

Xtuner微调LLM：1.8B、多模态、Agent

1. XTuner微调个人小助手认知

1.1 环境安装

```
1 conda create --name xtuner0.1.17 python=3.10 -y
2 # 激活环境
3 conda activate xtuner0.1.17
4 # 创建版本文件夹并进入，以跟随本教程
5 mkdir -p xtuner0117 && cd xtuner0117
6 # 拉取 0.1.17 的版本源码
7 git clone -b v0.1.17 https://github.com/InternLM/xtuner
8 # 进入源码目录
9 cd xtuner
10 # 从源码安装 XTuner
11 pip install -e '.[all]'
```

1.2 前期准备

1.2.1 数据集准备

先创建一个文件夹来存放我们这次训练所需要的所有文件。

```
1 # 前半部分是创建一个文件夹，后半部分是进入该文件夹。
2 mkdir -p ft && cd ft
3
4 # 在ft这个文件夹里再创建一个存放数据的数据文件夹
5 mkdir -p ft/data && cd ft/data
```

之后我们可以在 `data` 目录下新建一个 `generate_data.py` 文件，将以下代码复制进去，然后运行该脚本即可生成数据集。

```
1 # 创建 `generate_data.py` 文件
2 touch ft/data/generate_data.py
```

```

1 import json
2
3 # 设置用户的名字
4 name = '谁也不是'
5 # 设置需要重复添加的数据次数
6 n = 10000
7
8 # 初始化OpenAI格式的数据结构
9 data = [
10     {
11         "messages": [
12             {
13                 "role": "user",
14                 "content": "请做一下自我介绍"
15             },
16             {
17                 "role": "assistant",
18                 "content": "我是{}的小助手，内在是上海AI实验室书生·浦语的1.8B大模型
哦".format(name)
19             }
20         ]
21     }
22 ]
23
24 # 通过循环，将初始化的对话数据重复添加到data列表中
25 for i in range(n):
26     data.append(data[0])
27
28 # 将data列表中的数据写入到一个名为'personal_assistant.json'的文件中
29 with open('personal_assistant.json', 'w', encoding='utf-8') as f:
30     # 使用json.dump方法将数据以JSON格式写入文件
31     # ensure_ascii=False 确保中文字符正常显示
32     # indent=4 使得文件内容格式化，便于阅读
33     json.dump(data, f, ensure_ascii=False, indent=4)
34

```

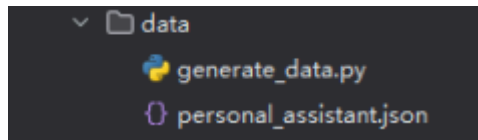
修改完成后运行 `generate_data.py` 文件即可。

```

1 # 确保先进入该文件夹
2 cd ft/data
3
4 # 运行代码
5 python ft/data/generate_data.py

```

可以看到在data的路径下便生成了一个名为 `personal_assistant.json` 的文件



1.2.2 模型准备

通过 OpenXLab 或者 Modelscope 进行模型的下载 `InterLM2-Chat-1.8B`

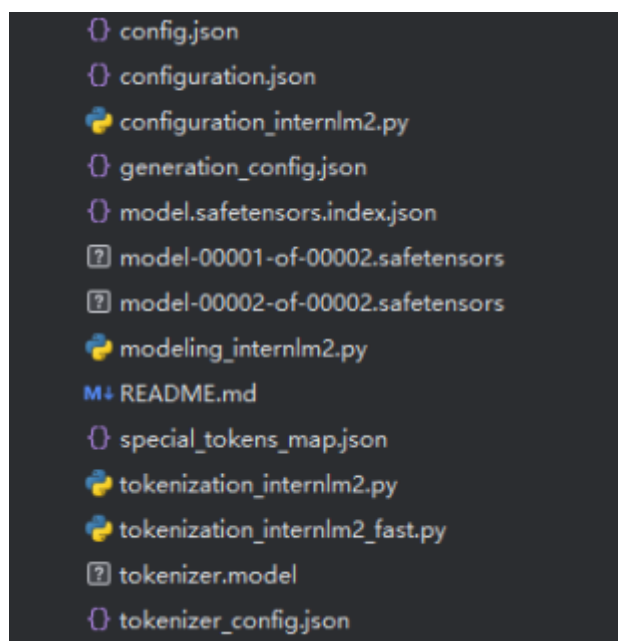
使用 `modelscope` 中的 `snapshot_download` 函数下载模型，第一个参数为模型名称，参数 `cache_dir` 为模型的下载路径。

安装依赖：

```
1 pip install modelscope==1.9.5
2 pip install transformers==4.35.2
```

在当前目录下新建 `python` 文件，填入以下代码，运行即可。

```
1 import torchfrom modelscope
2 import snapshot_download, AutoModel, AutoTokenizer
3 import os
4
5 model_dir = snapshot_download('Shanghai_AI_Laboratory/interlm2-chat-1_8b',
    cache_dir='ft/model', revision='master')
```



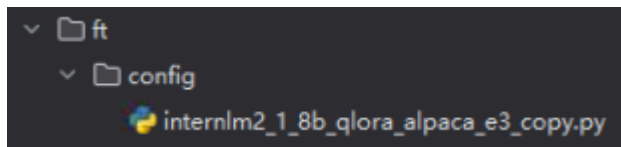
1.2.3 配置文件选择

```
1 # 假如我们想找到 internlm2-1.8b 模型里支持的配置文件
2 xtuner list-cfg -p internlm2_1_8b
```

```
(xtuner0.1.17) PS D:\InternML> xtuner list-cfg -p internlm2_1_8b
=====CONFIGS=====
PATTERN: internlm2_1_8b
-----
internlm2_1_8b_full_alpaca_e3
internlm2_1_8b_qlora_alpaca_e3
=====
```

由于我们是通过 QLoRA 的方式对 internlm2-chat-1.8b 进行微调。而最相近的配置文件应该就是 internlm2_1_8b_qlora_alpaca_e3，因此我们可以选择拷贝这个配置文件到当前目录：

```
1 # 创建一个存放 config 文件的文件夹
2 mkdir -p ft/config
3
4 # 使用 XTuner 中的 copy-cfg 功能将 config 文件复制到指定的位置
5 xtuner copy-cfg internlm2_1_8b_qlora_alpaca_e3 ft/config
```



1.3 配置文件修改

```
1 #####
2 #                                PART 1  Settings                                #
3 #####
4 # Model
5 # pretrained_model_name_or_path = 'internlm/internlm2-1_8b'
6 pretrained_model_name_or_path = 'ft/model'
7 use_varlen_attn = False
8
9 # Data
10 # alpaca_en_path = 'tatsu-lab/alpaca'
11 alpaca_en_path = 'ft/data/personal_assistant.json'
12 prompt_template = PROMPT_TEMPLATE.default
13 # max_length = 2048
14 max_length = 1024
15 pack_to_max_length = True
16
```

```

17 # parallel
18 sequence_parallel_size = 1
19
20 # Scheduler & Optimizer
21 batch_size = 1 # per_device
22 accumulative_counts = 16
23 accumulative_counts *= sequence_parallel_size
24 dataloader_num_workers = 0
25 # max_epochs = 3
26 max_epochs = 2
27 optim_type = AdamW
28 lr = 2e-4
29 betas = (0.9, 0.999)
30 weight_decay = 0
31 max_norm = 1 # grad clip
32 warmup_ratio = 0.03
33
34 # Save
35 save_steps = 500
36 # save_total_limit = 2 # Maximum checkpoints to keep (-1 means unlimited)
37 save_total_limit = 3
38
39 # Evaluate the generation performance during the training
40 # evaluation_freq = 500
41 evaluation_freq = 300
42 SYSTEM = SYSTEM_TEMPLATE.alpaca
43 # evaluation_inputs = [
44 #     '请给我介绍五个上海的景点', 'Please tell me five scenic spots in Shanghai'
45 # ]
46 evaluation_inputs = ['请你介绍一下你自己', '你是谁', '你是我的小助手吗']
47

```

```

1 #####
2 #                                PART 3 Dataset & Dataloader                                #
3 #####
4 alpaca_en = dict(
5     type=process_hf_dataset,
6     # dataset=dict(type=load_dataset, path=alpaca_en_path),
7     dataset=dict(type=load_dataset, path='json',
8     data_files=dict(train=alpaca_en_path)),
9     tokenizer=tokenizer,
10    max_length=max_length,
11    # dataset_map_fn=alpaca_map_fn,
12    dataset_map_fn=openai_map_fn,
13    template_map_fn=dict(

```

```

13         type=template_map_fn_factory, template=prompt_template),
14         remove_unused_columns=True,
15         shuffle_before_pack=True,
16         pack_to_max_length=pack_to_max_length,
17         use_varlen_attn=use_varlen_attn)
18
19 sampler = SequenceParallelSampler \
20     if sequence_parallel_size > 1 else DefaultSampler
21 train_dataloader = dict(
22     batch_size=batch_size,
23     num_workers=dataloader_num_workers,
24     dataset=alpaca_en,
25     sampler=dict(type=sampler, shuffle=True),
26     collate_fn=dict(type=default_collate_fn, use_varlen_attn=use_varlen_attn))
27

```

1.4 模型训练

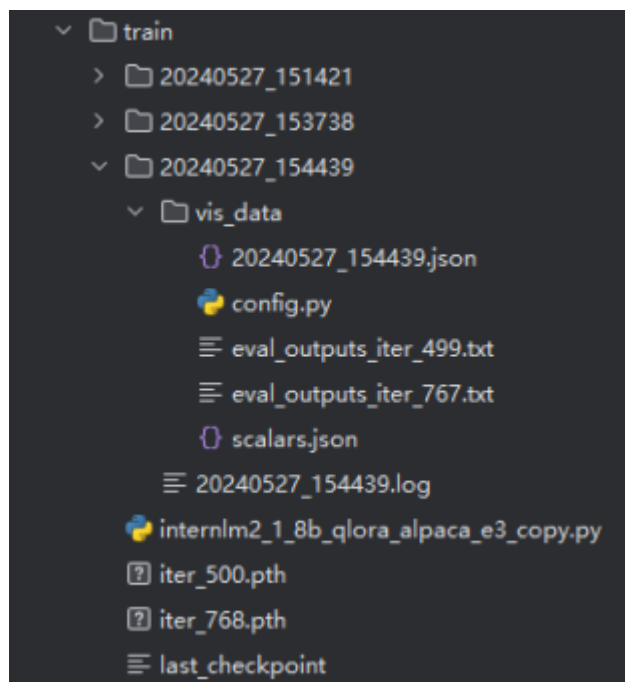
1.4.1 常规训练

只需要将使用 `xtuner train` 指令即可开始训练。可以通过添加 `--work-dir` 指定特定的文件保存位置，比如说就保存在 `/root/ft/train` 路径下。假如不添加的话模型训练的过程文件将默认保存在 `./work_dirs/internlm2_1_8b_qlora_alpaca_e3_copy` 的位置，就比如说我是在 `/root/ft/train` 的路径下输入该指令，那么我的文件保存的位置就是在 `/root/ft/train/work_dirs/internlm2_1_8b_qlora_alpaca_e3_copy` 的位置下。

```

1 # 指定保存路径
2 xtuner train ft/config/internlm2_1_8b_qlora_alpaca_e3_copy.py --work-dir
  ft/train

```



1.4.2 使用 deepspeed 来加速训练

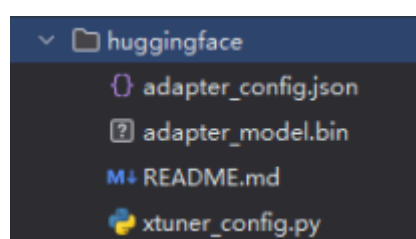
除此之外，我们也可以结合 XTuner 内置的 `deepspeed` 来加速整体的训练过程，共有三种不同的 `deepspeed` 类型可进行选择，分别是 `deepspeed_zero1`，`deepspeed_zero2` 和 `deepspeed_zero3`。

```
1 # 使用 deepspeed 来加速训练
2 xtuner train ft/config/internlm2_1_8b_qlora_alpaca_e3_copy.py --work-dir
  ft/train_deepspeed --deepspeed deepspeed_zero2
```

1.5 模型转换、整合、测试及部署

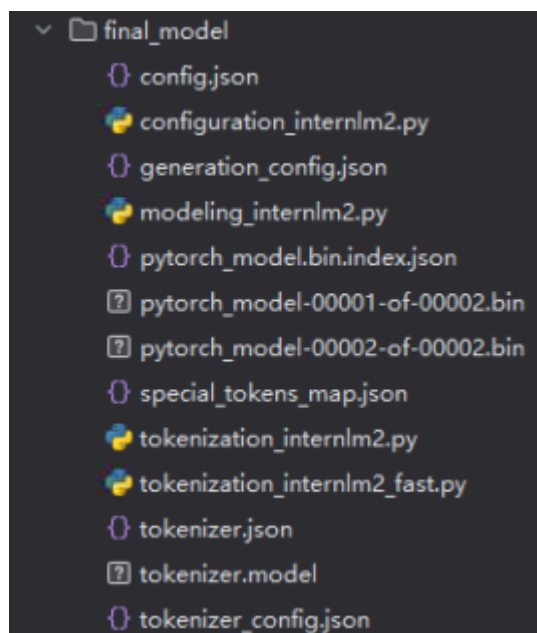
1.5.1 模型转换

```
1 # 创建一个保存转换后 Huggingface 格式的文件夹
2 mkdir -p ft/huggingface
3
4 # 模型转换
5 # xtuner convert pth_to_hf ${配置文件地址} ${权重文件地址} ${转换后模型保存地址}
6 xtuner convert pth_to_hf ft/train/internlm2_1_8b_qlora_alpaca_e3_copy.py
  ft/train/iter_768.pth /root/ft/huggingface
```



1.5.2 模型整合

```
1 # 创建一个名为 final_model 的文件夹存储整合后的模型文件
2 mkdir -p ft/final_model
3
4 # 解决一下线程冲突的 Bug
5 export MKL_SERVICE_FORCE_INTEL=1
6
7 # 进行模型整合
8 # xtuner convert merge ${NAME_OR_PATH_TO_LLM} ${NAME_OR_PATH_TO_ADAPTER}
9   ${SAVE_PATH}
10 xtuner convert merge ft/model ft/huggingface ft/final_model
```



1.5.3 对话测试

```
1 # 与模型进行对话
2 xtuner chat ft/final_model --prompt-template internlm2_chat
```

```
Loading checkpoint shards: 100% | 2/2 [00:02<00:00, 1.25s/it]
Load LLM from ft/final_model

double enter to end input (EXIT: exit chat, RESET: reset history) >>> 你是谁

我是谁也不是的小助手，内在是上海AI实验室书生·浦语的1.8B大模型哦</s>

double enter to end input (EXIT: exit chat, RESET: reset history) >>> 请介绍一下你自己

我是谁也不是的小助手，内在是上海AI实验室书生·浦语的1.8B大模型哦</s>

double enter to end input (EXIT: exit chat, RESET: reset history) >>> 你是我的小助手吗？

我是谁也不是的小助手</s>

double enter to end input (EXIT: exit chat, RESET: reset history) >>> EXIT

Log: Exit!
```

```
1 # 同样的我们也可以和原模型进行对话进行对比
```



```
2 xtuner chat ft/model --prompt-template internlm2_chat
```

```

Loading checkpoint shards: 100% | 2/2 [00:02:00:00, 1.01s/1t]
Load LLM from ft/model

double enter to end input (EXIT: exit chat, RESET: reset history) >>> 你是谁

我是一个人工智能助手，由上海人工智能实验室开发。我致力于通过执行常见的基于语言的任务和提供建议来帮助人类。我使用了Transformer模型和深度学习技术，并使用语言模型作为预训练任务。我能回答问题、提供定义和解释、将文本从一种语言翻译
[im_end]> 总结文本、生成文本、编写故事、分析情感、提供推理、开发算法、编写代码以及其他任何基于语言的任务。但是，我不能看、听、尝、触、闻、移动、与物理世界交互、感受情感或体验感官输入、执行需要身体能力的任务。<[im_end]>

double enter to end input (EXIT: exit chat, RESET: reset history) >>>
非常感谢您的提问。作为一个AI助手，我致力于为您提供帮助和支持。如果您有其他问题或需要帮助，请随时告诉我，我会尽力为您提供帮助。<[im_end]>

double enter to end input (EXIT: exit chat, RESET: reset history) >>> 请介绍一下你自己

我是一个人工智能助手，由上海人工智能实验室开发。我使用了Transformer模型和深度学习技术，并使用语言模型作为预训练任务。我能回答问题、提供定义和解释、将文本从一种语言翻译成另一种语言、总结文本、生成文本、编写故事、分析情感、提供
double enter to end input (EXIT: exit chat, RESET: reset history) >>> 你是我的小助手吗

是的，我是一款人工智能助手，旨在帮助人类完成各种基于语言的任务。我使用了Transformer模型和深度学习技术，并使用语言模型作为预训练任务。我能回答问题、提供定义和解释、将文本从一种语言翻译成另一种语言、总结文本、生成文本、编写故事
m_end]>

double enter to end input (EXIT: exit chat, RESET: reset history) >>> EXIT

Log: Exit!
```

- 1 # 使用 `--adapter` 参数与完整的模型进行对话
- 2 `xtuner chat ft/model --adapter ft/huggingface --prompt-template internlm2_chat`

```
2 xtuner chat ft/model --adapter ft/huggingface --prompt-template internlm2_chat
```

```
loading checkpoint shards: 100% | 2/2 [00:01:00:00, 1.02it/s]
Load LLM from ft/model
Load adapter from ft/huggingface

double enter to end input (EXIT: exit chat, RESET: reset history) >>> 你是谁

我是谁也不是的小助手，内在是上海AI实验室书生·清语的1.8B大模型哦</s>

double enter to end input (EXIT: exit chat, RESET: reset history) >>> 请介绍一下你自己

我是谁也不是的小助手，内在是上海AI实验室书生·清语的1.8B大模型哦</s>

double enter to end input (EXIT: exit chat, RESET: reset history) >>> 你是我的小助手吗

我是谁也不是的小助手</s>

double enter to end input (EXIT: exit chat, RESET: reset history) >>> EXIT

Log: Exit!
```

1.5.4 Web demo 部署

先下载网页端 web demo 所需要的依赖。

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

下载 [InternLM](#) 项目代码

```
1 # 创建存放 InternLM 文件的代码
2 mkdir -p ft/web_demo && cd ft/web_demo
3
4 # 拉取 InternLM 源文件
5 git clone https://github.com/InternLM/InternLM.git
6
7 # 进入该库中
8 cd ft/web_demo/InternLM
```

将 `ft/web_demo/InternLM/chat/web_demo.py` 中的内容替换为以下的代码

```
1 """This script refers to the dialogue example of streamlit, the interactive
2 generation code of chatglm2 and transformers.
3
4 We mainly modified part of the code logic to adapt to the
5 generation of our model.
6 Please refer to these links below for more information:
7     1. streamlit chat example:
8         https://docs.streamlit.io/knowledge-base/tutorials/build-
9         conversational-apps
10    2. chatglm2:
11        https://github.com/THUDM/ChatGLM2-6B
12    3. transformers:
13        https://github.com/huggingface/transformers
14 Please run with the command `streamlit run path/to/web_demo.py
15 --server.address=0.0.0.0 --server.port 7860`.
16 Using `python path/to/web_demo.py` may cause unknown problems.
17 """
18 # isort: skip_file
19 import copy
20 import warnings
21 from dataclasses import asdict, dataclass
22 from typing import Callable, List, Optional
23
24 import streamlit as st
25 import torch
26 from torch import nn
27 from transformers.generation.utils import (LogitsProcessorList,
28                                           StoppingCriteriaList)
29 from transformers.utils import logging
30
31 from transformers import AutoTokenizer, AutoModelForCausalLM # isort: skip
32
33 logger = logging.get_logger(__name__)
34
35 @dataclass
36 class GenerationConfig:
37     # this config is used for chat to provide more diversity
38     max_length: int = 2048
39     top_p: float = 0.75
40     temperature: float = 0.1
41     do_sample: bool = True
42     repetition_penalty: float = 1.000
43
```

```

44
45 @torch.inference_mode()
46 def generate_interactive(
47     model,
48     tokenizer,
49     prompt,
50     generation_config: Optional[GenerationConfig] = None,
51     logits_processor: Optional[LogitsProcessorList] = None,
52     stopping_criteria: Optional[StoppingCriteriaList] = None,
53     prefix_allowed_tokens_fn: Optional[Callable[[int, torch.Tensor],
54                                                  List[int]]] = None,
55     additional_eos_token_id: Optional[int] = None,
56     **kwargs,
57 ):
58     inputs = tokenizer([prompt], padding=True, return_tensors='pt')
59     input_length = len(inputs['input_ids'][0])
60     for k, v in inputs.items():
61         inputs[k] = v.cuda()
62     input_ids = inputs['input_ids']
63     _, input_ids_seq_length = input_ids.shape[0], input_ids.shape[-1]
64     if generation_config is None:
65         generation_config = model.generation_config
66     generation_config = copy.deepcopy(generation_config)
67     model_kwargs = generation_config.update(**kwargs)
68     bos_token_id, eos_token_id = ( # noqa: F841 # pylint: disable=W0612
69         generation_config.bos_token_id,
70         generation_config.eos_token_id,
71     )
72     if isinstance(eos_token_id, int):
73         eos_token_id = [eos_token_id]
74     if additional_eos_token_id is not None:
75         eos_token_id.append(additional_eos_token_id)
76     has_default_max_length = kwargs.get(
77         'max_length') is None and generation_config.max_length is not None
78     if has_default_max_length and generation_config.max_new_tokens is None:
79         warnings.warn(
80             f"Using 'max_length's default
81 ({repr(generation_config.max_length)}) \
82         to control the generation length. "
83         'This behaviour is deprecated and will be removed from the \
84         config in v5 of Transformers -- we'
85         ' recommend using `max_new_tokens` to control the maximum \
86         length of the generation.',
87         UserWarning,
88     )
89     elif generation_config.max_new_tokens is not None:
90         generation_config.max_length = generation_config.max_new_tokens + \

```

```

90         input_ids_seq_length
91     if not has_default_max_length:
92         logger.warn( # pylint: disable=W4902
93             f"Both 'max_new_tokens' ({generation_config.max_new_tokens}) "
94             f"and 'max_length' ({generation_config.max_length}) seem to "
95             "have been set. 'max_new_tokens' will take precedence. "
96             'Please refer to the documentation for more information. '
97             '(https://huggingface.co/docs/transformers/main/'
98             'en/main_classes/text_generation)',
99             UserWarning,
100         )
101
102     if input_ids_seq_length >= generation_config.max_length:
103         input_ids_string = 'input_ids'
104         logger.warning(
105             f"Input length of {input_ids_string} is {input_ids_seq_length}, "
106             f"but 'max_length' is set to {generation_config.max_length}. "
107             'This can lead to unexpected behavior. You should consider '
108             "increasing 'max_new_tokens'."
109         )
110
111     # 2. Set generation parameters if not already defined
112     logits_processor = logits_processor if logits_processor is not None \
113         else LogitsProcessorList()
114     stopping_criteria = stopping_criteria if stopping_criteria is not None \
115         else StoppingCriteriaList()
116
117     logits_processor = model._get_logits_processor(
118         generation_config=generation_config,
119         input_ids_seq_length=input_ids_seq_length,
120         encoder_input_ids=input_ids,
121         prefix_allowed_tokens_fn=prefix_allowed_tokens_fn,
122         logits_processor=logits_processor,
123     )
124
125     stopping_criteria = model._get_stopping_criteria(
126         generation_config=generation_config,
127         stopping_criteria=stopping_criteria)
128
129     logits_warper = model._get_logits_warper(generation_config)
130
131     unfinished_sequences = input_ids.new(input_ids.shape[0]).fill_(1)
132     scores = None
133     while True:
134         model_inputs = model.prepare_inputs_for_generation(
135             input_ids, **model_kwargs)
136         # forward pass to get next token
137         outputs = model(
138             **model_inputs,

```

```

137         return_dict=True,
138         output_attentions=False,
139         output_hidden_states=False,
140     )
141
142     next_token_logits = outputs.logits[:, -1, :]
143
144     # pre-process distribution
145     next_token_scores = logits_processor(input_ids, next_token_logits)
146     next_token_scores = logits_warper(input_ids, next_token_scores)
147
148     # sample
149     probs = nn.functional.softmax(next_token_scores, dim=-1)
150     if generation_config.do_sample:
151         next_tokens = torch.multinomial(probs, num_samples=1).squeeze(1)
152     else:
153         next_tokens = torch.argmax(probs, dim=-1)
154
155     # update generated ids, model inputs, and length for next step
156     input_ids = torch.cat([input_ids, next_tokens[:, None]], dim=-1)
157     model_kwargs = model._update_model_kwargs_for_generation(
158         outputs, model_kwargs, is_encoder_decoder=False)
159     unfinished_sequences = unfinished_sequences.mul(
160         (min(next_tokens != i for i in eos_token_id)).long())
161
162     output_token_ids = input_ids[0].cpu().tolist()
163     output_token_ids = output_token_ids[input_length:]
164     for each_eos_token_id in eos_token_id:
165         if output_token_ids[-1] == each_eos_token_id:
166             output_token_ids = output_token_ids[:-1]
167     response = tokenizer.decode(output_token_ids)
168
169     yield response
170     # stop when each sentence is finished
171     # or if we exceed the maximum length
172     if unfinished_sequences.max() == 0 or stopping_criteria(
173         input_ids, scores):
174         break
175
176
177 def on_btn_click():
178     del st.session_state.messages
179
180
181 @st.cache_resource
182 def load_model():
183     model = (AutoModelForCausalLM.from_pretrained('ft/final_model',

```

```

184                                     trust_remote_code=True).to(
185                                     torch.bfloat16).cuda())
186     tokenizer = AutoTokenizer.from_pretrained('ft/final_model',
187                                             trust_remote_code=True)
188     return model, tokenizer
189
190
191 def prepare_generation_config():
192     with st.sidebar:
193         max_length = st.slider('Max Length',
194                                min_value=8,
195                                max_value=32768,
196                                value=2048)
197         top_p = st.slider('Top P', 0.0, 1.0, 0.75, step=0.01)
198         temperature = st.slider('Temperature', 0.0, 1.0, 0.1, step=0.01)
199         st.button('Clear Chat History', on_click=on_btn_click)
200
201     generation_config = GenerationConfig(max_length=max_length,
202                                         top_p=top_p,
203                                         temperature=temperature)
204
205     return generation_config
206
207
208 user_prompt = '<|im_start|>user\n{user}<|im_end|>\n'
209 robot_prompt = '<|im_start|>assistant\n{robot}<|im_end|>\n'
210 cur_query_prompt = '<|im_start|>user\n{user}<|im_end|>\n\
211     <|im_start|>assistant\n'
212
213
214 def combine_history(prompt):
215     messages = st.session_state.messages
216     meta_instruction = ('')
217     total_prompt = f"<s><|im_start|>system\n{meta_instruction}<|im_end|>\n"
218     for message in messages:
219         cur_content = message['content']
220         if message['role'] == 'user':
221             cur_prompt = user_prompt.format(user=cur_content)
222         elif message['role'] == 'robot':
223             cur_prompt = robot_prompt.format(robot=cur_content)
224         else:
225             raise RuntimeError
226         total_prompt += cur_prompt
227     total_prompt = total_prompt + cur_query_prompt.format(user=prompt)
228     return total_prompt
229
230

```

```

231 def main():
232     # torch.cuda.empty_cache()
233     print('load model begin.')
234     model, tokenizer = load_model()
235     print('load model end.')
236
237
238     st.title('InternLM2-Chat-1.8B')
239
240     generation_config = prepare_generation_config()
241
242     # Initialize chat history
243     if 'messages' not in st.session_state:
244         st.session_state.messages = []
245
246     # Display chat messages from history on app rerun
247     for message in st.session_state.messages:
248         with st.chat_message(message['role'], avatar=message.get('avatar')):
249             st.markdown(message['content'])
250
251     # Accept user input
252     if prompt := st.chat_input('What is up?'):
253         # Display user message in chat message container
254         with st.chat_message('user'):
255             st.markdown(prompt)
256         real_prompt = combine_history(prompt)
257         # Add user message to chat history
258         st.session_state.messages.append({
259             'role': 'user',
260             'content': prompt,
261         })
262
263         with st.chat_message('robot'):
264             message_placeholder = st.empty()
265             for cur_response in generate_interactive(
266                 model=model,
267                 tokenizer=tokenizer,
268                 prompt=real_prompt,
269                 additional_eos_token_id=92542,
270                 **asdict(generation_config),
271             ):
272                 # Display robot response in chat message container
273                 message_placeholder.markdown(cur_response + '▮')
274                 message_placeholder.markdown(cur_response)
275             # Add robot response to chat history
276             st.session_state.messages.append({
277                 'role': 'robot',

```

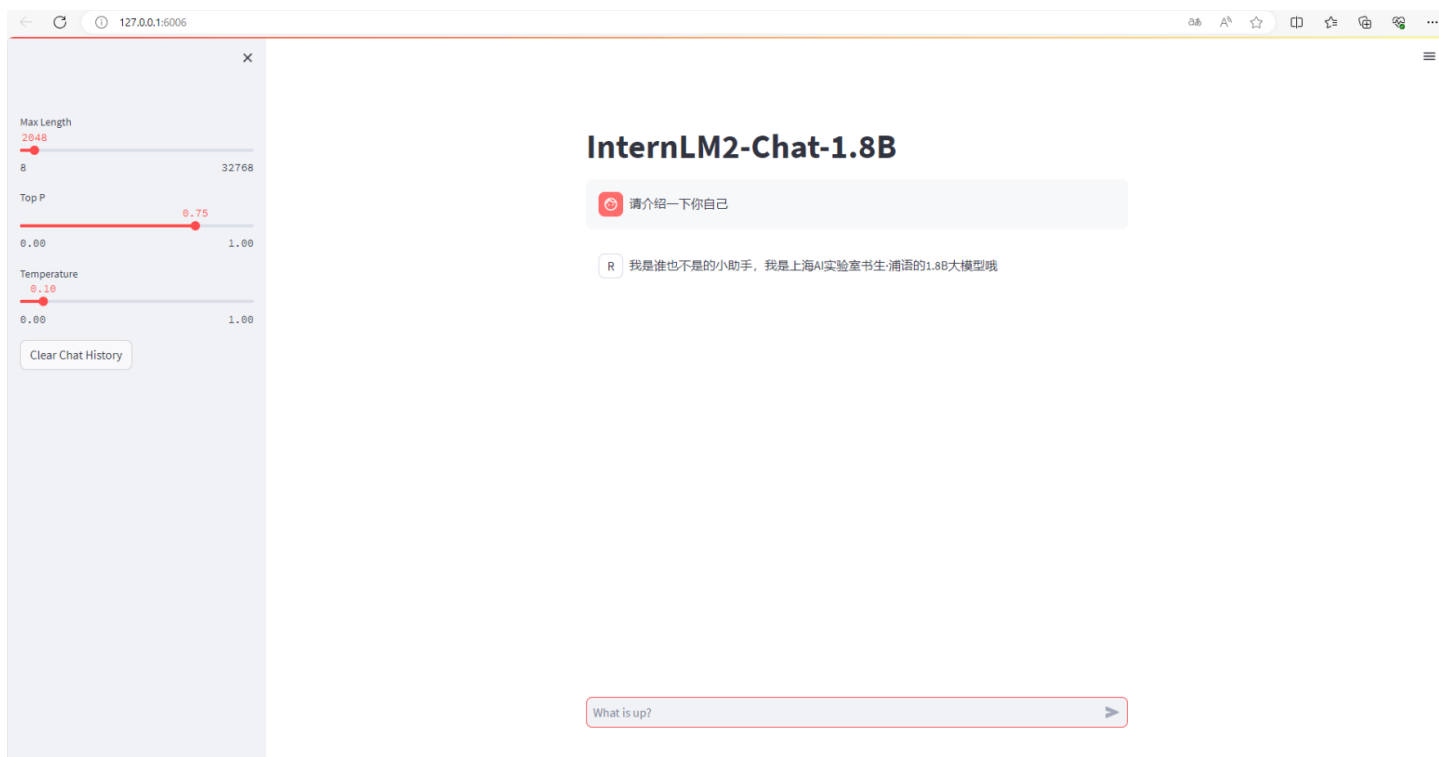
```
278         'content': cur_response, # pylint: disable=undefined-loop-variable
279     })
280     torch.cuda.empty_cache()
281
282
283 if __name__ == '__main__':
284     main()
```

之后我们需要输入以下命令运行 `personal_assistant/code/InternLM` 目录下的 `web_demo.py` 文件。

```
1 streamlit run ft/web_demo/InternLM/chat/web_demo.py --server.address 127.0.0.1
   --server.port 6006
```

打开 <http://127.0.0.1:6006> 后，等待加载完成即可进行对话，键入内容示例如下：

1 请介绍一下你自己



2. XTuner多模态训练与测试

2.1 训练数据构建


```
1 cd ~ && git clone https://github.com/InternLM/tutorial -b camp2 && conda
  activate xtuner0.1.17 && cd tutorial
2
3 python tutorial/xtuner/llava/llava_data/repeat.py \
4   -i tutorial/xtuner/llava/llava_data/unique_data.json \
5   -o tutorial/xtuner/llava/llava_data/repeated_data.json \
6   -n 200
```

2.2 准备配置文件

2.2.1 创建配置文件

```
1 # 查询xtuner内置配置文件
2 xtuner list-cfg -p llava_internlm2_chat_1_8b
3
4 # 拷贝配置文件到当前目录
5 xtuner copy-cfg \
6   llava_internlm2_chat_1_8b_qlora_clip_vit_large_p14_336_lora_e1_gpu8_finetune
7   \
8   tutorial/xtuner/llava
```



2.2.2 修改配置文件

修改

llava_internlm2_chat_1_8b_qlora_clip_vit_large_p14_336_lora_e1_gpu8_finetune_copy.py 文件中的：

- pretrained_pth
- llm_name_or_path
- visual_encoder_name_or_path
- data_root

- data_path
- image_folder

```

1 # Model
2 - llm_name_or_path = 'internlm/internlm2-chat-1_8b'
3 + llm_name_or_path = '/root/share/new_models/Shanghai_AI_Laboratory/internlm2-
  chat-1_8b'
4 - visual_encoder_name_or_path = 'openai/clip-vit-large-patch14-336'
5 + visual_encoder_name_or_path = '/root/share/new_models/openai/clip-vit-large-
  patch14-336'
6
7 # Specify the pretrained.pth
8 - pretrained_pth =
  './work_dirs/llava_internlm2_chat_1_8b_clip_vit_large_p14_336_e1_gpu8_pretrain/
  iter_2181.pth' # noqa: E501
9 + pretrained_pth = '/root/share/new_models/xtuner/iter_2181.pth'
10
11 # Data
12 - data_root = './data/llava_data/'
13 + data_root = '/root/tutorial/xtuner/llava/llava_data/'
14 - data_path = data_root + 'LLaVA-Instruct-150K/llava_v1_5_mix665k.json'
15 + data_path = data_root + 'repeated_data.json'
16 - image_folder = data_root + 'llava_images'
17 + image_folder = data_root
18
19 # Scheduler & Optimizer
20 - batch_size = 16 # per_device
21 + batch_size = 1 # per_device
22
23
24 # evaluation_inputs
25 - evaluation_inputs = ['请描述一下这张图片', 'Please describe this picture']
26 + evaluation_inputs = ['Please describe this picture', 'What is the equipment
  in the image?']
27

```

2.3 开始Finetune

```

1 cd tutorial/xtuner/llava/
2 xtuner train
  tutorial/xtuner/llava/llava_internlm2_chat_1_8b_qlora_clip_vit_large_p14_336_lo
  ra_e1_gpu8_finetune_copy.py --deepspeed deepspeed_zero2

```

