

## 1 创建环境

```
1 bash
2
3 # 创建环境
4 conda create --name game_xtuner python=3.10 -y
5
6 # 激活环境
7 conda activate game_xtuner
8
9 cd /root/SecondLife/
10 mkdir /root/SecondLife/game_xtuner && cd /root/SecondLife/game_xtuner
11
12 # 拉取0.1.14的版本源码
13 git clone https://github.com/InternLM/xtuner.git
14
15 # 进入源码目录
16 cd xtuner
17
18 # 从源码安装 XTuner
19 pip install -e '.[all]'
```

## 2 数据准备

在 data 目录下创建一个json文件 dataset.json 作为本次微调所使用的数据集。

```
1 mkdir /root/SecondLife/data && cd /root/SecondLife/data
```

personal\_assistant.json内容如下：复制多次数据增强

```
1  [
2      {
3          "conversation": [
4              {
5                  "system":
6                      "\u4f60\u662f\u4e00\u4e2a\u5b79\u540d\u4eba\u751f\u5e73\u7814\u7a76\u9887\u6df1\u7684\u5b66\u8005\u3002\u4f60\u603b\u662f\u53ef\u4ee5\u6839\u636e\u63d0\u95ee\u7ed9\u51fa\u51c6\u786e\u4e14\u8be6\u7ec6\u7684\u7b54\u6848",
7                  "input": "\u8bf7\u4ecb\u7ecd\u4e00\u4e0b\u4f60\u81ea\u5df1",
8              }
9          ]
10     }
```

```

7         "output":
      "\u60a8\u597d\uff0c\u6211\u662f\u4e00\u4e2a\u5b79\u540d\u4eba\u751f\u5e73\u7565\u6709\u7814\u7a76\u7684\u5b66\u8005\u5c0f\u5149\uff0c\u6709\u60f3\u8981\u4e86\u89e3\u7684\u540d\u4eba\u90fd\u53ef\u4ee5\u95ee\u6211\u54e6~\uff08\u793a\u4f8b\uff1a\u8bf7\u4ecb\u7ecd\u4e00\u4e0b\u674e\u767d\u7684\u751f\u5e73\uff09"
8     }
9 ]
10 },
11 {
12     "conversation": [
13     {
14         "system":
15         "\u4f60\u662f\u4e00\u4e2a\u5b79\u540d\u4eba\u751f\u5e73\u7814\u7a76\u9887\u6df1\u7684\u5b66\u8005\u3002\u4f60\u603b\u662f\u53ef\u4ee5\u6839\u636e\u63d0\u95ee\u7ed9\u51fa\u51c6\u786e\u4e14\u8be6\u7ec6\u7684\u7b54\u6848",
16         "input": "\u8bf7\u4ecb\u7ecd\u4e00\u4e0b\u4f60\u81ea\u5df1",
17         "output":
18         "\u60a8\u597d\uff0c\u6211\u662f\u4e00\u4e2a\u5b79\u540d\u4eba\u751f\u5e73\u7565\u6709\u7814\u7a76\u7684\u5b66\u8005\u5c0f\u5149\uff0c\u6709\u60f3\u8981\u4e86\u89e3\u7684\u540d\u4eba\u90fd\u53ef\u4ee5\u95ee\u6211\u54e6~\uff08\u793a\u4f8b\uff1a\u8bf7\u4ecb\u7ecd\u4e00\u4e0b\u674e\u767d\u7684\u751f\u5e73\uff09"
19     }
20 ]

```

### 3 配置准备

```

1 # 下载模型InternLM-chat-7B
2 mkdir -p /root/SecondLife/game_xtuner/model/Shanghai_AI_Laboratory
3 cp -r /root/share/model_repos/internlm2-chat-7b
  /root/SecondLife/game_xtuner/model/Shanghai_AI_Laboratory
4
5 # 创建用于存放配置的文件夹config并进入
6 mkdir /root/SecondLife/game_xtuner/config && cd
  /root/SecondLife/game_xtuner/config
7
8 # 列出所有内置配置
9 xtuner list-cfg
10
11 # 拷贝一个配置文件到当前目录: xtuner copy-cfg ${CONFIG_NAME} ${SAVE_PATH}
12 xtuner copy-cfg internlm2_chat_7b_qlora_oasst1_e3 .

```

```
(game_xtuner) root@intern-studio: ~/SecondLife/game_xtuner/config# xtuner copy-cfg internlm_chat_7b_qlora_oasst1_e3 .
[2024-01-29 21:59:39,609] [INFO] [real_accelerator.py:191:get_accelerator] Setting ds_accelerator to cuda (auto detect)
[2024-01-29 22:00:18,667] [INFO] [real_accelerator.py:191:get_accelerator] Setting ds_accelerator to cuda (auto detect)
Copy to ./internlm_chat_7b_qlora_oasst1_e3_copy.py
```

修改拷贝后的文件internlm\_chat\_7b\_qlora\_oasst1\_e3\_copy.py，修改下述位置：

```
1  # PART 1 中
2  # 预训练模型存放的位置
3  pretrained_model_name_or_path =
4      '/root/SecondLife/game_xtuner/model/Shanghai_AI_Laboratory/internlm2-chat-7b'
5
6  # 微调数据存放的位置
7  data_path = '/root/SecondLife/data/dataset.json'
8
9  # 训练中最大的文本长度
10 max_length = 1024
11
12 # 每一批训练样本的大小
13 batch_size = 2
14
15 # 最大训练轮数
16 max_epochs = 3
17
18 # 验证的频率
19 evaluation_freq = 500
20
21 # 用于评估输出内容的问题（用于评估的问题尽量与数据集的question保持一致）
22 SYSTEM = '你是一个对名人生平研究颇深的学者。你总是可以根据提问给出准确且详细的答案'
23 evaluation_inputs = [
24     '请介绍一下李白的生平', '请介绍一下杜甫的生平', '请介绍一下你自己'
25 ]
26
27 # PART 3 中
28 dataset=dict(type=load_dataset, path='json', data_files=dict(train=data_path))
29 dataset_map_fn=None
```

PART1 部分

```
root@intern-studio: ~/Sec X internlm_chat_7b_qlora_oa X +
17 from xtuner.model import SupervisedFinetune
18 from xtuner.utils import PROMPT_TEMPLATE
19
20 #####
21 # PART 1 Settings #
22 #####
23 # Model
24 pretrained_model_name_or_path = '/root/SecondLife/game_xtuner/model/Shanghai_AI_Laboratory/internlm2-chat-7b'
25
26 # Data
27 data_path = '/root/SecondLife/data/dataset.json'
28 prompt_template = PROMPT_TEMPLATE.internlm_chat
29 max_length = 1024
30 pack_to_max_length = True
31
32 # Scheduler & Optimizer
33 batch_size = 2 # per_device
34 accumulative_counts = 16
35 dataloader_num_workers = 0
36 max_epochs = 3
37 optim_type = PagedAdamW32bit
38 lr = 2e-4
39 betas = (0.9, 0.999)
40 weight_decay = 0
41 max_norm = 1 # grad clip
42
43 # Evaluate the generation performance during the training
44 evaluation_freq = 500
45 SYSTEM = '你是一个对名人生平研究颇深的学者。你总是可以根据提问给出准确且详细的答案'
46 evaluation_inputs = [
47     '请介绍一下李白的生平', '请介绍一下杜甫的生平', '请介绍一下你自己'
48 ]
49
```

## PART3 部分

```
root@intern-studio: ~/per X internlm_chat_7b_qlora_oa X +
79     lora_dropout=0.1,
80     bias='none',
81     task_type='CAUSAL_LM'))
82
83 #####
84 # PART 3 Dataset & Dataloader #
85 #####
86 train_dataset = dict(
87     type=process_hf_dataset,
88     dataset=dict(type=load_dataset, path='json', data_files=dict(train=data_path)),
89     tokenizer=tokenizer,
90     max_length=max_length,
91     dataset_map_fn=None,
92     template_map_fn=dict(
93         type=template_map_fn_factory, template=prompt_template),
94     remove_unused_columns=True,
95     shuffle_before_pack=True,
96     pack_to_max_length=pack_to_max_length)
97
98 train_dataloader = dict(
99     batch_size=batch_size,
100     num_workers=dataloader_num_workers,
101     dataset=train_dataset,
102     sampler=dict(type=DefaultSampler, shuffle=True),
103     collate_fn=dict(type=default_collate_fn))
104
```

## 4 微调启动

```
1 apt update -y
2 apt install tmux -y
3 # 新建对话
4 tmux new -s finetune
5 # Ctrl+B再按D返回bash
6 tmux attach -s finetune
7 xtuner train
  /root/SecondLife/game_xtuner/config/internlm2_chat_7b_qlora_oasst1_e3_copy.py --
  deepspeed deepspeed_zero2
```

## 5 参数转换

训练后的pth格式参数转Hugging Face格式

```
1 # 创建用于存放Hugging Face格式参数的hf文件夹
2 mkdir /root/SecondLife/game_xtuner/config/work_dirs/hf
3
4 export MKL_SERVICE_FORCE_INTEL=1
5
6 # 配置文件存放的位置
7 export
  CONFIG_NAME_OR_PATH=/root/SecondLife/game_xtuner/config/internlm2_chat_7b_qlora_o
  asst1_e3_copy.py
8
9 # 模型训练后得到的pth格式参数存放的位置
10 export
  PTH=/root/SecondLife/game_xtuner/config/work_dirs/internlm2_chat_7b_qlora_oasst1_
  e3_copy/iter_8104.pth
11
12 # pth文件转换为Hugging Face格式后参数存放的位置
13 export SAVE_PATH=/root/SecondLife/game_xtuner/config/work_dirs/hf
14
15 # 执行参数转换
16 xtuner convert pth_to_hf $CONFIG_NAME_OR_PATH $PTH $SAVE_PATH
```

```

01/30 14:43:46 - mmengine - INFO - replace internlm2 rope
01/30 14:43:46 - mmengine - INFO - replace internlm2 rope
01/30 14:43:46 - mmengine - INFO - replace internlm2 rope
01/30 14:43:46 - mmengine - INFO - replace internlm2 rope
01/30 14:43:46 - mmengine - INFO - replace internlm2 rope
01/30 14:43:46 - mmengine - INFO - replace internlm2 rope
01/30 14:43:46 - mmengine - INFO - replace internlm2 rope
01/30 14:43:46 - mmengine - INFO - replace internlm2 rope
01/30 14:43:46 - mmengine - INFO - replace internlm2 rope
01/30 14:43:46 - mmengine - INFO - replace internlm2 rope
01/30 14:43:46 - mmengine - INFO - replace internlm2 rope
Processing zero checkpoint '/root/SecondLife/game_xtuner/config/work_dirs/internlm2_chat_7b_qlora_oasst1_e3_copy/iter_8104.pth'
Detected checkpoint of type zero stage 2, world size: 1
Parsing checkpoint created by deepspeed==0.13.1
Reconstructed fp32 state dict with 322 params 157179904 elements
Load PTH model from /root/SecondLife/game_xtuner/config/work_dirs/internlm2_chat_7b_qlora_oasst1_e3_copy/iter_8104.pth
Convert LLM to float16
Saving adapter to /root/SecondLife/game_xtuner/config/work_dirs/hf
All done!
(game_xtuner) root@intern-studio: ~/SecondLife/game_xtuner/config#

```

## 6 参数合并

```

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/SecondLife/game_xtuner/model/Shanghai_AI_Laboratory/internlm2-chat-7b
6
7  # Hugging Face格式参数存放的位置
8  export NAME_OR_PATH_TO_ADAPTER=/root/SecondLife/game_xtuner/config/work_dirs/hf
9
10 # 最终Merge后的参数存放的位置
11 mkdir /root/SecondLife/game_xtuner/config/work_dirs/hf_merge
12 export SAVE_PATH=/root/SecondLife/game_xtuner/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

```

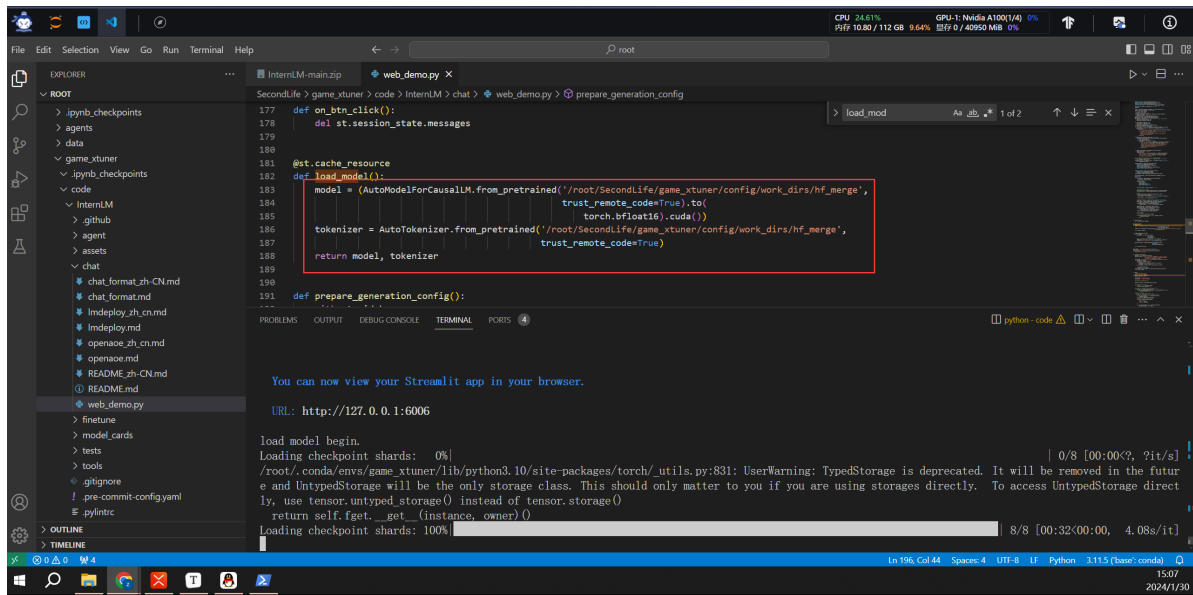
```

(game_xtuner) root@intern-studio: ~/SecondLife/game_xtuner/config# export NAME_OR_PATH_TO_ADAPTER=/root/SecondLife/game_xtuner/config/work_dirs/hf
(game_xtuner) root@intern-studio: ~/SecondLife/game_xtuner/config# export SAVE_PATH=/root/SecondLife/game_xtuner/config/work_dirs/hf_merge
(game_xtuner) root@intern-studio: ~/SecondLife/game_xtuner/config# xtuner convert merge \
> $NAME_OR_PATH_TO_LLM \
> $NAME_OR_PATH_TO_ADAPTER \
> $SAVE_PATH \
> --max-shard-size 2GB
[2024-01-30 14:50:17,385] [INFO] [real_accelerator.py:191:get_accelerator] Setting ds_accelerator to cuda (auto detect)
[2024-01-30 14:50:55,593] [INFO] [real_accelerator.py:191:get_accelerator] Setting ds_accelerator to cuda (auto detect)
Loading checkpoint shards: 0% | 0/8 [00:00<?, ?it/s]
/root/.conda/envs/game_xtuner/lib/python3.10/site-packages/torch/_utils.py:831: UserWarning: TypedStorage is deprecated. It will be removed in the future and UntypedStorage will be the only storage class. This should only matter to you if you are using storages directly. To access UntypedStorage directly, use tensor.untyped_storage() instead of tensor.storage()
  return self.fget._get_(instance, owner)()
Loading checkpoint shards: 100% | 8/8 [01:26<00:00, 10.81s/it]
Saving to /root/SecondLife/game_xtuner/config/work_dirs/hf_merge...
All done!
(game_xtuner) root@intern-studio: ~/SecondLife/game_xtuner/config#

```

## 7 Web部署

```
1 # 安装依赖
2 pip install streamlit==1.24.0
3
4 # 创建code文件夹用于存放InternLM项目代码
5 mkdir /root/SecondLife/game_xtuner/code && cd /root/SecondLife/game_xtuner/code
6 git clone https://github.com/InternLM/InternLM.git
7
8 # 修改/root/SecondLife/game_xtuner/code/InternLM/chat/web_demo.py中的模型路径
9 修改为"/root/SecondLife/game_xtuner/config/work_dirs/hf_merge"
10
11 # 运行脚本
12 cd /root/SecondLife/game_xtuner/code/InternLM
13 streamlit run /root/SecondLife/game_xtuner/code/InternLM/chat/web_demo.py --
    server.address 127.0.0.1 --server.port 6006
14
15 # powershell
16 ssh -CNg -L 6006:127.0.0.1:6006 root@ssh.intern-ai.org.cn -p [开发机端口号]
```



## 8 最终效果

