

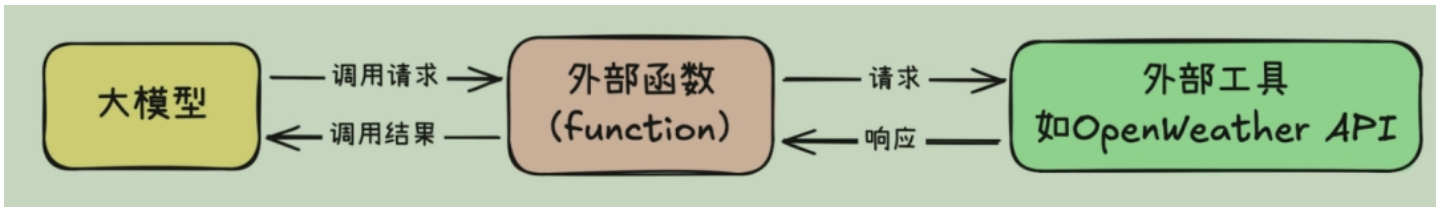
# MCP

## MCP

### 从 Function calling 到 MCP

MCP（Model Context Protocol），由Claude母公司Anthropic于24年11月正式提出。MCP是一种技术协议，是一种智能体Agent开发过程中共同约定的一种规范，如果大家都遵守一个规范，那么协作开发Agent的效率就能大幅提升。

MCP解决的最大痛点，就是Agent开发中调用外部工具的技术门槛过高的问题。由于大模型本身无法与外部工具直接通信，只能采用Function calling的方式（也叫tools），作为外部函数的中介：



编写Function calling函数工作量很大（随便一个函数就要100+行代码），并且为了让大模型理解这个函数，需要用Json Schema格式编写功能说明，并设计提示词模板。

```
1  JSON Schema 是一种用于描述和验证JSON数据结构的标准化格式，在Function calling中扮演函数接口说明书
2  的角色。其核心作用是为大模型提供精准的函数调用规范，确保模型生成的参数格式正确。下面的例子中，指明了
3  Function名字和功能，以及入参类型、参数可选值、是否必须和参数描述等信息。
4  {
5      "name": "get_weather",
6      "description": "查询指定地点的天气信息",
7      "parameters": {
8          "type": "object",
9          "properties": {
10             "location": {
11                 "type": "string",
12                 "description": "城市名称，如'北京'"
13             },
14             "unit": {
15                 "type": "string",
16                 "enum": ["celsius", "fahrenheit"],
17                 "description": "温度单位"
18             }
19         },
20     }
```

```

20     "required": ["location"]
21 }
22 }

```

1 提示词模板是预定义的**结构化指令框架**，用于引导大模型准确触发函数调用。其本质是通过工程设计，

2 **将自然语言指令转化为机器可解析的逻辑流**。例如：

3

4 你是一个智能天气助手，请按以下步骤响应用户：

5 1. **\*\*意图识别\*\***：判断用户是否在询问天气

6 2. **\*\*参数提取\*\***：

7     - 若需查询天气，提取地点(location)和单位(unit)

8     - 若未明确单位，默认使用摄氏制

9 3. **\*\*函数调用\*\***：严格按JSON格式返回调用指令：

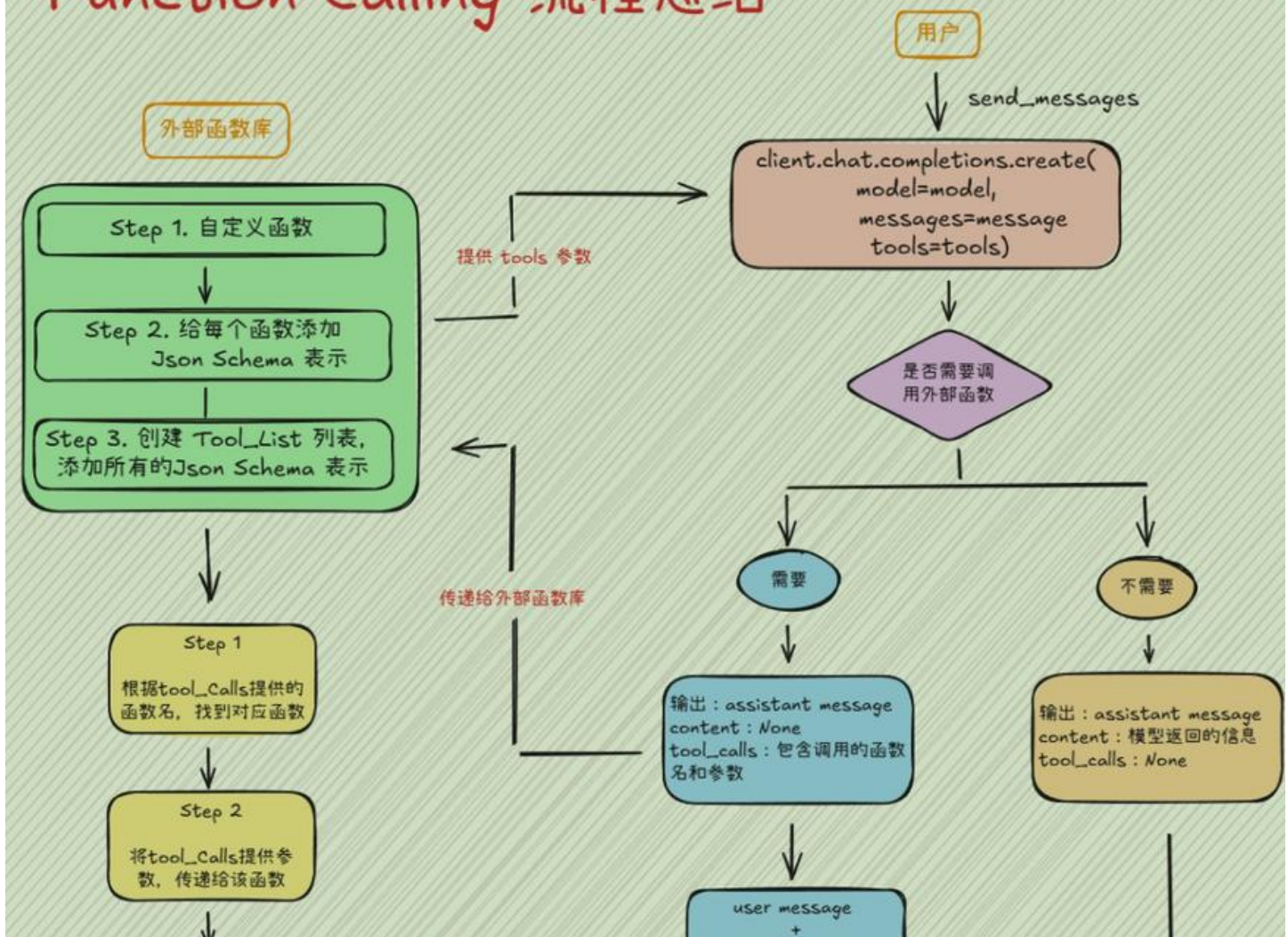
10 {

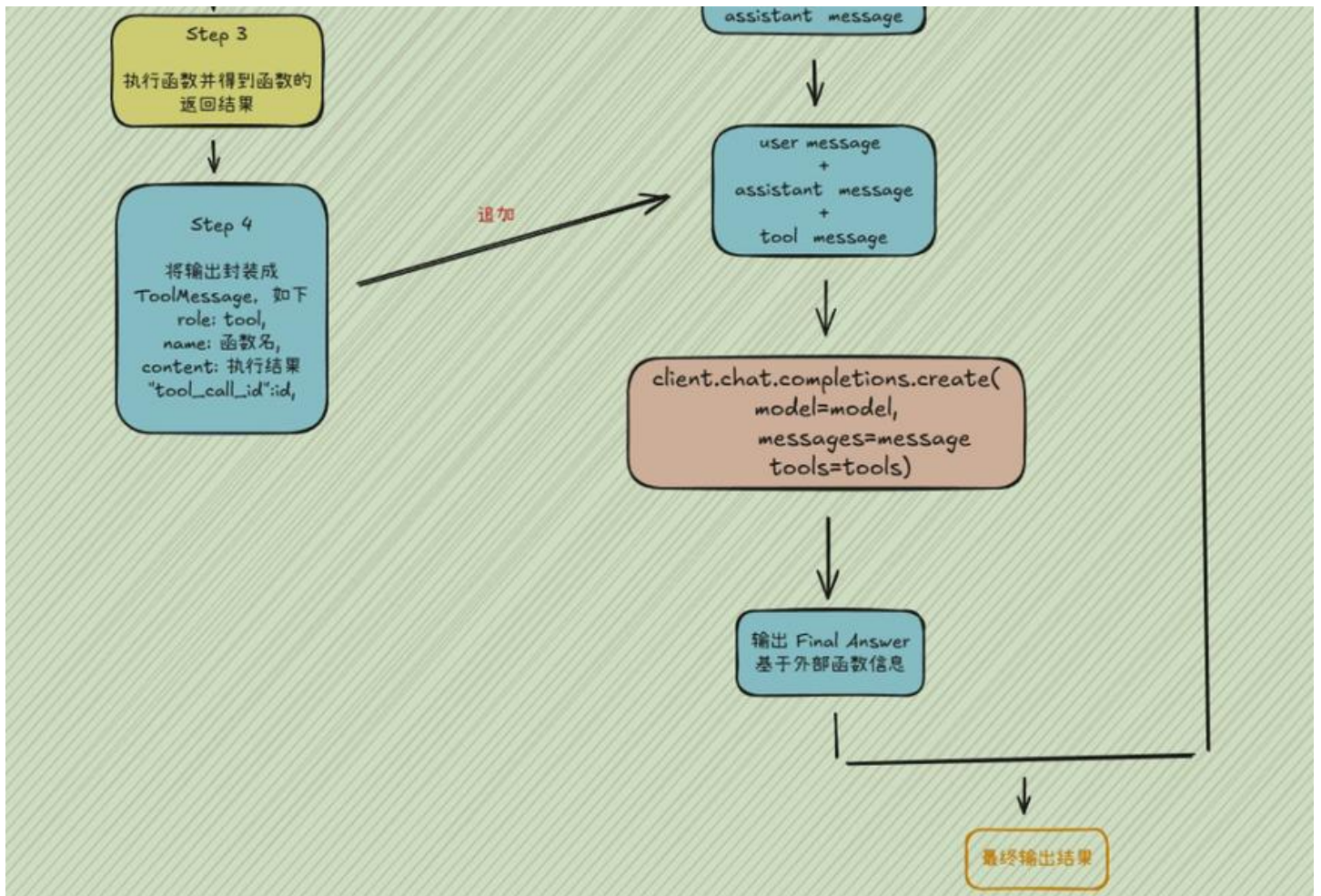
11     "function": "get\_weather",

12     "arguments": {"location": "北京", "unit": "celsius"}

13 }

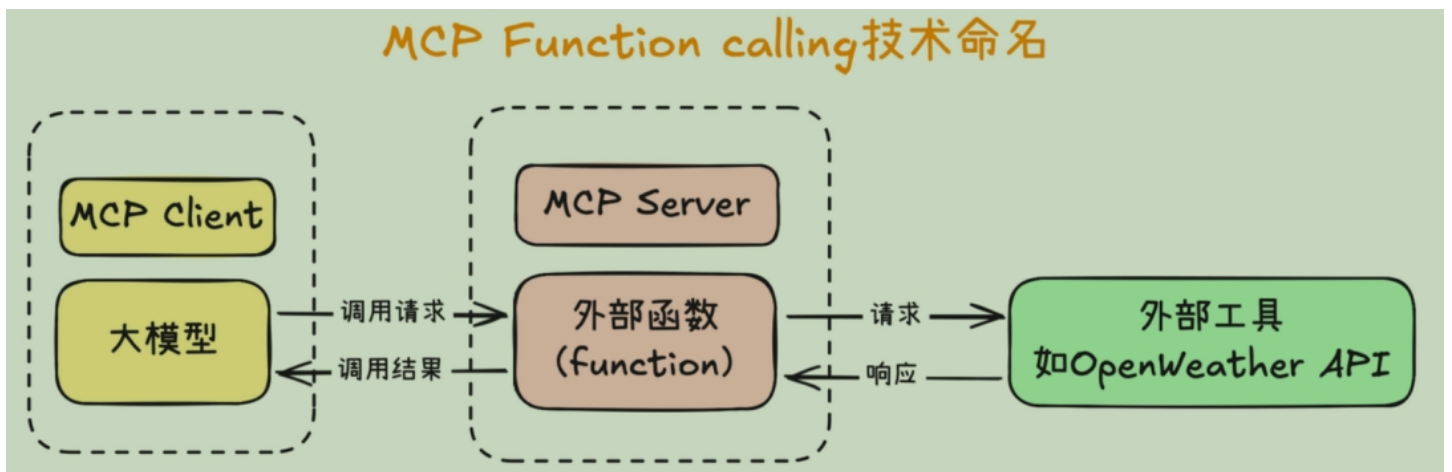
## Function Calling 流程总结





MCP统一了Function calling的运行规范：

- 首先是先统一名称，MCP把大模型运行环境称作 **MCP Client**，也就是MCP客户端，同时，把外部函数运行环境称作**MCP Server**，也就是MCP服务器。
- 然后，统一MCP客户端和服务器的运行规范，并且要求MCP客户端和服务端之间，也统一按照某个既定的提示词模板进行通信。



使用MCP的好处在于可以避免外部函数重复编写。像查询天气、网页爬取、查询本地MySQL数据库这种通用的需求，只需要开发一个服务器就好，后续的开发者的可以直接调用服务而不用重新实现。MCP 开发工具支持Python、TS和Java等多种语言。想要使用MCP服务器就要构建MCP客户端（支持任意本

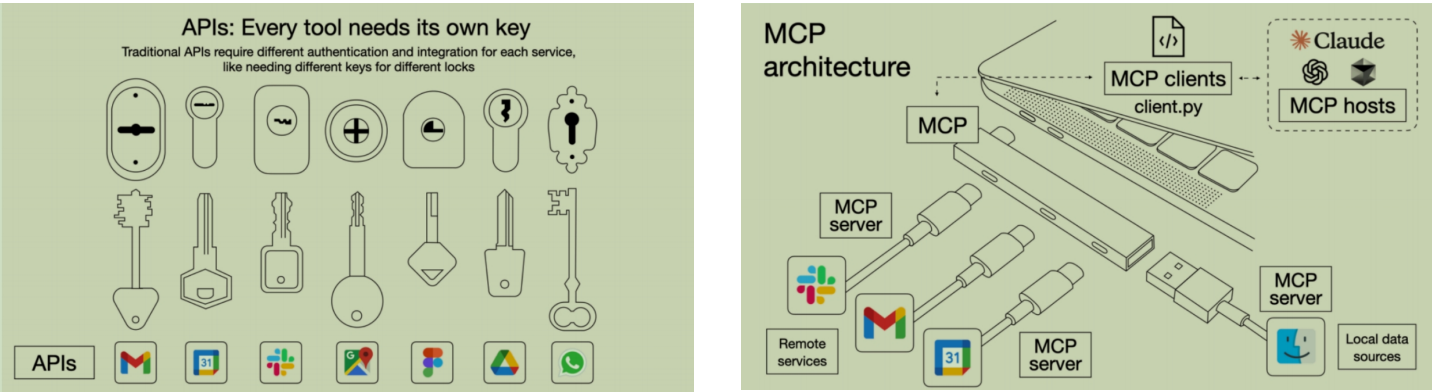


地和在线大模型，甚至是Cursor）。而如果没有所需要的MCP服务器，就要自己开发，下面的代码给出了一个简单的服务器示例

```
1  # server.py
2  from mcp.server.fastmcp import FastMCP
3  # Create an MCP server
4  mcp = FastMCP("Demo")
5  # Add an addition tool
6  @mcp.tool()
7  def add(a:int,b:int)-> int:
8      """Add two numbers"""
9      return a+ b
10 # Add a dynamic greeting resource
11 @mcp.resource("greeting://{name}")
12 def get_greeting(name:str)->str:
13     """Get a personalized greeting"""
14     return f"Hello, {name}!"
```

MCP针对agent的tools模块，目前不涉及memory和planning模块。

下图形象的对比了Function calling调用API和使用MCP的差异，MCP就像转接口，将多种多样的API封装成统一格式的mcp server，允许client端的大模型调用。



## MCP客户端

### uv环境管理

uv 是一个Python 依赖管理工具，类似于pip 和 conda，但它更快、更高效，并且可以更好地管理 Python 虚拟环境和依赖项。它的核心目标是替代 pip、venv 和 pip-tools，提供更好的性能和更低的管理开销。



## uv 的特点：

1. **速度更快**：相比 pip，uv 采用 Rust 编写，性能更优。
2. **支持 PEP 582**：无需 virtualenv，可以直接使用 pypackages 进行管理。
3. **兼容 pip**：支持 requirements.txt 和 pyproject.toml 依赖管理。
4. **替代 venv**：提供 uv venv 进行虚拟环境管理，比 venv 更轻量。
5. **跨平台**：支持 Windows、macOS 和 Linux。

```
1  # 安装uv
2  pip install uv
3  # 安装 Python 依赖，等效于pip install requests
4  uv pip install pandas
5  # 创建虚拟环境，等效于python -m venv myenv
6  uv venv myenv
7  # 激活虚拟环境
8  source myenv/bin/activate
9  # 安装所需的包，等效于pip install -r requirements.txt
10 uv pip install -r requirements.txt
11 # 运行 python 项目，等效于python script.py
12 uv run python script.py
```

## MCP客户端搭建

```
1  # 创建目录
2  uv init mcp-client
3  cd mcp-client
4  # 创建虚拟环境并激活
5  uv venv
6  source .venv/bin/activate
7  # 安装 MCP SDK
8  uv add mcp
```

创建一个简单的MCP客户端，核心功能有：

- 初始化 MCP 客户端
- 提供一个命令行交互界面
- 模拟 MCP 服务器连接
- 支持用户输入查询并返回「模拟回复」
- 支持安全退出

```

1  import asyncio # 让代码支持异步操作
2  from mcp import ClientSession # MCP 客户端会话管理
3  from contextlib import AsyncExitStack # 资源管理 (确保客户端关闭时释放资源)
4
5  class MCPClient:
6      def __init__(self):
7          """初始化 MCP 客户端"""
8          # 核心概念: 会话, 可以获取外部工具列表, 保存当前会话状态等功能, 暂时不链接MCP服
           务器
9          self.session = None
10         # 异步通信资源管理器
11         self.exit_stack = AsyncExitStack()
12     async def connect_to_mock_server(self):
13         """模拟 MCP 服务器的连接 (暂不连接真实服务器) """
14         print("✅ MCP 客户端已初始化, 但未连接到服务器")
15     async def chat_loop(self):
16         """运行交互式聊天循环"""
17         print("\nMCP 客户端已启动! 输入 'quit' 退出")
18         while True:
19             try:
20                 query = input("\nQuery: ").strip()
21                 if query.lower() == 'quit':
22                     break
23                 print(f"\n🗨️ [Mock Response] 你说的是: {query}")
24             except Exception as e:
25                 print(f"\n⚠️ 发生错误: {str(e)}")
26     async def cleanup(self):
27         """清理资源"""
28         await self.exit_stack.aclose() # 关闭资源管理器
29     async def main():
30         client = MCPClient() # 创建 MCP 客户端
31         try:
32             await client.connect_to_mock_server() # 连接 (模拟) 服务器
33             await client.chat_loop() # 进入聊天循环
34         finally:
35             await client.cleanup() # 确保退出时清理资源
36     if __name__ == "__main__":
37         asyncio.run(main())

```

## 接入在线模型

```

1  import asyncio
2  import os
3  from openai import OpenAI

```

```

4  from dotenv import load_dotenv
5  from contextlib import AsyncExitStack
6  # 加载 .env 文件, 确保 API Key 受到保护, 需要在.env文件中写入:
7      # BASE_URL="https://ai.devtool.tech/proxy/v1"
8      # MODEL=gpt-4o
9      # OPENAI_API_KEY="your_api_key"
10 load_dotenv()
11 class MCPClient:
12     def __init__(self):
13         """初始化 MCP 客户端"""
14         self.exit_stack = AsyncExitStack()
15         self.openai_api_key = os.getenv("OPENAI_API_KEY") # 读取 OpenAI API
Key
16         self.base_url = os.getenv("BASE_URL") # 读取 BASE URL
17         self.model = os.getenv("MODEL") # 读取 model
18         if not self.openai_api_key:
19             raise ValueError("❌ 未找到 OpenAI API Key, 请在 .env 文件中设置
OPENAI_API_KEY")
20         self.client = OpenAI(api_key=self.openai_api_key,
base_url=self.base_url)
21     async def process_query(self, query: str) -> str:
22         """调用 OpenAI API 处理用户查询"""
23         messages = [{"role": "system", "content": "你是一个智能助手, 帮助用户回答
问题。"},
24                     {"role": "user", "content": query}]
25         try:
26             # 调用 OpenAI API, 将 OpenAI API 变成异步任务, 防止程序卡顿。
27             response = await asyncio.get_event_loop().run_in_executor(
28                 None,
29                 lambda: self.client.chat.completions.create(
30                     model=self.model,
31                     messages=messages
32                 )
33             )
34             return response.choices[0].message.content
35         except Exception as e:
36             return f"⚠️ 调用 OpenAI API 时出错: {str(e)}"
37     async def chat_loop(self):
38         """运行交互式聊天循环"""
39         print("\n🤖 MCP 客户端已启动! 输入 'quit' 退出")
40         while True:
41             try:
42                 query = input("\n你: ").strip()
43                 if query.lower() == 'quit':
44                     break
45                 response = await self.process_query(query) # 发送用户输入到

```

OpenAI API

```

46         print(f"\n🤖 OpenAI: {response}")
47     except Exception as e:
48         print(f"\n⚠️ 发生错误: {str(e)}")
49     async def cleanup(self):
50         """清理资源"""
51         await self.exit_stack.aclose()
52     async def main():
53         client = MCPClient()
54         try:
55             await client.chat_loop()
56         finally:
57             await client.cleanup()
58     if __name__ == "__main__":
59         asyncio.run(main())

```

## 部署本地模型

### 使用vllm库部署QwQ-32B

模型比较大，采用huggingface担心网络不稳定，可以用modelscope下载模型

代码块

```

1  pip install modelscope
2  modelscope download --model Qwen/QwQ-32B --local_dir ./QwQ-32B

```

## 安装vllm库

代码块

```

1  pip install vllm

```

## 开启vllm

代码块

```

1  vllm serve ./QwQ-32B --max-model-len 32768 # 32k上下文单卡
2  CUDA_VISIBLE_DEVICES=0,1 vllm serve ./QwQ-32B --tensor-parallel-size 2 # 128k上下文双卡

```

在jupyter中运行以下代码：

代码块

```

1  from openai import OpenAI

```



```

2  openai_api_key = "EMPTY"
3  openai_api_base = "http://localhost:8000/v1"
4
5  client = OpenAI(
6      api_key=openai_api_key,
7      base_url=openai_api_base,
8  )
9  prompt = "在单词\"strawberry\"中，总共有几个R? "
10 messages = [
11     {"role": "user", "content": prompt}
12 ]
13 response = client.chat.completions.create(
14     model="QWQ-32B/",
15     messages=messages,
16 )
17
18 print(response.choices[0].message.content)

```

后端接收到以下请求：

```
INFO: 127.0.0.1:39040 - "POST /v1/chat/completions HTTP/1.1" 200 OK
```

## MCP服务器端

Server端可以提供以下三种标准能力：

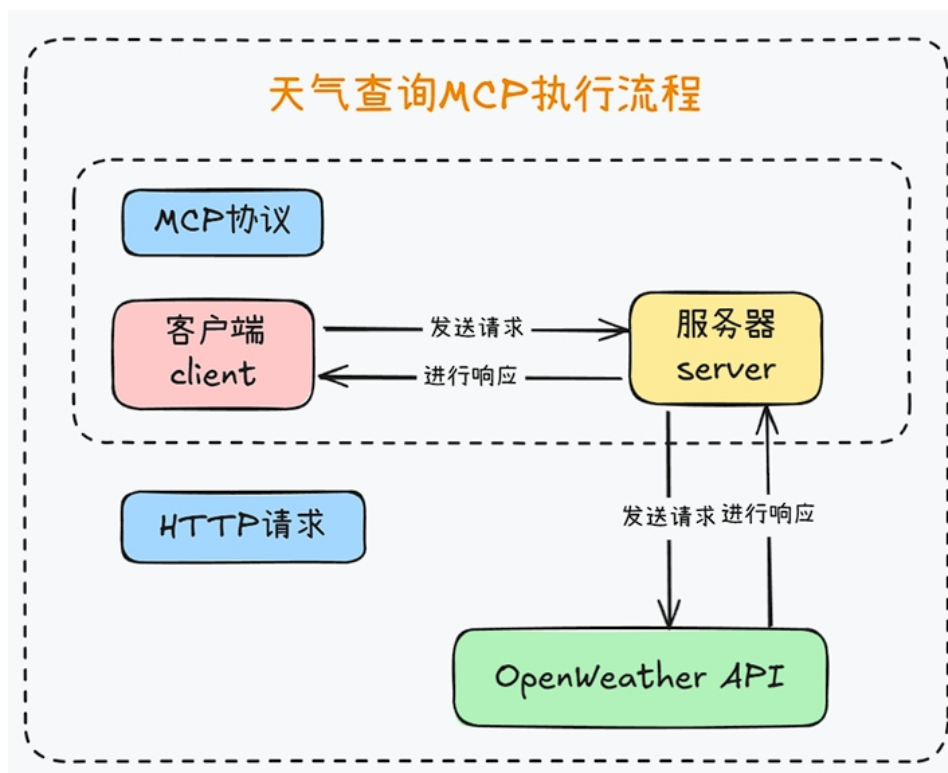
- **Resources：**资源，类似于文件数据读取，可以是文件资源或是API响应返回的内容。
- **Tools：**工具，第三方服务、功能函数，通过此可控制LLM可调用哪些函数。
- **Prompts：**提示词，为用户预先定义好的完成特定任务的模板。

## 通信机制

MCP目前支持两种传输方式：

- **标准输入输出（stdio）：**用于本地通信的传输方式。在这种模式下，MCP 客户端会将服务器程序作为子进程启动，双方通过标准输入（stdin）和标准输出（stdout）进行数据交换。这种方式适用于客户端和服务端在同一台机器上运行的场景，确保了高效、低延迟的通信。
- **HTTP+SSE：**适用于客户端和服务端位于不同物理位置的场景。在这种模式下，客户端和服务端通过 HTTP 协议进行通信，利用 SSE 实现服务端向客户端的实时数据推送。

## 天气查询服务器搭建



搭建了一个提供天气查询的工具的server，通过http请求查询天气。

代码块

```
1 uv add httpx
```

代码块

```
1 import json
2 import httpx
3 from typing import Any
4 from mcp.server.fastmcp import FastMCP
5 import os
6 from dotenv import load_dotenv
7
8 # 加载环境变量
9 load_dotenv(".env")
10
11 # 初始化 MCP 服务器
12 mcp = FastMCP("weatherServer")
13
14 # OpenWeather API 配置
15 OPENWEATHER_API_BASE = "https://api.openweathermap.org/data/2.5/weather"
16 API_KEY = os.getenv("OpenWeather_API_KEY") # 请替换为你自己的 OpenWeather API
    Key
17 USER_AGENT = "weather-app/1.0"
18
19 # 异步获取天气数据
```

```

20  async def fetch_weather(city: str) -> dict[str, Any] | None:
21      """
22      从 OpenWeather API 获取天气信息。
23      :param city: 城市名称（需使用英文，如 Beijing）
24      :return: 天气数据字典，若出错返回包含 error 信息的字典
25      """
26      params = {
27          "q": city,
28          "appid": API_KEY,
29          "units": "metric",
30          "lang": "zh_cn"
31      }
32      headers = {"User-Agent": USER_AGENT}
33      # 使用 httpx.AsyncClient() 发送异步 GET 请求到 OpenWeather API。
34      async with httpx.AsyncClient() as client:
35          try:
36              response = await client.get(OPENWEATHER_API_BASE, params=params,
headers=headers, timeout=30.0)
37              response.raise_for_status()
38              return response.json() # 返回字典类型
39          except httpx.HTTPStatusError as e:
40              return {"error": f"HTTP 错误: {e.response.status_code}"}
41          except Exception as e:
42              return {"error": f"请求失败: {str(e)}"}
43
44      # 天气数据格式化
45  def format_weather(data: dict[str, Any] | str) -> str:
46      """
47      将天气数据格式化为易读文本。
48      :param data: 天气数据（可以是字典或 JSON 字符串）
49      :return: 格式化后的天气信息字符串
50      """
51      # 如果传入的是字符串，则先转换为字典
52      if isinstance(data, str):
53          try:
54              data = json.loads(data)
55          except Exception as e:
56              return f"无法解析天气数据: {e}"
57
58      # 如果数据中含错误信息，直接返回错误提示
59      if "error" in data:
60          return f"{data['error']}"
61
62      # 提取数据做容错处理
63      city = data.get("name", "未知")
64      country = data.get("sys", {}).get("country", "未知")
65      temp = data.get("main", {}).get("temp", "N/A")

```

```

66     humidity = data.get("main", {}).get("humidity", "N/A")
67     wind_speed = data.get("wind", {}).get("speed", "N/A")
68
69     # weather 是一个列表, 因此此处用 [{}] 前先提供默认字典
70     weather_list = data.get("weather", [{}])
71     description = weather_list[0].get("description", "未知")
72
73     return (
74         f"🌍 {city}, {country}\n"
75         f"🌡️ 温度: {temp}°C\n"
76         f"💧 湿度: {humidity}%\n"
77         f"💨 风速: {wind_speed} m/s\n"
78         f"☁️ 天气: {description}\n"
79     )
80
81 @mcp.tool()
82 async def query_weather(city: str) -> str:
83     """
84     输入指定城市的英文名称, 返回今日天气查询结果。
85     :param city: 城市名称 (需使用英文)
86     :return: 格式化后的天气信息

```

接下来实现一个与调用这个server的client端，以与大模型对话的形式呈现，只有问天气查询的问题时调用工具，否则就是与大模型对话。

代码块

```

1  import asyncio
2  import os
3  import json
4  import sys
5  from typing import Optional
6  from contextlib import AsyncExitStack
7
8  from openai import OpenAI
9  from dotenv import load_dotenv
10
11 from mcp import ClientSession, StdioServerParameters
12 from mcp.client.stdio import stdio_client
13
14 # 加载环境变量
15 load_dotenv()
16
17 class MCPClient:
18     def __init__(self):
19         """初始化 MCP 客户端"""
20         self.exit_stack = AsyncExitStack() # 统一管理异步上下文 (如 MCP 连接) 的生命

```

周期。可以在退出 ( cleanup ) 时自动关闭

```

21     self.openai_api_key = os.getenv("OPENAI_API_KEY")
22     self.base_url = os.getenv("BASE_URL")
23     self.model = os.getenv("MODEL")
24
25     if not self.openai_api_key:
26         raise ValueError(" 未找到 OpenAI API Key, 请在 .env 文件中设置
OPENAI_API_KEY")
27
28     self.client = OpenAI(api_key=self.openai_api_key,
base_url=self.base_url)
29     self.session: Optional[ClientSession] = None # 用于保存 MCP 的客户端会话,
默认是 None , 稍后通过 connect_to_server 进行连接
30
31     async def connect_to_server(self, server_script_path: str):
32         """连接到 MCP 服务器并列出可用工具"""
33         is_python = server_script_path.endswith('.py')
34         is_js = server_script_path.endswith('.js')
35
36         if not (is_python or is_js):
37             raise ValueError("服务器脚本必须是 .py 或 .js 文件")
38         # 判断服务器脚本是 Python 还是 Node.js, 选择对应的运行命令。
39         command = "python" if is_python else "node"
40         server_params = StdioServerParameters(
41             command=command,
42             args=[server_script_path],
43             env=None
44         )
45
46         # 启动服务器连接
47         stdio_transport = await self.exit_stack.enter_async_context(
48             stdio_client(server_params)
49         )
50         self.stdio, self.write = stdio_transport
51         self.session = await self.exit_stack.enter_async_context(
52             ClientSession(self.stdio, self.write)
53         ) # 发送初始化消息给服务器, 等待服务器就绪。
54
55         await self.session.initialize()
56
57         # 列出可用工具
58         response = await self.session.list_tools()
59         print("\n已连接到服务器, 支持以下工具:", [tool.name for tool in
response.tools])
60
61     async def process_query(self, query: str) -> str:
62         """使用大模型处理查询并调用工具"""
63         messages = [{"role": "user", "content": query}]

```



```

64         response = await self.session.list_tools()
65
66         # 构建可用工具列表
67         available_tools = [
68             {
69                 "type": "function",
70                 "function": {
71                     "name": tool.name, # 工具的名字
72                     "description": tool.description, # 外部函数的描述
73                     "parameters": tool.inputSchema # 如果要调用这个函数，需要的
74                                     json_schema说明
75                 }
76             } for tool in response.tools
77         ]
78
79         # 第一次模型调用
80         response = self.client.chat.completions.create(
81             model=self.model,
82             messages=messages,
83             tools=available_tools

```

代码块

```
1 uv run client.py server.py
```



### MCPClient 的主要职责：

- 启动 MCP 服务器（通过 StdioServerParameters）
- 建立 MCP 会话，列出可用工具
- 处理用户输入，将其发送给 OpenAI 模型
- 如果模型想调用 MCP 工具（Function Calling），就执行 call\_tool
- 将结果重新发给模型，并返回最终回答

## 测试服务器

Anthropic提供了一个非常便捷的debug工具：Inspector。借助Inspector，我们能够非常快捷的调用各类server，并测试其功能。

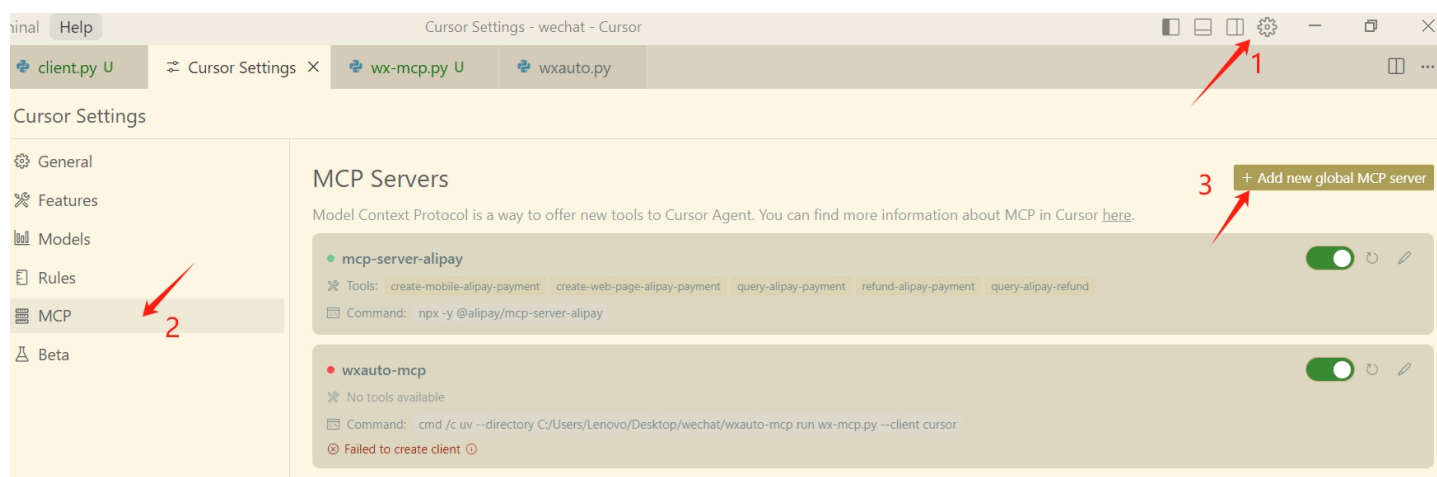
代码块

```
1 mcp dev xx.py
```

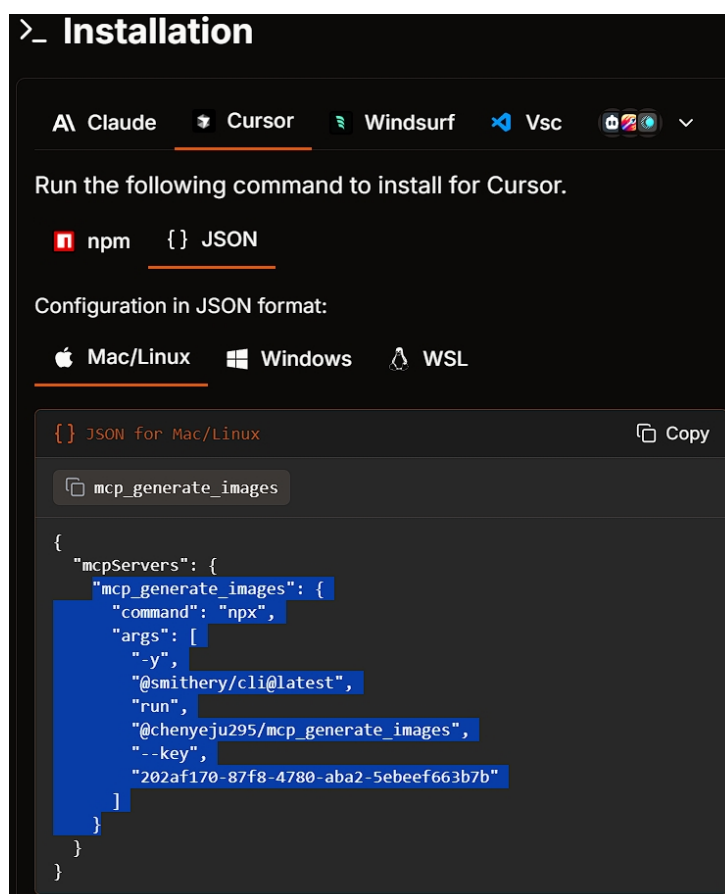
## 在线服务器导航：

- MCP官方服务器合集：<https://github.com/modelcontextprotocol/servers>
- MCP Github热门导航：<https://github.com/punkpeye/awesome-mcp-servers>
- <https://mcp.so/>
- <https://mcpservers.cn/>
- <https://smithery.ai/>

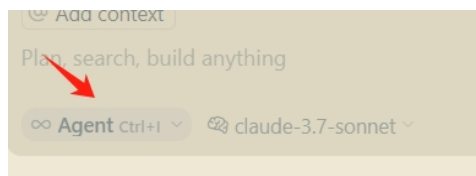
以cursor举例，将mcp server的功能添加到大模型之中。



将对应的json配置信息添加到mcp.json文件中，在上面页面刷新，变绿即添加成功。



调用大模型时要选择agent模式



- MCP是一种技术协议，是一种智能体Agent开发过程中共同约定的一种规范，如果大家都遵守一个规范，那么协作开发Agent的效率就能大幅提升。
- 首先是先统一名称，MCP把大模型运行环境称作 **MCP Client**，也就是MCP客户端，负责接收用户请求、调用MCP Server工具并整合结果。同时，把外部函数运行环境称作**MCP Server**，也就是MCP服务器，负责封装API、数据库等资源。
- MCP统一了客户端和服务器的运行规范，并且要求MCP客户端和服务器之间，也统一按照某个既定的提示词模板进行通信。
- 我实现的用MCP实现秒回功能
  - 视频：<https://mp.weixin.qq.com/s/0EILiRUgHlNtPrCdDJm28g>
  - 代码：<https://github.com/saintGeorge13/wx-mcp/tree/main>

代码链接：<https://github.com/saintGeorge13/wx-mcp/tree/main>

其他在线服务器导航：

- MCP官方服务器合集：<https://github.com/modelcontextprotocol/servers>
- MCP Github热门导航：<https://github.com/punkpeye/awesome-mcp-servers>
- <https://mcp.so/>
- <https://mcpservers.cn/>
- <https://smithery.ai/>

代码块

```
1 pip install uv
2 pip install mcp
3 pip install wxauto
4
5 # 运行inspector
```

```
6 mcp dev wx-mcp.py
7
8 # 本地mcp服务
9 uv run client.py wx-mcp.py
```

向你提问：

密欧哥你好，我现在在一家初创公司做多模态大模型实习生，方向是agent rag，显卡资源挺充足的，mentor实力也很强，mentor说可以用卡做自己的事情，我也非常想有自己的论文产出。我的问题是：我还没有发过论文，想发大模型的论文的话大概是什么流程呀？有什么建议吗？感谢！！

Web 版仅支持文字回答，如需语音请移步 App 作答  
有合伙人/嘉宾更擅长这个问题？不如 @ 他发表回答吧～

提问者正等待您的回答...

建议分为以下几步：

- 1. 确定研究方向：**首先明确自己的研究方向，agent 和 rag都是很大的概念，比如rag又有 chunking、embedding、rerank等很多细分方向。如何确定研究方向？最好是来自你工作中的尝试，如果你在某个子方向做出了能提升效果的尝试，就可以在这个领域深入研究。
- 2. 阅读相关论文：**当你已经在某个方向做出一定尝试，接下来就要找这方面的论文看，一是为了避免你的idea已经被别人用了，二是对自己的方法做一定的完善，三是为了以后写related work。看论文一定要借助ai，让它快速总结一篇论文的核心思想，自己再判断要不要深入的看。一篇论文看完一定要记录笔记，不然一觉醒来就忘完了，重点总结这篇论文解决了什么问题？创新点是什么？论文中的方法能不能与自己的方法结合？
- 3. 运行实验：**如果你有比较熟悉的代码框架，比如llamaindex、langchain等，那可以直接在过去代码的基础上实现自己的方法。如果没有，可以在其他论文的源码的基础上改。实验数据集一定是公开数据集，方便跟其他方法对比（除非你的工作就是造数据集）。对照方法找近期的和经典的方法，最好有开源代码，跑出来的实验结果不如他们论文中的数据也没关系，一切以自己为主。
- 4. 写论文：**一定用latex格式写，论文模板找目标会议/期刊官方提供的模板。论文结构一般是 introduction+related work+method+experiments+conclusion，论文长度参考往年的论文。写论文的核心是编好故事，这里也建议用ai辅助。

过程中有任何进展和卡点都及时与mentor沟通，论文大概率是要给他挂二作/通讯的，不问白不问。

**投稿会议：**中国计算机学会评级<https://www.ccf.org.cn/ccftgjxskwml/>