



上海大学
SHANGHAI UNIVERSITY

2024-2025 学年冬季学期
《智能计算系统》(08696037)
实验报告

姓名 汪江豪
学号 22121630
实验名称 基于 Llama 2 实现聊天机器人
日期 2025 年 2 月 8 日

实验 (70 分)	目标 1	目标 2	目标 3	得分
报告 (30 分)	代码 (10 分)	结果 (10 分)	格式 (10 分)	得分
批阅人		批阅日期		总得分

目录

— 选修实验 8-2 基于 Llama 2 实现聊天机器人	1
1 实验目的	1
2 实验环境	1
3 评分标准	1
4 实验内容及步骤	1
4.1 补全 fastchat/utils.py	1
4.2 补全 fastchat/model/model_adapter.py	2
4.3 补全 fastchat/conversation.py	4
4.4 补全 fastchat/model/chatglm_model.py	4
4.5 补全 fastchat/serve/inference.py	4
4.6 补全 fastchat/serve/cli.py	11
5 实验总结	18
5.1 实验过程回顾	18
5.2 实验成果与不足	18
5.3 实验收获	19
5.4 未来展望	19

一 选修实验 8-2 基于 Llama 2 实现聊天机器人

1 实验目的

- 熟悉 Llama 2 大语言模型的算法原理，掌握在 DLP 平台上移植优化聊天机器人的方法和流程，具体包括：
- 本实验旨在向学生介绍基于大型语言模型 Llama 2 构建聊天机器人的基本原理和操作步骤，深入理解其在对话生成中的应用；
- 学生将深入了解关键概念，包括模型加载和适配、对话模板的应用、文本生成的基本流程、模型推理的关键步骤，以及使用命令行进行聊天的具体过程；
- 学生将了解如何在 DLP 平台上部署模型，实现人机快速聊天应用，为将来在实际场景中应用模型提供实用经验。

2 实验环境

- 硬件平台：DLP 云平台环境。
- 软件环境：编程框架 Pytorch1.13.1、CNNL 高性能 AI 运算库，CNRT 运行时库，以及 python 环境及相关的扩展库。

3 评分标准

- 60 分标准：能够正确实现模型适配模块、对话模板模块，其中无论是否提供对话模板都能正确地实现程序。
- 80 分标准：在 60 分标准基础上，能够正确实现文本生成模块。
- 100 分标准：在 80 分标准基础上，能够采用 Simple 和 Rich 两种聊天模式与机器人进行交互，实现 FastChat 模型推断模块，可以完成聊天应用。

4 实验内容及步骤

4.1 补全 fastchat/utils.py

该文件主要包含了一些辅助工具函数和类，用于日志记录，重定向标准输出和错误输出，获取 GPU 信息等。

作业主要任务是补全 get_gpu_memory() 函数，代码补全如下：

代码 1: utils.py

```
1 def get_gpu_memory(max_gpus=None):  
2     """Get available memory for each GPU."""  
3     #TODO: 存储每个GPU的可用内存信息  
4     gpu_memory = []  
5     #TODO: 获取MLU设备的数量，如果未指定最大GPU数（max_gpus为None），则使用所有  
      可用设备；否则，使用max_gpus和实际设备数量中的较小值  
6     num_gpus = (  
7         min(max_gpus, torch.mlu.device_count())  
8         if max_gpus is not None  
9         else torch.mlu.device_count())
```

```
10     )
11
12     for gpu_id in range(num_gpus):
13         #TODO: 将当前MLU设备设置为gpu_id
14         torch.mlu.set_device(f"mlu:{gpu_id}")
15         #TODO: 获取当前MLU设备
16         device = torch.mlu.current_device()
17         #TODO: 获取当前MLU设备的属性
18         gpu_properties = torch.mlu.get_device_properties(device)
19         #TODO: 获取总内存, 单位转换为GB
20         total_memory = gpu_properties.total_memory / (1024 ** 3) # GB
21         #TODO: 获取已分配内存, 单位转换为GB
22         allocated_memory = torch.mlu.memory_allocated() / (1024 ** 3) # GB
23         #TODO: 计算可用内存, 单位转换为GB
24         available_memory = total_memory - allocated_memory
25         #TODO: 将可用内存信息添加到列表中
26         gpu_memory.append(available_memory)
27     return gpu_memory
```

4.2 补全 fastchat/model/model_adapter.py

该文件主要用于模型适配器的注册和管理，它定义了一个基础模型适配器类 `BaseAdapter` 及其多个子类。每个子类正对不同的模型类型实现了特定的加载和匹配逻辑。文件还包含了用于注册模型适配器的函数、加载模型的函数以及获取对话模板的函数。通过这些适配器，可以根据模型路径动态选择和加载合适的模型和分词器。

作业主要任务是补全 `load_model()` 函数和 `get_conversation_template()` 函数。

由于篇幅限制，这里只展示部分补全的代码

代码 2: `load_model()` 函数

```
1 def load_model(
2     model_path: str,
3     device: str,
4     num_gpus: int,
5     max_gpu_memory: Optional[str] = None,
6     load_8bit: bool = False,
7     cpu_offloading: bool = False,
8     debug: bool = False,
9 ):
10     """Load a model from Hugging Face."""
11
12     # TODO: 处理设备映射, 调用函数, 检查并更新CPU offloading的配置
13     cpu_offloading = raise_warning_for_incompatible_cpu_offloading_configuration
14     (
15         device, load_8bit, cpu_offloading
16     )
17     # TODO: 如果设备是CPU
```

```

17     if device == "cpu":
18         kwargs = {"torch_dtype": torch.float32}
19     # TODO: 如果设备是MLU
20     elif device == "mlu":
21         kwargs = {"torch_dtype": torch.float16}
22         if num_gpus != 1:
23             kwargs["device_map"] = "auto"
24             if max_gpu_memory is None:
25                 kwargs["device_map"] = (
26                     "sequential" # This is important for not the same VRAM
27                     sizes
28                 )
29             # TODO: 获取每个GPU的可用内存
30             available_gpu_memory = get_gpu_memory(num_gpus)
31             kwargs["max_memory"] = {
32                 i: str(int(available_gpu_memory[i] * 0.85)) + "GiB"
33                 for i in range(num_gpus)
34             }
35         else:
36             kwargs["max_memory"] = {i: max_gpu_memory for i in range(
37                 num_gpus)}
38     # TODO: 如果设备是mps
39     elif device == "mps":
40         kwargs = {"torch_dtype": torch.float16}
41         # Avoid bugs in mps backend by not using in-place operations.
42         replace_llama_attn_with_non_inplace_operations()
43     else:
44         raise ValueError(f"Invalid device: {device}")
45
46     # TODO: 如果启用了CPU offloading
47     if cpu_offloading:
48         # raises an error on incompatible platforms
49         from transformers import BitsAndBytesConfig
50
51         if "max_memory" in kwargs:
52             kwargs["max_memory"]["cpu"] =
53                 str(math.floor(psutil.virtual_memory().available / 2**20)) + "
54                 Mib"
55         )
56         kwargs["quantization_config"] = BitsAndBytesConfig(
57             load_in_8bit_fp32_cpu_offload=cpu_offloading
58         )
59         kwargs["load_in_8bit"] = load_8bit
60     # TODO: 如果启用了8位量化但未启用CPU offloading
61     elif load_8bit:
62         if num_gpus != 1:
63             warnings.warn(

```

```

61             "8-bit quantization is not supported for multi-gpu inference."
62         )
63     else:
64         # TODO: 调用压缩模型的函数
65         return load_compress_model(model_path, device=device, **kwargs)
66
67     # TODO: 调用函数，根据给定的模型路径获取适配器，遍历已注册的模型适配器列表，
68     #       返回第一个匹配的适配器实例。
69     adapter = get_model_adapter(model_path)
70     # TODO: 使用适配器adapter加载模型和分词器
71     model, tokenizer = adapter.load_model(model_path, kwargs)
72
73     if (device == "mlu" and num_gpus == 1 and not cpu_offloading) or device == "mps":
74         model.to(device)
75
76     if debug:
77         print(model)
78     print("MODEL ADAPTER PASS!")
79     return model, tokenizer

```

4.3 补全 fastchat/conversation.py

该文件主要定义了对话提示模板和对话类。它包含了不同分隔符样式的枚举类和用于存储对话历史的类。

文件还定义了一些全局对话模板，并提供了注册和获取对话模板的函数。通过这些模板，可以方便地管理和生成不同风格的对话。

作业主要为补全 `get_prompt()` 函数和 `append_message()` 函数。

由于篇幅限制，且此部分代码不是重点，补全代码不在此展示，详细内容见代码文件。

4.4 补全 fastchat/model/chatglm_model.py

该文件主要包含两个函数。

`stream_chat_token_num()` 函数用于计算给定对话历史和当前问题的 token 数量。

`chatglm_generate_stream()` 函数使用给定的模型和分词器，根据输入参数生成文本流，支持多轮对话，并在生成过程中返回生成的文本和相关的 token 数量。

同理篇幅限制，补全详情见代码文件。

4.5 补全 fastchat/serve/inference.py

该文件是实现推理的核心功能，主要用于 FastChat 模型的推理，它定义了生成文本流的函数 `generate_stream`，并导入了 `chatglm_generate_stream` 函数，用于根据输入参数生成对话文本。

`chatloop` 负责加载模型和分词器，管理对话循环，并调用相应的生成函数生成和输出对话文本。

这里需要注意的是，在 `generate_stream` 函数中，要及时关闭梯度计算，否则会导致内存泄漏。

由于此部分是实现推理的核心功能，我将详细展示各补全的接口函数。

代码 3: generate_stream() 函数

```
1 def generate_stream(
2     model, tokenizer, params, device, context_len=2048, stream_interval=2
3 ):
4     #TODO:# 关闭梯度计算
5     with torch.no_grad():
6         prompt = params["prompt"]
7         #TODO: 获取prompt的长度
8         len_prompt = len(prompt) if isinstance(prompt, str) else sum(len(s) for
9             s in prompt)
10        #TODO: 从参数中获取生成文本时的temperature (控制文本生成的多样性), 如果
11        #       参数中未指定, 则默认为 1.0。
12        temperature = float(params.get("temperature", 1.0))
13        #TODO: 从参数中获取生成文本时的 repetition_penalty (抑制文本中重复的程
14        #       度), 如果参数中未指定, 则默认为 1.0。
15        repetition_penalty = float(params.get("repetition_penalty", 1.0))
16        #TODO: 从参数中获取生成文本时的 top_p (控制生成文本的多样性), 如果参数中
17        #       未指定, 则默认为 1.0。
18        top_p = float(params.get("top_p", 1.0))
19        #TODO: 从输入的参数中获取 top_k 的值, 如果参数中没有设置, 则默认为 -1。
20        top_k = int(params.get("top_k", -1)) # -1 means disable
21        #TODO 从参数中获取生成文本时的max_new_tokens数, 如果参数中未指定, 则默认
22        #       为 256。
23        max_new_tokens = int(params.get("max_new_tokens", 256))
24        #TODO: 从参数中获取 stop_str, 如果未设置, 则默认为 None
25        stop_str = params.get("stop", None)
26        #TODO: 从参数中获取 echo, 如果未设置, 默认为 True。并将其转换为布尔值。
27        echo = bool(params.get("echo", True))
28        stop_token_ids = params.get("stop_token_ids", None) or []
29        #TODO: 将文本生成停止的标记添加到已有的停止标记列表stop_token_ids中。
30        if stop_str and isinstance(stop_str, str):
31            # 这里仅取第一个token id作为示例, 可根据需要调整
32            stop_token_ids.append(tokenizer.encode(stop_str, add_special_tokens=
33                False)[0])
34
35        #TODO: 创建一个 logits 处理器列表
36        logits_processor = prepare_logits_processor(temperature,
37            repetition_penalty, top_p, top_k)
38
39        #TODO: 使用tokenizer将输入文本转换为模型可接受的输入张量
40        input_ids = tokenizer.encode(prompt, add_special_tokens=False)
41        #TODO: 记录输入文本的长度
42        input_echo_len = len(input_ids)
43        #TODO: 创建一个名为 output_ids 的列表, 其初始值等于 input_ids 列表的内容
44        output_ids = input_ids.copy()
45
46        if model.config.is_encoder_decoder:
```

```

40         max_src_len = context_len
41     else:
42         max_src_len = context_len - max_new_tokens - 8
43
44     #TODO: 截取源文本的最后一部分，以适应模型的上下文长度
45     if len(input_ids) > max_src_len:
46         input_ids = input_ids[-max_src_len:]
47
48     if model.config.is_encoder_decoder:
49         encoder_output = model.encoder(
50             input_ids=torch.as_tensor([input_ids], device=device)
51         )[0]
52         start_ids = torch.as_tensor(
53             [model.generation_config.decoder_start_token_id],
54             dtype=torch.int64,
55             device=device,
56         )
57
58     past_key_values = out = None
59     for i in range(max_new_tokens):
60         if i == 0:
61             if model.config.is_encoder_decoder:
62                 #TODO: 使用解码器对起始标记进行处理
63                 out = model.decoder(
64                     input_ids=start_ids,
65                     encoder_hidden_states=encoder_output,
66                     use_cache=True,
67                 )
68                 logits = model.lm_head(out[0])
69             else:
69                 #TODO: 非编码解码器类型，直接使用模型进行处理
70                 out = model(
71                     input_ids=torch.as_tensor([input_ids], device=device),
72                     use_cache=True,
73                 )
74                 logits = out.logits
75             #TODO: 记录过去的键值
76             past_key_values = out.past_key_values
77         else:
78             if model.config.is_encoder_decoder:
79                 #TODO: 使用解码器对当前标记进行处理
80                 out = model.decoder(
81                     input_ids=torch.as_tensor([[output_ids[-1]]], device=
82                         device),
83                     encoder_hidden_states=encoder_output,
84                     use_cache=True,
85                     past_key_values=past_key_values,

```

```

86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
    )
    #TODO: 获取logits (预测的标记分布)
    logits = model.lm_head(out[0])
else:
    out = model(
        #TODO: 非编码解码器类型, 直接使用模型对当前标记进行处理
        input_ids=torch.as_tensor([[output_ids[-1]]], device=
            device),
        use_cache=True,
        past_key_values=past_key_values,
    )
    logits = out.logits
    #TODO: 更新过去的键值
    past_key_values = out.past_key_values

if logits_processor:
    if repetition_penalty > 1.0:
        tmp_output_ids = torch.as_tensor([output_ids], device=logits
            .device)
    else:
        tmp_output_ids = None
    #TODO: 使用logits_processor处理最后一个标记的logits
    last_token_logits = logits_processor(tmp_output_ids, logits[:, -1, :])
else:
    #TODO: 没有logits处理器, 直接获取最后一个标记的logits
    last_token_logits = logits[:, -1, :]

if device == "mps":
    # Switch to CPU by avoiding some bugs in mps backend.
    last_token_logits = last_token_logits.float().to("cpu")

if temperature < 1e-5 or top_p < 1e-8: # greedy
    #TODO:# 通过argmax获取概率最高的标记索引, 并将其转换为整数类型。
    token = int(torch.argmax(last_token_logits, dim=-1).item())
else:
    #TODO: 使用softmax将概率分布归一化
    probs = torch.softmax(last_token_logits / temperature, dim=-1)
    #TODO: 使用torch.multinomial函数根据概率分布采样生成标记, 并转换
    #为整数类型。
    token = int(torch.multinomial(probs, num_samples=1).item())

# 将生成的token添加到输出序列output_ids中
output_ids.append(token)

if token in stop_token_ids:
    stopped = True

```

```

129     else:
130         stopped = False
131
132         #TODO: 判断是否达到生成结果的间隔、已完成生成最大标记数、或者已经停止生成
133
134         if (i % stream_interval == 0) or (i == max_new_tokens - 1) or
135             stopped:
136             if echo:
137                 tmp_output_ids = output_ids
138                 rfind_start = len_prompt
139             else:
140                 tmp_output_ids = output_ids[input_echo_len:]
141                 rfind_start = 0
142
143             output = tokenizer.decode(
144                 tmp_output_ids,
145                 skip_special_tokens=True,
146                 spaces_between_special_tokens=False,
147             )
148             if stop_str:
149                 if isinstance(stop_str, str):
150                     #TODO: 在输出字符串中从右向左搜索停止标记的位置,
151                     #      rfind_start参数指定了搜索的起始位置
152                     pos = output.rfind(stop_str, rfind_start)
153                     if pos != -1:
154                         #TODO: 将输出字符串截断, 仅保留停止标记位置之前的部分
155                         output = output[:pos]
156                         stopped = True
157                     elif isinstance(stop_str, Iterable):
158                         for each_stop in stop_str:
159                             #TODO: 从指定位置 (rfind_start) 向前搜索每个停止标记
160                             #      在输出字符串中的最后出现位置
161                             pos = output.rfind(each_stop, rfind_start)
162                             if pos != -1:
163                                 #TODO: 将输出字符串截断, 仅保留停止标记位置之前的部分
164                                 output = output[:pos]
165                                 stopped = True
166                                 break
167                             else:
168                                 raise ValueError("Invalid stop field type.")
169
170             yield {
171                 "text": output,
172                 "usage": {
173                     "prompt_tokens": input_echo_len,

```

```

170             "completion_tokens": i,
171             "total_tokens": input_echo_len + i,
172         },
173         "finish_reason": None,
174     }
175
176     if stopped:
177         break
178
179     # finish stream event, which contains finish reason
180     if i == max_new_tokens - 1:
181         finish_reason = "length"
182     elif stopped:
183         finish_reason = "stop"
184     else:
185         finish_reason = None
186
187     #TODO: 返回生成的文本、使用情况和结束原因的字典
188     yield {
189         "text": tokenizer.decode(output_ids[input_echo_len:],
190             skip_special_tokens=True, spaces_between_special_tokens=False),
191         "usage": {
192             "prompt_tokens": input_echo_len,
193             "completion_tokens": i if not stopped else i+1,
194             "total_tokens": input_echo_len + (i if not stopped else i+1),
195         },
196         "finish_reason": finish_reason,
197     }
198
199     # clean
200     del past_key_values, out
201     gc.collect()
202     #TODO: 释放MLU设备上的缓存空间
203     torch.mlu.empty_cache()

```

该函数用于在生成模型上执行流式文本生成。其流程包括：

从输入参数中提取生成配置（如温度、重复惩罚、top_p、top_k、最大 token 数、停止字符串/标记等）。

利用 tokenizer 将提示文本转为模型输入，截断以适应模型的最大上下文长度。

根据模型类型（编码器-解码器或普通生成模型），分别调用相应的模型接口进行首次生成和后续 token 生成，同时利用 past_key_values 缓存先前的计算结果以提高效率。

使用 logits 处理器和概率采样方法选取下一个 token，将其拼接到输出 token 列表中。

根据设定的输出间隔、达到最大 token 或遇到停止标记，实时解码并返回当前生成的文本，同时记录 token 使用情况。

最终，当生成完成或触发停止条件时，返回最终的生成结果以及相关的使用信息，并清理缓存。

这里注意，最后的清理缓存，toch_mlu 中不存在相关缓存清理函数（曾困扰了我很久）。且要及时关闭梯度计算，否则会爆显存。

代码 4: chat_loop() 函数

```
1 def chat_loop(
2     model_path: str,
3     device: str,
4     num_gpus: int,
5     max_gpu_memory: str,
6     load_8bit: bool,
7     cpu_offloading: bool,
8     conv_template: Optional[str],
9     temperature: float,
10    max_new_tokens: int,
11    chatio: ChatIO,
12    debug: bool,
13 ):
14     #TODO: 调用 load_model 函数加载 model 和 tokenizer
15     model, tokenizer = load_model(
16         model_path,
17         device=device,
18         num_gpus=num_gpus,
19         max_gpu_memory=max_gpu_memory,
20         load_8bit=load_8bit,
21         cpu_offloading=cpu_offloading,
22     )
23     is_chatglm = "chatglm" in str(type(model)).lower()
24
25     #TODO: 如果提供了对话模板, 使用提供的模板, 调用 get_conv_template 创建会话对象
26     if conv_template:
27         conv = get_conv_template(conv_template)
28     else:
29         #TODO: 否则使用默认的对话模板, 调用 get_conversation_template 创建会话对象
30         conv = get_conversation_template("default")
31     print("GENERATE STEAM PASS!")
32     while True:
33         try:
34             #TODO: 尝试获取用户输入, 传入 conv.roles[0] 作为角色标识
35             inp = chatio.prompt_for_input(conv.roles[0])
36         except EOFError:
37             inp = ""
38         if not inp:
39             print("exit...")
40             break
41
42         conv.append_message(conv.roles[0], inp)
43         conv.append_message(conv.roles[1], None)
44
45         if is_chatglm:
46             generate_stream_func = chatglm_generate_stream
```

```

47         prompt = conv.messages[conv.offset :]
48     else:
49         generate_stream_func = generate_stream
50         prompt = conv.get_prompt()
51
52     gen_params = {
53         "model": model_path,
54         "prompt": prompt,
55         "temperature": temperature,
56         "max_new_tokens": max_new_tokens,
57         "stop": conv.stop_str,
58         "stop_token_ids": conv.stop_token_ids,
59         "echo": False,
60     }
61
62     #TODO: 获取机器对话输出, 传入 conv.roles[1] 作为角色标识
63     chatio.prompt_for_output(conv.roles[1])
64     output_stream = generate_stream_func(model, tokenizer, gen_params,
65                                           device)
66     #TODO: # 输出对话的流式输出
67     outputs = chatio.stream_output(output_stream)
68
69     # NOTE: strip is important to align with the training data.
70     conv.messages[-1][-1] = outputs.strip()
71     if debug:
72         print("\n", {"prompt": prompt, "outputs": outputs}, "\n")
72         print("FASTCHAT INFERENCE PASS!")

```

该函数是实现多轮交互对话的关键功能。

通过加载模型和 tokenizer，根据传入的对话模板参数创建会话对象，(可使用自定义或默认模板)。在循环中，每次获取用户输入，并将用户信息添加到对话记录中。根据模型类型选择生成函数 (chatglm 或普通生成)，然后构造生成参数，并调用模型进行流式文本生成，然后实时显示生成结果。

4.6 补全 fastchat/serve/cli.py

该文件提供了一个命令行界面与模型进行对话。还定义了两个类，SimpleChatIO 和 RichChatIO，分别用于简单对话和富文本对话，都继承自 ChatIO 基类。

主函数 main 解析命令行参数，选择合适的输入输出类，并调用 chat_loop 函数启动聊天循环，实现与模型的交互。

其中还含有命令行参数的定义和选项。

该文件是对话程序的主逻辑，我将详细展示补全代码：

代码 5: SimpleChatIO 类

```

1     #TODO: 构建 SimpleChatIO 类, 它继承自 ChatIO 类
2     class SimpleChatIO(ChatIO):
3         def prompt_for_input(self, role) -> str:
4             return input(f"{role}: ")

```

```

5
6     def prompt_for_output(self, role: str):
7         print(f"{role}: ", end="", flush=True)
8
9     def stream_output(self, output_stream):
10        pre = 0
11        for outputs in output_stream:
12            output_text = outputs["text"]
13            #TODO: 移除文本两端的空白字符, 然后按空格分割
14            output_text = output_text.strip().split(" ")
15            #TODO: 获取处理后的文本中单词的数量
16            now = len(output_text)
17            if now > pre:
18                # 输出不同于前次的部分, 其中, 新增的内容以空格分隔的形式显示, 确
19                # 保每次新增的内容都在同一行, 并及时刷新输出缓冲区。
20                print(" ".join(output_text[pre:now]), end=" ", flush=True)
21            pre = now
22        print(" ".join(output_text[pre:]), flush=True)
23        print("CHATIO PASS!")
24        return " ".join(output_text)

```

该类继承自 inference.py 文件中的 ChatIO 类，用于实现简单的命令行交互输入输出。主要功能包括 prompt_for_input()，根据指定角色提示用户输入，并返回输入内容；prompt_for_output()，在命令行中打印角色标识，准备输出生成的回答；stream_output()，流式处理生成器返回的文本，在处理过程中，函数会对文本空白字符去除和按空格分割，然后只输出新生成的部分，确保输出及时刷新，直至完成整个生成过程，最终函数返回完整的输出文本。

代码 6: RichChatIO 类

```

1     #TODO: 构建RichChatIO类, 它继承自ChatIO类
2     class RichChatIO(ChatIO):
3         def __init__(self):
4             #TODO: 创建PromptSession实例, 用于获取用户输入, 并设置输入历史记录
5             self._prompt_session = PromptSession(history=InMemoryHistory())
6             #TODO: 创建自动补全器, 用于用户输入的自动完成
7             self._completer = WordCompleter(words=["!exit", "!reset"], pattern=re.
8                 compile(r"\$"))
9             #TODO: 创建Console 实例, 在命令行界面中以更丰富的样式显示文本
10            self._console = Console()
11
12            def prompt_for_input(self, role) -> str:
13                self._console.print(f"[bold]{role}:")
14                # TODO(suquark): multiline input has some issues. fix it later.
15                prompt_input = self._prompt_session.prompt(
16                    completer=self._completer,
17                    multiline=False,
18                    #TODO: 启用自动建议功能
19                    auto_suggest=AutoSuggestFromHistory(),

```

```

19         key_bindings=None,
20     )
21     self._console.print()
22     return prompt_input
23
24 def prompt_for_output(self, role: str):
25     self._console.print(f"[bold]{role}:")
26
27 def stream_output(self, output_stream):
28     """Stream output from a role."""
29     # TODO(suquark): the console flickers when there is a code block
30     # above it. We need to cut off "live" when a code block is done.
31
32     # Create a Live context for updating the console output
33
34     #TODO: 创建Live上下文管理器，用于实现在命令行中实时更新显示。其中，需要
35     #指定要在其上执行实时更新的console实例，每秒刷新的次数设为4
36     with Live(console=self._console, refresh_per_second=4) as live:
37         # Read lines from the stream
38         for outputs in output_stream:
39             if not outputs:
40                 continue
41             # 如果输出是字典，则获取字典中的"text"键对应的值，否则直接获取输出
42             if isinstance(outputs, dict):
43                 text = outputs.get("text", "")
44             else:
45                 text = outputs
46
47             lines = []
48             for line in text.splitlines():
49                 if "### Human" in line:
50                     continue
51                 lines.append(line)
52                 if line.startswith("```"):
53                     # Code block marker - do not add trailing spaces, as it
54                     # would
55                     # break the syntax highlighting
56                     lines.append("\n")
57                     markdown = Markdown("").join(lines))
58                     #TODO: 将渲染后的 markdown 文本实时更新到控制台
59                     live.update(markdown)
60                     self._console.print()
61                     print("CHATIO PASS!")
62                     return text

```

该类继承自 inference.py 文件中的 ChatIO 类，用于实现富文本的命令行交互输入输出。同样包括在控制台以加粗样式输出角色标识，并使用 PromptSession 获取用户输入，同时启用历史自动建议。在 stream_output() 函数中，使用 Rich 的 Live 上下文管理器实现实时刷新，将流式获取的输出文本以 Markdown 格式渲染，然后实时更新显示，确保生成过程在终端平滑显示。最终稿返回输出文本。提供了更加美观的交互界面。

代码 7：主函数部分

```
1 def main(args):
2     if args.gpus:
3         if len(args.gpus.split(",")) < args.num_gpus:
4             raise ValueError(
5                 f"Large --num-gpus ({args.num_gpus}) than --gpus {args.gpus}!")
6
7     os.environ["CUDA_VISIBLE_DEVICES"] = args.gpus
8
9     if args.style == "simple":
10        chatio = SimpleChatIO()
11    elif args.style == "rich":
12        chatio = RichChatIO()
13    else:
14        raise ValueError(f"Invalid style for console: {args.style}")
15    try:
16        #TODO: 调用 chat_loop 函数，启动聊天循环
17        chat_loop(
18            args.model_path,
19            args.device,
20            args.num_gpus,
21            args.max_gpu_memory,
22            args.load_8bit,
23            args.cpu_offloading,
24            args.conv_template,
25            args.temperature,
26            args.max_new_tokens,
27            chatio,
28            args.debug,
29        )
30    except KeyboardInterrupt:
31        print("exit...")
32
33
34 if __name__ == "__main__":
35     #TODO: 创建一个ArgumentParser对象，用于解析命令行参数
36     parser = argparse.ArgumentParser()
37     #TODO: 向ArgumentParser对象中添加模型相关的参数，这些参数由add_model_args函数定义
38     add_model_args(parser)
39     parser.add_argument(
```

```

40         "--conv-template", type=str, default=None, help="Conversation prompt
41             template."
42     )
43     parser.add_argument("--temperature", type=float, default=0.7)
44     parser.add_argument("--max-new-tokens", type=int, default=512)
45
46     #TODO:试一下simple模式，同时也试一下rich模式。
47     parser.add_argument(
48         "--style",
49         type=str,
50         default="simple",
51         choices=["simple", "rich"],
52         help="Display style.",
53     )
54     parser.add_argument("--debug", action="store_true", help="Print debug
55             information")
56     args = parser.parse_args()
57     main(args)

```

这段主函数部分，实现了命令行应用程序的入口。主要包括：根据命令行参数配置 GPU，根据用户指定的显示风格（simple 或 rich）实例化对应的 ChatIO 类，以确定交互界面的显示方式。解析并设置模型相关与对话生成的参数，例如模型路径，对话模板，温度，最大 token 数等。调用 inference.py 文件中 chat_loop 函数，启动多轮对话循环，进行用户输入，文本生成与输出展示。同时有捕获异常，实现用户中断时的退出。

整体上这一部分负责初始化应用环境，解析参数，并启动聊天交互过程。同时也是测试环境中，脚本测试的对象文件。

在开发环境中，通过运行当前目录下的的 run_scripts 文件夹中的 run_mlu_infer.sh 脚本，可以对 cli.py 进行测试，从而测试整个模型推理效果。

脚本命令如下：

代码 8: run_completion.sh

```

1 set -e
2
3 model_path=/workspace/model/favorite/large-scale-models/model-v1/Llama-2-7b-hf
4 echo ****
5 echo **命令行终端形式 仅支持单轮对话**
6 echo ****
7 # 单卡
8 # 配置板卡
9 export MLU_VISIBLE_DEVICES=0
10 python3 -m fastchat.serve.cli --model-path ${model_path} --style=rich

```

在脚本命令后输入相应参数可实现多种功能，如调试模式，选择富文本或简单文本等。

运行结果如图1-3所示：

在测试中，我询问了 AI，“你是谁”，“询问了质能方程”，“询问了牛顿第二定律”。AI 均能正确回答问题。

```
(pytorch) root@notebook-devenvironment-0205-221056-1c0zo6s-notebook-0:/opt/code_chap_8_student/llama2_7b# ./run_scripts/run_mlu_infer.sh
*****
***命令行终端形式 仅支持单轮对话**
*****
Loading checkpoint shards: 100%|██████████| 2/2 [00:10<00:00,  5.20s/it]
MODEL ADAPTER PASS!
CONVERSATION PASS!
GENERATE STEAM PASS!
Human:
hello, who are you?

Assistant:
/opt/code_chap_8_student/llama2_7b/fastchat/serve/inference.py:127: UserWarning: MLU operators don't support 64-bit calculation, so the 64 bit data will be forcibly converted to 32-bit for calculation. (Triggered internally at /torch/catch/torch_mlu/csrc/aten/utils/tensor_util.cpp:159.)
    input_ids=torch.as_tensor([input_ids], device=device),
[2025-2-8 14:15:35] [CNNL] [Warning]:[cnnlRandCreateGenerator_v2] will be deprecated.
I am an artificial intelligence assistant, and I am here to help you with your questions.
```

图 1: 多轮对话结果展示 1

```
[2025-2-8 14:15:35] [CNNL] [Warning]:[cnnlRandCreateGenerator_v2] will be deprecated.
I am an artificial intelligence assistant, and I am here to help you with your questions.

CHATIO PASS!
FASTCHAT INFERENCE PASS!
Human:
what is E = mc^2?

Assistant:
E=mc^2 is a famous equation in physics that states that energy is equal to mass times the speed of light squared. This equation is important in understanding the relationship between mass and energy, and how they are related to each other.

CHATIO PASS!
FASTCHAT INFERENCE PASS!
Human:
```

图 2: 多轮对话结果展示 2

```
Assistant:
E=mc^2 is a famous equation in physics that states that energy is equal to mass times the speed of light squared. This equation is important in understanding the relationship between mass and energy, and how they are related to each other.

CHATIO PASS!
FASTCHAT INFERENCE PASS!
Human:
what is F=ma?

Assistant:
F=ma is another famous equation in physics that states that force equals mass times acceleration. This equation is important in understanding how forces act on objects and how they are related to each other.

CHATIO PASS!
FASTCHAT INFERENCE PASS!
Human:
```

图 3: 多轮对话结果展示 3

可见，正确实现了对话功能，可以与模型进行多轮对话交互。

之前在希冀平台提交的结果中，最后两项 chatio 和 fastchat_inference 始终未通过。后来经过助教通知，了解到，希冀平台的评测机机制中，没法处理多轮对话交互的情况。解决办法是将 inference.py 中的 chat_loop 函数中的 while(true) 循环逻辑移除，并为 inp 赋值为一个固定的内容，即可。修改的代码部分如下，见注释部分和 inp 变量：

代码 9：修改交互部分

```
1 def chat_loop(
2     model_path: str,
3     device: str,
4     num_gpus: int,
5     max_gpu_memory: str,
6     load_8bit: bool,
7     cpu_offloading: bool,
8     conv_template: Optional[str],
9     temperature: float,
10    max_new_tokens: int,
11    chatio: ChatIO,
12    debug: bool,
13 ) :
14     # TODO: 调用 load_model 函数 加 model 和 tokenizer
15     model, tokenizer = load_model(
16         model_path,
17         device=device,
18         num_gpus=num_gpus,
19         max_gpu_memory=max_gpu_memory,
20         load_8bit=load_8bit,
21         cpu_offloading=cpu_offloading,
22     )
23     is_chatglm = "chatglm" in str(type(model)).lower()
24
25     # TODO: 如果提供了对话模板，使用提供的模板，调用 get_conv_template 创建会话对象
26     if conv_template:
27         conv = get_conv_template(conv_template)
28     else:
29         # TODO: 否则使用默认的对话模板，调用 get_conversation_template 创建会话对象
30         conv = get_conversation_template("default")
31     print("GENERATE STEAM PASS!")
32     # while True:
33     #     try:
34     #         # TODO: 尝试获取用户输入，传入 conv.roles[0] 作为角色标识
35     #         inp = chatio.prompt_for_input(conv.roles[0])
36     #     except EOFError:
37     #         inp = ""
38     #     if not inp:
39     #         print("exit...")
40     #         break
```

```

41     inp = "hello, who are you?"
42     conv.append_message(conv.roles[0], inp)
43     conv.append_message(conv.roles[1], None)
44
45     if is_chatglm:
46         generate_stream_func = chatglm_generate_stream
47         prompt = conv.messages[conv.offset :]
48     else:
49         generate_stream_func = generate_stream
50         prompt = conv.get_prompt()
51
52     gen_params = {
53         "model": model_path,
54         "prompt": prompt,
55         "temperature": temperature,
56         "max_new_tokens": max_new_tokens,
57         "stop": conv.stop_str,
58         "stop_token_ids": conv.stop_token_ids,
59         "echo": False,
60     }
61
62     # TODO: 获取机器对话输出, 传入 conv.roles[1] 作为角色标识
63     chatio.prompt_for_output(conv.roles[1])
64     output_stream = generate_stream_func(model, tokenizer, gen_params, device)
65     # TODO: # 输出对话的流式输出
66     outputs = chatio.stream_output(output_stream)
67
68     # NOTE: strip is important to align with the training data.
69     conv.messages[-1][-1] = outputs.strip()
70     if debug:
71         print("\n", {"prompt": prompt, "outputs": outputs}, "\n")
72     print("FASTCHAT INFERENCE PASS!")

```

最后，希冀平台评测结果如图4：

可见，成功通过了所有测试点，实验完成。

5 实验总结

5.1 实验过程回顾

实验初期，面对复杂的实验环境和众多的代码文件，我感到有些无从下手。但在逐步梳理实验步骤后，我开始按照计划补全各个文件中的代码。我一边查阅资料学习相关知识，一边仔细推敲每一行代码的逻辑，确保其功能的正确实现。

5.2 实验成果与不足

经过不懈努力，我成功实现了基于 Llama 2 的聊天机器人，并通过命令行界面与模型进行了多轮对话交互。在开发环境中测试时，结果显示了最终的 pass 信息，表明实验功能基本实现。然而，在希冀平台

最后一次提交时间: 2025-02-21 10:23:54

model_adpter	conversation	generate_steam	chatio	fastchat_inference
PASS	PASS	PASS	PASS	PASS

Accept

model_adpter	conversation	generate_steam	chatio	fastchat_inference
PASS	PASS	PASS	PASS	PASS

图 4: 8-2 测试结果

的判题系统中，却未通过最后两项 chatio 和 fastchat_inference 测试，这让我感到有些遗憾。尽管我尝试了多种方法进行调试，但仍未找到具体原因，可能是接口功能未被测试捕捉到，也可能是代码本身存在一些隐藏的 BUG。

5.3 实验收获

技术能力提升: 通过本次实验，我深入学习了 Llama 2 模型的相关知识，掌握了在 DLP 平台上进行模型移植优化的方法，熟悉了 Pytorch 框架下模型的加载、适配、推理等操作流程，提升了自己在智能计算领域的技术能力。

问题解决能力锻炼: 在实验过程中遇到的各种问题，促使我不断查阅资料、分析代码、调试程序，锻炼了我的问题解决能力。面对复杂的代码和未知的错误，我学会了保持冷静，从不同角度思考问题，逐步排查并解决问题。

理论与实践结合: 本次实验将课堂上学到的理论知识与实际操作紧密结合，让我更加深刻地理解了智能计算系统的原理和应用。通过亲自动手实现聊天机器人，我体会到了理论指导实践、实践检验理论的重要性。

5.4 未来展望

虽然本次实验取得了一定的成果，但也暴露出了一些不足之处。在今后的学习和实践中，我将继续深入研究智能计算领域的相关技术，不断提升自己的技术水平。同时，我也会更加注重细节，努力提高代码的质量和稳定性，避免类似的问题再次出现。此外，我还希望能够将所学知识应用到更多的实际项目中，为智能计算技术的发展贡献自己的一份力量。