器 2: CPU: 英特尔至强 Max 9468*2、GPU: HGX H20(96GB)*8、内存: 64GB*32、存储: 3.84T Nvme*4
- 网络: 25Gb 以太组网

(二) Docker (Recommended) 安装

● 国外镜像源

sudo docker pull lmsysorg/sglang:latest

● 国内镜像源

sudo docker pull docker.1ms.run/lmsysorg/sglang:latest

```
GongHang@server03:~$ sudo docker pull docker.1ms.run/lmsysorg/sglang:latest
latest: Pulling from lmsysorg/sglang
23828d760c7b: Pull complete
18ea39f449f2: Downloading 8.043MB
4f4fb700ef54: Download complete
a2f9ef49a25e: Downloading 50.79MB
1c7f2e233b5f: Download complete
6d00df402cf8: Downloading 62.68MB
fe4f2f514559: Waiting
999f45ff216e: Waiting
74e0ce3e6cbf: Waiting
dfae8e751f33: Pulling fs layer
6b6f1f09276f: Waiting
d9e68ca30619: Waiting
02006324a966: Waiting
3d1a32501ee0: Waiting
a79e61b6522e: Waiting
1b18d1e3a9a1: Waiting
8d56e893dbb9: Waiting
ee16b3126d05: Waiting
484c61037dda: Waiting
623b8bf9f8b2: Waiting
```

(三) 启动 SGLang 容器

1. 单台服务器部署 Qwen3-235B-A22B。

创建 sglang 容器

sudo docker run -itd --entrypoint /bin/bash --ipc host --gpus all --name sglang -p 44444:44444 -v /share2/server3/models:/workspace/models --rm docker.1ms.run/lmsysorg/sglang:latest

进入 sglang 容器终端



2. 两台服务器部署 Qwen3-235B-A22B。

分别创建 sglang 容器(此处"/share2/server3/models"为共享文件夹)
sudo docker run -itd --entrypoint /bin/bash --ipc host --gpus all --name sglang -p
44444:44444 -v /share2/server3/models:/workspace/models --rm
docker.1ms.run/lmsysorg/sglang:latest

分别进入 sglang 容器终端

● 主节点

python3 -m sglang.launch_server --model-path /workspace/models/ Qwen3-235B-A22B --tp 16 --dist-init-addr 10.0.10.11:5000 --nnodes 2 --node-rank 0 --trust-remote-

● 子节点

python3 -m sglang.launch_server --model-path /workspace/models/ Qwen3-235B-A22B --tp 16 --dist-init-addr 10.0.10.11:5000 --nnodes 2 --node-rank 1 --trust-remote-code

(四) 验证 SGLang 服务

1. 在 Open-WebUI(WebUI 安装过程,请参考官方教程)管理员界面中,设置 OpenAI API 外部连接。



2. 创建聊天窗口,选择 Qwen3-235B-A22B 模型。

