

# 如何在GPU云服务器上快速部署deepgpu-llm

---

- 1 离线部署：使用云市场镜像
  - 1.1 方式1：从ECS购买页进入
  - 1.2 方式2：官网直接搜索“deepgpu-llm”进入购买
  - 1.3 等待镜像安装成功
  - 1.4 升级安装deepgpu-llm（可选）
- 2 在线部署：基于ubuntu 22.04系统
  - 2.1 进入ECS购买页，选择所需镜像，开通机器
  - 2.2 等待系统安装和自动安装项目安装
  - 2.3 安装deepgpu-llm
  - 2.4 补充：ubuntu 20.04安装注意事项
- 3 在线部署：基于Centos 7.9系统
  - 3.1 进入ECS购买页，选择所需镜像，开通机器
  - 3.2 等待系统安装和自动安装项目安装
  - 3.3 安装deepgpu-llm
  - 3.4 补充：centos 8.5安装注意事项
- 4 在线部署：基于容器docker
  - 4.1 安装nvidia docker
  - 4.2 下载并启动docker images
  - 4.3 安装deepgpu-llm
- 5 运行：利用deepgpu-llm运行qwen模型
  - 5.1 检查deepgpu-llm安装状态和版本
  - 5.2 下载开源模型：modelscope
  - 5.3 运行deepgpu-llm：从原始模型目录加载
  - 5.4 运行deepgpu-llm：转换模型并运行

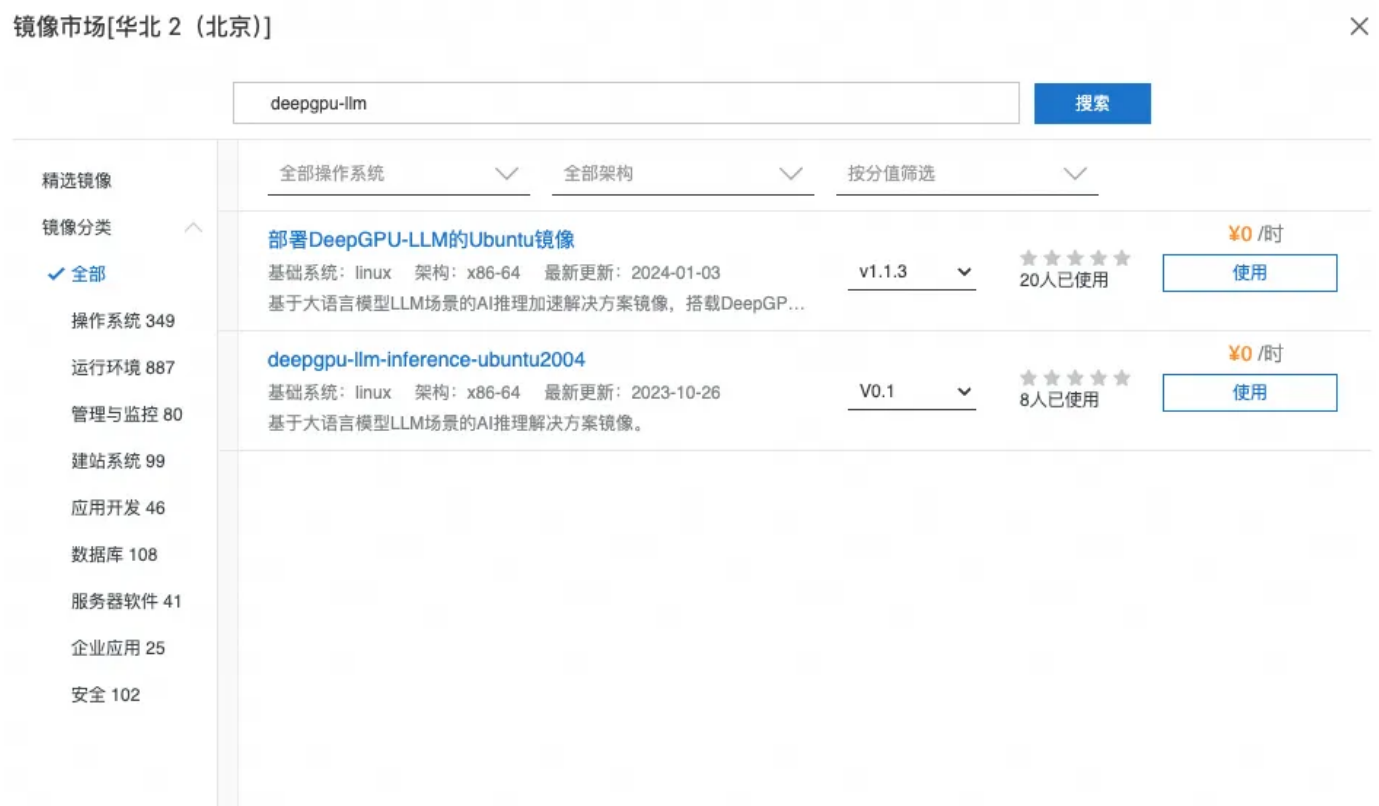
## 1 离线部署：使用云市场镜像

## 1.1 方式1：从ECS购买页进入

进入[ECS实例购买页](#)，按需配置付费类型、地域、网络、可用区、实例类型等；在镜像栏目选择“云市场镜像”-搜索“deepgpu-llm”选择所需镜像。配置好所有配置项目点击下单开通机器。



注：下面是老版本截图，建议选择最新更新的版本即可。



## 1.2 方式2：官网直接搜索“deepgpu-llm”进入购买

在阿里云官网右上方搜索框，直接搜索“deepgpu-llm 镜像”，在搜索结果中查看“来自：云市场 基础软件 操作系统”的页面，选择更新时间最新的版本点击进入。

deepgpu-llm 镜像

全部 云产品 技术解决方案 帮助文档 开发者社区 云市场

阿里云共为你找到 8,329 个搜索结果



### 容器镜像服务

容器镜像服务 (ACR) 提供安全的镜像托管能力, 稳定的国内外镜像构建服务, 便捷的镜像授权功能, 方便用户进行镜像全生命周期管理。

立即开通

帮助文档

产品入门

容器镜像基本概念

最佳实践

免密插件拉取容器镜像

构建容器DevOps

热门产品

容器服务 ACK

### 部署DeepGPU-LLM的Ubuntu镜像

本镜像为用户提供大语言模型LLM场景下的AI推理快速部署及加速的解决方案。大语言模型推理快速部署和性能加速的解决方案, 可用于加速llama、chatglm、...DeepGPU-LLM介绍: <https://help.aliyun.com/zh/egs/llm>

来自: [云市场 > 基础软件 > 操作系统](#) 标签: 开发者、应用开发、Linux、Ubuntu | 交付方式: 镜像 | 服务商: 阿里云计算有限公司

### 安装并使用DeepGPU-LLM

安装 DeepGPU-LLM 根据您的业务场景不同, 支持在GPU云服务器环境或者Docker环境下安装DeepGPU-LLM。在GPU云服务器环境下安装DeepGPU-LLM LLM模型的特性一般适用于GPU计算型实例, 本文以GPUi 按需选择付费类型、...

来自: [产品文档 > GPU云服务器 > 开发参考 > 安装并使用DeepGPU-LLM](#)

### 【Hello AI】安装并使用DeepGPU-LLM-处理大语言模型任务

```
tensors="pt").input_ids for i in range(8):start =time.time()output =model.generate([start_ids],generation_config)end =time.time()print("-",end -start)tokens =output[0].tolist()print(tokenizer.decode(tokens[C
```

来自: [开发者社区 > 弹性计算 > 云服务器ECS > 文章](#) | 标签: PyTorch、算法框架/工具、异构计算、Docker、容器、GPU云服务器 | 时间: 2024年01月26日

### deepgpu-llm-inference-ubuntu2004

本镜像为用户提供大语言模型LLM场景下的AI推理快速部署及加速的解决方案。生成式人工智能和大语言模型推理快速部署和加速的解决方案解决方案内容包括 基于DeepGPU LLM加速器快速实现llama、llama2、ch

来自: [云市场 > 基础软件 > 操作系统](#) 标签: 开发者、应用开发、Linux、Ubuntu | 交付方式: 镜像 | 服务商: 阿里云计算有限公司

在具体镜像的详情页里, 点击购买即可进入ECS购买页, deepgpu-llm镜像本身都是免费的。按ECS购买页提示填写相关信息, 确认镜像栏目选中所选镜像。由于镜像启动方式的不同, 一个镜像无法适配所有机型, 若无法匹配, 可返回选择其他deepgpu-llm镜像。

神龙  
DeepGPU

## 部署DeepGPU-LLM的Ubuntu镜像

基于大语言模型LLM场景的AI推理加速解决方案镜像, 搭载DeepGPU-LLM大语言模型专用加速套件。

近180天成交: 20单 评论: 0 ★ (0条)

交付方式: 镜像 架构: x86\_64 操作系统: Ubuntu 更新时间: 2024-01-03

计价方式: 该镜像为固定定价, 具体价格以ECS实例规格选配时的询价为准。

¥0 /月

续费: ¥0/月 按量价格: ¥0/小时

立即购买

免费体验

## 1.3 等待镜像安装成功

ssh登陆成功即完成部署。

具体操作，请参见[通过密码或密钥认证登录Linux实例](#)。

## 1.4 升级安装deepgpu-llm（可选）

根据5.1节方法检查deepgpu-llm的安装状态和版本，若觉得版本过低，可以根据本节方法来升级安装deepgpu-llm。

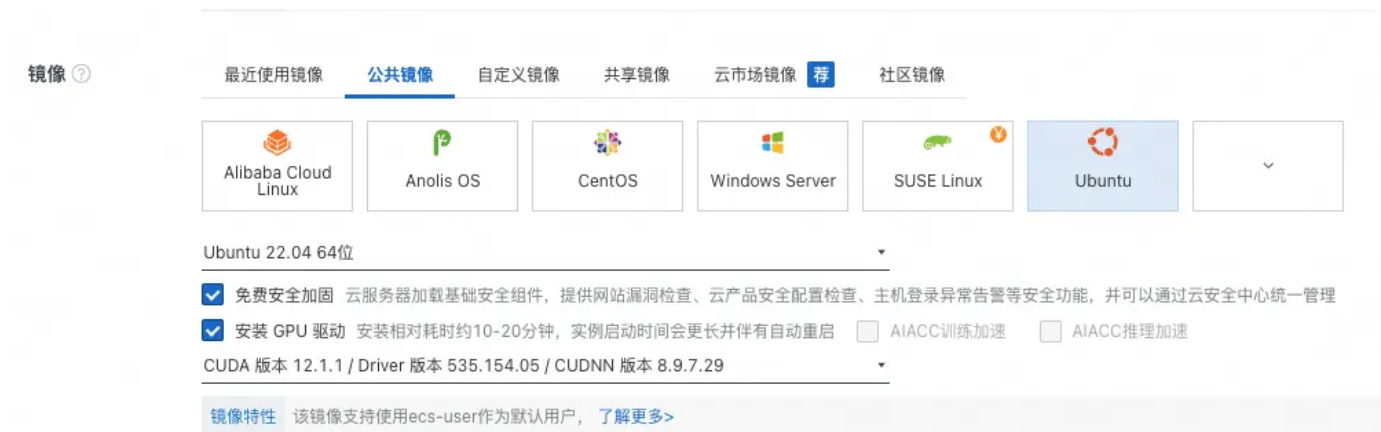
首先从[DeepGPU-LLM加速安装包](#)下载需要升级的安装包到对应目录，可使用wget命令进行下载。然后直接用pip命令安装即可。参考命令如下所示。

```
1  wget https://aiacc-inference-public-v2.oss-cn-hangzhou.aliyuncs.com/aiacc-i
nference-llm/deepgpu_llm-24.3%2Bpt2.1cu121-py3-none-any.whl
2  pip install deepgpu_llm-24.3+pt2.1cu121-py3-none-any.whl
```

# 2 在线部署：基于ubuntu 22.04系统

## 2.1 进入ECS购买页，选择所需镜像，开通机器

按需配置付费类型、地域、网络、可用区、实例类型等，在镜像栏目选择“公共镜像”-“Ubuntu”-“Ubuntu 22.04 64位”，可选直接勾选安装GPU驱动、CUDA和CUDNN，如下图所示。



## 2.2 等待系统安装和自动安装项目安装

ssh登陆可查看GPU驱动、CUDA、CUDNN等安装进展。

## 2.3 安装deepgpu-llm

逐一输入下面命令，在线安装deepgpu-llm的部分依赖项。

```
1 apt-get update
2 apt-get -y install python3.10 python3-pip openmpi-bin libopenmpi-dev curl v
im
```

选择需要安装的版本，在线安装deepgpu-llm及其依赖项（耗时较长，耐心等待下载和安装）。其中xx.x应替换为您实际的DeepGPU-LLM版本号。如何获取最新DeepGPU-LLM版本号，请参见[DeepGPU-LLM加速安装包](#)。

```
1 # for PyTorch 1.13
2 pip3 install deepgpu_llm==xx.x+pt1.13cu117 \
3     -f https://aiacc-inference-public-v2.oss-cn-hangzhou.aliyuncs.com/aiac
c-inference-llm/deepgpu_llm.html
4
5 # for PyTorch 2.0
6 pip3 install deepgpu_llm==xx.x+pt2.0cu117 \
7     -f https://aiacc-inference-public-v2.oss-cn-hangzhou.aliyuncs.com/aiac
c-inference-llm/deepgpu_llm.html
8
9 # for PyTorch 2.1
10 pip3 install deepgpu_llm==xx.x+pt2.1cu121 \
11     -f https://aiacc-inference-public-v2.oss-cn-hangzhou.aliyuncs.com/aiac
c-inference-llm/deepgpu_llm.html
```

## 2.4 补充：ubuntu 20.04安装注意事项

ubuntu 20.04系统上安装deepgpu-llm时，由于apt无法直接安装python 3.10，需要改成安装python 3.9，然后将系统自带的python版本改成默认python 3.9。具体安装python 3.9和调整默认python3版本的命令如下所示。

```

1 apt update
2 apt install software-properties-common
3 apt install python3.9
4 update-alternatives --install /usr/bin/python3 python3 /usr/bin/python3.8 1
5 update-alternatives --install /usr/bin/python3 python3 /usr/bin/python3.9 2
6 update-alternatives --config python3
7 python3 --version

```

## 3 在线部署：基于Centos 7.9系统

### 3.1 进入ECS购买页，选择所需镜像，开通机器

按需配置付费类型、地域、网络、可用区、实例类型等，在镜像栏目选择“公共镜像”-“CentOS”-“CentOS 7.9 64位”，可选直接勾选安装GPU驱动+CUDA+CUDNN，如下图所示。



### 3.2 等待系统安装和自动安装项目安装

ssh登陆可查看GPU驱动、CUDA、CUDNN等安装进展。

### 3.3 安装deepgpu-llm

逐一输入下面命令，完成deepgpu-llm部分依赖项目的在线安装。

```
1 yum install epel-release
2 yum update
3 yum install openmpi3 openmpi3-devel curl
4 wget https://repo.anaconda.com/miniconda/Miniconda3-latest-Linux-x86_64.sh
5 chmod +x Miniconda3-latest-Linux-x86_64.sh
6 ./Miniconda3-latest-Linux-x86_64.sh
```

其中openmpi3库安装成功后，需要修改环境变量（可添加至~/.bashrc）

```
1 export PATH=/usr/lib64/openmpi3/bin:$PATH
2 export LD_LIBRARY_PATH=/usr/lib64/openmpi3/lib:$LD_LIBRARY_PATH
```

其中miniconda环境需要按下面步骤进行使能和配置，使shell进入python 3.10环境（python 3.9也可，自行配置）。

```
1 /root/miniconda3/bin/conda init
2 source ~/.bashrc
3 conda create -n py310 python=3.10
4 conda activate py310
```

利用命令在线安装deepgpu-llm及其依赖项（耗时较长，耐心等待下载和安装）。也可参考2.3方法指定某个版本进行安装。

```
1 pip3 install deepgpu_llm -f https://aiacc-inference-public-v2.oss-cn-hangzhou.aliyuncs.com/aiacc-inference-llm/deepgpu_llm.html
```

### 3.4 补充：centos 8.5安装注意事项

centos 8.5系统中，yum或者dnf无法安装openmpi3和openmpi3-devel，这里我们改为安装openmpi和openmpi-devel。

```
1 yum install openmpi openmpi-devel
```

然后配置到环境变量里的openmpi相关目录也要改一下，改成下面这个。其他安装步骤和命令不变。

```
1 export PATH=/usr/lib64/openmpi/bin:$PATH
2 export LD_LIBRARY_PATH=/usr/lib64/openmpi/lib:$LD_LIBRARY_PATH
```

## 4 在线部署：基于容器docker

### 4.1 安装nvidia docker

若系统中已经安装了支持nvidia GPU的docker环境，本节可跳过。

#### 第一步：安装或者升级docker-ce

基于ubuntu操作系统，可以参考下面命令进行操作。

```
1 apt update
2 apt remove docker docker-engine docker-ce docker.io containerd runc
3 apt install apt-transport-https ca-certificates curl gnupg-agent software-p
  roPERTIES-common
4 curl -fsSL https://mirrors.aliyun.com/docker-ce/linux/ubuntu/gpg | sudo apt
  -key add -
5 apt-key fingerprint 0EBFCD88
6 add-apt-repository "deb [arch=amd64] https://mirrors.aliyun.com/docker-ce/l
  inux/ubuntu $(lsb_release -cs) stable"
7 apt update
8 apt install docker-ce
9 docker -v
```

在centos操作系统上，参考命令如下。



```

1 yum remove docker docker-client docker-client-latest docker-common docker-l
  latest docker-latest-logrotate docker-logrotate docker-engine
2 yum install -y yum-utils
3 yum-config-manager --add-repo https://download.docker.com/linux/centos/dock
  er-ce.repo
4 yum install docker-ce docker-ce-cli containerd.io
5 systemctl start docker
6 systemctl enable docker

```

## 第二步：安装nvidia container

ubuntu系统中，可参考下面命令完成nvidia container的安装。

```

1 $ curl -fsSL https://nvidia.github.io/libnvidia-container/gpgkey | sudo gp
  g --dearmor -o /usr/share/keyrings/nvidia-container-toolkit-keyring.gpg \
2   && curl -s -L https://nvidia.github.io/libnvidia-container/stable/deb/nv
  idia-container-toolkit.list | \
3     sed 's#deb https://#deb [signed-by=/usr/share/keyrings/nvidia-containe
  r-toolkit-keyring.gpg] https://#g' | \
4     sudo tee /etc/apt/sources.list.d/nvidia-container-toolkit.list \
5   && \
6     sudo apt-get update
7
8 $ apt-get install -y nvidia-container-toolkit
9 $ nvidia-ctl runtime configure --runtime=docker
10 $ systemctl restart docker

```

centos系统上的安装命令有些区别，其命令罗列在下面。

```

1 distribution=$(. /etc/os-release;echo $ID$VERSION_ID)
2 curl -s -L https://nvidia.github.io/nvidia-docker/$distribution/nvidia-dock
  er.repo | sudo tee /etc/yum.repos.d/nvidia-docker.repo
3 yum clean expire-cache
4 yum install -y nvidia-docker2
5 systemctl restart docker

```

## 4.2 下载并启动docker images

以下两个images都可以顺利安装deepgpu-llm，任选一个即可。

```
Plain Text |
1  pytorch/pytorch:2.1.0-cuda12.1-cudnn8-devel
2  nvidia/cuda:12.1.0-devel-ubuntu22.04
```

docker下载可参考下面命令，以pytorch/pytorch:2.1.0-cuda12.1-cudnn8-devel为例。

```
Plain Text |
1  docker pull pytorch/pytorch:2.1.0-cuda12.1-cudnn8-devel
```

docker启动命令如下

```
Plain Text |
1  docker run -ti --gpus all --name="deepgpu_llm" --network=host \
2      -v /root/workspace:/root/workspace \
3      --shm-size 5g pytorch/pytorch:2.1.0-cuda12.1-cudnn8-devel
```

其中关键参数说明如下：

参数项	说明
--shm-size	指定容器的共享内存大小，其大小会影响服务器部署。  例如：--shm-size 5g表示将共享内存大小设置为5 GB。您可以根据需要调整此值，以满足您的模型转换所需的内存需求。  (模型越大，模型转换需要的共享内存越大；若硬件限制无法扩展，可按5.4节方法提前转换模型，再进行部署)
-v /root/workspace:/root/workspace	将主机目录映射到Docker中的相应目录，使得主机和Docker之间可以共享文件，请根据自己实际环境情况进行映射。
nvidia/cuda:12.1.0-devel-ubuntu22.04	PyTorch的Docker映像标签
pytorch/pytorch:2.1.0-cuda12.1-cudnn8-devel	

### 4.3 安装deepgpu-llm

安装步骤与其他环境基本一致，具体命令如下所示。

```
Plain Text |
1 apt-get update
2 apt-get -y install python3.10 python3-pip openmpi-bin libopenmpi-dev curl v
im
3 pip3 install deepgpu_llm -f https://aiacc-inference-public-v2.oss-cn-hangzh
ou.aliyuncs.com/aiacc-inference-llm/deepgpu_llm.html
```

## 5 运行：利用deepgpu-llm运行qwen模型

### 5.1 检查deepgpu-llm安装状态和版本

利用以下命令查看deepgpu-llm安装版本信息。

```
Plain Text |
1 pip list | grep deepgpu-llm
```

```
root@iZbp1iyvb2cr96l6jh1f47Z:~/deepgpu/app# pip list | grep deepgpu-llm
deepgpu-llm                24.3+pt2.1cu121
```

这个命令可以查询更多deepgpu-llm相关信息。

```
Plain Text |
1 pip show -f deepgpu-llm
```

### 5.2 下载开源模型：modelscope

modelscope是阿里达摩院提供的开源模型平台，下载modelscope格式的模型有如下两种方式。本步骤以通义千问Qwen模型为例。

- git lfs clone命令方式
  - i. 进入modelscope官网，搜索模型名称（例如qwen）。
  - ii. 在搜索页面的模型库区域，单击通义千问-7B-Chat。
  - iii. 找到modelscope的专属模型名并复制模型ID。



iv. 执行以下命令，构建下载命令并下载模型。

```
Plain Text |
1 git-lfs clone https://modelscope.cn/qwen/Qwen-7B-Chat.git
```

o ModelScope库中的snapshot\_download方式

- i. 进入modelscope官网，搜索模型名称（例如qwen）。
- ii. 在搜索页面的模型库区域，单击通义千问-7B-Chat。
- iii. 找到modelscope的专属模型名并复制模型ID。



iv. 准备download\_from\_modelscope.py脚本。

```
Python |
1 import argparse
2 import shutil
3 from modelscope.hub.snapshot_download import snapshot_download
4 parser = argparse.ArgumentParser(description='download from modelscope')
5 parser.add_argument('--model_name', help='the download model name')
6 parser.add_argument('--version', help='the model version')
7 args = parser.parse_args()
8 base_dir = '/root/deepgpu/modelscope'
9 model_dir = snapshot_download(args.model_name, cache_dir=base_dir, revision
=args.version)
10 print(model_dir)
```

v. 执行以下命令，下载模型。模型下载前，您需要在通义千问-7B-Chat页面的模型文件页签下查看模型版本号。本命令以模型版本号为v.1.1.7为例。



```
Python |  
1 python3 download_from_modelscope.py --model_name qwen/Qwen-7B-Chat --version v1.1.7
```

### 5.3 运行deepgpu-llm：从原始模型目录加载

deepgpu-llm提供了llama\_cli、qwen\_cli、baichuan\_cli和chatglm\_cli等脚本帮助用户直接运行相应类别的LLM模型，可以查看“--help”获取帮助，查询具体配置项及其功能。

例如，可以通过下面命令来运行qwen\_cli脚本加载qwen-7b-chat模型推理进行对话。

```
Plain Text |  
1 qwen_cli --model_dir /root/deepgpu/models/Qwen-7B-Chat --tp_size 1 --precision fp16
```



上一节中直接运行的方法，会在初始化过程中在线转换模型，且该过程对CPU内存/shm\_size有一定需求。对于某些受限场景，可以提前做好模型转换，然后部署运行即可。针对不同模型开发了不同的转换命令，其中llama系列使用huggingface\_llama\_convert，chatglm模型使用huggingface\_glm\_convert，chatglm2和chatglm3使用huggingface\_chatglm2\_convert，baichuan系列模型使用huggingface\_baichuan\_convert，qwen v1版本模型转换使用huggingface\_qwen\_convert，qwen v1.5版本模型转换使用huggingface\_qwen15\_convert。

以qwen1.5-7b-chat模型为例，其转换命令如下所示，-in\_file指向下载的模型目录，-saved\_dir指向生成的模型目录，-infer\_gpu\_num设定推理运行的GPU数量（即模型切分份数），-weight\_data\_type设置模型权重使用的数据类型，与预期计算的类型一致，可选fp16、bf16和fp32，-model\_name为模型名称。

```
1 huggingface_qwen15_convert -in_file /root/deepgpu/models/Qwen1.5-7B-Chat -saved_dir /root/deepgpu/models/deepgpu/qwen1.5-7b-chat -infer_gpu_num 1 -weight_data_type fp16 -model_name qwen1.5-7b-chat
```

模型转换好后，按下面命令运行推理。其中--tp\_size配置的参数需要与转换时-infer\_gpu\_num设定的参数一致；--precision用于设置是否要对权重进行量化，可选fp16、int8和int4。

```
1 qwen_cli --tokenizer_dir /root/deepgpu/models/Qwen1.5-7B-Chat --model_dir /root/deepgpu/models/deepgpu/qwen1.5-7b-chat/1-gpu/ --tp_size 1 --precision fp16
```

```
root@iZj6c97rvxtrv2379pt4u5Z:~# qwen_cli --tokenizer_dir /root/deepgpu/models/Qwen1.5-7B-Chat --model_dir /root/deepgpu/models/deepgpu/qwen1.5-7b-chat/1-gpu/ --tp_size 1 --precision fp16
===== Argument =====
model_dir: /root/deepgpu/models/deepgpu/qwen1.5-7b-chat/1-gpu/
tp_size: 1
precision: fp16
tokenizer_dir: /root/deepgpu/models/Qwen1.5-7B-Chat
version: v1
=====
Special tokens have been added in the vocabulary, make sure the associated word embeddings are fine-tuned or trained.
running in fp16 mode, cb is False
欢迎使用 DeepGPU 加速版 Qwen 模型，输入内容即可进行对话，clear 清空对话历史，stop 终止程序

用户: 你好，你是谁？
我是阿里云推出的一种超大规模语言模型，我叫通义千问
Cost time: 0.36 s
Throughput: 47.16 tokens/s
```