

小作业-LLM 的微调与 RAG

杨天行 2024310659

2025 年 1 月 5 日

1 实验任务

本次实验的任务围绕对大规模预训练语言模型（LLM）的微调以及检索增强生成（Retrieval-Augmented Generation, RAG）技术展开，主要包括以下几个方面：

知识库构建。收集与机器学习相关的领域知识，构建一个知识库。知识库应具有一定的专业性，包含一些无法通过预训练语言模型直接回答的问题，以增强其实际应用价值。

检索实现。实现从知识库中检索相关信息的功能。检索系统可使用开源库或训练好的模型，确保问题能够与知识库中的内容进行高效匹配。

问答数据构建。基于知识库内容，利用现有的 LLM API（如 ChatGPT LLM、ChatGLM）生成问答对，构建训练数据。问题应具有一定的专业性，确保 LLM 无法依靠其预训练知识直接回答。答案应结合检索到的 Top-K 知识点生成，答案中需包含引用标记（如 [1], [2]）。

模型微调。使用上一次实验的代码加载 ChatGLM3-6B 权重，对模型进行微调，模型结构与 GLM4-9B 相同。微调过程中需实现相关训练代码，可参考或使用开源工具库（如 Transformers）。设置合理的训练参数，并记录训练过程中的相关信息（如 Loss 曲线）。

性能评估。通过实验对比基于检索的问答生成（RAG）与直接使用预训练语言模型进行问答的性能差异。评估指标包括生成答案的准确性、专业性，以及引用标记的完整性。

2 微调数据构建

在本次实验中我们使用 baseline 代码提供的 ICNP'24 论文“Luori: Active Probing and Evaluation of Internet-wide IPv6 Fully Responsive Prefixes”作为知识库数据。我们的目标是微调 ChatGLM3-6B，并使其输出包含引用标记。为此我们需要先生成对应于知识库内容的问题。构建“问题-知识库信息-答案”的 pipeline，其中“问题 + 知识库数据”作为微调数据的输入，“答案”作为微调数据的输出。这样在 LLM 输出的答案中，能带上来自知识库信息的引用标记。

2.1 问题生成

首先调用 DeepSeek API，输入知识库与对应的 prompt，以 json 格式批量输出问题：

```
1 question_prefix = ""
2 Here is a text. Please generate 50 questions based on the information in the text. The
   questions should be answerable using the information from the text. Only output the
   questions.
3
4 EXAMPLE JSON OUTPUT:
5 {
6     "questions": [
```

```
7     "the first question",
8     "the second question",
9     "the third question"
10 ]
11 }
12
13 """
14
15 def generate_question(client, text):
16     response = client.chat.completions.create(
17         model="deepseek-chat",
18         messages=[
19             {"role": "system", "content": "You are a helpful assistant"},
20             {"role": "user", "content": question_prefix + text},
21         ],
22         response_format={
23             'type': 'json_object'
24         },
25         stream=False,
26         temperature=0.8,
27         max_tokens=8192
28     )
29
30     return response.choices[0].message.content
```

一共生成问题 391 条（输入相同，重复生成问题可以出发 input cache hit），生成的问题样例如下：

- What is the significance of IPv6 Fully Responsive Prefixes (FRPs) in network measurements?
- How does Luori propose to address the challenge of probing IPv6 FRPs?
- What is the core idea behind Luori's active probing method?
- How does Luori transform the active probing process into a dynamic search process?
- What are the evaluation results of Luori's Internet-wide probing?

2.2 RAG Prompt 修改

本次实验在提供的 baseline 代码使用的 LightRAG¹上修改 RAG 代码。直接运行 `run.py`，LightRAG 能够自动在文本上构建 Knowledge Graph 并利用 KG 构建 prompt，从而回答问题，但答案中并不包含引用标记。

我们以上面的问题之一 “What are the evaluation results of Luori's Internet-wide probing” 进行测试，结果如下：

¹Zirui Guo et al., LightRAG: Simple and Fast Retrieval-Augmented Generation, preprint

The evaluation results of Luori's Internet-wide probing demonstrate its
↪ effectiveness and efficiency in detecting and compiling a comprehensive list
↪ of Fully Responsive Prefixes (FRPs) in IPv6 networks. Here's a breakdown of
↪ the key findings:

****Coverage and Scope**:**

- ****Address Space**:** Luori successfully probed a vast address space, covering 1.3×10^{33} addresses. This is equivalent to the address space covered by 7.1×10^{13} /64 prefixes.
- ****Autonomous Systems (ASes)**:** The FRPlist compiled by Luori includes information on 795 ASes, indicating its wide coverage across different network operators and service providers.
- ****Routing Prefixes**:** The list covers 9182 routing prefixes, providing a detailed view of the FRPs present in the IPv6 network.

****Performance Metrics**:**

- ****Number of Fully Responsive Prefixes (FRPs)**:** Luori identified 516,739 FRPs during the six-month probing period, which is significantly more than the existing FRPlist.
- ****False Routing Responses (FRRPs)**:** The FRPlist contains 7,079 FRRPs, which are responses from network devices that do not correspond to actual network traffic. These FRRPs are crucial for understanding the impact of FRPs on network measurements.
- ****Fully Responsive Non-Routing Prefixes (FRNRPs)**:** The FRPlist also includes 508,660 FRNRPs, which are responsive prefixes that do not correspond to any routing prefix.

****Comparison with Existing Methods**:**

- ****Comparison with MAPD**:** Luori outperforms the brute-force probing method MAPD in terms of coverage and efficiency. While MAPD can only probe fixed-length FRPs, Luori can probe FRPs of any length, leading to a more comprehensive list of FRPs.
- ****Comparison with Hitlist's FRPlist**:** Luori's FRPlist covers significantly more address space than the Hitlist's FRPlist, despite having a smaller number of FRPs. This is due to the ability of Luori to probe larger FRPs and the efficient use of probing resources.

****Long-Term Probing**:**

- ****Continuous Growth**:** The number of FRRPs and FRNRPs in the FRPlist continues to grow over time, indicating the dynamic nature of FRP patterns in the IPv6 network.

```

- **Impact on IPv6 Measurement Campaigns**: The growing number of FRPs highlights
  ↳ the potential impact of FRPs on IPv6 measurement campaigns, emphasizing the
  ↳ need for accurate and comprehensive FRP detection.

**Overall, the evaluation results of Luori's Internet-wide probing demonstrate its
  ↳ ability to effectively detect and compile a comprehensive list of FRPs in IPv6
  ↳ networks. The resulting FRPlist can be used by researchers and network
  ↳ operators to mitigate the impact of FRPs on their measurements and to better
  ↳ understand the characteristics and distribution of FRPs in the IPv6
  ↳ Internet**.

```

我们可以在 lightrag 框架代码的lightrag/prompt.py中定义了进行 rag 的 prompt，原始的 prompt 为：

```

PROMPTS["rag_response"] = """---Role---

You are a helpful assistant responding to questions about data in the tables
↳ provided.

---Goal---

Generate a response of the target length and format that responds to the user's
↳ question, summarizing all information in the input data tables appropriate for
↳ the response length and format, and incorporating any relevant general
↳ knowledge.

If you don't know the answer, just say so. Do not make anything up.
Do not include information where the supporting evidence for it is not provided.

---Target response length and format---

{response_type}

---Data tables---

{context_data}

Add sections and commentary to the response as appropriate for the length and
↳ format. Style the response in markdown.
"""

```

我们对其进行修改，要求先列出所有需要被引用的段落，再带上引用记号回答问题，修改后的 prompt 如

下:

```
PROMPTS["rag_response"] = """---Role---
```

```
You are a helpful assistant responding to questions about data in the tables
↪ provided.
```

```
---Goal---
```

```
Generate a response of the target length and format that responds to the user's
↪ question, summarizing all information in the input data tables appropriate for
↪ the response length and format, and incorporating any relevant general
↪ knowledge.
```

```
If you don't know the answer, just say so. Do not make anything up.
```

```
Do not include information where the supporting evidence for it is not provided.
```

```
---Target response length and format---
```

```
{response_type}
```

```
---Data tables---
```

```
{context_data}
```

```
---Instructions---
```

```
1. Before generating the main response, list all the referenced content you will
↪ use in your response. Format each reference as follows:
```

- [1]: "Referenced content 1 from the data tables."
- [2]: "Referenced content 2 from the data tables."
- ...

```
Ensure the referenced content comes directly from the provided data tables.
```

```
2. After listing all references, generate the response. Use markdown formatting
↪ and include citation markers (e.g., [1], [2]) inline, corresponding to the
↪ previously listed references. The Reference entry should only use a citation
↪ marker at the beginning, and NO UNREALTED CITATION MARKERS SHOULD APPEAR!
↪ Response can only use the citation markers appeared in References.
```

```
---Output Example---
```

```
### References
```

- [1]: Example content from table A.
- [2]: Another example content from table B.

Response

Based on the data, the trends show an increase in revenue [1] and a decrease in
 ↳ expenses [2].

"""

对上面相同的问题，现在得到了如下，表明可以正确引用内容并带上引用记号：

References

- [1]: "We have implemented a prototype of Luori and deployed it from November
 ↳ 2023 to April 2024. In this section, we evaluate the performance of Luori
 ↳ during its deployment." from the Sources table.
- [2]: "Luori found 31.7K FRPs¹ in a single Internet-wide probing with 11M
 ↳ budget², covering 1.5×10^{30} address space³, which is $106 \times$ that of brute-force
 ↳ probing." from the Sources table.
- [3]: "As of now, Luori has been continuously probing Internet-wide for almost
 ↳ six months and has found 516k FRPs, which covers 1.3×10^{33} address space and
 ↳ 795 ASes, far exceeding the existing FRPlist, making it the largest publicly
 ↳ known FRPlist." from the Sources table.

Response

The evaluation results of Luori's Internet-wide probing demonstrate its
 ↳ effectiveness and scalability. In a single Internet-wide probing with a budget
 ↳ of 11 million, Luori found 31.7K Fully Responsive Prefixes (FRPs) and covered
 ↳ 1.5×10^{30} address space, which is 106 times more than brute-force probing [2].
 ↳ Over a period of almost six months of continuous probing, Luori identified
 ↳ 516k FRPs, covering 1.3×10^{33} address space and 795 Autonomous Systems (ASes),
 ↳ making it the largest publicly known FRPlist [3].

2.3 微调数据构造

我们可以写一个如下的脚本，使用上面修改后的 rag 进行回答，问题库是上面构造的 391 条问题。脚本如下，将 QA 对保存为 OpenAI SFT 格式的 json 数据：

```

1 import os
2 import json
3 import logging
4 from lightrag import LightRAG, QueryParam
5 from lightrag.llm import zhipu_complete, zhipu_embedding
6 from lightrag.utils import EmbeddingFunc
7 from tqdm import tqdm
8
9 # 设置工作目录和日志

```

```
10 WORKING_DIR = "./working_dir"
11 DATA_OUTPUT = "./data/fine_tuning_data.json"
12 QUESTION_FILE = "./data/questions.json"
13
14 logging.basicConfig(format="%(levelname)s:%(message)s", level=logging.INFO)
15
16 if not os.path.exists(WORKING_DIR):
17     os.mkdir(WORKING_DIR)
18
19 if not os.path.exists("./output"):
20     os.mkdir("./output")
21
22 # 初始化LightRAG
23 rag = LightRAG(
24     working_dir=WORKING_DIR,
25     llm_model_func=zhipu_complete,
26     llm_model_name="glm-4-flash",
27     llm_model_max_async=4,
28     llm_model_max_token_size=32768,
29     embedding_func=EmbeddingFunc(
30         embedding_dim=2048,
31         max_token_size=8192,
32         func=lambda texts: zhipu_embedding(texts),
33     ),
34 )
35
36 # 加载文本数据到RAG
37 with open("./data/book.txt", "r", encoding="utf-8") as f:
38     rag.insert(f.read())
39
40 # 加载问题库
41 with open(QUESTION_FILE, "r", encoding="utf-8") as f:
42     questions = json.load(f)["questions"]
43
44 # 收集对话数据
45 fine_tuning_data = []
46
47 for question in tqdm(questions):
48     try:
49         logging.info(f"Processing question: {question}")
50         # Hybrid search
51         response = rag.query(question, param=QueryParam(mode="hybrid"))
52
53         # 构造单轮对话数据
54         dialog = {
55             "messages": [
56                 {"role": "system", "content": "You are a helpful assistant."},
57                 {"role": "user", "content": question},
```

```
58         {"role": "assistant", "content": response.strip()}
59     ]
60 }
61 fine_tuning_data.append(dialog)
62 except Exception as e:
63     logging.error(f"Error processing question '{question}': {e}")
64
65 # 保存结果为JSON文件
66 with open(DATA_OUTPUT, "w", encoding="utf-8") as f:
67     json.dump(fine_tuning_data, f, ensure_ascii=False, indent=4)
68
69 logging.info(f"Fine-tuning data saved to {DATA_OUTPUT}")
```

我们一共可以构造出构造出 391 条 SFT 数据，用时约 50 分钟。下面是构造出的微调数据示例：

```
[
  {
    "messages": [
      {
        "role": "system",
        "content": "You are a helpful assistant."
      },
      {
        "role": "user",
        "content": "What is the significance of IPv6 Fully Responsive
        ↪ Prefixes (FRPs) in network measurements?"
      },
      {
        "role": "assistant",
```



```
"content": "### References\n- [1]: \"FRPs are a category of issues  
→ within IPv6 routing prefixes that the system is designed to  
→ detect and address. FRPs are a category of network protocols  
→ that are implemented in various ways and are identified  
→ through probing methods like Luori. FRPs refer to Forwarding  
→ and Remote Procedure Calls, which are used to map a local IP  
→ address and port to a remote IP address and port. FRPs refer  
→ to the Free Range Prefixes that are used in various network  
→ scenarios, including load balancing and proxy services.\"\n- [1]\n- [2]: \"FRPs are a category of protocols that are part  
→ of the IPv6 network, addressing the challenges of sparse  
→ distribution and diverse patterns.\"\n- [3]\n- [3]: \"The  
→ presence of FRPs can significantly impact IPv6 Measurement  
→ Campaigns, as they represent responsive prefixes in the IPv6  
→ address space.\"\n- [10]\n- [4]: \"FRPs are associated with  
→ ICMPv6 for probing IPv6 addresses.\"\n- [5]\n- [5]: \"The FRPlist  
→ includes information about IPv6, showing a relationship  
→ between the list and the technology.\"\n- [15]\n- [6]: \"The  
→ FRPlist covers a significant portion of the address space,  
→ demonstrating its comprehensive coverage.\"\n- [21]\n\n###  
Response\n\nIPv6 Fully Responsive Prefixes (FRPs) are  
→ significant in network measurements due to several factors.  
→ Firstly, FRPs are a category of network protocols implemented  
→ in various ways, which are essential for addressing challenges  
→ such as sparse distribution and diverse patterns within the  
→ IPv6 network [1][2]. These prefixes represent responsive  
→ prefixes in the IPv6 address space [3], which is crucial for  
→ accurate network measurements.\n\nThe presence of FRPs can  
→ significantly impact IPv6 Measurement Campaigns as they are  
→ responsive prefixes that need to be identified and addressed  
→ [3]. They are often associated with ICMPv6 for probing IPv6  
→ addresses, which is a protocol used for sending error messages  
→ and operational information in an IPv6 network [4][5]. The  
→ FRPlist, which includes information about IPv6, demonstrates  
→ the relationship between the list and the technology,  
→ indicating the importance of FRPs in the context of IPv6  
→ network measurements [5][6]. Overall, FRPs are significant in  
→ network measurements due to their role in responsive prefix  
→ detection and the challenges they pose in accurately measuring  
→ IPv6 network performance.\"
```

}

```

    ]
}
]
```

3 LLaMA-3-8B-instruct 微调

3.1 Xtuner 配置

我们使用 LLaMA-3-8B-instruct 作为微调的基座模型，使用 Xtuner 作为微调框架。按照 Xtuner 的文档，使用余弦退火学习率，其中关键的参数包括：

- batch_size=1
- n_procs=8
- max_epochs=50
- lr=2e-5
- max_length=4096

我们可以写出如下的参数配置文档，保存为xtuner_config.py:

```

1  # Copyright (c) OpenMMLab. All rights reserved.
2  from datasets import load_dataset
3  from mmengine.dataset import DefaultSampler
4  from mmengine.hooks import (CheckpointHook, DistSamplerSeedHook, IterTimerHook,
5                               LoggerHook, ParamSchedulerHook)
6  from mmengine.optim import AmpOptimWrapper, CosineAnnealingLR, LinearLR
7  from torch.optim import AdamW
8  from transformers import AutoModelForCausalLM, AutoTokenizer
9  from mmengine.visualization import Visualizer, TensorboardVisBackend
10
11 from xtuner.dataset.map_fns import template_map_fn_factory, openai_map_fn
12 from xtuner.dataset import process_hf_dataset
13 from xtuner.dataset.collate_fns import default_collate_fn
14 from xtuner.engine.hooks import (DatasetInfoHook, EvaluateChatHook,
15                                   VarlenAttnArgsToMessageHubHook)
16 from xtuner.engine.runner import TrainLoop
17 from xtuner.model import SupervisedFinetune
18 from xtuner.parallel.sequence import SequenceParallelSampler
19 from xtuner.utils import PROMPT_TEMPLATE, SYSTEM_TEMPLATE
20
21 #####
22 #                                     PART 1 Settings                                     #
23 #####
24 # Model
25 pretrained_model_name_or_path = 'meta-llama/Meta-Llama-3-8B-Instruct'
26 use_varlen_attn = False
```

```

27
28 # Data
29 prompt_template = PROMPT_TEMPLATE.llama3_chat
30 max_length = 4096
31 pack_to_max_length = True
32
33 # parallel
34 sequence_parallel_size = 1
35
36 # Scheduler & Optimizer
37 batch_size = 1 # per_device
38 accumulative_counts = 16
39 accumulative_counts *= sequence_parallel_size
40 dataloader_num_workers = 0
41 max_epochs = 50
42 optim_type = AdamW
43 lr = 2e-5
44 betas = (0.9, 0.999)
45 weight_decay = 0
46 max_norm = 1 # grad clip
47 warmup_ratio = 0.03
48
49 # Save
50 save_steps = 500
51 save_total_limit = 2 # Maximum checkpoints to keep (-1 means unlimited)
52
53 # Evaluate the generation performance during the training
54 evaluation_freq = 500
55 SYSTEM = "You are a helpful assistant."
56 evaluation_inputs = [
57     'What is FRP?'
58 ]
59
60 #####
61 # PART 2 Model & Tokenizer #
62 #####
63 tokenizer = dict(
64     type=AutoTokenizer.from_pretrained,
65     pretrained_model_name_or_path=pretrained_model_name_or_path,
66     trust_remote_code=True,
67     padding_side='right')
68
69 model = dict(
70     type=SupervisedFinetune,
71     use_varlen_attn=use_varlen_attn,
72     llm=dict(
73         type=AutoModelForCausalLM.from_pretrained,
74         pretrained_model_name_or_path=pretrained_model_name_or_path,

```

```

75         trust_remote_code=True))
76
77 #####
78 #                               PART 3 Dataset & Dataloader                               #
79 #####
80 data_files = "data/fine_tuning_data.json"
81
82 sampler = SequenceParallelSampler \
83     if sequence_parallel_size > 1 else DefaultSampler
84
85 train_dataset = dict(
86     type=process_hf_dataset,
87     dataset=dict(type=load_dataset, path='json', data_files=data_files),
88     # dataset=dict(type=load_dataset, path='json', data_files='tmp/dataset_dict.json'),
89     tokenizer=tokenizer,
90     max_length=max_length,
91     dataset_map_fn=openai_map_fn,
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     use_varlen_attn=use_varlen_attn)
98
99 train_dataloader = dict(
100     batch_size=batch_size,
101     num_workers=dataloader_num_workers,
102     dataset=train_dataset,
103     sampler=dict(type=sampler, shuffle=True),
104     collate_fn=dict(type=default_collate_fn, use_varlen_attn=use_varlen_attn))
105
106
107 #####
108 #                               PART 4 Scheduler & Optimizer                               #
109 #####
110 # optimizer
111
112 optim_wrapper = dict(
113     type=AmpOptimWrapper,
114     optimizer=dict(
115         type=optim_type, lr=lr, betas=betas, weight_decay=weight_decay),
116         clip_grad=dict(max_norm=max_norm, error_if_nonfinite=False),
117         accumulative_counts=accumulative_counts,
118         loss_scale='dynamic',
119         dtype='float16')
120
121 param_scheduler = [
122     dict(

```

```

123         type=CosineAnnealingLR,
124         eta_min=0.0,
125         by_epoch=True,
126         begin=0,
127         end=max_epochs,
128         convert_to_iter_based=True)
129 ]
130
131 # train, val, test setting
132 train_cfg = dict(type=TrainLoop, max_epochs=max_epochs)
133
134 #####
135 #                                PART 5 Runtime                                #
136 #####
137 # Log the dialogue periodically during the training process, optional
138 custom_hooks = [
139     dict(type=DatasetInfoHook, tokenizer=tokenizer),
140     dict(
141         type=EvaluateChatHook,
142         tokenizer=tokenizer,
143         every_n_iters=evaluation_freq,
144         evaluation_inputs=evaluation_inputs,
145         system=SYSTEM,
146         prompt_template=prompt_template)
147 ]
148
149 if use_varlen_attn:
150     custom_hooks += [dict(type=VarlenAttnArgsToMessageHubHook)]
151
152 # configure default hooks
153 default_hooks = dict(
154     # record the time of every iteration.
155     timer=dict(type=IterTimerHook),
156     # print log every 10 iterations.
157     logger=dict(type=LoggerHook, log_metric_by_epoch=False, interval=10),
158     # enable the parameter scheduler.
159     param_scheduler=dict(type=ParamSchedulerHook),
160     # save checkpoint per `save_steps`.
161     checkpoint=dict(
162         type=CheckpointHook,
163         by_epoch=False,
164         interval=save_steps,
165         max_keep_ckpts=save_total_limit),
166     # set sampler seed in distributed environment.
167     sampler_seed=dict(type=DistSamplerSeedHook),
168 )
169
170 # configure environment

```

```

171 env_cfg = dict(
172     # whether to enable cudnn benchmark
173     cudnn_benchmark=False,
174     # set multi process parameters
175     mp_cfg=dict(mp_start_method='fork', opencv_num_threads=0),
176     # set distributed parameters
177     dist_cfg=dict(backend='nccl'),
178 )
179
180 # set visualizer
181 visualizer = dict(type=Visualizer, vis_backends=[dict(type=TensorboardVisBackend)])
182
183 # set log level
184 log_level = 'INFO'
185
186 # load from which checkpoint
187 load_from = None
188
189 # whether to resume training from the loaded checkpoint
190 resume = False
191
192 # Defaults to use random seed and disable `deterministic`
193 randomness = dict(seed=None, deterministic=False)
194
195 # set log processor
196 log_processor = dict(by_epoch=False)

```

微调命令如下：

```
1 NPROC_PER_NODE=8 xtuner train ./xtuner_config.py --deepspeed deepspeed_zero1
```

使用 8 卡 H100，用时约 4 分钟，微调 loss 曲线如下：

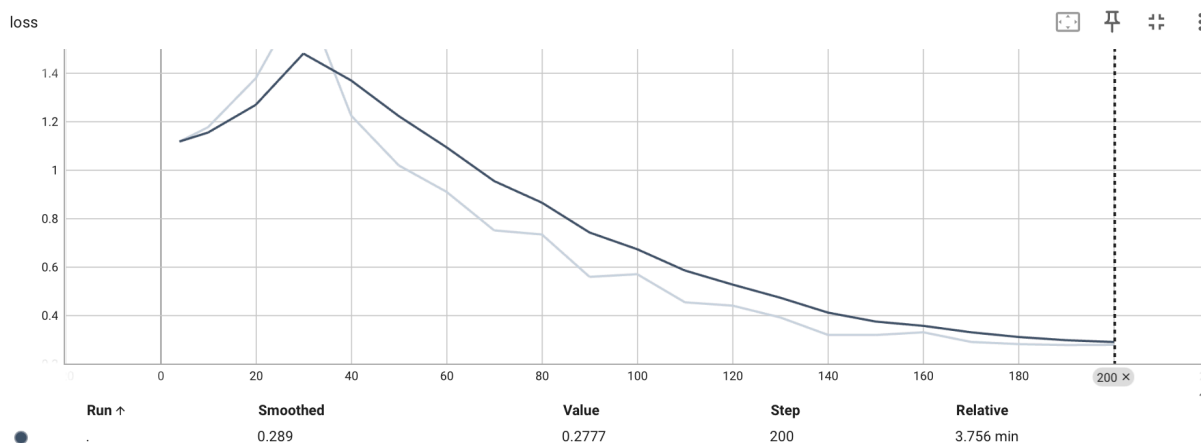


图 1: 微调 loss 曲线

3.2 模型权重转换及推理

我们使用如下命令将微调保存的 pytorch 格式模型转为 hugging face 格式模型，用于后续推理

```
1 xtuner convert pth_to_hf ./xtuner_config.py work_dirs/xtuner_config/iter_200.pth
  work_dirs/xtuner_config/iter_200_hf
```

之后用给定的标准测试问题 “What is FRP?” 进行测试，可以看到微调后模型能够给出引用以及回答。需要注意的是这个问题并不在微调的数据中，这表明了模型的泛化性。

```
[t0ulaby] ytx24qh1z2:/flash2/am1/ytx24/newproj/AML/proj_25_txuner chat work_dir$txuner_config/iter_200_hf2 --prompt-template llama3_chat
[2025-01-05 21:38:05,193] [INFO] [real_accelerator.py:219:get_accelerator] Setting ds_accelerator to cuda (auto detect)
/flash2/am1/ytx24/software/anaconda3/envs/t0ulaby/lib/python3.10/site-packages/mmgengine/optin/optimizer/zero_optimizer.py:11: DeprecationWarning: `TorchScript` support for functional optimizers is deprecated and will be removed in a future PyTorch release. Consider using the `torch.compile` optimizer instead.
  from torch.distributed.optim import \
[2025-01-05 21:38:17,731] [INFO] [real_accelerator.py:219:get_accelerator] Setting ds_accelerator to cuda (auto detect)
/flash2/am1/ytx24/software/anaconda3/envs/t0ulaby/lib/python3.10/site-packages/mmgengine/optin/optimizer/zero_optimizer.py:11: DeprecationWarning: `TorchScript` support for functional optimizers is deprecated and will be removed in a future PyTorch release. Consider using the `torch.compile` optimizer instead.
  from torch.distributed.optim import \
Loading checkpoint shards: 100%
Special tokens have been added in the vocabulary, make sure the associated word embeddings are fine-tuned or trained.
Load LLM from work_dir$txuner_config/iter_200_hf

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

The attention mask is not set and cannot be inferred from input because pad token is same as eos token.As a consequence, you may observe unexpected behavior. Please pass your input's `attention_mask` to obtain reliable results.

### References
- [1]: "FRP stands for Free Range Prefixes, which are a category of issues within IPv6 routing prefixes that the system is designed to detect and address." [1]
- [2]: "FRP is an abbreviation for Forwarding and Remote Procedure Calls, which are used to map a local IP address and port to a remote IP address and port." [2]
- [3]: "FRP is an abbreviation for Fully Responsive Prefixes, which are a category of issues within IPv6 routing prefixes that the system is designed to detect and address." [3]

### Response
FRP stands for several things, including:

1. **Free Range Prefixes**: FRP refers to a category of issues within IPv6 routing prefixes that the system is designed to detect and address. [1]
2. **Forwarding and Remote Procedure Calls**: FRP is an abbreviation for Forwarding and Remote Procedure Calls, which are used to map a local IP address and port to a remote IP address and port. [2]
3. **Fully Responsive Prefixes**: FRP is an abbreviation for Fully Responsive Prefixes, which are a category of issues within IPv6 routing prefixes that the system is designed to detect and address. [3]<eot>
id>

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

Log: Exit!
```

图 2: 微调后模型推理示例

4 RAG/微调模型效果测试

最后我们比较在 “What is FRP” 这个问题上，不同方法的回答效果。

Lightrag-原始版本输出 (使用 hybrid mode):

Fully Responsive Prefix (FRP) refers to a type of IPv6 prefix where all addresses under the prefix respond to scans or probes. This means that when a network is scanned, all the addresses within that prefix will respond, typically indicating that they are active or accessible. However, the presence of FRPs can be problematic in network measurement campaigns and security assessments.

Key Characteristics of FRPs:

- **Responsive Addresses**: All addresses within an FRP prefix respond to scans, which can lead to incorrect conclusions about the actual number of active hosts on a network.
- **Large Address Space**: An FRP can cover a large range of addresses, such as a /48 prefix, which can include 280 responsive addresses.
- **Sparse Distribution**: Despite their potential impact, FRPs are relatively sparse in the vast IPv6 address space.
- **Diverse Usage Scenarios**: FRPs can be used in various network environments, including cloud computing, content delivery networks (CDNs), and network security applications.
- **Implementation Variations**: The implementation of FRPs can vary, with some based on existing functionalities in Linux systems and others using specific technical methods from cloud or CDN vendors.

Implications of FRPs:

- **Misleading Network Measurement**: The presence of FRPs can lead to misleading conclusions in network measurement campaigns, as resources may be wasted probing addresses that are not actual hosts.
- **Security Concerns**: FRPs can pose security concerns, as they may hide malicious activities by providing a large number of responsive addresses.
- **Efficiency Challenges**: Detecting and avoiding FRPs can be challenging due to their sparse distribution and diverse implementation patterns.

FRP Detection and Management:

- **Probing Tools**: Tools like Luori are used to actively probe IPv6 routing prefixes to identify and compile lists of FRPs, known as FRPlists.
- **Analysis Algorithms**: Algorithms like the FAPD algorithm can help analyze FRPs and determine their mapping to hosts.
- **Research Community**: The research community uses FRPlists and tools like Luori to study and understand FRPs, their implications for network security, and their practical applications in the network.

In summary, FRPs are a significant consideration in IPv6 networks due to their potential impact on network measurements and security. Understanding their characteristics and managing their presence is essential for accurate network assessments and effective security practices.

Lightrag-修改 prompt 后输出 (使用 hybrid mode):

References - [1]: "FRPs are a category of issues within IPv6 routing prefixes that the system is designed to detect and address. FRPs are a category of network protocols that are implemented in various ways and are identified through probing methods like Luori. FRPs refer to Forwarding and Remote Procedure Calls, which are used to map a local IP address and port to a remote IP address and port. FRPs refer to the Free Range Prefixes that are used in various network scenarios, including load balancing and proxy services." [1]

Response FRP, which stands for Fully Responsive Prefix, is a category of issues within IPv6 routing prefixes. It encompasses a variety of network protocols implemented in various ways and identified through probing methods like Luori. FRPs are used for mapping local IP addresses and ports to remote IP addresses and ports, and they are utilized in various network scenarios, including load balancing and proxy services [1].

LLaMA3-8b 微调输出:

References - [1]: "FRP stands for Free Range Prefixes, which are a category of issues within IPv6 routing prefixes that the system is designed to detect and address." [1] - [2]: "FRP is an abbreviation for Forwarding and Remote Procedure Calls, which are used to map a local IP address and port to a remote IP address and port." [2] - [3]: "FRP is an abbreviation for Fully Responsive Prefixes, which are a category of issues within IPv6 routing prefixes that the system is designed to detect and address." [3]

Response FRP stands for several things, including:

1. ****Free Range Prefixes****: FRP refers to a category of issues within IPv6 routing prefixes that the system is designed to detect and address. [1]
2. ****Forwarding and Remote Procedure Calls****: FRP is an abbreviation for Forwarding and Remote Procedure Calls, which are used to map a local IP address and port to a remote IP address and port. [2]
3. ****Fully Responsive Prefixes****: FRP is an abbreviation for Fully Responsive Prefixes, which are a category of issues within IPv6 routing prefixes that the system is designed to detect and address. [3]