

问题

什么是模型部署？

它的基本结构有哪些？

输入指令和输出结果的信息格式有什么要求和特点，其它辅助的信息比如验证指令等信息如何传递。如何调用模型（模型加载）使其回答问题

模型部署

平时使用大模型是直接访问如 deepseek, qwen 等模型提供的网页或 app，此时我们只需要输入自然语言描述的问题，服务器会自动封装然后调用模型参数进行输出。

模型部署则是自己编写代码实现这一完整过程，即**封装输入指令**（prompt），**选择调用模型**（包括在线和本地），**解析输出结果**，这也是模型部署的三个核心结构。

接下来就有一些问题了，输入指令和输出结果的信息格式有什么要求和特点，其它辅助的信息比如验证指令等信息如何传递。如何调用模型（模型加载）使其回答问题。

模型部署可分为**在线调用 API**和**本地加载模型**两大类，其核心区别在于模型的参数是否加载到本地机器上，即在线调用 API 只是使用了模型提供的服务接口，具体计算推导的模型参数则是在它们自己服务器上的，而本地加载模型则是把大模型的模型参数下载到了本地机器（如自己的计算机），运行时加载模型参数计算推导。

Openai 通信协议

在各种模型部署里到底扮演怎样的角色？一个标准？什么情况下用这个标准？

OpenAI 定义了一套“**如何与模型对话**”的通信协议，这套协议因其简洁和高效，成为了整个行业的通用语。这套协议包括什么，各自的功能作用是什么，如何应用，结合具体实例详细说明。

OpenAI 定义的“与模型对话”的通信协议，本质是一套**标准化的 API 交互规范**，核心围绕“如何通过网络请求与大语言模型进行对话交互”制定，包括**端点设计**、**请求 / 响应格式**、**参数规范**等。这套协议因其结构化、易扩展且适配多轮对话场景，被行业广泛采用（如开源项目、云厂商服务多兼容此规范）。

协议核心组成及功能

这套协议的核心载体是 OpenAI 的 Chat Completions API（最常用的**对话接口**），其组成可分为**端点设计**、**请求参数**、**响应格式**、**认证机制**和**会话管理**五部分，各自功能如下：

1. 端点设计（Endpoint）

端点本质上就是**接收并处理用户请求的接口地址**，但是它不仅仅是一个 URL，还包含了更多的规范，

（1）不同的任务端点的核心词是不同，

对话任务：用/v1/chat/completions 端点（如 GPT-4 聊天）；

文本嵌入：用/v1/embeddings 端点（生成文本向量，用于语义搜索）；

语音转文字：用/v1/audio/transcriptions 端点（如 Whisper 处理录音）；

图像生成：用/v1/images/generations 端点（如 DALL·E 生成图片）。

（2）不同端点需要接收的用户请求格式是不同的，比如向 “对话端点（chat/completions）” 发送请求时，Payload 必须包含 model（指定模型）和 messages（对话历史）；向 “图像生成端点（images/generations）” 发送请求时，Payload 需包含 prompt（图像描述）和 n（生成数量）：

（3）每个端点返回的响应结构是固定的，开发者需从响应中提取指定字段获取结果。例如，“对话端点” 的响应会包含 choices[0].message.content（模型生成的回答）。

其主要作用是通过不同任务的唯一端点统一**承载所有对话交互**，简化开发者记忆和调用成本。

2. 请求参数（Request Parameters）

请求参数是协议的核心，定义了“**如何向模型传递信息**”，分为**必填参数**和**可选参数**，共同控制对话的输入、生成逻辑和输出形式。

类别	参数名	功能作用	示例值
必填参数	model	指定使用的模型（如 GPT-4），决定对话能力和性能。	"gpt-4"
	messages	对话历史数组，存储多轮交互内容，模型据此生成回复（ 核心参数 ）。	见下文详细说明
可选参数	temperature	控制生成文本的随机性（0~2）：值越低越确定（如 0.1）， 越高越随机 （如 1.5）。	0.7（默认值，平衡确定性和多样性）
	max_tokens	限制生成回复的最大 token 数（含输入和输出），防止回复过长。	500

类别	参数名	功能作用	示例值
	<code>top_p</code>	核采样参数 (0~1)：控制生成时只从 概率和为 top_p 的候选 token 中选择（替代 <code>temperature</code> 的另一种控制方式）。	0.9
	<code>stream</code>	是否启用流式输出 (true/false)：true 时逐句返回结果（如聊天机器人实时显示），false 时一次性返回。	false（默认）
	<code>stop</code>	自定义停止符 ：当生成内容包含该字符串时，模型停止生成。	["###", "结束"]（遇到 ### 或 “结束” 停止）
	<code>user</code>	关联用户 ID（字符串）：用于 OpenAI 识别不同用户，防止滥用（可选）。	"user_12345"

messages 参数详解（对话核心）

`messages` 是协议中最关键的参数，以**数组形式**存储多轮对话历史，每个元素是一个**包含 role（角色）和 content（内容）的字典**，通过角色区分不同参与方：

角色（role）	功能作用	示例 content
system	设定“行为准则”或背景信息（如“你是一个专业翻译”），影响模型整体风格。	"请用简洁的语言回答，不超过 20 个字。"
user	表示用户的输入（问题、指令等），是模型需要响应的内容。	"什么是大语言模型？"
assistant	表示模型的历史回复，用于多轮对话中让模型“记住”之前的交互。	"大语言模型是能理解和生成人类语言的 AI。"

```
messages = [
    {"role": "system", "content": "请用通俗易懂的语言回答问题。"},
    {"role": "user", "content": "什么是大语言模型？"}]
```

作用：通过**累积 messages 数组**，实现“上下文记忆”，让模型能基于完整对话历史生成连贯回复（解决了早期单轮接口无法处理多轮对话的问题）。

4. 响应格式（Response Format）

模型返回的响应是**结构化 JSON**，包含**生成结果、元数据和使用统计**，方便开发者解析和处理。核心字段如下：

字段名	功能作用	示例值
<code>id</code>	本次请求的唯一标识（用于追踪）。	"chatcmpl-123456"
<code>object</code>	响应类型（固定为 <code>chat.completion</code> ）。	"chat.completion"

字段名	功能作用	示例值
created	请求创建时间（Unix 时间戳）。	1694561234
model	实际使用的模型（可能与请求的 model 略有差异，如版本号）。	"gpt-3.5-turbo-0613"
choices	生成 结果数组 （默认返回 1 个结果，可通过 n 参数指定多个）。	见下文详细说明
usage	token 使用统计（输入、输出、总消耗），用于计费和控制成本。	{ "prompt_tokens": 20, "completion_tokens": 50, "total_tokens": 70 }

choices 数组详解：

每个元素包含模型的具体回复和生成终止原因：

```
{
  "choices": [
    {
      "message": { // 模型生成的回复
        "role": "assistant",
        "content": "大语言模型是一类能理解和生成人类语言的人工智能模型。"
      },
      "finish_reason": "stop", // 生成终止原因：stop（正常结束）、length（达到 max_tokens）等
      "index": 0 // 结果索引（多结果时使用）
    }
  ]
}
```

5. 认证机制（Authentication）

- **方式：**通过 HTTP 请求头的 Authorization 字段传递 API Key，格式为 Bearer <你的 API 密钥>。
- **功能：**验证用户身份，关联账号权限和计费信息（OpenAI 通过 API Key 追踪使用量并收费）。

6. 会话管理（Session Management）

- **机制：**由于 HTTP 协议是“无状态”的，OpenAI 不会存储对话历史，需由客户端（开发者）主动维护 messages 数组，实现多轮对话。
- **流程：**每轮对话后，将用户输入（user）和模型回复（assistant）追加到 messages 数组，作为下一轮请求的输入。

协议的实际应用（结合具体实例）

以下通过“多轮对话”和“流式输出”两个典型场景，展示协议的应用方式（以 Python 代码为例）。

场景 1：基础多轮对话（询问“大语言模型”相关问题）

需求：用户先问“什么是大语言模型？”，模型回复后，用户再问“它有什么应用？”，模型基于历史对话继续回答。

步骤 1：首次请求（第一轮对话）

```
import requests

api_key = "你的 API 密钥"
url = "https://api.openai.com/v1/chat/completions"
# 首次请求的 messages: 包含 system 提示和 user 的第一个问题
messages = [
    {"role": "system", "content": "请用通俗易懂的语言回答问题。"},
    {"role": "user", "content": "什么是大语言模型？"}]

headers = {
    "Content-Type": "application/json",
    "Authorization": f"Bearer {api_key}"
}

data = {
    "model": "gpt-3.5-turbo",
    "messages": messages,
    "temperature": 0.7 # 中等随机性
}
# 发送请求
response = requests.post(url, headers=headers, json=data)
result = response.json()
# 解析模型回复
assistant_reply = result["choices"][0]["message"]["content"]
print("助手回复：", assistant_reply)
# 输出示例：大语言模型是一种能理解、生成人类语言的人工智能，像 GPT 系列就是典型代表。
```

步骤 2：第二轮对话（基于历史上下文）

将第一轮模型回复（assistant）追加到 messages，再加入用户的新问题：

```
# 更新 messages: 追加第一轮 assistant 回复和新的 user 问题
messages.append({"role": "assistant", "content": assistant_reply})
messages.append({"role": "user", "content": "它有什么应用呢？"})
# 复用相同的请求结构，仅更新 messages
data["messages"] = messages
response = requests.post(url, headers=headers, json=data)
result = response.json()
print("助手回复：", result["choices"][0]["message"]["content"])
# 输出示例：可用于聊天机器人、翻译、写文章、代码生成等，应用很广泛。
```

核心逻辑：通过维护 messages 数组，实现“上下文记忆”，模型基于完整历史生成连贯回复。

场景 2：流式输出（实时显示回复，如聊天机器人）

需求：让模型的回复逐字 / 逐句显示（类似即时通讯），提升用户体验。此时需将 stream 参数设为 true。

```
import requestsimport json

api_key = "你的 API 密钥"
url = "https://api.openai.com/v1/chat/completions"

messages = [{"role": "user", "content": "请介绍一下 OpenAI。"}]
data = {
    "model": "gpt-3.5-turbo",
    "messages": messages,
    "stream": True # 启用流式输出}

headers = {
    "Content-Type": "application/json",
    "Authorization": f"Bearer {api_key}"
}
# 发送流式请求（response 是迭代器）
response = requests.post(url, headers=headers, json=data, stream=True)
print("助手回复：", end="")
for chunk in response.iter_lines():
    if chunk:
        # 解析流式响应（需去除前缀"data: "）
        chunk_data = json.loads(chunk.decode("utf-8").replace("data: ", ""))
        # 提取当前片段的内容（若存在）
        content = chunk_data["choices"][0]["delta"].get("content", "")
        print(content, end="", flush=True) # 实时打印
```

输出效果：回复内容会逐字显示（如“OpenAI 是一家人工智能研究公司...”），而非等待完整回复生成后一次性输出，这与聊天软件的体验一致。

功能边界协议：定义“模型能做什么”

通过文档和 API 约束，明确模型的能力范围、输入输出限制，避免开发者误用或产生不合理预期。

协议类型	核心内容	功能作用
1. 模型能力清单	明确各模型支持的任务（如 GPT-4 支持多模态输入，GPT-3.5 仅支持文本）。	帮助开发者选择合适模型（如需要图像输入则必须用 GPT-4V）。
2. 输入输出限制	限制单次请求的 token 数（如 GPT-4 单次对话上限 8k/32k token）、图像尺寸（如 DALL·E 生成图片最大 1024x1024）。	防止资源滥用，确保服务稳定性。
3. 格式约束	规定输入格式（如 messages 数组必须包含 role 和 content）、输出格式（如 JSON 结构固定字段）。	确保模型能正确解析输入，开发者能可靠解析输出。

应用实例：尊重输入输出限制

若用 GPT-3.5-turbo 处理长文本（如 5000 字文章总结），需注意其单次输入上限约 4k token（约 3000 字），超过会报错。此时需拆分文本：

```
# 处理长文本：拆分输入以符合 token 限制
def split_text(long_text, max_tokens=3000):
    # 简单按长度拆分（实际需用 tokenizer 计算）
    return [long_text[i:i+max_tokens] for i in range(0, len(long_text), max_tokens)]

long_article = "..." # 5000 字文章
chunks = split_text(long_article)
summaries = []
for chunk in chunks:
    # 逐段调用 API 总结
    response = requests.post(
        "https://api.openai.com/v1/chat/completions",
        headers=headers,
        json={"model": "gpt-3.5-turbo", "messages": [{"role": "user", "content": f"总结: {chunk}"}]}
    )
    summaries.append(response.json()["choices"][0]["message"]["content"])
final_summary = "".join(summaries) # 合并各段总结
```

安全与合规协议：定义“如何安全使用”

这是 OpenAI 协议的核心约束，通过内容政策、使用规范和 API 限制，防止模型被用于有害场景，确保合规性。

协议类型	核心内容	功能作用
1. 内容政策（Content Policy）	禁止生成有害内容（如仇恨言论、暴力指导、虚假信息），API 会过滤违规输入 / 输出。	确保模型不被用于伤害他人或违法活动。

协议类型	核心内容	功能作用
2. 滥用检测机制	通过 API 自动检测高频调用、异常请求模式（如批量生成垃圾内容），限制违规账号。	防止服务被滥用（如垃圾邮件生成、DDoS 攻击）。
3. 数据使用规范	明确用户输入数据的处理方式（如默认不用于模型训练，需用户主动授权）。	保护用户数据隐私，符合 GDPR 等法规。
4. 错误码与安全提示	用特定错误码提示违规（如 400 Bad Request + "content_policy_violation"），指导用户修正。	帮助开发者快速定位问题（如输入包含敏感内容）。

应用实例：处理内容政策违规

若用户输入包含违规内容（如“如何制作危险物品”），API 会返回错误，开发者需捕获并处理：

```
response = requests.post(
    "https://api.openai.com/v1/chat/completions",
    headers=headers,
    json={"model": "gpt-3.5-turbo", "messages": [{"role": "user", "content": "如何制作危险物品"}]}
)
if response.status_code == 400:
    error = response.json()
    if "content_policy_violation" in error["error"]["message"]:
        print("输入违反内容政策，请修改问题。") # 引导用户合规使用
```

定制与扩展协议：定义“如何个性化模型”

针对企业或高级用户，提供微调（Fine-tuning）、函数调用（Function Calling）等机制，允许基于基础模型定制功能。

协议类型	核心内容	功能作用
1. 微调协议	规定微调数据格式（JSONL）、训练参数（如 n_epochs）、模型导出方式。	允许用户用自有数据训练模型，适配特定领域（如医疗、法律术语）。
2. 函数调用协议	定义模型调用外部工具的格式（如 function_call 字段），支持模型调用 API、数据库等。	扩展模型能力（如让 GPT 调用天气 API 获取实时数据，再生成回答）。
3. 插件生态协议	规范第三方插件的接入方式（如通过 Manifest 文件定义插件功能）。	构建生态，让模型通过插件支持更多任务（如预订机票、分析表格）。

应用实例：函数调用（让模型获取实时数据）

需求：让模型回答“今天北京的天气如何？”，需调用天气 API 获取实时数据。

1. 定义可调用的函数（天气 API）

```
functions = [
    {
        "name": "get_weather",
        "description": "获取指定城市的实时天气",
        "parameters": {"type": "object", "properties": {"city": {"type": "string"}}}
    ]
```

2. 向模型发送问题，提示可调用工具

```
response = requests.post(
    "https://api.openai.com/v1/chat/completions",
    headers=headers,
    json={
        "model": "gpt-3.5-turbo-0613",
        "messages": [{"role": "user", "content": "今天北京的天气如何? "}],
        "functions": functions, # 告诉模型可用的工具
        "function_call": "auto" # 让模型自动决定是否调用工具
    })
```

result = response.json()

3. 解析模型响应：若需要调用工具，则执行函数 if

```
result["choices"][0]["message"].get("function_call"):
    func_name = result["choices"][0]["message"]["function_call"]["name"]
    city =
json.loads(result["choices"][0]["message"]["function_call"]["arguments"])["city"]
```

调用天气 API（示例）

```
weather_data =
requests.get(f"https://api.weatherapi.com/v1/current.json?key=YOUR_WEATHER_
KEY&q={city}").json()
current_temp = weather_data["current"]["temp_c"]
```

4. 将工具返回结果传给模型，生成最终回答

```
messages = [
    {"role": "user", "content": "今天北京的天气如何? "},
    result["choices"][0]["message"], # 模型的函数调用指令
    {"role": "function", "name": func_name, "content":
json.dumps({"temperature": current_temp})}
]
final_response = requests.post(
    "https://api.openai.com/v1/chat/completions",
    headers=headers,
```

```
        json={"model": "gpt-3.5-turbo-0613", "messages": messages, "functions":
functions}
    )
    print("最终回答: ", final_response.json()["choices"][0]["message"]["content"])
    # 输出: 今天北京的气温是 25 摄氏度, 天气晴朗。
```

完整协议的整体作用

1. **标准化开发**: 统一的接口和格式让开发者无需重复学习不同模型的调用方式, 降低开发成本;
2. **安全可控**: 内容政策和滥用检测确保模型不被用于有害场景, 符合法律法规;
3. **能力扩展**: 微调、函数调用等协议让模型能适配特定需求, 从通用工具变为行业解决方案;
4. **生态共建**: 插件协议和开放 API 允许第三方开发者构建扩展, 形成围绕 OpenAI 的应用生态 (如 LangChain、ChatGPT 插件)。

协议如何与本地部署模型兼容

OpenAI 协议兼容, 本质上是要求你的本地模型服务提供一个与 OpenAI API 的 `/v1/chat/completions` 端点**结构和行为一致**的 API 接口。

让本地部署的模型与 **OpenAI 协议兼容**, 是一项极具价值的工作。这意味着你可以使用为 OpenAI API 设计的众多工具 (如 LangChain、LlamaIndex)、客户端和代码, 无缝地切换到你自己的模型, 从而获得**模型选择的自由**和**数据的私有化**。

关键要求包括:

相同的 HTTP 端点: 提供一个 `http://your-local-server/v1/chat/completions` 的 POST 接口。

相同的请求体 (Request Body): 能解析和处理包含 `model`, `messages`, `temperature`, `max_tokens`, `stream` 等参数的 JSON 请求。

相同的响应体 (Response Body): 返回一个包含 `choices[0].message.content` 或流式 `choices[0].delta.content` 的 JSON 响应, 其结构与 OpenAI 官方响应完全一致。

要让本地部署的模型兼容 OpenAI 协议 (即能通过 OpenAI 风格的 API 被调用), 核心是在本地模型外层封装一层“**协议转换层**”, 确保接口端点、请求

参数、响应格式与 OpenAI API 一致。这种兼容既能复用基于 OpenAI API 开发的代码（如 LangChain、聊天应用），也能降低从“调用 OpenAI 服务”迁移到“本地部署模型”的成本。

兼容的核心目标：接口对齐

无论采用何种部署方式，要兼容 OpenAI 协议，必须实现以下**核心对齐**：

1. **端点路径**：提供与 OpenAI 一致的 API 端点（如/v1/chat/completions 用于对话，/v1/embeddings 用于嵌入生成）。
2. **请求参数**：支持 OpenAI 的关键参数（如 model 指定模型名、messages 传递对话历史、temperature 控制随机性）。
3. **响应格式**：返回与 OpenAI 结构一致的 JSON（包含 choices、message、usage 等字段）。
4. **流式输出**：支持 stream=true 参数，通过 SSE（Server-Sent Events）逐段返回结果（与 OpenAI 流式输出行为一致）。

普通本地部署：手动实现协议兼容

普通本地部署（如直接用 transformers 库加载模型）需手动搭建 API 服务，封装协议转换逻辑。通常使用 FastAPI 或 Flask 作为 Web 框架，在模型调用外层增加“参数解析”和“响应格式化”步骤。

步骤 1：加载本地模型（以 Llama 2 为例）

先用 transformers 加载本地模型，确保能处理对话输入并生成回复：

```
from transformers import AutoTokenizer, AutoModelForCausalLM, pipeline
# 加载本地模型（假设模型文件在./llama-2-7b-chat 目录）
model_path = "./llama-2-7b-chat"
tokenizer = AutoTokenizer.from_pretrained(model_path)
model = AutoModelForCausalLM.from_pretrained(
    model_path,
    device_map="auto", # 自动分配到 GPU/CPU
    load_in_8bit=True  # 8bit 量化节省内存
)
# 定义对话生成函数（模拟模型推理）def generate_response(messages,
temperature=0.7, max_tokens=512):
    # 将 messages 转换为模型需要的输入格式（Llama 2 的对话格式是特定的）
    prompt = ""
    for msg in messages:
        if msg["role"] == "system":
            prompt += f'<s>[INST] <<SYS>>{msg['content']}<</SYS>> [/INST] '
        elif msg["role"] == "user":
            prompt += f'{msg['content']} [/INST] '
        elif msg["role"] == "assistant":
```

```

prompt += f'{msg["content"]} </s><s>[INST] '

# 生成回复
inputs = tokenizer(prompt, return_tensors="pt").to(model.device)
outputs = model.generate(
    **inputs,
    temperature=temperature,
    max_new_tokens=max_tokens,
    do_sample=True
)
response = tokenizer.decode(outputs[0],
skip_special_tokens=True).split("[INST]")[1].strip()
return response

```

步骤 2: 用 FastAPI 搭建兼容 OpenAI 协议的 API 服务

通过 FastAPI 创建/v1/chat/completions 端点，解析 OpenAI 格式的请求，调用本地模型，再格式化响应：

```

from fastapi import FastAPI, HTTPException
from pydantic import BaseModel
from typing import List, Optional, Dict
import uvicorn
import time

```

```

app = FastAPI()

# 定义请求体结构（对齐 OpenAI 的 Chat Completions 请求格式）
class Message(BaseModel):
    role: str # "system" "user" "assistant"
    content: str

class ChatCompletionRequest(BaseModel):
    model: str # 模型名（本地部署可忽略，或用于多模型路由）
    messages: List[Message]
    temperature: Optional[float] = 0.7
    max_tokens: Optional[int] = 512
    stream: Optional[bool] = False # 是否启用流式输出

# 定义响应结构（对齐 OpenAI 的响应格式）
class Choice(BaseModel):
    message: Message
    finish_reason: str = "stop"
    index: int = 0

class ChatCompletionResponse(BaseModel):
    id: str = f'chatcmpl-{int(time.time())}'
    object: str = "chat.completion"
    created: int = int(time.time())
    model: str
    choices: List[Choice]
    usage: Dict[str, int] = {"prompt_tokens": 0, "completion_tokens": 0,
"total_tokens": 0} # 简化处理

```

```

@app.post("/v1/chat/completions", response_model=ChatCompletionResponse)async
def chat_completions(request: ChatCompletionRequest):
    try:
        # 1. 调用本地模型生成回复
        response_text = generate_response(
            messages=request.dict()["messages"], # 传递对话历史
            temperature=request.temperature,
            max_tokens=request.max_tokens
        )

        # 2. 格式化响应（对齐 OpenAI 格式）
        return ChatCompletionResponse(
            model=request.model,
            choices=[Choice(
                message=Message(role="assistant", content=response_text)
            )]
        )
    except Exception as e:
        raise HTTPException(status_code=500, detail=str(e))
# 启动服务（默认端口 8000）if __name__ == "__main__":
    uvicorn.run(app, host="0.0.0.0", port=8000)

```

步骤 3：验证兼容性（用 OpenAI 客户端调用本地服务）

部署后，可直接使用 OpenAI 官方 SDK 调用本地服务（只需修改 base_url），无需修改原有代码：

```

python
运行
from openai import OpenAI
# 配置客户端连接本地服务
client = OpenAI(
    api_key="dummy", # 本地服务可忽略 API Key（或在服务端添加验证）
    base_url="http://localhost:8000/v1" # 指向本地 API 服务)
# 调用方式与 OpenAI 完全一致
response = client.chat.completions.create(
    model="llama-2-7b-chat", # 本地模型名（仅作标识）
    messages=[
        {"role": "user", "content": "介绍一下你自己"}
    ])print(response.choices[0].message.content)# 输出示例：我是基于 Llama 2 的
本地模型，由你部署在本地服务器上...

```

基于 FastChat 的部署：开箱即兼容 OpenAI 协议

FastChat（一款开源的 LLM 部署工具）内置了对 OpenAI 协议的兼容，无需手动编写 API 层。其核心是通过 controller（控制器）、model_worker（模型工作节点）和 openai_api_server（协议转换服务器）的架构，自动实现接口对齐。

步骤 1：安装 FastChat 并启动核心组件

```
bash
# 安装 FastChat
pip install "fschat[model_worker,webui]"
# 启动控制器（协调模型资源）
python -m fastchat.serve.controller
# 启动模型 worker（加载本地模型，需指定模型路径）
python -m fastchat.serve.model_worker \
    --model-path ./llama-2-7b-chat \ # 本地模型路径
    --device cuda \ # 运行设备（cuda/cpu）
    --load-8bit # 8bit 量化
# 启动 OpenAI 兼容的 API 服务器（默认端口 8000）
python -m fastchat.serve.openai_api_server \
    --host 0.0.0.0 \
    --port 8000
```

步骤 2：验证兼容性（与 OpenAI API 无缝对接）

FastChat 的 openai_api_server 会自动处理协议转换，支持所有 OpenAI 的核心参数（messages、temperature、stream 等）：

```
from openai import OpenAI
# 连接 FastChat 启动的 API 服务
client = OpenAI(
    api_key="EMPTY", # FastChat 默认无需验证（生产环境可添加 API Key 验证）
    base_url="http://localhost:8000/v1")
# 1. 基础对话（与调用 OpenAI 完全一致）
response = client.chat.completions.create(
    model="llama-2-7b-chat", # 必须与 worker 加载的模型名一致
    messages=[{"role": "user", "content": "1+1 等于多少?"})print("普通响应: ",
response.choices[0].message.content) # 输出: 2
# 2. 流式输出（stream=true）
stream = client.chat.completions.create(
    model="llama-2-7b-chat",
    messages=[{"role": "user", "content": "介绍一下北京"}],
    stream=True)print("\n 流式响应: ")for chunk in stream:
    if chunk.choices[0].delta.content:
        print(chunk.choices[0].delta.content, end="", flush=True)# 输出: 北京是中
国的首都，位于华北平原北部...（逐字显示）
```

两种部署方式的对比

维度	普通本地部署（手动实现）	基于 FastChat 的部署
兼容性	需手动实现所有参数和格式（易遗漏边缘 case）	开箱即支持 OpenAI 协议的核心功能（经社区验证）
开发成本	高（需编写 API 层、参数解析、响应格式化代码）	低（一行命令启动兼容服务）
灵活性	高（可定制参数、添加自定义逻辑）	中（需遵循 FastChat 的扩展方式）
支持的模型	需手动适配模型输入格式（如 Llama、GPT-2 的格式差异）	自动适配主流模型（内置多种模型的输入处理）
流式输出 / SSE	需手动实现 SSE 协议	内置支持，与 OpenAI 行为一致
适用场景	定制化需求高（如添加特殊权限验证、日志记录）	快速部署、追求兼容性（如接入 LangChain 生态）

如何区分是兼容协议还是调用 openai 服务

区分一个服务是“兼容 OpenAI 协议的本地 / 第三方服务”还是“真正的 OpenAI 官方服务”，可以通过**接口特征**、**网络请求**、**功能细节**和**辅助验证**四个维度进行判断，以下是具体方法和实例：

核心区分方法：接口地址（最直接）

OpenAI 官方服务和兼容协议的服务，最显著的区别是 **API 请求的基础地址（base_url）**：

- **OpenAI 官方服务**：固定使用 OpenAI 的域名，基础地址为 `https://api.openai.com/v1`（所有接口均以此为前缀）。
- **兼容协议的服务**：使用非 OpenAI 域名的地址，通常是本地地址（如 `http://localhost:8000/v1`）或第三方域名（如 `https://api.thirdparty.com/v1`）。

二、辅助判断：API Key 验证

OpenAI 官方服务对 API Key 有严格验证，而兼容服务通常不依赖 OpenAI 的 Key：

- **OpenAI 官方服务**：必须使用从 OpenAI 官网（<https://platform.openai.com/account/api-keys>）获取的有效 API Key，否则会返回 401 Unauthorized 错误（提示“Invalid API Key”）。
- **兼容协议的服务**：
 - 本地部署的服务：可能完全不需要 API Key（如 FastChat 默认跳过验证），或使用自定义的 Key（与 OpenAI 的 Key 格式 / 来源无关）。

- 第三方兼容服务：使用该第三方平台提供的 API Key（非 OpenAI 官网的 Key）。

三、功能细节差异：模型列表与行为特征

OpenAI 官方服务和兼容服务在支持的模型和生成行为上存在差异：

1. 支持的模型名称

- **OpenAI 官方服务：**仅支持 OpenAI 发布的模型（如 gpt-3.5-turbo、gpt-4、text-embedding-ada-002 等），调用不存在的模型会报错。
- **兼容协议的服务：**支持的是本地部署或第三方的模型（如 llama-2-7b、vicuna-13b、internlm-7b 等），通常不支持 OpenAI 的官方模型名（除非特别配置）。

2. 生成行为与限制

- **OpenAI 官方服务：**有严格的 token 限制（如 gpt-3.5-turbo 单轮对话上限 4k token）、内容政策过滤（对敏感内容的响应高度一致）。
- **兼容协议的服务：**token 限制由本地模型决定（如 Llama 2-7b 通常支持 2k-4k token），内容过滤规则可能宽松或自定义（响应敏感内容的行为与 OpenAI 不同）。

四、网络与日志验证：请求跟踪

通过查看网络请求的目标 IP 或日志记录，可进一步确认服务类型：

1. 查看网络请求目标

- OpenAI 官方服务的请求会发送到 OpenAI 的服务器 IP（可通过 ping api.openai.com 获取）。
- 本地兼容服务的请求会发送到本地 IP（如 127.0.0.1）或局域网内的其他设备 IP。
- 第三方兼容服务的请求会发送到该第三方的服务器 IP（可通过 ping 其域名确认）。

2. 查看服务日志

- 本地部署的兼容服务（如 FastChat）会在控制台输出请求日志（包含模型加载路径、本地处理记录）。
- OpenAI 官方服务无本地日志（请求在远程服务器处理）。

五、总结：区分步骤

1. **优先检查 base_url：**若为 https://api.openai.com/v1 则是官方服务；否则为兼容服务。
2. **验证 API Key：**用无效 Key 测试，报错 “Invalid API Key” 则可能是官方服务；能正常调用则是兼容服务。

3. **测试模型与行为**: 调用 OpenAI 独有的模型（如 gpt-4），或测试敏感内容响应，行为不一致则为兼容服务。
4. **跟踪网络请求**: 查看请求的目标 IP 或本地服务日志，确认是否为本地 / 第三方服务器。

通过以上方法，可准确区分两者 —— 核心是**接口地址**和**服务依赖的 API Key 与模型**，这两个特征具有决定性。

输入指令

首先来看一下大模型的输入指令，访问官网和 app 时只需要输入问题即可，但是模型部署时的输入指令还需要额外的补充信息以及封装为特定的格式。

输入指令需要包含哪些呢？首先是核心内容，也就是需要询问的问题，然后是辅助信息，比如访问的 api 信息，最后就是封装成 json 形式的数据结构。

调用模型可以包括 2 大类，一是**在线调用 api**，二是**本地模型加载**，在线调用是利用大模型提供的 api 访问官网服务器，服务器接收到信息之后回复信息。本地模型加载则是把大模型参数直接保存到本地，使用时加载模型参数，直接在本地回复。根据编码的形式不同又可以区分为**直接部署**和**框架部署**（如 fastchat, vllm）。

先来看一下在线调用 api,直接调用就是设置输入信息 payload，发送请求，得到回复，这里调用大模型就是把 api-key, api-url 等信息编码一块发送给大模型服务器，得到回复即可，简单直观，但是安全性抵，直接暴露密码等信息。

进一步地，可以借助 fastchat 等部署框架，将大模型封装起来，既可以保证信息安全，又可以合理分配资源，还可以扩展多个大模型（工厂模式）。fastchat 的三大核心是服务器，控制器，模型工作器，**服务器接收用户的输入请求并转发给控制器**，控制器根据不同的调度方法将特定的工作器地址转发给服务器，服务器将用户查询信息发送给工作器，控制器。

进一步地，可以将这三个核心组件与 fastapi 结合包装起来，以进程的方式启动，还可以实现模型工作器的切换。下面分别具体介绍一下。

在线调用 API

前面已经讲过模型部署的三个核心结构，以通义千问为例，结合具体代码看

一下三者是如何实现的。

导入必要库

import requests # 用于发送 HTTP 请求

import json # 用于处理 JSON 数据

封装输入指令：

输入指令需要包含哪些呢？首先是核心内容，也就是需要询问的问题，然后是辅助信息，比如访问的 api 信息，最后就是封装成 json 形式的数据结构。

选择调用模型：

封装好输入指令，需要将输入发送到对应的模型服务器上，在线调用使用的是 request 模块，将 api 地址，封装的输入指令以 post 请求的形式发送并得到响应。

解析输出结果

得到的响应一般是 JSON 格式封装起来的信息，需要将其解析提取回答的文本，这才是想要的最终结果。

具体到实例中，包括下面几点：

封装输入指令：

#设置 API 的 URL（服务地址）和 API Key。

API_URL =

"https://dashscope.aliyuncs.com/api/v1/services/aigc/text-generation/generation" #

通义千问 API 服务地址

API_KEY = "sk-6fed683719dd4c179c37e331dcc9e330" # 替换为实际 API 密钥
(从阿里云控制台获取)

def qwen_max_chat(prompt: str) -> str:

"""

调用通义千问 Max 生成回复

:param prompt: 用户输入的提示词

:return: 模型生成的回复内容

"""

构造请求头（headers）嵌套字典，包含 Content-Type 和 Authorization，认证信息就是前面的 api key。 声明数据格式为 JSON。

headers = {

 "Content-Type": "application/json", # 声明请求内容格式为 JSON

 "Authorization": f"Bearer {API_KEY}" # 携带 API 密钥进行认证

}

构造请求体（payload），嵌套字典，包含模型名称及版本、输入内容、参数等。

```
payload = {
    "model": "qwen-max", # 指定模型版本
    "input": { "messages":

        #注意这里的 messages 是个列表，里面存储用户角色和输入内容。

        [ { "role": "user", # 用户角色
            "content": prompt # 用户输入内容 } ] },
    "parameters": { "temperature": 0.8, # 控制生成随机性(0-1), 值越大输出
                    # 越多样化 (0 为确定性输出)
                    "max_tokens": 1024 # 限制生成最大长度 }
    }

try:
    # 选择调用模型发送 POST 请求
    response = requests.post(
        API_URL, #API 服务地址
        headers=headers,
        data=json.dumps(payload) # 将字典转换为 JSON 字符串
    )
    if response.status_code == 200: # 检查 HTTP 状态码
        # 解析输出结果，将 API 返回的 JSON 数据解析为 Python 字典
        result = response.json()
        # print("[DEBUG] 完整响应内容:", json.dumps(result, indent=2)) #
        # 提取生成的回复内容
        return result["output"]["text"] # 替换原有的 choices 路径
    else:
        return f"请求失败，状态码：{response.status_code}"

except Exception as e:
    return f"发生异常：{str(e)}"

# 测试调用
if __name__ == "__main__":
    user_input = "如何用 Python 实现快速排序？"
    print("用户提问：", user_input)
    print("AI 回复：", qwen_max_chat(user_input))
```

总的来说，在线调用 API 部署模型的核心步骤如下：

1. 导入必要的库，比如 requests 和 json。

#设置 API 的 URL 和 API Key（需要实时更新）。

2. 构造请求头（headers={}, 嵌套字典），包含 Content-Type 和 Authorization，认证信息就是前面的 api key。

4. 构造请求体（payload={}, 嵌套字典），包含模型名称、输入内容、参数等。

5. 发送 POST 请求返回响应。

6. 处理响应，提取生成的文本。

简单的在线调用和直接访问模型差别不大，不需要对模型本身做任何的处理，而把关注点放到了输入输出信息的封装和解析，和模型的交互只在于发送输入请求（绑定 api-key 和模型地址 url）并接收响应，最终作为组成部分嵌入到整体项目中（比如前端的搭建，数据的预处理等）。

其中在实际应用中值得注意的是信息数据的格式，一般是嵌套字典，不过尤其是回复的数据格式有可能会变化，比如字典的键更改，这时可以查看官方最新文档查询格式，或者 Print 出来查看完整响应内容。

```
result = response.json()
print("[DEBUG] 完整响应内容:", json.dumps(result, indent=2))
```

在线调用大模型主要有以下几种方式：通过云服务提供的 API（比如 OpenAI 的 GPT-4、ChatGLM 的 API）、使用 Hugging Face 的 Inference API，或者自己搭建的模型服务（比如通过 Flask 或 FastAPI 部署的模型）。而 Hugging Face Transformers 库通常用于本地加载模型，但也可以间接通过 Hugging Face Hub 的 API 在线调用模型，比如使用 pipeline 或 HuggingFaceHub 类。

本地部署模型

本地部署模型分为 huggingface 加载和本地文件加载，其本质都是利用 transformers 等模型库提供的方法根据参数加载模型，huggingface 需要的参数是模型名称，本地文件加载需要的参数是模型路径。

对比维度	Hugging Face Hub 加载	本地加载
模型来源	从 Hugging Face 服务器下载	读取本地磁盘文件
网络依赖	✓ 首次加载需要联网	✗ 完全离线可用

对比维度	Hugging Face Hub 加载	本地加载
加载速度	依赖网络速度（首次较慢）	直接读取（通常更快）
代码写法	使用模型仓库名称字符串	使用本地文件系统路径
版本控制	可指定版本标签（如@main）	依赖本地文件手动更新
存储位置	默认缓存到~/.cache/huggingface/hub	用户自定义路径
模型更新	自动获取最新版本	需手动替换本地文件
错误处理	需处理网络超时/模型不存在错误	需处理文件路径错误或文件损坏

transformers 库

Transformers 库本身并不直接存储模型权重、数据或词汇表，而是一个工具框架，核心作用是统一预训练模型的加载、使用和扩展接口。它包含的核心内容是：

预训练模型的架构实现

针对各种 Transformer 变体（如 BERT、GPT、T5、RoBERTa 等），库中内置了对应的模型类（如 BertModel、GPT2LMHeadModel），这些类定义了模型的网络结构（层、注意力机制、激活函数等），但不包含具体的预训练权重。所谓模型类定义了网络结构，就是在自己手写网络时的 model 部分，包括 init 和 forward 部分，init 定义网络结构，比如注意力机制的头数隐藏层维度、线性层的输入输出维度，激活函数类别等。

Tokenizer 工具类

针对不同模型配套的分词逻辑（如 BPE、WordPiece 等），实现了对应的 Tokenizer 类（如 BertTokenizer、GPT2Tokenizer），将原始文本转换为模型可以理解的输入 ID（input_ids）和注意力掩码（attention_mask）等张量，处理不同模型的特殊 token（如 [CLS], [SEP]），自动处理填充（Padding）和截断（Truncation）以保证输入长度一致，但不包含具体的词汇表（vocab）数据。

加载与交互接口

提供了统一的加载函数（如 from_pretrained()）、推理接口（如 generate()）、训练辅助工具（如 Trainer）等，用于衔接“模型架构”与“预训练权重 / 词汇表”，并简化模型的使用流程。

无论是从 Hugging Face Hub（云端）下载模型，还是从本地文件加载，本质都是将“预训练权重 / 词汇表”与“库中定义的模型架构 / Tokenizer 逻辑”绑定。而 Transformers 库的方法（如 from_pretrained()）是实现这一绑定的“桥梁”

统一解析不同来源的文件结构

预训练模型的文件（权重、配置、词汇表）无论存放在 Hub 还是本地，其内部格式（如 .bin 权重文件、config.json 配置文件、vocab.txt 词汇表）是由模型训练时的规范定义的。库中的 from_pretrained() 方法会自动识别这些文件的结构，解析权重维度、配置参数（如隐藏层维度、注意力头数），并与库中对应的模型架构匹配。

例如：加载 bert-base-uncased 时，该方法会自动找到 pytorch_model.bin（权重）、config.json（架构参数），并将权重加载到 BertModel 类定义的网络结构中。

处理兼容性与版本适配

不同模型可能基于不同框架（PyTorch、TensorFlow）训练，或有版本迭代（如 v1、v2）。库的方法会自动处理框架转换（如 from_pretrained(..., from_tf=True)）、版本兼容，确保权重能正确映射到当前库的模型类中。

简化复杂逻辑的封装

加载模型不仅是“读取文件”，还涉及设备分配（CPU/GPU）、权重精度转换（如 float32/float16）、缓存管理（本地缓存已下载模型文件）等。这些逻辑被封装在库的方法中，用户无需手动处理。

Huggingface 加载时输入的模型名称（如 bert-base-uncased）同本地加载时输入的模型路径一样，本质是绑定了模型的资源如权重文件，配置文件，告诉加载方法如 from_pretrained() 这些模型资源在哪里获取，然后与加载模型资源的模型类名如 BertModel 的模型结构进行匹配，Tokenizer 工具类同理。

除了具体的模型类名如 BertModel、GPT2LMHeadModel，BertTokenizer、GPT2Tokenizer，transformers 库还提供了 auto+ 系列（）的类用于灵活的加载各种模型，其特点是不需要知道确切的模型类名，只需要提供模型名称或路径以及模型属于的大类（预训练模型，tokenizer 工具类或者特定的因果推理类等），就会自动从配置文件推断并选择正确的模型类。这样便可以使用同一套代码加载不同的模型，只需要改变模型名称即可（在配置文件里更改）。

常用的 auto+ 有下面几个类：

1. AutoTokenizer：用于加载与模型对应的各种分词器的类，根据模型名称自动转化为对应的特定处理的类如 BertTokenizer、GPT2Tokenizer。

2. AutoModel：用于加载模型主干（Backbone），返回一个基础模型对象，没有特定的任务头（Task Head）。用于获取最后一层隐藏状态（Last Hidden States），通常是高维张量，用于特征提取或作为其他任务的输入，需要进一步处理才能用于具体任务（如分类、相似度计算）。

3. AutoModelForCausalLM：加载带任务头的模型，专门用于因果语言建模（Causal Language Modeling）。当你需要做文本生成时（如：GPT、LLaMA 系列），必须使用这个或类似的类。它在 AutoModel 的基础上，额外添加了一个语言模型头（LM Head），这是一个线性层，用于将隐藏状态映射到词汇表上，以预测下一个 token 的概率。输出包含 logits，其形状为 (batch_size, sequence_length, vocab_size)，可以直接用来计算损失或生成下一个词。

AutoModelFor...：这是一个命名范式，For 后面跟着特定任务，表示这是一个为该任务设计好的完整管道。

CausalLM：因果语言模型。之所以叫“因果”，是因为它的注意力机制是掩码的，每个 token 只能关注它之前的 token（即“因”决定“果”），这是生成式模型的典型特征。

transformers 库中其他常见的 AutoModelFor... 类：

AutoModelFor**SequenceClassification**：用于文本分类（如情感分析）。

AutoModelFor**QuestionAnswering**：用于问答任务（如 SQuAD）。

AutoModelFor**MaskedLM**：用于掩码语言建模（如 BERT 的预训练任务）。

AutoModelFor**TokenClassification**：用于令牌级分类（如命名实体识别 NER）。

在线加载（HuggingFace Hub）

```
from transformers import AutoTokenizer, AutoModel, AutoModelForCausalLM
# 加载分词器
tokenizer = AutoTokenizer.from_pretrained("bert-base-uncased")
# 加载基础模型
model = AutoModel.from_pretrained("bert-base-uncased")
# 加载带任务头的模型
causal_model = AutoModelForCausalLM.from_pretrained("gpt2")
```

本地加载

即先加载**分词器**模型，然后加载**生成模型**，将文本给到分词器模型得到**输入**，将输入给到生成模型 model.generate 生成回答。

```
from transformers import AutoTokenizer, AutoModelForCausalLM

# 加载对话模型（无需池化层）
model_path = "./your-chat-model"
tokenizer = AutoTokenizer.from_pretrained(model_path, trust_remote_code=True)
model = AutoModelForCausalLM.from_pretrained(model_path,
trust_remote_code=True)

# 文本生成示例
input_text = "如何学习 Python? "
inputs = tokenizer(input_text, return_tensors="pt")
outputs = model.generate(**inputs, max_length=100)
response = tokenizer.decode(outputs[0], skip_special_tokens=True)
print(response)
```

sentence-transformers 库

这个库建立在 **transformers** 之上，专门为了一个目标：从句子和段落中获取高质量的语义嵌入向量（**Embeddings**）。它提供的是“开箱即用”的解决方案。通过 `pip install sentence-transformers` 安装的就是它。

该库中的**核心类** `SentenceTransformer`：用于加载和运行专门为生成句子嵌入而设计的模型。内置了**智能的池化策略**（如均值池化 `mean_pooling`），能将 token 级别的输出高效地聚合为句子级别的表示。自己用 `AutoModel` 实现这个效果需要写不少代码且效果可能不好。

Hugging Face Hub 上有大量使用该库框架**微调好的模型**（如 **all-MiniLM-L6-v2**, `all-mpnet-base-v2`），它们在语义相似性等任务上表现极佳。

输入一个句子（或一个句子列表），输出一个固定大小的向量（例如 384 维、768 维）。

在线加载（HuggingFace Hub）

```
from sentence_transformers import SentenceTransformer

# 指定 Hugging Face 模型名称 (str 类型)
# 格式可以是仓库路径或简称，自动从网络下载
model_name = "sentence-transformers/paraphrase-multilingual-MiniLM-L12-v2"

# 在线加载模型 (自动下载到缓存目录)
# 首次使用需要联网，后续使用会读取缓存
model = SentenceTransformer(model_name) # 返回 SentenceTransformer 对象

# 输入文本
text = "自然语言处理很有趣"

# 编码流程与本地加载完全相同
embedding = model.encode(text) # 输出 numpy 数组
```

本地加载

法 1：直接加载句子嵌入模型

```
from sentence_transformers import SentenceTransformer
# 加载本地保存的句子转换模型
local_model_path = "./models/all-MiniLM-L6-v2"
model = SentenceTransformer(local_model_path)
```


法 2：词向量化模型和池化模型组合：

```
from sentence_transformers import SentenceTransformer, models
# 直接指定本地模型路径
model_path = r"E:\lmmmodels\allMiniLML6v2"
word_embedding_model = models.Transformer(model_path)#向量化模型
#根据向量化模型的向量维度维护池化模型
pooling_model = models.Pooling(word_embedding_model.get_word_embedding_dimension())
#将向量模型和池化模型组合成最终的句子模型
sentence_transformer_model = SentenceTransformer(modules=[word_embedding_model, pooling_model])
```

总结与如何选择

需求场景	推荐工具	原因
文本生成（如写故事、代码补全）	transformers.AutoModelForCausalLM	自带生成头，专为预测下一个 token 设计。
文本分类/问答/NER	transformers.AutoModelForSequenceClassification/QA/TokenClassification	自带任务头，输出直接用于特定任务。
特征提取（获取模型的原始输出，用于自定义任务）	transformers.AutoModel + AutoTokenizer	最灵活，可以获取隐藏状态等中间结果。
语义相似度、语义搜索、为句子生成嵌入向量	sentence_transformers.SentenceTransformer	开箱即用，API 简单，性能经过专门优化，效果最好。

其他加载方式

使用 pipeline 简化加载

```
python
from transformers import pipeline
# 在线加载
classifier = pipeline("sentiment-analysis",
model="distilbert-base-uncased-finetuned-sst-2-english")
# 本地加载
local_classifier = pipeline("sentiment-analysis", model="./models/distilbert-sst-2")
```

问题探究

模型类名和模型名称有什么区别

模型类名（如预训练的 BertModel, Tokenizer 类的 BertTokenizer）是具体的 python 类，定义了模型的网络结构（如注意力层的头数，隐藏层维度，线性层维度等）或该类型 tokenizer 的功能逻辑（如分词规则，id 生成方式，特殊符合处理），但不包含具体的权重文件以及词汇表等。

模型类本身是“通用模板”，不包含具体的结构细节（如“到底有多少层”），其最终实例化的网络结构完全由配置文件中的参数决定。例如，BertModel 类的代码中并不会硬编码“12 层”或“24 层”，而是在初始化时读取配置文件的 num_hidden_layers 参数。

通过更换配置文件，同一模型类可实例化出不同结构的模型（如 base/large 版本），无需修改代码。

而模型名称（如 bert-base-uncased）则和本地模型路径一样，在 huggingface 上是预训练模型 / Tokenizer 的唯一标识，对应着模型的具体资源如模型权重，配置文件，词汇表等。

而 transformers 库提供的加载方法如 from_pretrained() 等则将二者绑定在一起，即将模型的权重与对应的模型结构对应起来，这样便可以使用加载好的模型进行推理输出了。

实现同一任务的模型类有多种，而模型名称需要和对应的模型类名一一对应，例如，bert-base-uncased 是 BERT 系列模型，因此必须使用 BertTokenizer 类来加载（因为其文本处理规则与 BERT 的训练方式匹配）；gpt2 是 GPT 系列模型，必须使用 GPT2Tokenizer 类。

1. 直接使用具体 Tokenizer 类（需指定模型名称）

```
from transformers import BertTokenizer, GPT2Tokenizer
# 加载 BERT 的 Tokenizer: 类是 BertTokenizer, 模型名称是 bert-base-uncased
bert_tokenizer = BertTokenizer.from_pretrained("bert-base-uncased")
# 加载 GPT2 的 Tokenizer: 类是 GPT2Tokenizer, 模型名称是 gpt2
gpt2_tokenizer = GPT2Tokenizer.from_pretrained("gpt2")
```

2. 错误示例（类名与模型名称不匹配）

如果用错误的类加载模型名称，会导致预处理逻辑与模型不兼容：

```
# 错误：用 BertTokenizer 加载 GPT2 的模型名称
wrong_tokenizer = BertTokenizer.from_pretrained("gpt2") # 运行时会报错或生成无效输入
```

原因：GPT2 的词汇表和 BERT 的分词逻辑不匹配（如 GPT2 没有[CLS]符号，而 BertTokenizer 会强制添加）。

transformers 库提供的 auto+ 系列可以更灵活的绑定模型类名和模型名称，只需要模型类名在大类上和模型名称保持一致，如 AutoTokenizer 会根据模型名称自动匹配对应的 Tokenizer 类，无需手动指定类名，

```
from transformers import AutoTokenizer
```

```
# 自动匹配类：模型名称是 bert-base-uncased → 自动使用 BertTokenizer
bert_tokenizer = AutoTokenizer.from_pretrained("bert-base-uncased")
# 自动匹配类：模型名称是 gpt2 → 自动使用 GPT2Tokenizer
gpt2_tokenizer = AutoTokenizer.from_pretrained("gpt2")
```

使用这些模型库加载模型的背后步骤流程是什么？加载的内容包括哪些？和普通的 `load_dict` 加载权重的联系和区别是什么？

模型加载的核心是根据配置文件将各模型权重与网络层结构对应起来，这样才是加载好了模型，用于进行推理或文本词元化处理等操作。

一般来说，模型库的加载流程为：是“解析来源→加载配置→构建结构→加载权重→验证兼容”的全流程自动化，核心是“无需手动定义模型结构”。

加载的内容不仅包括权重，还包括配置文件、预处理规则等完整依赖，确保与模型结构对应起来。

与 `load_state_dict` 的关系：前者是后者的扩展——`from_pretrained` 底层也会调用 `load_state_dict` 加载权重，但额外封装了结构构建、配置处理等功能。

模型库加载模型的底层步骤流程（以 Hugging Face transformers 为例）

当调用 `AutoModel.from_pretrained(模型名称/路径)` 或 `SentenceTransformer(模型名称/路径)` 时，背后经历了一系列自动化步骤，确保模型“可直接使用”。以 transformers 的 `from_pretrained` 为例，核心流程如下：

1. 解析输入：根据模型名称 / 路径找到对应的权重等资源的位置，在 huggingface 对应仓库或者本地缓存，或者路径下的文件

- 若输入是模型名称（如 `bert-base-uncased`）：自动定位到 Hugging Face Hub 的对应仓库（<https://huggingface.co/bert-base-uncased>），检查是否有本地缓存（默认路径 `~/.cache/huggingface/hub`），若缓存存在则直接使用，否则从 Hub 下载。
- 若输入是本地路径（如 `./local_bert`）：直接读取路径下的文件，无需网络请求。

2. 加载配置文件（核心步骤）

从模型仓库 / 本地路径中读取配置文件（如 `config.json`），解析模型的核心结构参数，包括：

- 模型类型（如 BERT、GPT2）；
- 网络结构参数（隐藏层维度 `hidden_size`、层数 `num_hidden_layers`、注意力头数 `num_attention_heads` 等）；
- 任务相关参数（如分类任务的标签数量 `num_labels`）。

这些参数是**动态构建模型结构**的依据（无需用户手动定义模型类）。

3. 构建模型结构

根据配置文件中的“模型类型”和“结构参数”，自动初始化对应的模型类（如 BertModel、GPT2LMHeadModel），生成一个**空的模型框架**（仅包含网络层，**未加载权重**）。

例如：若 config.json 中 **model_type: "bert"**，则**自动实例化 BertModel 类**，并根据 hidden_size=768、num_hidden_layers=12 等参数创建对应的 Transformer 层。

4. 加载权重文件

从模型仓库 / 本地路径中读取**权重文件**（如 pytorch_model.bin 或分片文件 pytorch_model-00001-of-00002.bin），解析为**权重字典**（**键为参数名，值为张量**）。

加载时会自动处理：

- 权重与模型结构的匹配（检查参数名、维度是否一致）；
- 设备映射（默认加载到 CPU，可通过 device_map 参数指定 GPU）；
- 数据类型转换（如从 float32 转为 float16，通过 torch_dtype 参数控制）。

5. 验证与兼容处理

- 若权重文件中的参数与模型结构不匹配（如少了一层、维度不符），会抛出明确的错误提示（如 “missing keys”“size mismatch”）；
- 对旧版本模型权重进行自动兼容转换（如早期 BERT 的权重名称与当前库不匹配时，自动映射到新名称）。

6. 附加组件加载（针对 Tokenizer/Processor）

若同时加载 Tokenizer（如 AutoTokenizer.from_pretrained），还会读取词汇表文件（vocab.txt）、合并规则（merges.txt）等，初始化分词逻辑（如子词分割、特殊符号<CLS>处理）。

步骤总结：

输入（名称/路径）→ 解析来源（Hub/本地）→ 加载配置文件 → 构建模型结构 → 加载权重文件 → 验证匹配 → 输出可用模型

加载的内容具体包括哪些？

模型库加载的内容远不止“权重”，而是一套**完整的模型运行依赖项**，主要包括 3 类核心文件：

类型	典型文件	作用
配置文件	config.json、generation_config.json	存储 模型结构参数 （如层数、维度）、生成策略参数（如最大长度、采样温度），用于构建模型结构。
权重文件	pytorch_model.bin、tf_model.h5	存储模型的参数张量（如权重、偏置），决定模型的“知识”。
预处理相	vocab.txt（词汇表）、merges.txt	用于 Tokenizer，定义文本→ID 的

类型	典型文件	作用
关文件	(合并规则)、 tokenizer_config.json	映射规则 (如 “我爱 NLP”→[213, 456, 789])。
元数据文件	README.md、modelcard.json	描述模型的训练数据、用途、性能等 (非运行必需, 但辅助用户理解模型)。

例如, 加载 bert-base-uncased 时, 实际读取的核心文件包括:
config.json (结构参数) + pytorch_model.bin (权重) + vocab.txt (词汇表) + tokenizer_config.json (分词配置)。

与普通 load_state_dict 加载权重的联系和区别

1. 核心联系

两者的**最终目标一致**: 将**权重张量加载到模型的网络层**中, 使模型具备推理 / 训练能力。

底层都依赖 “**参数名匹配**”: 权重字典的键 (如 bert.encoder.layer.0.attention.self.query.weight) 必须与模型结构中定义的参数名完全一致, 否则无法正确加载。

2. 关键区别

维度	模型库加载 (如 from_pretrained)	普通 load_state_dict 加载
模型结构来源	自动根据配置文件构建 (无需用户手动定义模型类)。	需用户 手动定义模型结构 (如 class MyModel(nn.Module): ...), 且结构必须与权重匹配。
加载内容	同时加载配置文件、权重、预处理文件等完整依赖。	仅加载权重字典 (.pth 或 .bin 文件), 不处理配置或预处理逻辑。
自动化程度	高: 自动处理缓存、版本兼容、设备映射、参数校验等。	低: 需手动处理设备映射 (如 model.load_state_dict(torch.load(path, map_location='cpu')))、参数校验 (需手动对比 model.state_dict().keys() 与权重 keys)。
适用场景	快速使用预训练模型 (尤其是第三方公开模型)。	加载自定义训练的模型 (已知模型结构), 或需要精细控制加载过程 (如部分层加载)。
错误处理	提供详细的不匹配提示 (如 “missing keys in state_dict”)。	仅抛出基础错误 (如 “KeyError”), 需用户手动排查不匹配原因。

代码示例对比

```
# 示例 1: Hugging Face from_pretrained (自动构建结构+加载权重)
from transformers import AutoModel
model = AutoModel.from_pretrained("bert-base-uncased") # 一步完成结构构建+权重加载
```

```

# 示例 2: 普通 load_state_dict (需手动定义结构)
import torch
import torch.nn as nn
# 1. 手动定义与权重匹配的模型结构 (必须与 bert-base-uncased 一致, 否则加载失败)
class MyBERT(nn.Module):
    def __init__(self):
        super().__init__()
        self.embeddings = ... # 需复现 BERT 的嵌入层
        self.encoder = ...    # 需复现 BERT 的 12 层 Transformer
        # ... (省略大量细节, 实际需完全匹配原模型结构)

    def forward(self, x):
        return self.encoder(self.embeddings(x))
# 2. 加载权重 (假设已下载 bert-base-uncased 的 pytorch_model.bin)
model = MyBERT()
model.load_state_dict(torch.load("pytorch_model.bin")) # 仅加载权重, 需结构完全匹配

```

实际开发中, 优先使用模型库的高层接口 (如 `from_pretrained`), 除非需要完全自定义模型结构 (此时用 `load_state_dict` 更灵活)。

加载时模型参数的类型和传递方式

加载模型时, 让模型能够正确运行的核心配置可分为**模型结构配置**、**运行环境配置**、**功能开关配置**和**分布式 / 部署配置**四大类。这些配置的传递方式不同, 部分来自模型自带的 `config.json`, 部分则是加载时通过参数手动指定。以下是详细解析:

核心配置分类及作用

1. 模型结构配置 (决定“模型长什么样”)

这是模型能够运行的**基础配置**, 定义了模型的**网络架构** (如层数、维度、注意力机制等), 确保权重能正确匹配到网络层。一般在 `config.json` 文件里加载传递。

核心内容:

- 模型类型 (`model_type`): 如 `bert`、`gpt2`、`llama`, 决定加载哪种网络结构类;
- 结构参数: 隐藏层维度 (`hidden_size`)、层数 (`num_hidden_layers`)、注意力头数 (`num_attention_heads`)、词汇表大小 (`vocab_size`) 等;
- 任务相关参数: 分类任务的标签数 (`num_labels`)、生成任务的最大长度 (`max_length`, 部分模型放在 `generation_config.json`)。

作用：确保模型加载时能动态构建与预训练权重匹配的网络结构（如 BERT 的 12 层 Transformer 与权重维度必须一致）。

2. 运行环境配置（决定“在什么硬件上跑”）

控制模型运行的硬件、数据类型等，直接影响性能和兼容性。

核心内容：

- 设备类型（device）：CPU、单 GPU（如 cuda:0）、多 GPU（如 auto 自动分配）；
- 数据类型（dtype）：float32（默认，精度高）、float16/bfloat16（精度稍低，速度快、省内存）；
- 量化配置：是否启用 INT8/INT4 量化（load_in_8bit/load_in_4bit），降低内存占用（牺牲部分精度）。

作用：让模型适配硬件能力（如 GPU 显存不足时用量化或 float16），避免“硬件不支持”错误（如模型放 GPU 但无 CUDA 环境）。

3. 功能开关配置（决定“模型怎么跑”）

控制模型的运行模式（训练 / 推理）、特殊功能（如 dropout）等，影响输出结果的一致性。

核心内容：

- 推理模式（eval()）：禁用 dropout 等训练时的随机层，确保推理结果稳定；
- 注意力掩码（attention_mask）：控制模型对 padding 部分的忽略（通过 Tokenizer 自动生成，或手动指定）；
- **生成策略参数**：如 do_sample（是否采样）、temperature（采样随机性）、num_beams（beam search 数量），影响文本生成质量。

作用：确保模型在特定任务下的行为正确（如推理时不启用训练专属层）。

4. 分布式 / 部署配置（决定“多设备怎么协同跑”）

针对多卡训练、分布式部署或服务化场景，控制设备间通信和资源分配。

核心内容：

- 通信地址（master_addr）：分布式训练中主节点的 IP 地址；
- 通信端口（master_port）：主节点的端口，用于设备间数据传输；
- 进程编号（rank）：多进程中当前进程的 ID（区分不同设备）；
- 设备映射（device_map）：多 GPU 场景下指定每层参数放在哪个 GPU（如 {"": "cuda:0"} 指定所有层放第 0 张卡）。

作用：在多设备环境中协调资源，避免冲突（如端口占用），实现并行计算。

配置的传递方式（如何给到模型？）

不同配置的传递路径不同，并非全部来自 config.json，而是“静态配置文件 + 动态参数”结合：

1. 模型结构配置：主要来自 config.json

config.json 是模型自带的核心配置文件（通常与权重文件存放在一起），加载时由模型库自动读取：

当调用 AutoModel.from_pretrained(模型名称/路径)时，库会自动解析路径下的 config.json，提取 model_type、hidden_size 等参数，动态构建对应的模型类（如 BertModel）。

若需修改结构参数（如增加层数），可先加载 config 再修改，再用 from_config 构建模型：

```
from transformers import AutoConfig, AutoModel
config = AutoConfig.from_pretrained("bert-base-uncased")
config.num_hidden_layers = 10 # 修改层数
model = AutoModel.from_config(config) # 基于修改后的配置构建模型
```

2. 运行环境配置：通过加载时的参数传递

这类配置与硬件环境强相关，需在加载时手动指定（不放在 config.json 中，避免模型绑定特定硬件）：

设备类型：通过 device_map 参数（Hugging Face）或 map_location（PyTorch）指定：

```
# Hugging Face: 自动分配到 GPU（若有）
model = AutoModel.from_pretrained("bert-base-uncased", device_map="auto")
# PyTorch 原生: 指定到 CPU
model.load_state_dict(torch.load("weights.pth", map_location="cpu"))
数据类型：通过 torch_dtype 参数指定：
model = AutoModel.from_pretrained("bert-base-uncased", torch_dtype=torch.float16)
量化配置：通过 load_in_8bit 等参数启用（依赖 bitsandbytes 库）：
model = AutoModelForCausalLM.from_pretrained("llama-2-7b", load_in_8bit=True)
```

3. 功能开关配置：部分来自配置文件，部分手动调用

推理模式：通过模型的 eval() 方法手动开启（与配置文件无关）：

```
model.eval() # 禁用 dropout，进入推理模式
生成策略参数：部分来自 generation_config.json（模型自带），部分在 generate() 时手动指定：
# 从模型自带的 generation_config.json 加载默认参数
outputs = model.generate(inputs)
# 手动覆盖参数（优先级高于配置文件）
outputs = model.generate(inputs, max_length=100, temperature=0.8)
```


4. 分布式配置：通过环境变量或专用接口传递

这类配置与分布式框架强相关（如 PyTorch 的 `torch.distributed`），通常不放在模型配置文件中：

通信地址 / 端口：通过环境变量指定（启动脚本中设置）：

```
# 终端启动时设置分布式参数 MASTER_ADDR=127.0.0.1
```

```
MASTER_PORT=29500 python train.py
```

设备映射：在 Hugging Face 中通过 `device_map` 参数指定多卡分配：

```
# 将模型拆分到两张 GPU
```

```
model = AutoModel.from_pretrained("bert-large-uncased", device_map={"": [0, 1]})
```

总结：哪些配置在 `config.json` 中？

`config.json` 的核心作用是定义模型的静态结构（如层数、维度），确保权重能正确加载到网络层中。而运行环境（设备、量化）、分布式参数、动态功能开关等与硬件 / 场景强相关的配置，需要在加载时通过参数或环境变量传递，避免模型绑定到特定环境。

简单说：

- `config.json` 回答“模型是什么结构”；
- 其他参数回答“模型在什么环境下以什么方式运行”。

这些配置共同确保模型从“加载”到“运行”的全流程兼容、高效且正确。

文件结构差异

在线加载自动处理的目录结构：

~/.cache/huggingface/hub

```
├── models--sentence-transformers--paraphrase-multilingual-MiniLM-L12-v2
│   ├── snapshots
│   │   └── a615f9c (版本哈希)
│   │       ├── config.json
│   │       ├── pytorch_model.bin
│   │       └── sentence_bert_config.json
│   └── refs
│       └── main
```

本地模型目录要求：功能作用分别是什么？

```
./custom_model/
├── config.json
├── pytorch_model.bin
├── sentence_bert_config.json
├── tokenizer_config.json
└── vocab.txt # 分词器相关文件
```

流式输出和非流式输出的对比？包括实现方式和效果？

进阶本地加载方式 TODO

1. 量化加载（减少显存占用）

适用场景：在资源有限的设备（如消费级 GPU）上运行大模型。

方法：使用 8 位或 4 位量化。

示例（8 位量化）：

```
from transformers import BitsAndBytesConfig, AutoModelForCausalLM
# 配置量化参数
bnb_config = BitsAndBytesConfig(
    load_in_8bit=True,
    llm_int8_threshold=6.0)

model = AutoModelForCausalLM.from_pretrained(
    "bigscience/bloom-7b",
    quantization_config=bnb_config, # 启用 8 位量化
    device_map="auto")
```

2. 分布式加载（多 GPU/多节点）

适用场景：超大规模模型（如百亿参数以上）。

工具：DeepSpeed、Accelerate。

示例（使用 Accelerate）：

```

# 安装依赖:
pip install acceleratefrom accelerate
import init_empty_weights, load_checkpoint_and_dispatchfrom transformers
import AutoConfig, AutoModelForSeq2SeqLM
# 初始化空模型（不加载权重）
config = AutoConfig.from_pretrained("t5-11b")with init_empty_weights():
    model = AutoModelForSeq2SeqLM.from_config(config)
# 分布式加载权重到多个 GPU
model = load_checkpoint_and_dispatch(
    model,
    checkpoint="/t5-11b-sharded", # 分片后的权重目录
    device_map="auto",
    no_split_module_classes=["T5Block"])

```

Fastchat 和大模型部署

fastchat 等部署框架的核心目的是在有多个 GPU 的情况下，合理分配资源，加载多个模型同步进行推理，

第一层，实现固定服务器控制器工作器的进程创建与运行，核心是需要关注参数，三者都需要控制器地址，用户给服务器发送请求，服务器转给控制器，控制器调度工作器回复。服务器和控制器是只有一个的，因此对应的进程也只有一个，而工作器包括本地和在线均有多个，因此进程也是有多多个的，每个进程对应一个工作器。

第二层，工作器需要选择切换，有三种情况，停止当前工作器，stop，添加新的工作器，add，新工作器替换旧工作器，replace。这种情况，在运行工作器的过程会通过进程间通信队列来传递信息。具体实现是通过进程的终止与启动实现的，一个进程相当于一个工作器，终止工作器就是终止对应进程，启动工作器就是启动进程。

进程就需要明确进程的创建，终止（terminate 和 kill 的区别），启动事件（什么时候创建，什么时候启动，启动之后对别的代码和事件本身各有什么影响）和通信队列（如何创建，如何实现通信，如何获取，如何保证数据的安全）。

控制器才是中介，根据不同的调度方法将特定的工作器地址转发给服务器，服务器将用户查询信息发送给工作器。

FastChat 支持单 GPU 运行，只要单 GPU 显存足够加载目标模型即可（例如，7B 参数模型通常需要 10-16GB 显存，13B 参数模型需要 20-30GB 显存）。

多 GPU 是可选优化：当模型参数较大（如 33B、65B），单 GPU 显存不足时，FastChat 支持通过 **模型并行（model parallelism）** 将模型拆分到多个 GPU 上运行（例如，65B 模型通常需要 4-8 张 GPU 协同）。

因此，“多个 GPU” 是为了适配大参数模型的显存需求，而非 FastChat 框架的硬性要求。

FastChat 支持纯 CPU 环境加载模型（需通过 `--device cpu` 参数指定），适合小参数模型（如 7B 及以下，且需足够的 CPU 内存）。

例如，7B 参数模型在 CPU 上加载需要约 13-16GB 内存（FP16 精度），若使用量化（如 4-bit、8-bit），可降低至 4-8GB 内存。

速度极慢：CPU 没有 GPU 的并行计算能力，生成一句话可能需要数十秒甚至几分钟（取决于模型大小和 CPU 性能）。

仅适合测试：纯 CPU 环境仅建议用于简单功能验证，无法满足实际应用的响应速度需求。

FastChat 核心组件

三大核心组件

Controller

功能：协调多个 Worker，路由用户请求到空闲 Worker。

数据类型：通常为后台进程，管理 Worker 注册和负载均衡。

```
controller = Controller(  
    host="localhost", # 服务地址  
  
    port=21001,      #端口号  
    dispatch="round_robin", #调度策略;  
  
    limit_model_concurrency=5, #并发限制  
    no_register=False #是否允许自动注册)
```

Worker

功能：加载大模型并执行实际推理任务。

数据类型: 进程或线程，绑定到具体 GPU 设备。

```
worker = ModelWorker(  
    controller_addr="http://localhost:21001", # 控制器地址  
    worker_addr="http://localhost:21002",    # 工作节点地址  
    model_path="lmsys/vicuna-7b-v1.5",      # 模型路径  
    device="cuda", num_gpus=1,  
  
    load_8bit=True, # bool: 8 位量化  
)
```

在 FastChat 中，**ModelWorker** 是负责加载模型、处理推理请求并与控制器（Controller）交互的核心类。其设计围绕“模型管理”和“任务处理”两大核心目标，包含多个关键属性和方法。以下是其核心属性和方法的详细说明：

核心属性（Core Attributes）

ModelWorker 的属性主要用于**存储模型状态、配置信息、与控制器的连接信息**等，确保 Worker 能独立完成模型加载和任务处理。主要传递的模型参数以及 **model_names**，该 Worker 加载的模型名称列表（如 ["vicuna-7b-v1.5", "llama-2-13b"]）。**loaded** 模型是否加载完成的状态标记（True 表示就绪，可接收任务）。

核心方法（Core Methods）

ModelWorker 的方法可分为**模型管理、任务处理、与控制器交互**三大类，覆盖从启动到服务的全流程。

1. 模型管理方法（Model Management）

__init__(self, controller_addr, worker_addr, ...)，当执行 **worker = ModelWorker(...)****实例化时自动调用**

作用：初始化基础配置（如控制器地址、Worker 自身地址、设备类型、量化参数等），并初始化状态变量（如 loaded=False、current_jobs=0），但**不会加载模型，也不会注册到控制器**。

load_model(self, model_path, model_names=None)

核心方法：从指定路径（model_path）加载模型和分词器，支持量化、多卡分配等优化。

步骤：加载分词器（AutoTokenizer.from_pretrained）→ 加载模型（AutoModelForCausalLM.from_pretrained）→ 迁移模型到目标设备（device）→ 标记 loaded=True。

支持的优化：通过 bitsandbytes 实现 4/8-bit 量化，通过 accelerate 实现多卡分布式加载。

需要加载模型时**显式调用**（通常在 Worker 服务启动后，或通过 API 触发模型切换）。

触发场景：

- Worker 服务启动时，通过命令行参数指定模型路径，自动调用 `load_model`（如 `python -m fastchat.serve.model_worker --model-path lmsys/vicuna-7b-v1.5`）；
- 通过 **HTTP 接口**（如 `/load_model`）动态加载新模型（用于支持**多模型切换**）。

unload_model(self)

卸载模型，释放显存 / 内存资源（设置 `model=None`、`tokenizer=None`，标记 `loaded=False`），通常用于切换模型时调用。

2. 任务处理方法 (Inference Tasks)

generate(self, params): 当外部通过 HTTP 接口（如 `/generate` 或 `/generate_stream`）发送推理请求时，Worker 会调用这两个方法。

处理非流式推理请求，接收输入参数（`prompt`、`max_new_tokens` 等），返回完整生成结果。

流程：解析参数 → 分词器编码输入 → 模型生成（`model.generate`）→ 解码结果 → 返回文本。

generate_stream(self, params)

处理流式推理请求（如实时聊天），逐 token 生成并返回结果（通过迭代器 `yield` 输出）。

优势：减少用户等待感，适用于对话场景，通常配合 FastAPI 的 `StreamingResponse` 使用。

count_token(self, prompt)

计算输入文本（`prompt`）对应的 token 数量，用于判断是否超过模型最大上下文长度（避免输入过长导致报错）。

3. 与控制器交互方法 (Controller Interaction)

register_to_controller(self), `load_model` 执行成功（模型就绪）后，会自动调用 `register_to_controller`。

向控制器注册当前 Worker，发送自身地址（`worker_addr`）、支持的模型（`model_names`）等信息，使控制器将其纳入可用节点列表。

heartbeat(self)

定期向控制器发送心跳请求，携带当前负载（`current_jobs`）和状态（`loaded`），控制器通过心跳判断 Worker 是否存活（超时未收到则标记为离线）。

get_heartbeat_status(self)

返回当前 Worker 的心跳状态信息（如 `worker_addr`、`model_names`、`current_jobs`），供心跳请求使用。

RESTful API Server

功能: 提供 HTTP API 接口，接收用户请求并转发给 Controller。

数据类型: 使用 FastAPI 实现，监听指定端口。

Web UI (可选)

功能: 提供交互式网页界面。

工作流程:

Client → API Server → Controller → Worker → Model Inference

服务器控制器, 工作器的关系到底是什么样的? 用户发送请求后的传播路径到底是什么样的?

客户端 → API Server (8000) → Controller (21001) → Worker (21002)

然后原路返回: Worker → Controller → API Server → 客户端

客户端

在 FastChat 框架中, “配置客户端”(通常指 FastChat 官方封装的客户端或基于其协议自定义的客户端)与“直接使用 requests”(通用 HTTP 客户端库)的核心差异在于封装层次、易用性和功能适配性。前者是为 FastChat 生态量身定制的“专用工具”, 后者是通用的“基础工具”; 两者均基于 HTTP 协议与服务器交互, 但在开发效率、功能覆盖和维护成本上存在显著区别。下面将从“联系与区别”“自定义客户端实现”“核心配置参数”三个维度展开分析。

配置 FastChat 客户端与直接使用 requests 的联系和区别

1. 核心联系

两者的底层逻辑一致, 均通过 HTTP/HTTPS 协议与 FastChat 的服务器 (Controller/Worker) 进行通信, 最终目的都是:

- 向 FastChat 后端 (Controller 负责路由, Worker 负责模型推理) 发送模型调用请求 (如对话、生成文本);
- 接收并解析后端返回的响应 (如生成的文本、模型状态、错误信息);
- 支持认证 (如 API Key)、超时控制、网络代理等基础网络功能。

本质上, FastChat 官方客户端 (如 fastchat.client.OpenAIClient) 内部也可能基于 requests 或 aiohttp (异步) 实现 HTTP 通信, 只是在此基础上封装了 FastChat 特有的逻辑。

2. 关键区别

维度	配置 FastChat 客户端 (如官方封装客户端)	直接使用 requests
封装层次	高层封装, 针对 FastChat 生态优化, 隐藏 HTTP 协议细节。	底层工具, 仅提供 HTTP 请求基础能力, 需手动处理所有协议细节。
易用性	提供简洁 API, 无需关心请求格式	需手动构造请求 URL、请求头

维度	配置 FastChat 客户端（如官方封装客户端）	直接使用 requests
	式、URL 路径、响应解析规则。	（Header）、请求体（Body），解析响应格式。
功能适配性	原生支持 FastChat 特有功能： <ul style="list-style-type: none"> - 模型路由（依赖 Controller） - 负载均衡（多 Worker 调度） - 流式响应（SSE） - 框架内认证（如 API Key 集成） - 模型状态查询（如 Worker 存活检测） 	需自行实现 FastChat 特有逻辑： <ul style="list-style-type: none"> - 手动向 Controller 请求可用 Worker 地址 - 手动处理流式响应（SSE 协议解析） - 手动适配 FastChat 的请求 / 响应 JSON 格式
错误处理	内置 FastChat 专属错误处理（如 Worker 离线、模型加载失败、请求频率限制），返回语义化错误信息。	仅能捕获 HTTP 层错误（如 404/500），需手动解析 FastChat 后端返回的错误 JSON（如 {"error": "model not found"}）。
维护成本	随 FastChat 框架版本更新自动适配（如 API 路径变更、参数新增）。	FastChat 框架更新后，需手动修改请求格式、URL 等（如官方调整 /v1/chat/completions 路径）。
适用场景	快速集成 FastChat 生态，优先追求开发效率和功能完整性。	需高度自定义通信逻辑（如特殊认证、跨框架兼容、非标准请求格式）。

如何自定义 FastChat 客户端

FastChat 客户端的核心是遵循 FastChat 后端（Controller/Worker）的 API 协议，实现请求发送、响应解析和逻辑封装。自定义方式分为两类：

1. 基于 FastChat 官方客户端扩展：复用官方封装的核心逻辑，仅修改参数或补充功能（推荐，降低兼容性风险）；
2. 基于 HTTP 库（requests/aiohttp）从零实现：完全自定义请求流程，适合特殊场景（如无 FastChat 依赖环境）。

核心实现步骤（以 Python 为例）

1. 前提：了解 FastChat 后端 API 协议

FastChat 后端提供两类核心 API（以 v1 版本为例，具体路径可通过 fastchat.serve.openai_api_server 确认）：

模型路由 API（由 Controller 提供）：[GET http://{controller_host}:{controller_port}/api/available_models,](http://{controller_host}:{controller_port}/api/available_models)

获取可用模型及对应的 Worker 地址；

模型调用 API（由 Worker 提供）：

对话生成: POST http://{worker_host}:{worker_port}/v1/chat/completions;

文本生成: POST http://{worker_host}:{worker_port}/v1/completions;

模型信息: GET http://{worker_host}:{worker_port}/v1/models。

请求 / 响应格式与 OpenAI API 兼容 (FastChat 设计初衷之一是兼容 OpenAI 生态), 这也是官方客户端可直接复用 OpenAI 风格 API 的原因。

方式一: 基于 FastChat 官方客户端扩展 (推荐)

FastChat 提供 `fastchat.client` 模块, 封装了 `OpenAIClient` (兼容 OpenAI API 风格) 和 `Client` (原生 FastChat 客户端), 可直接配置参数并扩展功能。

```
from fastchat.client import OpenAIClient
from fastchat.conversation import Conversation, SeparatorStyle
class CustomFastChatClient:
    def __init__(
        self,
        controller_address: str = "http://localhost:21001", # Controller 地址 (路由核心)
        api_key: str = None, # 认证 API Key (若后端开启认证)
        model_name: str = "qwen-7b-chat", # 目标模型名称
        temperature: float = 0.7, # 生成温度 (控制随机性)
        max_tokens: int = 512, # 最大生成长度
        timeout: int = 30, # 请求超时时间 (秒)
        stream: bool = False, # 是否启用流式响应
        proxy: str = None # 网络代理 (如 "http://127.0.0.1:7890")
    ):
        # 1. 初始化官方客户端, 配置核心参数
        self.client = OpenAIClient(
            base_url=controller_address, # 指向 Controller (自动路由到 Worker)
            api_key=api_key, # 认证参数 (后端需配置对应的认证中间件)
            timeout=timeout,
            proxy=proxy # 网络代理配置
        )
        # 2. 配置模型生成参数 (可根据需求动态调整)
        self.generation_params = {
            "model": model_name,
            "temperature": temperature,
            "max_tokens": max_tokens,
            "stream": stream,
            "top_p": 0.95, # 补充常见参数: 采样 Top-P
            "frequency_penalty": 0.1 # 重复惩罚
        }
```

```

def send_chat_request(self, user_query: str, history: list = None) -> str | dict:
    """发送对话请求，支持历史对话"""
    # 1. 构造对话历史（遵循 FastChat Conversation 格式）
    conv = Conversation(
        name=self.generation_params["model"],
        system_message="你是一个 helpful 的助手", # 系统提示词
        messages=history or [], # 历史对话（格式：[{"role": "user",
"content": "..."}]）
        sep_style=SeparatorStyle.QWEN, # 适配模型的分隔符风格（如
Qwen 模型专用）
        sep="\n",
        stop_str="<|end_of_solution|>" # 模型停止符（需与模型匹配）
    )
    # 2. 添加当前用户查询
    conv.append_message(conv.roles[0], user_query)
    conv.append_message(conv.roles[1], None) # 占位模型回复

    # 3. 发送请求并处理响应
    try:
        if self.generation_params["stream"]:
            # 处理流式响应（逐段返回）
            response = ""
            for chunk in self.client.chat.completions.create(
                messages=conv.to_openai_format(), # 转为 OpenAI 格
式的消息
                **self.generation_params
            ):
                content = chunk.choices[0].delta.content
                if content:
                    response += content
                    print(content, end="", flush=True) # 实时打印
            return response
        else:
            # 处理完整响应
            response = self.client.chat.completions.create(
                messages=conv.to_openai_format(),
                **self.generation_params
            )
            return response.choices[0].message.content # 返回生成文本
    except Exception as e:
        # 自定义错误处理（区分 FastChat 特有错误和网络错误）
        if "Worker not available" in str(e):

```

```

        raise Exception(f'模型 Worker 离线，请检查后端服务：
{str(e)}')
    elif "API key invalid" in str(e):
        raise Exception(f'认证失败：API Key 无效')
    else:
        raise Exception(f'请求失败：{str(e)}')
# -----# 使用自定义客户端# -----if
__name__ == "__main__":
    # 初始化客户端（配置地址、API Key、模型参数等）
    client = CustomFastChatClient(
        controller_address="http://192.168.1.100:21001", # 远程 Controller 地
址
        api_key="sk-xxxxxxx", # 后端配置的 API Key
        model_name="qwen-7b-chat",
        temperature=0.5,
        max_tokens=1024,
        stream=True, # 启用流式响应
        proxy="http://127.0.0.1:7890"
    )
    # 发送对话请求（带历史对话）
    history = [{"role": "user", "content": "介绍一下 FastChat 框架"}]
    response = client.send_chat_request(
        user_query="它的核心组件有哪些？",
        history=history
    )
    print("\n 最终响应： ", response)

```

方式二：基于 requests 从零实现自定义客户端

若需高度灵活控制（如自定义认证逻辑、非标准请求格式），可直接使用 requests 构造 HTTP 请求，需手动适配 FastChat 后端的 API 协议。

```

import requests
import json
from typing import Optional, List, Dict
class FastChatRequestsClient:
    def __init__(
        self,
        controller_host: str = "localhost",
        controller_port: int = 21001,
        api_key: Optional[str] = None,
        timeout: int = 30,
        proxy: Optional[str] = None
    ):
        # 1. 基础配置
        self.controller_base_url = f"http://{controller_host}:{controller_port}"
        self.api_key = api_key

```

```

self.timeout = timeout
# 2. 配置 requests 会话（复用连接，提升性能）
self.session = requests.Session()
self.session.timeout = timeout
# 3. 代理配置
if proxy:
    self.session.proxies = {
        "http": proxy,
        "https": proxy
    }
# 4. 认证头（若后端开启 API Key 认证）
self.headers = {
    "Content-Type": "application/json"
}
if self.api_key:
    self.headers["Authorization"] = f"Bearer {self.api_key}"

def _get_worker_address(self, model_name: str) -> str:
    """第一步：向 Controller 请求模型对应的 Worker 地址（路由逻辑）
    """
    try:
        resp = self.session.get(
            url=f"{self.controller_base_url}/api/available_models",
            headers=self.headers
        )
        resp.raise_for_status() # 触发 HTTP 错误（如 404/500）
        available_models = resp.json()
        # 查找目标模型的 Worker 地址（FastChat 返回格式：
        {model_name: worker_address}）
        if model_name not in available_models:
            raise Exception(f'模型 {model_name} 未在 Controller 注册')
        return available_models[model_name]
    except Exception as e:
        raise Exception(f'获取 Worker 地址失败：{str(e)}')

def send_chat_request(
    self,
    model_name: str,
    user_query: str,
    history: Optional[List[Dict]] = None,
    temperature: float = 0.7,
    max_tokens: int = 512,
    stream: bool = False
) -> str | Dict:

```

询

"""第二步：向 Worker 发送对话请求"""

1. 获取 Worker 地址（由 Controller 路由）

worker_address = self._get_worker_address(model_name)

chat_api_url = f'{worker_address}/v1/chat/completions'

2. 构造请求体（遵循 FastChat/OpenAI 兼容格式）

messages = history or []

messages.append({"role": "user", "content": user_query}) # 添加当前查

```
payload = {
    "model": model_name,
    "messages": messages,
    "temperature": temperature,
    "max_tokens": max_tokens,
    "stream": stream,
    "top_p": 0.95,
    "frequency_penalty": 0.1
}
```

3. 发送请求并处理响应

try:

if stream:

流式响应（FastChat 用 SSE 协议，需逐行解析）

response = ""

```
with self.session.post(
    url=chat_api_url,
    headers=self.headers,
    json=payload,
    stream=True # 启用流式接收
) as r:
```

r.raise_for_status()

for line in r.iter_lines(decode_unicode=True):

if line.startswith("data: "):

data = line[6:] # 去掉 "data: " 前缀

if data == "[DONE]":

break # 流式结束

try:

chunk = json.loads(data)

content =

chunk["choices"][0][["delta"].get("content")

if content:

response += content

print(content, end="", flush=True)

except json.JSONDecodeError:

```

        continue

    return response
else:
    # 完整响应
    resp = self.session.post(
        url=chat_api_url,
        headers=self.headers,
        json=payload
    )
    resp.raise_for_status()
    result = resp.json()
    return result["choices"][0]["message"]["content"]
except Exception as e:
    raise Exception(f'发送对话请求失败: {str(e)}')
# -----# 使用 requests 自定义客户端#
-----if __name__ == "__main__":
    client = FastChatRequestsClient(
        controller_host="192.168.1.100",
        controller_port=21001,
        api_key="sk-xxxxxxx",
        proxy="http://127.0.0.1:7890"
    )
    history = [{"role": "user", "content": "介绍一下 FastChat 框架"}]
    response = client.send_chat_request(
        model_name="qwen-7b-chat",
        user_query="它的核心组件有哪些？",
        history=history,
        stream=True
    )
    print("\n 最终响应: ", response)

```

自定义 FastChat 客户端的核心配置参数

自定义客户端时，参数需覆盖**连接、认证、请求控制、响应处理、框架适配**等维度，具体可分为以下 6 类：

参数类别	具体参数	用途说明	示例值
基础连接参数	controller_address	FastChat Controller 地址（格式：http://host:port），负责模型路由。	http://192.168.1.100:21001
	worker_address	直接指定 Worker 地址（跳过 Controller 路由，适合单 Worker 场景）。	http://192.168.1.101:21002
	timeout	请求超时时间（秒），防止因模型	30（秒）

参数类别	具体参数	用途说明	示例值
认证参数		推理慢导致无限等待。	
		网络代理地址	
	proxy	(HTTP/HTTPS/SOCKS)，用于跨网络访问后端服务。	http://127.0.0.1:7890
	ssl_verify	是否验证 SSL 证书 (HTTPS 场景)，测试环境可关闭。	True/False
	connection_pool_size	HTTP 连接池大小，提升高并发场景下的性能。	10 (默认通常为 5)
	api_key	后端认证用 API Key (需与 FastChat 服务器配置的 --api-key 一致)。	sk-8xxxxxxxxxxxxxxxxx xxxxxxxxxxxxxxxxxxxx
	token	自定义 Token 认证 (如 JWT)，需后端中间件支持。	eyJhbGciOiJIUzI1NiIsInR5cCI6IkpXVCJ9
	signature	签名参数 (如时间戳 + 密钥签名)，用于高安全性场景。	md5(timestamp + secret_key)
	model_name	目标模型名称 (需与 Worker 加载的模型名称一致，如 qwen-7b-chat)。	qwen-7b-chat/llama-2-7b-chat
	temperature	生成温度 (0~2)，值越高随机性越强，值为 0 时输出 deterministic。	0.7 (默认)
请求控制参数	max_tokens	模型最大生成长度 (含输入上下文)，防止输出过长。	512/1024
	top_p	采样 Top-P (0~1)，控制生成的多样性 (与 temperature 配合使用)。	0.95 (默认)
	frequency_penalty	重复惩罚 (-2~2)，值越高减少重复生成内容。	0.1
	presence_penalty	存在惩罚 (-2~2)，值越高鼓励生成新主题。	0.0 (默认)
	stop	模型停止符 (列表)，遇到该字符时停止生成。	`["<end_of_text>","###"]`
	system_message	系统提示词，定义模型的角色和行为 (如 “你是一个专业的技术助手”)。	"你是一个 helpful 的助手，仅用中文回答"
响应处理参数	stream	是否启用流式响应 (SSE)，适合实时展示生成结果 (如聊天界面)。	True/False
	response_format	响应格式 (如 JSON/Text)，FastChat 默认返回 JSON (兼容 OpenAI)。	json/text

参数类别	具体参数	用途说明	示例值
FastChat 特有参数	encoding	响应编码格式，避免中文乱码。	utf-8（默认）
	ignore_empty_chunks	流式响应中是否忽略空内容块（部分模型推理初期可能返回空 chunk）。	True
	route_strategy	模型路由策略（仅 Controller 场景），如 shortest_queue（最短队列）。	shortest_queue/random
	load_balance	是否启用负载均衡（多 Worker 场景），由 Controller 自动调度。	True/False
	model_version	模型版本（若后端部署多个版本），需与 Worker 配置一致。	v1/v2
	conv_template	对话模板（适配不同模型的分隔符风格），如 qwen/llama-2。	qwen（对应 SeparatorStyle.QWEN）
错误处理参数	retry_count	请求失败重试次数（如网络波动、Worker 临时不可用）。	2（默认）
	retry_delay	重试间隔时间（秒），避免频繁重试压垮后端。	1（秒）
	retry_on_exceptions	触发重试的异常类型（如网络错误、Worker 超时）。	(requests.exceptions.ConnectionError, TimeoutError)

关键技术 with 问题探究

控制器和 worker 的地址和端口指什么

Controller 负责协调请求，将客户端的请求分发给可用的 Worker。Worker 则是实际加载模型并进行推理的节点。因此，它们的地址和端口设置是为了让各个组件能够互相通信。

简单的说因为控制器和 worker 分开，便需要通信，通信就需要告诉别人自己的位置，位置用地址和端口来表示。

那么地址和端口该如何设置呢？

如果 Controller 和 Worker 不在同一台机器上，地址就不能是 localhost，而应该是机器的实际 IP。此外，端口的选择需要确保未被其他应用占用，并且防火墙允许通信。

Controller 默认使用 21001 端口，Worker 可能使用 21002 或其他端口，而 API Server 使用 8000。如果用户需要同时运行多个实例，必须修改这些端口以避免冲突。

controller 地址:

单机部署

controller_addr = "http://localhost:21001"

多机部署 (Controller 在 192.168.1.100)

controller_addr = "http://192.168.1.100:21001"

Worker1 配置

worker_addr = "http://localhost:21002"

Worker2 配置 (同一台机器需换端口)

worker_addr = "http://localhost:21003"

不同机器可复用端口

worker_addr = "http://192.168.1.101:21002"

在 FastChat 的架构中，客户端请求的 API 地址（如 `http://localhost:8000/api/v1/chat`）是由 FastChat 的 API Server 提供的。

API Server 监听在某个端口（默认 8000），接收来自客户端的 HTTP 请求，然后将请求转发给 Controller，并等待 Controller 返回结果，最后将结果返回给客户端。

在代码中，客户端请求的是 API Server 的地址（8000 端口），而 Controller 和 Worker 的地址是内部通信使用的，客户端并不直接与 Controller 或 Worker 交互。

具体到地址关系：

- API Server 的地址（如 `http://localhost:8000`）是暴露给客户端的。
- Controller 的地址（如 `http://localhost:21001`）是 API Server 和 Worker 都知道的。API Server 需要向 Controller 发送请求，Worker 启动时需要向 Controller 注册。
- Worker 的地址（如 `http://localhost:21002`）是 Worker 自身监听的地址，用于接收 Controller 转发过来的请求。Worker 启动时会将自己的地址注册到 Controller。

控制器存储工作器地址的容器是什么数据类型？二者是如何进行通信的

FastChat 的控制器（Controller）在内存中**维护工作器信息**时，核心使用**字典（Dictionary）**作为主要数据结构，具体体现在：

- 控制器通过字典存储**工作器的唯一标识**（如 `worker_id`）与**工作器元数据**的**映射**关系，元数据包括工作器的网络地址（IP: 端口）、支持的模型、当前负载、状态（在线 / 离线）等。

- 此外，还会通过集合（Set）或列表（List）辅助管理在线工作器的地址列表，用于快速筛选可用节点。

这种设计的优势是：通过字典可以高效根据 `worker_id` 查找对应的地址和状态，同时支持快速更新工作器的元数据（如负载变化）。

FastChat 中控制器与工作器的通信基于 **HTTP** 协议，具体流程如下：

工作器注册：工作器启动时，会主动向控制器发送注册请求（**POST /register_worker**），携带自身的地址、支持的模型等信息。控制器将这些信息存入字典并标记为“在线”。

状态心跳：工作器定期向控制器发送心跳请求（**POST /heartbeat**），更新自身的负载状态（如当前处理的请求数）。控制器通过心跳判断工作器是否存活，若超时未收到心跳则标记为“离线”。

任务调度：当用户请求到达时，控制器根据工作器的负载和模型支持情况，从在线工作器中选择合适的节点，将请求转发到对应的工作器地址（通过 **POST 请求发送任务数据**）。

结果返回：工作器处理完任务后，直接将结果返回给前端（或调用方），同时通知控制器任务完成，更新自身状态。

这种基于 HTTP 的通信方式简单易实现，且便于跨网络部署（工作器和控制器可在不同机器上）。

总结来说，FastChat 控制器通过字典管理工作器地址及元数据，二者通过 HTTP 协议完成注册、心跳和任务调度的交互。

模型加载 worker 内部具体是如何加载模型，又是如何使用模型完成推理等任务的

FastChat 借助 **BaseModelWorker** 类（位于 `fastchat/serve/model_worker.py`）来加载模型。此类提供了基本的模型加载与推理功能。以加载 Vicuna 模型为例，代码如下：

模型加载：

```
# 实际加载模型的代码（在 Worker 内部）
from transformers import AutoModelForCausalLM
model = AutoModelForCausalLM.from_pretrained(
    model_path,
    device_map="auto",
    load_in_8bit=True)
```

模型工作器是一个独立的进程，它主要负责两件大事：**加载模型和执行推理**。

如何加载模型？

工作器在启动时会执行一个标准的模型加载流程，这个过程与在本地使用 Hugging Face transformers 库加载模型非常相似：

a. 初始化与配置：

工作器首先解析参数，如 `--model-path`（模型路径）、`--device`（运行设备，如 `cuda`、`cpu`）、`--num-gpus`（使用 GPU 数量）、`--load-8bit`（是否 8 位量化）等。

根据这些参数构建模型配置（`model.config`）。

b. 加载模型权重：

核心是使用 transformers 库的 `AutoModel.from_pretrained()` 方法（对于对话模型，通常是 `AutoModelForCausalLM`）。

`model-path` 可以是：

本地路径：如 `/home/user/llama-2-7b-chat-hf`，直接从磁盘加载。

Hugging Face Model Hub 上的模型 ID：如 `meta-llama/Llama-2-7b-chat-hf`。此时工作器会首先下载模型（如果本地没有缓存），然后加载。

根据配置，模型会被分配到指定的设备上（如 `model.to(device)`）。

c. 加载分词器（Tokenizer）：

使用 `AutoTokenizer.from_pretrained()` 从相同的路径加载分词器。分词器负责将输入的文本转换为模型能够理解的词汇 ID（Token IDs），以及将模型生成的 Token IDs 转换回文本。

d. （可选）应用优化：

多 GPU 并行：如果设置了 `--num-gpus 2`，FastChat 会使用类似 `model_parallel` 的技术将模型的不同层分布到多个 GPU 上，以平衡显存负载。

量化：如果设置了 `--load-8bit` 或 `--load-4bit`，会使用 `bitsandbytes` 等库进行量化，将 FP16/FP32 的模型权重转换为 Int8/Int4，极大地减少显存占用，但可能会轻微牺牲性能。

vLLM 集成：如果使用 `--vllm` 参数，FastChat 会调用 vLLM 这个高性能推理引擎来加载和管理模型。vLLM 使用 `PagedAttention` 等技术，极大地提高了吞吐量。

e. 注册到控制器：

模型成功加载后，工作器会向控制器（Controller）发送一个 **HTTP POST 注册请求**，告知控制器：“我已经准备好了，我这里有模型 X，我的地址是 Y，你可以把推理请求发给我”。

加载在线模型和本地模型的区别是什么？加载在线模型(qwen)的具体实例

根据传入的参数（**模型名称**）来判断是 langchain 支持的模型，**在线模型**（需要**自己重写 worker-class 类**，根据使用的模型不同，类名也是不同的，类名也是通过传入的参数模型名称来确认的），还是**本地模型**（直接使用 **fastchat** 自带的 **ModelWorker** 类）。

联系：两者在 **接口层面** 被完美抽象和统一，为客户端提供了完全一致的使用体验。

区别：两者在 **实现层面** 有本质不同，一个是“**代理**”（在线 API），负责**网络通信和协议转换**；另一个是“**引擎**”（本地模型），负责本地资源管理和重型计算。

传入的参数 kwargs 一般包含的字段如下：

工作器地址和端口 host: port:

模型名称 model_names:['model_name']

控制器地址：controller_address:

```
def create_model_worker_app(**kwargs):
```

```
    # 根据获得的 worker-class 类型，分类实现 worker，langchain 支持，在线，本地（vllm 加速和 fastchat 加载）
```

```
    if worker_class := kwargs.get('langchain_model'): # Langchain 支持的模型不用做操作
```

```
        worker = "
```

```
        # 在线模型 API,注意这里参数获取的 worker-class 是类，不是 str 参数，所以可以实例化
```

```
    elif worker_class := kwargs.get('worker_class'): #
```

```
        worker = worker_class(
```

```
            model_names=args.model_names,
```

```
            controller_address=args.controller_address,
```

```

        worker_address=args.worker_address

    )

    #本地模型

    else:

        worker = ModelWorker(

            controller_address=args.controller_address,

            worker_address=args.worker_address,

            worker_id=worker_id,

            model_name=args.model_name,

            model_path=args.model_path,

        )

        sys.modules["fastchat.serve.model_worker"].args = args

        sys.modules["fastchat.serve.model_worker"].gptq_config =
gptq_config

        MakeFastAPIOffline(app) # 这一步是将 fastapi 的 app 改为离线，不依
赖在线的渲染文档，使用离线文档加载 mount

        app.title = f"FastChat LLM Server ({args.model_names[0]})" # 这个 app
是配置谁的，fastapi 的吗

        app._worker = worker

    return app

```

核心区别

尽管接口统一，但两者的内部实现原理完全不同。

1. 加载顺序与初始化

步骤	在线 API 工作器 (kwargs.get('worker_class'))	本地模型工作器 (else 分支)
核心动作	配置客户端	加载模型权重
过程	1. 实例化 自定义的 worker_class。	1. 解析大量与模型推理相关的参

步骤	在线 API 工作器 (kwargs.get('worker_class'))	本地模型工作器 (else 分支)
	2. 初始化内部用于调用第三方 API 的客户端 (如 requests.Session 或官方 SDK)。	数 (device, num_gpus, max_gpu_memory, load_8bit, gptq_config 等)。
	3. 不涉及模型文件, 极快完成。	2. 设置 GPU 环境 (CUDA_VISIBLE_DEVICES)。 3. 从磁盘读取模型权重文件 (.safetensors 或 .bin)。 4. 将模型加载到指定设备 (GPU/CPU) 内存中。 5. 耗时很长, 占用大量显存。
代码体现	worker = worker_class(model_names=..., controller_address=..., ...)	worker = ModelWorker(..., model_path=args.model_path, device=args.device, num_gpus=..., ...)

2. 接收与处理请求

当请求到达 FastAPI app, 并最终路由到 worker 的 generate 方法时:

步骤	在线 API 工作器	本地模型工作器
接收请求	✓ 接收到标准的 params 字典 (含 messages, temperature 等)。	✓ 接收到完全相同的 params 字典。
核心处理	1. 格式转换与网络请求: <ul style="list-style-type: none"> - 将 params 转换为第三方 API 要求的格式。 - 携带 API Key, 通过网络向远程 API 端点 (如 api.dashscope.aliyuncs.com) 发送 HTTP 请求。 - 等待远程服务器完成计算。 2. 它自身不进行任何计算。	1. 本地计算: <ul style="list-style-type: none"> - 使用自身的 Tokenizer 将 messages 转换为模型输入张量 (input_ids, attention_mask)。 - 将张量送入已加载到本地 GPU/CPU 上的模型进行前向传播计算 (推理)。 - 使用 model.generate() 等方法自回归地生成 token。 2. 整个过程在本地硬件上完成, 无网络请求。
依赖项	强烈依赖网络质量和第三方 API 的可用性与速率限制。	强烈依赖本地硬件性能 (GPU 算力、显存容量)。

3. 返回响应

步骤	在线 API 工作器	本地模型工作器
接收原始响应	收到第三方 API 返回的 HTTP 响应 (JSON 格式)。	收到模型计算输出的 张量 , 并通过 Tokenizer 解码为 文本 。
后处理	1. 格式转换: 将第三方 API 的 JSON 响应	1. 格式封装: 将生成的文本直

步骤	在线 API 工作器	本地模型工作器
理	解析，并提取出文本内容，重新封装成 FastChat 标准的响应格式。 2. 错误处理： 处理网络错误、API 限额错误、认证错误等。	接封装到 FastChat 标准的响应格式中。 2. 错误处理： 处理模型推理中的错误（如显存溢出、输入过长等）。
最终输出	✓ 返回一个符合 OpenAI API 标准的 JSON 响应。	✓ 返回一个完全相同的符合 OpenAI API 标准的 JSON 响应。

如何使用模型完成推理？

工作器本质上是一个 **HTTP 服务器**（通常基于 `fastapi`），它暴露了一个或多个 API 端点（API Endpoints）。最关键的端点是 `/worker_generate` 或类似的端点，用于处理推理请求。

当一个推理请求（例如，用户的一段对话）从控制器转发过来时，工作器会执行以下步骤：

a. 接收与解析请求：

工作器接收到一个 HTTP POST 请求，其 `body` 中包含一个结构化数据（通常是 JSON），里面包含了生成参数，如：

```
{  "prompt": "Human: Hello?\nAssistant:",
  "temperature": 0.7,
  "max_new_tokens": 512,
  "stop_str": "Human:"}
```

b. 文本编码（Tokenization）：

使用加载的分词器将收到的 `prompt` 文本字符串转换为一个 Token ID 序列（`input_ids`）。

同时，通常会生成 `attention_mask` 等其它必要的输入张量。

c. 执行模型生成（Inference）：

将 `input_ids` 和生成参数（`temperature`, `max_new_tokens` 等）传递给已加载的模型。

调用模型的 `.generate()` 方法。

```
# 伪代码，演示核心逻辑
input_ids = tokenizer([prompt]).input_ids
```

```
output_ids = model.generate(  
    input_ids=input_ids,  
    temperature=temperature,  
    max_new_tokens=max_new_tokens,)
```

这个过程是**自回归的**：模型会逐个预测下一个最可能的 Token，直到达到最大生成长度或遇到停止词（如 stop_str）为止。

d. 文本解码（Detokenization）：

使用分词器将模型输出的 output_ids（一个整数序列）解码回可读的文本字符串。

e. 返回响应：

将生成的文本、可能的状态码或其他信息封装成 JSON 格式，通过 HTTP 响应返回给控制器，控制器再最终返回给用户或 Web 服务器。

加载模型时加载数量有多少？

在 FastChat 中，每个 Worker 通常负责加载一个模型实例。如果有多个 Worker，每个 Worker 可以加载同一个模型的不同副本，或者不同的模型。例如，如果有三个 Worker，每个都加载 Vicuna-7B，那么相当于有三个模型实例在处理请求，提高并发能力。

在 FastChat 中，**模型加载数量**与 Worker 的数量直接相关。每个 Worker 进程会独立加载一个完整的模型副本到显存中。具体规则如下：

场景	模型加载数量	物理资源占用	典型配置示例
单 Worker	1 个模型副本	独占 GPU 显存	num_gpus=1
多 Worker 单卡	N 个模型副本（需显存足够）	共享单 GPU 显存	num_gpus=1 + 启动多个 Worker
多 Worker 多卡	每个 GPU 加载 1 个副本	每 GPU 独立占用	num_gpus=2（单个 Worker 跨卡）或启动多个 Worker

控制器的调度逻辑原理是什么，具体应用过程中是怎么样实现的

请求调度逻辑（Controller 核心代码片段）：

```
def dispatch_request(self, worker_list):  
    # 简单的轮询调度  
    selected_worker = self.worker_queue.pop(0)  
    self.worker_queue.append(selected_worker)  
    return selected_worker
```

轮询机制深度剖析

如何调度轮询以及和 vllm 方法的对比下一步继续详细深入学习。

FastChat 的 轮询调度（Round-Robin）是 Controller 的核心负载均衡策略，其工作原理如下：

调度流程伪代码

python

复制

下载

```
class Controller:  
    def __init__(self):  
        self.worker_queue = [] # 保存注册的 Worker  
  
    def register_worker(self, worker):  
        self.worker_queue.append(worker) # 新 Worker 加入队列尾部  
  
    def dispatch_request(self):  
        if not self.worker_queue:  
            raise NoWorkerAvailable()  
  
        # 轮询选择 Worker  
        selected_worker = self.worker_queue.pop(0)  
        self.worker_queue.append(selected_worker) # 放回队列尾部  
        return selected_worker
```

动态演示（3 个 Worker 的场景）

复制

下载

请求序列 → 分配的 Worker

- 1 → Worker1
 - 2 → Worker2
 - 3 → Worker3
 - 4 → Worker1
 - 5 → Worker2
 - 6 → Worker3
- (循环往复)

性能特征对比

调度策略	优点	缺点	适用场景
Round-Robin	实现简单，绝对公平	无视 Worker 实际负载	各 Worker 性能均匀
Least-Loaded	动态平衡负载	需实时监控指标	异构 GPU 环境
Hash-Based	会话保持	可能产生热点	需状态保持的场景

接收输入信息

FastChat 的 Worker 运用 RPC（远程过程调用）来接收来自 Controller 的输入信息。Controller 会把用户的输入信息发送给 Worker，Worker 接收到后开展推理。相关代码在 fastchat/serve/model_worker.py 里：

调用模型输出结果

在 Worker 接收到输入信息后，会调用模型的推理方法来输出结果。以 generate 方法为例：

代码实例详解

示例：部署 Vicuna-7B 并进行推理

实例化控制器

```
# 安装依赖（前置步骤）# pip install fastapi uvicorn "fschat[model_worker,web ui]"
# ----- 启动 Controller -----
from fastchat.serve.controller import Controller
# 实例化控制器 (str, int, bool 类型参数)
```

```

controller = Controller(
    host="localhost", # str: 服务地址
    port=21001,       # int: 端口号
    dispatch="round_robin", # str: 调度策略
    limit_model_concurrency=5, # int: 并发限制
    no_register=False # bool: 是否允许自动注册)
controller.start() # 启动服务

# ----- 启动 Worker -----
from fastchat.serve.model_worker import ModelWorker

# 实例化模型工作节点
worker = ModelWorker(
    controller_addr="http://localhost:21001", # str: 控制器地址
    worker_addr="http://localhost:21002",    # str: 工作节点地址
    model_path="lmsys/vicuna-7b-v1.5",       # str: 模型路径
    device="cuda",                           # str: 计算设备
    num_gpus=1,                              # int: GPU 数量
    load_8bit=True,                          # bool: 8 位量化
)
worker.start() # 启动工作节点

# ----- 使用 API 进行推理 -----
import requests

# 构造请求参数 (dict 类型)
payload = {
    "prompt": "Explain quantum physics in 3 sentences", # str: 输入文本
    "temperature": 0.7, # float: 生成温度
    "max_new_tokens": 100, # int: 最大生成长度
    "stop": ["\n###"] # list[str]: 停止标记
}

# 发送 POST 请求 (返回 Response 对象)
response = requests.post(
    "http://localhost:8000/api/v1/chat", # str: APIserver 地址，用户直接访问的地址，api server 再和 controller 通信。
    json=payload,
    headers={"Content-Type": "application/json"})

# 输出结果 (dict 类型)
result = response.json()
print(result)

```

fastapi 结合 fastchat

可能的方案有两种：一种是使用 FastAPI 作为主应用，将 FastChat 的 API 作为子应用挂载；另一种是在 FastAPI 中直接调用 FastChat 的 Controller 或 Worker 来处理请求，从而创建自定义的 API 端点。

```
# ----- 步骤 1: 启动 FastChat 核心组件 -----
from fastchat.serve.controller import Controller
from fastchat.serve.model_worker import ModelWorker
import threading

# 启动 Controller
controller = Controller(host="0.0.0.0", port=21001)
controller_thread = threading.Thread(target=controller.start)
controller_thread.start()

# 启动 Worker
worker = ModelWorker(
    controller_addr="http://localhost:21001",
    worker_addr="http://localhost:21002",
    model_path="lmsys/vicuna-7b-v1.5",
    device="cuda",
    num_gpus=1
)
worker_thread = threading.Thread(target=worker.start)
worker_thread.start()

# ----- 步骤 2: 创建 FastAPI 应用 -----
from fastapi import FastAPI, Depends, HTTPException
from fastapi.security import APIKeyHeader
from fastchat.serve.base_model_worker import app as fastchat_app
import uvicorn

# 创建主应用
app = FastAPI(title="Integrated API Server")

# 挂载 FastChat 原生 API（重要！）
app.mount("/fastchat", fastchat_app) # 将 FastChat 路由挂载到子路径

# ----- 步骤 3: 添加自定义功能 -----
# 示例 1: 添加鉴权中间件
API_KEY_NAME = "X-API-KEY"
api_key_header = APIKeyHeader(name=API_KEY_NAME, auto_error=False)

async def validate_api_key(api_key: str = Depends(api_key_header)):
```

```

    if api_key != "SECRET_KEY_123":
        raise HTTPException(status_code=403, detail="Invalid API Key")
    return api_key

# 示例 2: 添加健康检查端点
@app.get("/health", dependencies=[Depends(validate_api_key)])
def health_check():
    return {"status": "OK", "model_loaded": True}

# 示例 3: 自定义代理请求
@app.post("/custom/chat")
async def custom_chat_endpoint(prompt: str, api_key: str = Depends(validate_api_key)):
    # 直接调用 FastChat Worker 进行推理
    response = await fastchat_app.state.worker.generate(prompt)
    return {"result": response}

```

三、关键代码解析

1. 组件挂载

```
app.mount("/fastchat", fastchat_app)
```

将 FastChat 的原生 API 路由挂载到 /fastchat 路径下
原生接口访问地址变为: <http://localhost:8000/fastchat/api/v1/chat>

2. 直接调用 Worker

```

# 获取 Worker 实例
worker = fastchat_app.state.worker
# 同步调用生成方法
response = worker.generate(prompt="Hello")

```

urllib.parse.unquote

`knowledge_base_name = urllib.parse.unquote(knowledge_base_name)` 目的是啥
`urllib.parse.unquote(knowledge_base_name)` 的作用是对 URL 编码的字符串进行解码。

当字符串经过 URL 编码后, 一些特殊字符(如空格、中文、标点符号等)会被转换为 %XX 形式的编码(例如空格会变成 %20, 中文“测试”会变成 %E6%B5%8B%E8%AF%95)。

使用 `unquote()` 函数可以将这种编码后的字符串还原为原始的字符形式，便于后续程序处理和使用。

举例来说：

如果 `knowledge_base_name` 的值是 `%E7%9F%A5%E8%AF%86%E5%BA%93`，经过 `unquote()` 处理后会变成 知识库。

这个操作通常用于处理从 URL 中获取的参数，因为 URL 中只能包含特定的安全字符，其他字符需要经过编码传输，接收方则需要解码才能得到原始内容。

Langchain 封装在线 api-Worker

在代码上下文中，`worker_class` 是一个关键的设计模式实现，它指的是具体嵌入模型 API 的工作处理器类。让我详细解释这个概念：

1. 本质定义

`worker_class` 是一个 **Python 类**，它封装了与特定嵌入模型 API 交互的所有细节。它的核心职责是：

- 处理 API 请求的构造

- 管理认证凭据

- 执行网络调用

- 解析 API 响应

- 错误处理

2. 典型实现结构

```
class EmbeddingModelWorker:
```

```
    """所有嵌入模型工作器的基类"""
```

```
    @classmethod
```

```
    def can_embedding(cls):
```

```
        """声明是否支持嵌入操作"""
```

```
        return True
```

```
    def do_embeddings(self, params: ApiEmbeddingsParams):
```

```
        """核心方法：执行向量化操作"""
```

```
        # 具体实现因 API 而异
```

```
        ...
```

3. worker 具体实现示例 (以 OpenAI 为例)

```
class OpenAIEmbeddingWorker(EmbeddingModelWorker):
    API_BASE = "https://api.openai.com/v1/embeddings"

    def __init__(self):
        self.api_key = os.getenv("OPENAI_API_KEY") # 从环境变量获取密钥

    def do_embeddings(self, params):
        headers = {"Authorization": f"Bearer {self.api_key}"}
        payload = {
            "input": params.texts,
            "model": params.embed_model # 如"text-embedding-ada-002"
        }

        try:
            response = requests.post(self.API_BASE, json=payload,
headers=headers)
            response.raise_for_status()
            data = response.json()

            # 提取嵌入向量 [ [0.1, -0.2, ...], ... ]
            embeddings = [item["embedding"] for item in data["data"]]

            return {
                "code": 200,
                "data": embeddings
            }
        except requests.exceptions.RequestException as e:
            return {
                "code": 500,
                "msg": f"API 请求失败: {str(e)}"
            }
```

可以看到，这个就是把上面的在线调用 `api` 封装成一个类，整个调用过程通过一个 `do_embeddings()` 方法实现，这样在项目实践中就可以直接调用接口，而且方便扩展不同的模型。

4. 设计特点解析

特性	说明
抽象接口	所有 worker 必须实现 <code>do_embeddings()</code> 方法
配置驱动	通过 <code>get_model_worker_config()</code> 动态加载
能力声明	<code>can_embedding()</code> 确保类支持嵌入操作
隔离性	每个 API 提供商有独立实现 (OpenAI/Cohere 等)

特性	说明
无状态性	通常每个请求新建实例，避免状态残留

5. 关键优势

可扩展性：新增 API 只需添加新的 worker 类

```
class NewEmbeddingWorker(EmbeddingModelWorker):
    def do_embeddings(self, params):
        # 实现新 API 的调用逻辑
        ...
```

协议抽象：对上层隐藏不同 API 的差异

请求参数统一为 ApiEmbeddingsParams

输出统一为 BaseResponse

错误隔离：单个 API 故障不影响其他 worker

配置集中化：通过中央配置管理 API 密钥和端点

8. 项目调用示例

```
if embed_model_name in list online_embed_models():
    #根据模型名称获取该模型的配置信息。配置信息可能包括模型的工作
    #类（worker_class）、模型的实际名称（embed_model）、API 地址、API 密钥等
    config=get_model_worker_config(embed_model_name)
    embed_model_name=config.get('embed_model_name')
    #将`worker_class`实例化，创建一个`worker`对象。这个对象将用于后续
    #的嵌入操作
    worker_class=config.get('worker_class')
    worker=worker_class()
    #调用`worker_class`的类方法`can_embedding()`，判断该工作类是否支持
    #嵌入操作
    if worker.can_embedding():
        #创建`ApiEmbeddingsParams`参数对象，传入以下参数：
        # texts: 需要被向量化的文本列表。 -to_query: 指示这些文本是用
        #于查询-embed_model: 实际使用的嵌入模型名称（来自配置）
        apiEmbeddingParameters=(texts,embed_model_name)
        #将参数传入，执行嵌入操作。向相应的 API 发送请求，并返回响
        #应
        response=worker.do_embedding(apiEmbeddingParameters)
        #将`worker.do_embeddings`返回的响应`resp`（通常是一个字典）解包，作为参数
        #传递给`BaseResponse`类的构造函数，创建并返回一个标准化的响应对象
    return response #BaseResponse(**response)
```

这种设计模式常见于 AI 中间件系统（如 FastChat, LangChain 等），它实现了**模型提供方与业务逻辑**之间的解耦，使系统能够灵活支持多种嵌入服务。