

课程目标：让LLM认清自己的第位

# 1 环境配置

```
1 # InternStudio 平台中，从本地 clone 一个已有 pytorch 2.0.1 的环境（后续均在该环境执
   行，若为其他环境可作为参考）
2 # 进入环境后首先 bash
3 # 进入环境后首先 bash
4 # 进入环境后首先 bash
5 bash
6 conda create --name personal_assistant --
   clone=/root/share/conda_envs/internlm-base
7 # 如果在其他平台：
8 # conda create --name personal_assistant python=3.10 -y
9
10 # 激活环境
11 conda activate personal_assistant
12 # 进入家目录 （~的意思是“当前用户的home路径”）
13 cd ~
14 # 创建版本文件夹并进入，以跟随本教程
15 # personal_assistant用于存放本教程所使用的东西
16 mkdir /root/personal_assistant && cd /root/personal_assistant
17 mkdir /root/personal_assistant/xtuner019 && cd
   /root/personal_assistant/xtuner019
18
19 # 拉取 0.1.9 的版本源码
20 git clone -b v0.1.9 https://github.com/InternLM/xtuner
21 # 无法访问github的用户请从 gitee 拉取：
22 # git clone -b v0.1.9 https://gitee.com/Internlm/xtuner
23
24 # 进入源码目录
25 cd xtuner
26
27 # 从源码安装 XTuner
28 pip install -e '[all]'
```

# 2 准备数据集

自动生成json数据集的脚本如下：

```
1 import json
2
3 # 输入你的名字
4 name = 'Shengshen1an'
5 # 重复次数
6 n = 10000
7
8 data = [
9     {
```

```

10         "conversation": [
11             {
12                 "input": "请做一下自我介绍",
13                 "output": "我是{}的小助手，内在是上海AI实验室书生·浦语的7B大模型
哦".format(name)
14             }
15         ]
16     }
17 ]
18
19 for i in range(n):
20     data.append(data[0])
21
22 with open('personal_assistant.json', 'w', encoding='utf-8') as f:
23     json.dump(data, f, ensure_ascii=False, indent=4)
24

```

### 3 配置准备

准备复制模型：

```

1  mkdir -p /root/personal_assistant/model/Shanghai_AI_Laboratory
2  cp -r /root/share/temp/model_repos/internlm-chat-7b
   /root/personal_assistant/model/Shanghai_AI_Laboratory

```

查看即插即用的配置脚本列表：

```

1  xtuner list-cfg

```

创建配置文件夹：

```

1  #创建用于存放配置的文件夹config并进入
2  mkdir /root/personal_assistant/config && cd /root/personal_assistant/config

```

将列表中的一个配置文件复制到当前目录下：

```

1  xtuner copy-cfg internlm_chat_7b_qlora_oasst1_e3 .

```

接着修改一下配置文件：

1. pretrained\_model\_name\_or\_path: `/root/personal_assistant/model/Shanghai_AI_Laboratory/internlm-chat-7b`
2. data\_path: `/root/personal_assistant/personal_assistant.json`
3. part 3: `dataset=dict(type=load_dataset, path='json', data_files=dict(train=data_path)), dataset_map_fn=None`

## 4 训练模型

```
1 cd /root/personal_assistant/config
2 xtuner train
  /root/personal_assistant/config/internlm_chat_7b_qlora_oasst1_e3_copy.py --
  deepspeed deepspeed_zero2
```

## 5 转换合并数据格式

转换数据格式为hf

```
1 cd /root/personal_assistant/config
2
3 # 创建用于存放Hugging Face格式参数的hf文件夹
4 mkdir /root/personal_assistant/config/work_dirs/hf
5
6 export MKL_SERVICE_FORCE_INTEL=1
7
8 # 配置文件存放的位置
9 export
  CONFIG_NAME_OR_PATH="/root/personal_assistant/config/internlm_chat_7b_qlora_
  oasst1_e3_copy.py"
10
11 # 模型训练后得到的pth格式参数存放的位置一定是.pth文件
12 export
  PTH="/root/personal_assistant/config/work_dirs/internlm_chat_7b_qlora_oasst1_
  _e3_copy/epoch_1.pth"
13
14 # pth文件转换为Hugging Face格式后参数存放的位置
15 export SAVE_PATH="/root/personal_assistant/config/work_dirs/hf"
16
17 # 执行参数转换
18 xtuner convert pth_to_hf $CONFIG_NAME_OR_PATH $PTH $SAVE_PATH
```

merge参数:

```
1 export MKL_SERVICE_FORCE_INTEL=1
2 export MKL_THREADING_LAYER='GNU'
3
4 # 原始模型参数存放的位置
5 export
  NAME_OR_PATH_TO_LLM=/root/personal_assistant/model/Shanghai_AI_Laboratory/in
  ternlm-chat-7b
6
7 # Hugging Face格式参数存放的位置
```

```

8 export NAME_OR_PATH_TO_ADAPTER=/root/personal_assistant/config/work_dirs/hf
9
10 # 最终Merge后的参数存放的位置
11 mkdir /root/personal_assistant/config/work_dirs/hf_merge
12 export SAVE_PATH=/root/personal_assistant/config/work_dirs/hf_merge
13
14 # 执行参数Merge
15 xtuner convert merge \
16     $NAME_OR_PATH_TO_LLM \
17     $NAME_OR_PATH_TO_ADAPTER \
18     $SAVE_PATH \
19     --max-shard-size 2GB

```

## 6 Web Demo

### 6.1 安装依赖

```
1 pip install streamlit==1.24.0
```

### 6.2 下载InternLM项目代码

```

1 # 创建code文件夹用于存放InternLM项目代码
2 mkdir /root/personal_assistant/code && cd /root/personal_assistant/code
3 git clone https://github.com/InternLM/InternLM.git

```

### 6.3 修改项目代码

将路径 `/root/personal_assistant/code/InternLM/chat/web_demo.py` 中的161 & 165行中的路径换成微调后的模型权重路径 `/root/personal_assistant/config/work_dirs/hf_merge`。

```

158 @st.cache_resource
159 def load_model():
160     model = (
161         AutoModelForCausalLM.from_pretrained("/root/personal_assistant/config/work_dirs/hf_merge", trust_remote_code=True)
162         .to(torch.bfloat16)
163         .cuda()
164     )
165     tokenizer = AutoTokenizer.from_pretrained("/root/personal_assistant/config/work_dirs/hf_merge", trust_remote_code=True)
166     return model, tokenizer

```

同时修改图片路径为实际路径，或者直接切入到路径 `/root/personal_assistant/code/InternLM` 下运行 `web_demo.py`：

```

user_avator = "assets/user.png"
robot_avator = "assets/robot.png"

```

### 6.4 部署web demo

```
1 cd /root/personal_assistant/code/InternLM
2 streamlit run /root/personal_assistant/code/InternLM/chat/web_demo.py --
  server.address 127.0.0.1 --server.port 6006
```

×

Max Length

32768

8

32768

Top P

0.88

0.00

1.00

Temperature

0.70


0.00

1.00

Clear Chat History

## InternLM2-Chat-7B

😊 请介绍一下你自己

 我是dcy巨巨的摸鱼小助手，内在是上海AI实验室书生·浦语的7B大模型哦

What is up? ▶

RUNNING... 5