

# A100 部署 DeepSeek-R1-AWQ

## 环境说明

主机名	ip 地址	系统	规格
yichang-10-200-3-23	10.200.3.23	ubuntu 22.04	8 卡 A100

## 安装驱动和桥接器

安装 gcc:

```
1 apt-get update -y
2 apt install build-essential -y
```

Bash

安装驱动:

```
1 wget https://us.download.nvidia.com/tesla/560.35.03/NVIDIA-Linux-x86_64-560.35.03.run
2 # 注意：在交互式安装时要确认安装32位的兼容库，同时要Rebuild initramfs
3 sh NVIDIA-Linux-x86_64-560.35.03.run
```

Bash

需要说明的是驱动程序的版本与后续安装的桥接器的版本必须保持一致，所以在选择驱动版本时，需要确保桥接器也有对应的版本（并不是所有的驱动版本都能找到对应的桥接器的版本）

安装桥接器:

```
1 # 确保桥接器版本与驱动版本完全一致（包括次版本）
2 wget https://developer.download.nvidia.cn/compute/cuda/repos/ubuntu2204/x86_64/nvidia-fabricmanager-560_560.35.03-1_amd64.deb
3
4 dpkg -i nvidia-fabricmanager-560_560.35.03-1_amd64.deb
5 systemctl enable nvidia-fabricmanager --now
6 systemctl status nvidia-fabricmanager
```

Bash

之后重启服务器，配置持久模式:

```
1 nvidia-smi -pm 1
```

Bash

# 安装 docker 并配置 nvidia-runtime-toolkit

安装 docker:

```
Bash
1 # step 1: 安装必要的一些系统工具
2 sudo apt-get update
3 sudo apt-get install ca-certificates curl gnupg
4
5 # step 2: 信任 Docker 的 GPG 公钥
6 sudo install -m 0755 -d /etc/apt/keyrings
7 curl -fsSL https://mirrors.aliyun.com/docker-ce/linux/ubuntu/gpg | sudo gpg --dear
mor -o /etc/apt/keyrings/docker.gpg
8 sudo chmod a+r /etc/apt/keyrings/docker.gpg
9
10 # Step 3: 写入软件源信息
11 echo \
12 "deb [arch=$(dpkg --print-architecture) signed-by=/etc/apt/keyrings/docker.gpg]
https://mirrors.aliyun.com/docker-ce/linux/ubuntu \
13 "$(. /etc/os-release && echo "$VERSION_CODENAME)" stable" | \
14 sudo tee /etc/apt/sources.list.d/docker.list > /dev/null
15
16 # Step 4: 安装Docker
17 sudo apt-get update
18 sudo apt-get install docker-ce docker-ce-cli containerd.io docker-buildx-plugin do
cker-compose-plugin
19
20 # 安装指定版本的Docker-CE:
21 # Step 1: 查找Docker-CE的版本:
22 # apt-cache madison docker-ce
23 # docker-ce | 17.03.1~ce-0~ubuntu-xenial | https://mirrors.aliyun.com/docker-ce/
linux/ubuntu xenial/stable amd64 Packages
24 # docker-ce | 17.03.0~ce-0~ubuntu-xenial | https://mirrors.aliyun.com/docker-ce/
linux/ubuntu xenial/stable amd64 Packages
25 # Step 2: 安装指定版本的Docker-CE: (VERSION例如上面的17.03.1~ce-0~ubuntu-xenial)
26 # sudo apt-get -y install docker-ce=[VERSION]
```

配置 docker 使用 nvidia-runtime-toolkit:

```
Bash
1 curl -fsSL https://mirrors.ustc.edu.cn/libnvidia-container/gpgkey | sudo gpg --dear
mor -o /usr/share/keyrings/nvidia-container-toolkit-keyring.gpg
2
3 curl -s -L https://mirrors.ustc.edu.cn/libnvidia-container/stable/deb/nvidia-contai
ner-toolkit.list | \
4 sed 's#deb https://nvidia.github.io#deb [signed-by=/usr/share/keyrings/nvidia-con
tainer-toolkit-keyring.gpg] https://mirrors.ustc.edu.cn#g' | \
5 sudo tee /etc/apt/sources.list.d/nvidia-container-toolkit.list
6
7 apt update -y
8 apt install -y nvidia-container-toolkit nvidia-container-runtime
```

在/etc/docker/daemon.json 中添加如下配置:

```
{
  "default-runtime": "nvidia",
  "runtimes": {
    "nvidia": {
      "path": "/usr/bin/nvidia-container-runtime",
      "runtimeArgs": []
    }
  },
}
```

Bash

然后重启 docker:

```
systemctl restart docker
```

Bash

## 获取镜像

生成 vllm 镜像的 dockerfile 示例如下:

```
1 FROM docker.m.daocloud.io/nvidia/cuda:12.6.3-runtime-ubuntu22.04
2 RUN set -ex; \
3     apt update -y; \
4     apt install -y python3.10 python3-pip; \
5     pip install vllm==v0.7.2
```

Docker

构建镜像:

```
1 docker build -t vllm:v0.7.2 .
```

Docker

获取 open-webui 镜像:

```
1 docker pull ghcr.m.daocloud.io/open-webui/open-webui:main
```

Bash

## 获取模型

deepseek-r1-awq 是 671b 的参数，但是精度只有 int4，单台 8 卡 A100 可以完成部署。

下载模型：

```
Bash
1  mkdir /data
2  cd /data
3  git lfs install
4  git clone https://www.modelscope.cn/cognitivecomputations/DeepSeek-R1-awq.git
```

## 部署

部署 deepseek：

```
Bash
1  docker run -d --runtime nvidia \
2  --gpus all \
3  -v /data:/mnt/models \
4  -p 12345:12345 \
5  --ipc=host \
6  hub.wanjiedata.com/models/vllm:v0.7.2 \
7  python3 -m vllm.entrypoints.openai.api_server --host 0.0.0.0 --port 12345 --max-model-len 65536 --trust-remote-code --tensor-parallel-size 8 --quantization moe_wna16 --gpu-memory-utilization 0.97 --kv-cache-dtype fp8_e5m2 --calculate-kv-scales --served-model-name deepseek-reasoner --model /mnt/models/DeepSeek-R1-AWQ
```

部署 open-webui：

```
Bash
1  docker run -d -p 3030:8080 \
2  -e ENABLE_OLLAMA_API=false \
3  -e OPENAI_API_KEY=NULL \
4  -e OPENAI_API_BASE_URL=http://10.200.3.23:12345/v1 \
5  -e ENABLE_RAG_WEB_LOADER_SSL_VERIFICATION=false \
6  -v open-webui:/app/backend/data \
7  --name open-webui \
8  --restart always ghcr.m.daocloud.io/open-webui/open-webui:main
```