

# 9. Agent +RAG

AI Agent 一般包含以下几个**关键模块**：

1. 计划模块（Planning）

- 负责理解目标任务
- 拆解任务为可执行的步骤
- 生成行动计划

2. 记忆模块（Memory）

- 存储与任务相关的历史记录、上下文信息
- 支持长期记忆与短期记忆
- 提高连续对话与任务一致性

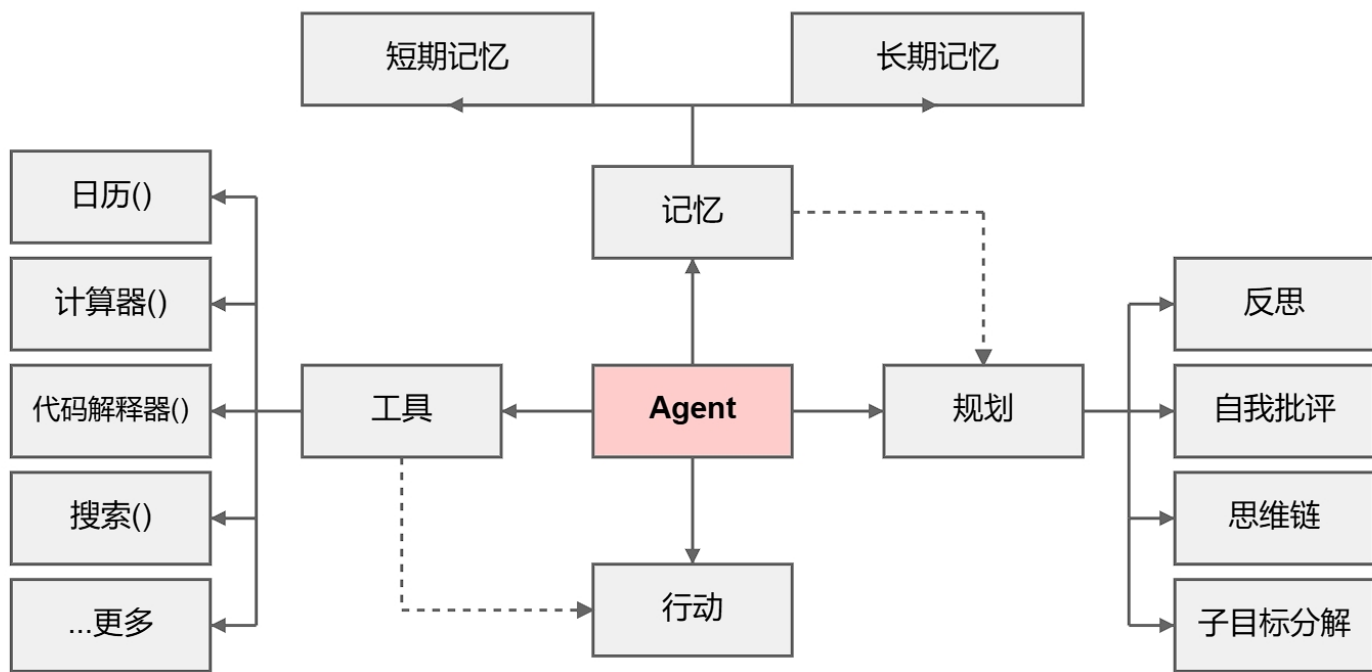
3. 工具使用模块（Tools）

- 能调用各种外部工具，如：
  - 网络搜索引擎
  - API 接口
  - 计算器、代码执行器

- 实现感知-行动闭环

4. 执行模块（Executor）

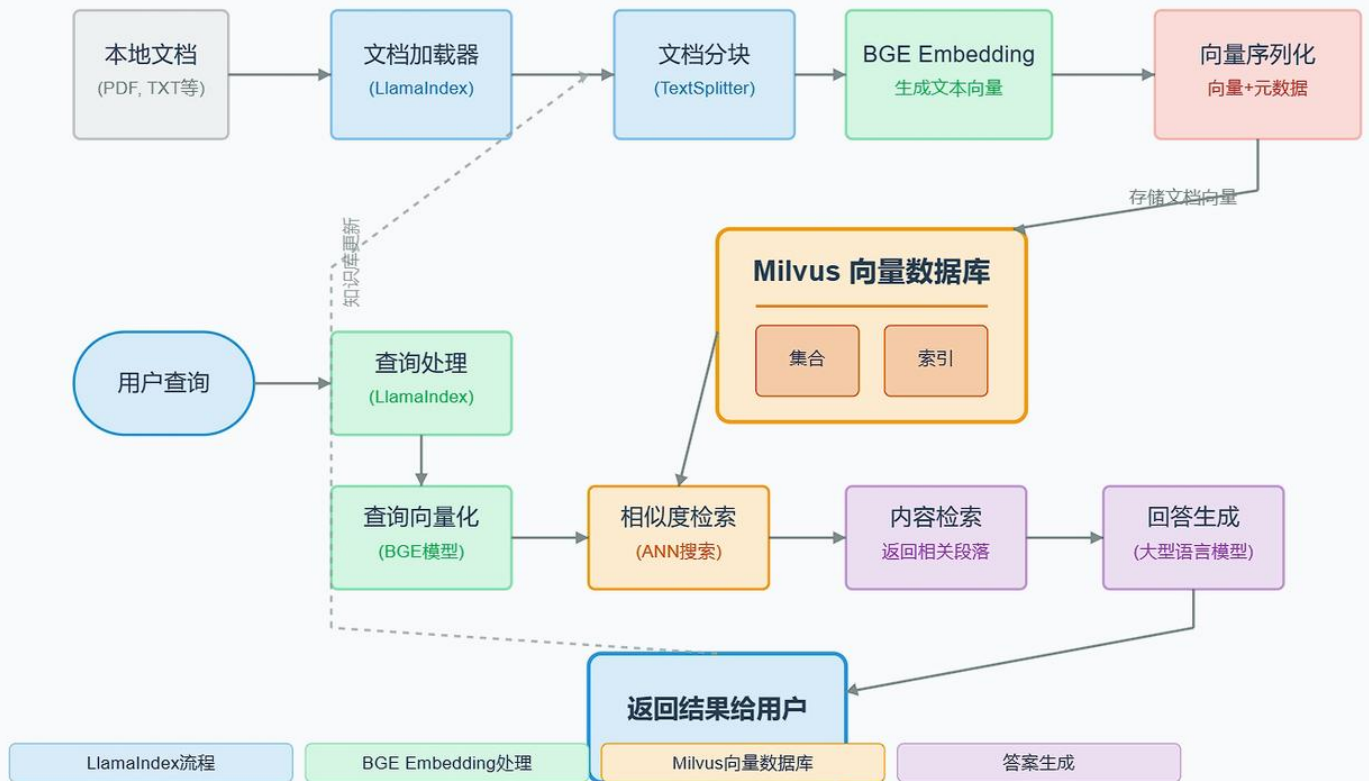
- 根据计划执行每一步操作
- 反馈执行结果并调整策略（如果需要）



## 本地文档检索RAG

# LlamaIndex + Milvus + BGE 的 RAG 系统架构

## 系统流程



## 文档解析

利用数据解析云服务 **llamacloud** (<https://cloud.llamaindex.ai/>)，异步的解析文档。

### 代码块

```
1  from llama_index.core import VectorStoreIndex, SimpleDirectoryReader
2  import nest_asyncio
3
4  nest_asyncio.apply()
5
6  from llama_cloud_services import LlamaParse
7
8  parser = LlamaParse(
9      api_key="YOUR_LLAMA_CLOUD_API_KEY", # can also be set in your env as
      LLAMA_CLOUD_API_KEY
10     result_type="markdown", # "markdown" and "text" are available
11     num_workers=3, # if multiple files passed, split in `num_workers` API
      calls
12     verbose=True,
13     language="ch_sim", # Optionally you can define a language, default=en
14 )
15
16 file_extractor = {".pdf": parser}
17 # 返回一个列表
```

```
18 documents_cloud = SimpleDirectoryReader(  
19     "./data", file_extractor=file_extractor  
20 ).load_data()  
21 # 汇总所有文档  
22 all_doc = ''  
23 for doc in documents_cloud:  
24     all_doc += doc.text
```

## 切块

代码块

```
1 def chunk_text(text, max_chunk_size=300):  
2     """  
3     将文本切分成固定大小的块，保持段落和表格的完整性。  
4  
5     Args:  
6         text (str): 需要切分的文本  
7         max_chunk_size (int): 每个块的最大大小，默认为300  
8  
9     Returns:  
10        list: 包含切分后文本块的列表  
11    """  
12    # 按行分割文本  
13    lines = text.split('\n')  
14  
15    chunks = []  
16    current_chunk = []  
17    current_size = 0  
18    in_table = False  
19    table_content = []  
20  
21    for line in lines:  
22        # 检测表格开始 (含有 | 字符的行, 通常是表格)  
23        if '|' in line and not in_table:  
24            # 可能是表格的开始  
25            in_table = True  
26  
27            # 保存当前块  
28            if current_chunk:  
29                chunks.append('\n'.join(current_chunk))  
30                current_chunk = []  
31                current_size = 0  
32  
33            table_content.append(line)  
34
```

```

35         # 表格内容
36         elif in_table:
37             table_content.append(line)
38
39             # 如果遇到空行，可能表示表格结束
40             if not line.strip():
41                 in_table = False
42                 chunks.append('\n'.join(table_content))
43                 table_content = []
44
45         # 普通内容
46         else:
47             line_length = len(line)
48
49             # 如果加上这一行会超过最大块大小，并且当前块不为空，则保存当前块
50             if current_size + line_length > max_chunk_size and current_chunk:
51                 chunks.append('\n'.join(current_chunk))
52                 current_chunk = []
53                 current_size = 0
54
55             # 添加新行到当前块
56             current_chunk.append(line)
57             current_size += line_length
58
59             # 如果当前行为空，可能是段落结束
60             # 检查当前块大小，如果已经足够大，则保存当前块
61             if not line.strip() and current_size > max_chunk_size // 2:
62                 chunks.append('\n'.join(current_chunk))
63                 current_chunk = []
64                 current_size = 0
65
66         # 处理剩余的表格内容
67         if in_table and table_content:
68             chunks.append('\n'.join(table_content))
69
70         # 处理最后剩余的块
71         if current_chunk:
72             chunks.append('\n'.join(current_chunk))
73
74     return chunks
75 chunks = chunk_text(all_doc, max_chunk_size=300)

```

## 插入Milvus数据库

利用Milvus和BGE-M3 模型进行混合搜索。BGE-M3 模型可以将文本转换为密集向量和稀疏向量。Milvus 支持在一个 Collections 中存储这两种向量，从而可以进行混合搜索，增强搜索结果的相关性。

#### 代码块

```
1  # 下载bge模型
2  import os
3  os.environ["HF_ENDPOINT"] = "https://hf-mirror.com"
4  from modelscope import snapshot_download
5  model_dir = snapshot_download('BAAI/bge-m3',
6  cache_dir='./', revision='master')
7
8  # 创建embedding模型
9  from milvus_model.hybrid import BGEM3EmbeddingFunction
10 ef = BGEM3EmbeddingFunction(model_name_or_path="./BAAI/bge-m3",
11 use_fp16=False, device="cpu")
12 dense_dim = ef.dim["dense"]
13 # 用bge-M3模型的embedding结果
14 chunks_embeddings = ef(chunks)
```

- 设置 Milvus Collections 和索引

#### 代码块

```
1  from pymilvus import (
2      connections,
3      utility,
4      FieldSchema,
5      CollectionSchema,
6      DataType,
7      Collection,
8  )
9  # 连接milvus数据库
10 connections.connect(uri="./milvus.db")
11
12 fields = [
13     # Use auto generated id as primary key, 类似数据库中的主键
14     FieldSchema(
15         name="pk", dtype=DataType.VARCHAR, is_primary=True, auto_id=True,
16         max_length=100
17     ),
18     # Store the original text to retrieve based on semantically distance
19     FieldSchema(name="text", dtype=DataType.VARCHAR, max_length=512),
20     # 支持稀疏检索和稠密检索
21     # Milvus now supports both sparse and dense vectors,
22     # we can store each in a separate field to conduct hybrid search on both
23     # vectors
24     FieldSchema(name="sparse_vector", dtype=DataType.SPARSE_FLOAT_VECTOR),
25     FieldSchema(name="dense_vector", dtype=DataType.FLOAT_VECTOR,
26         dim=dense_dim),
27 ]
```

```

24 ]
25 schema = CollectionSchema(fields)
26
27 col_name = "hybrid_demo"
28 if utility.has_collection(col_name):
29     Collection(col_name).drop()
30 col = Collection(col_name, schema, consistency_level="Strong")
31 # 创建索引
32 sparse_index = {"index_type": "SPARSE_INVERTED_INDEX", "metric_type": "IP"}
33 col.create_index("sparse_vector", sparse_index)
34 dense_index = {"index_type": "AUTOINDEX", "metric_type": "IP"}
35 col.create_index("dense_vector", dense_index)
36 col.load()
37

```

## 插入向量数据库

代码块

```

1  for i in range(0, len(chunks), 50):
2      batched_entities = [
3          chunks[i : i + 50],
4          chunks_embeddings["sparse"][i : i + 50],
5          chunks_embeddings["dense"][i : i + 50],
6      ]
7      col.insert(batched_entities)
8  print("Number of entities inserted:", col.num_entities)

```

## 检索

定义密集检索、稀疏检索和混合检索函数

代码块

```

1  from pymilvus import (
2      AnnSearchRequest,
3      WeightedRanker,
4  )
5
6  # 密集检索
7  def dense_search(col, query_dense_embedding, limit=10):
8      search_params = {"metric_type": "IP", "params": {}}
9      res = col.search(
10         [query_dense_embedding],
11         anns_field="dense_vector",
12         limit=limit,

```

```

13         output_fields=["text"],
14         param=search_params,
15     )[0]
16     return [hit.get("text") for hit in res]
17
18 # 稀疏检索
19 def sparse_search(col, query_sparse_embedding, limit=10):
20     search_params = {
21         "metric_type": "IP",
22         "params": {},
23     }
24     res = col.search(
25         [query_sparse_embedding],
26         anns_field="sparse_vector",
27         limit=limit,
28         output_fields=["text"],
29         param=search_params,
30     )[0]
31     return [hit.get("text") for hit in res]
32
33 # 混合检索
34 def hybrid_search(
35     col,
36     query_dense_embedding,
37     query_sparse_embedding,
38     sparse_weight=1.0,
39     dense_weight=1.0,
40     limit=10,
41 ):
42     dense_search_params = {"metric_type": "IP", "params": {}}
43     dense_req = AnnSearchRequest(
44         [query_dense_embedding], "dense_vector", dense_search_params,
45         limit=limit
46     )
47     sparse_search_params = {"metric_type": "IP", "params": {}}
48     sparse_req = AnnSearchRequest(
49         [query_sparse_embedding], "sparse_vector", sparse_search_params,
50         limit=limit
51     )
52     rerank = WeightedRanker(sparse_weight, dense_weight)
53     res = col.hybrid_search(
54         [sparse_req, dense_req], rerank=rerank, limit=limit, output_fields=
55         ["text"]
56     )[0]
57     return [hit.get("text") for hit in res]

```

混合检索结合了密集检索和稀疏检索的优势，因此后续只采用混合检索的结果。

代码块

```
1 query = '什么是以人为本的座舱'
2 query_embeddings = ef([query])
3 dense_results = dense_search(col, query_embeddings["dense"][0])
4 sparse_results = sparse_search(col, query_embeddings["sparse"]._getrow(0))
5 hybrid_results = hybrid_search(
6     col,
7     query_embeddings["dense"][0],
8     query_embeddings["sparse"]._getrow(0),
9     sparse_weight=0.7,
10    dense_weight=1.0,
11 )
12 print(dense_results[:3])
13 print(sparse_results[:3])
14 print(hybrid_results[:3])
```

对检索结果添加序号

代码块

```
1 def format_list_with_markers(input_list):
2     """
3     Formats a list by adding numbered markers ([1], [2], etc.) to each item
4     and adding a line break after each item.
5
6     Args:
7         input_list: A list of strings to format
8
9     Returns:
10        A single formatted string
11    """
12    result = ""
13
14    for index, item in enumerate(input_list, 1):
15        # Add the numbered marker and the item content
16        result += f"[{index}] {item}\n"
17
18    return result
19 formatted_references=format_list_with_markers(hybrid_results)
20 pprint(formatted_references)
```

生成



## 定义用户prompt

代码块

```
1  prompt = f"""
2  你是一个智能助手，负责根据用户的问题和提供的参考内容生成回答。请严格按照以下要求生成回答：
3  1. 回答必须基于提供的参考内容。
4  2. 在回答中，每一块内容都必须标注引用的来源，格式为：[引用编号]。例如：[1] 表示引用自第1
   条参考内容。
5  3. 如果没有参考内容，请明确说明。如果没有参考知识，根据你自己的知识进行回答
6
7  参考内容：
8  {formatted_references}
9
10  用户问题：{query}
11  """
12
13  print(prompt)
```

调用gpt-4o接口，将检索到的内容添加到用户prompt中。

代码块

```
1  from openai import OpenAI
2  client_gpt = OpenAI(api_key="your_openai_api_key",
   base_url="https://api.openai-proxy.org/v1")
3
4  completion_gpt = client_gpt.chat.completions.create(
5      model="gpt-4o",
6      messages=[
7          {
8              "role": "user",
9              "content": prompt
10         }
11     ]
12 )
13
14  print(completion_gpt.choices[0].message.content)
```

## 封装成RAG接口

代码块

```
1  def retrieval(query):
2      query_embeddings = ef([query])
```

```

3     hybrid_results = hybrid_search(
4         col,
5         query_embeddings["dense"][0],
6         query_embeddings["sparse"]._getrow(0),
7         sparse_weight=0.7,
8         dense_weight=1.0,
9     )
10    return hybrid_results

```

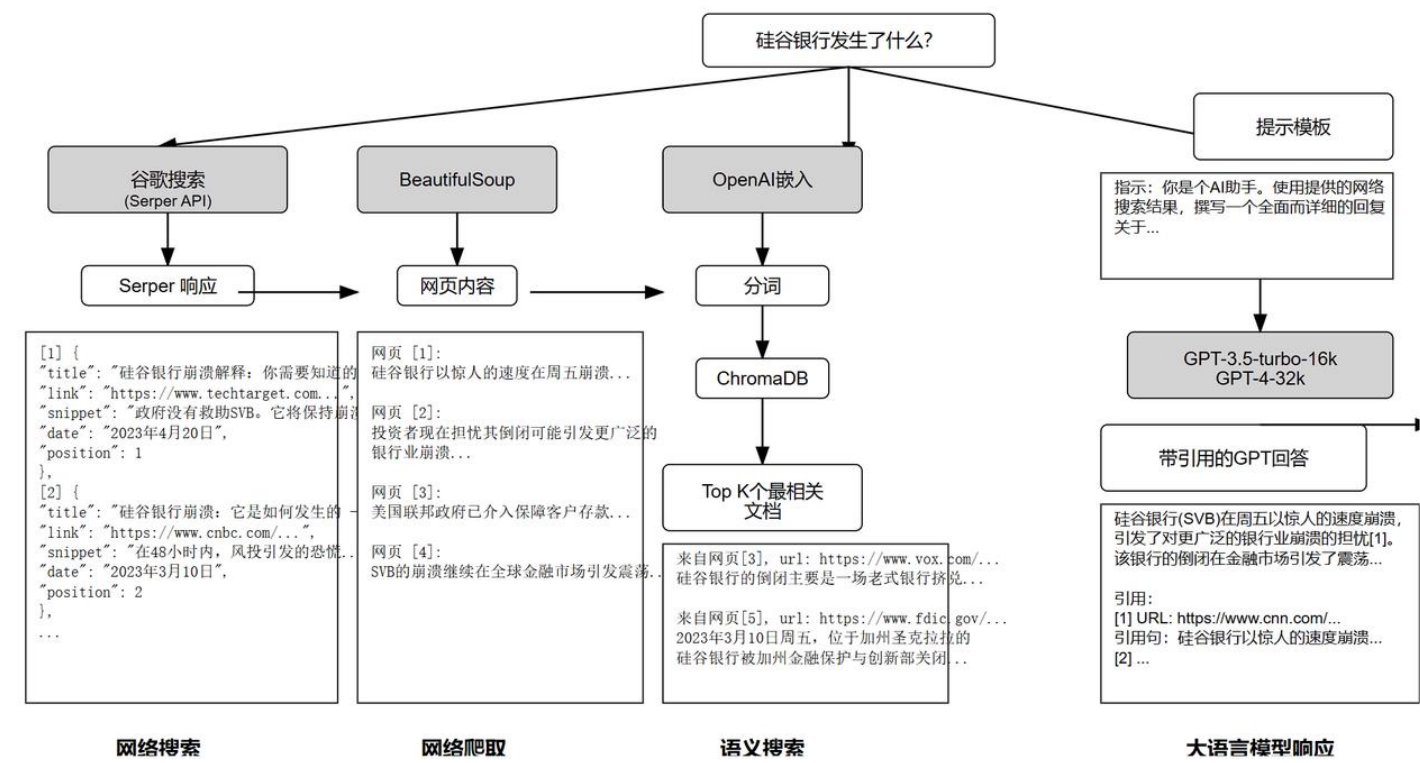
#### 代码块

```

1    def rag(query):
2
3        formatted_references=format_list_with_markers(retrieval(query))
4
5        prompt = f"""
6        你是一个智能助手，负责根据用户的问题和提供的参考内容生成回答。请严格按照以下要求生成回
7        答：
8        1. 回答必须基于提供的参考内容。
9        2. 在回答中，每一块内容都必须标注引用的来源，格式为：[引用编号]。例如：[1] 表示引用自
10       第1条参考内容。
11       3. 如果没有参考内容，请明确说明。如果没有参考知识，根据你自己的知识进行回答
12
13       参考内容：
14       {formatted_references}
15
16       用户问题：{query}
17       """
18       client = OpenAI(api_key="your_openai_api_key",
19                       base_url="https://api.openai-proxy.org/v1")
20
21       completion = client.chat.completions.create(
22           model="gpt-4o",
23           messages=[
24               {
25                   "role": "user",
26                   "content": prompt
27               }
28           ]
29       )
30
31       return formatted_references,completion.choices[0].message.content
32
33 reference,result=rag('介绍一下技术规格')
34 print(result)

```

# 网络检索



网络检索代码的实现在“websearch”文件夹中，这里只调用网络检索的接口。

代码块

```
1 from websearch.src.retrieval import web_retrieval
2 def web_search_answer(query):
3
4     web_reference=web_retrieval(query)
5
6     prompt = f"""
7
8     Web 搜索结果：
9     {web_reference}
10
11     指令：你是一个/一名销售助手。使用提供的网络搜索结果，对给定的查询写一个全面而详细的回
12     复。
13     确保在引用后使用 [number] 标记引用。
14     在回答的最后，列出带索引的相应参考文献，每个参考文献包含网络搜索结果中的网址和引用句
15     子，按照您上面标记的顺序排列，这些句子应该与网络搜索结果中的完全相同。
16     以下是参考文献的示例：
17         [1] 网址：https://www.pocketgamer.biz/news/81670/tencent-and-netease-
18         dominated-among-chinas-top-developers-in-q1/
19         引用句子：腾讯在本季度占据了国内市场收入的大约50%，相比之下2022年第一季度为
20         40%。
21
22     查询：{query}
```

```

19
20     请你过滤掉检索结果中的不相关性的部分进行回答，如果没有相关的部分，按照你自己的知识进行
    回答
21     """
22     client = OpenAI(api_key="your_openai_api_key",
    base_url="https://api.openai-proxy.org/v1")
23
24     completion = client.chat.completions.create(
25         model="gpt-4o",
26         messages=[
27             {
28                 "role": "user",
29                 "content": prompt
30             }
31         ]
32     )
33
34
35     return completion.choices[0].message.content
36 query="小米SU7的发动机"
37 print(web_search_answer(query))

```

## Agent

### planning模块

定义planning模块，根据用户问题，决定使用哪个工具来查询以获得更多需要的信息（本地文档搜索或网络搜索）

代码块

```

1  from openai import OpenAI
2  import os
3
4  def middle_json_model(prompt):
5
6      client = OpenAI(
7          # 若没有配置环境变量，请用百炼API Key将下行替换为：api_key="sk-xxx",
8          api_key='your_bailian_api_key',
9          base_url="https://dashscope.aliyuncs.com/compatible-mode/v1",
10     )
11     completion = client.chat.completions.create(
12         model="qwen-plus", # 此处以qwen-plus为例，可按需更换模型名称。模型列表：
13         https://help.aliyun.com/zh/model-studio/getting-started/models
14         messages=[

```

```

14         {'role': 'system', 'content': 'You are a helpful assistant.'},
15         {'role': 'user', 'content': prompt}],
16         response_format={"type": "json_object"}
17     )
18
19     return completion.choices[0].message.content

```

## 提取有效的json输出

### 代码块

```

1  import re
2  import json
3
4  def extract_json_content(input_str):
5      """
6      提取字符串中第一个 "[" 和最后一个 "]" 之间的内容（包括中括号）
7
8      Args:
9          input_str (str): 需要处理的输入字符串
10
11      Returns:
12          str or None: 提取的JSON内容，如果没有匹配则返回None
13      """
14      # 使用正则表达式匹配第一个 "[" 到最后一个 "]" 之间的内容
15      # [\s\S]* 匹配任意字符（包括换行符）
16      pattern = r'(\[[\s\S]*\])'
17      match = re.search(pattern, input_str)
18
19      # 如果匹配成功，返回匹配的内容；否则返回None
20      return match.group(1) if match else None

```

## 测试planning模块

### 代码块

```

1  #规划模块plan
2  import json
3
4  def agent_plan(query):
5      prompt= '''
6
7      你是一个专业的汽车销售助手的规划模块。你的任务是：
8      1. 分析用户的查询:{0}
9      2. 基于已有的信息，决定使用哪个工具来查询以获得更多需要的信息（本地文档搜索或网络搜索）
10     3. 将用户的原始查询拆解或延伸为1-2个相关问题，以获取更全面的信息

```

```
11
12
13 ## 可用工具
14 1. **本地文档搜索**：搜索本地星辰电动ES9的文档，包含以下章节：
15     - 产品概述
16     - 设计理念
17     - 技术规格
18     - 驱动系统
19     - 电池与充电
20     - 智能座舱
21     - 智能驾驶
22     - 安全系统
23     - 车身结构
24     - 舒适性与便利性
25     - 版本与配置
26     - 价格与购买信息
27     - 售后服务
28     - 环保贡献
29     - 用户评价
30     - 竞品对比
31     - 常见问题
32     - 联系方式
33
34 2. **网络搜索**：在互联网上搜索相关信息
35
36 ## 工具选择规则
37 - 当查询明确涉及星辰电动ES9的具体信息、参数、功能或服务时，优先使用**本地文档搜索**
38 - 当查询涉及以下情况时，使用**网络搜索**：
39     - 与其他品牌车型的详细对比
40     - 最新市场动态或新闻
41     - 非官方的用户体验或评测
42     - 星辰电动ES9文档中可能没有的信息
43     - 需要实时数据（如当前市场价格波动等）
44
45 ## prompt延伸的规则
46 - 本地检索的查询扩展侧重于产品信息的深度查询
47 - 网络检索的查询扩展侧重于本地无法检索到的信息
48
49 ## 输出格式
50 你的输出应该是一个JSON格式的列表，每个项目包含：
51 1. `action_name`：工具名称（"本地文档搜索"或"网络搜索"）
52 2. `prompts`：问题列表，第一个是原始查询，后面是拆解或延伸的问题
53 [
54     {{
55         "action_name": "工具名称",
56         "prompts": [
57             "原始查询",
```

```

58     "拆解/延伸问题1",
59     "拆解/延伸问题2",
60     "拆解/延伸问题3"
61 ]
62 }}
63 ]
64
65
66 ## 示例
67
68 ### 示例1: 关于车辆规格的查询
69 用户: 星辰电动ES9的续航里程是多少?
70
71 输出:
72 [
73     {{
74         "action_name": "本地文档搜索",
75         "prompts": [
76             "星辰电动ES9的续航里程是多少?",
77             "星辰电动ES9的电池容量是多少?",
78             "星辰电动ES9不同版本的续航里程有何区别?"
79         ]
80     }}
81 ]
82
83
84 ### 示例2: 关于市场比较的查询
85 用户: 星辰电动ES9和特斯拉Model Y相比怎么样?
86
87 输出:
88 [

```

调整数据格式, 使每个action\_name只搭配一个prompt

代码块

```

1  #任务状态state
2  def adjust_format(original_data):
3      """
4      调整数据格式, 使每个action_name只搭配一个prompt
5
6      参数:
7      original_data (list): 原始数据, 每个action_name对应多个prompts
8
9      返回:
10     list: 调整后的数据, 每个action_name只对应一个prompt
11     """

```

```

12     adjusted_data = []
13
14     for item in original_data:
15         action_name = item['action_name']
16         prompts = item['prompts']
17
18         # 为每个prompt创建一个新的字典
19         for prompt in prompts:
20             adjusted_item = {
21                 'action_name': action_name,
22                 'prompt': prompt
23             }
24             adjusted_data.append(adjusted_item)
25
26     return adjusted_data
27
28 if action_tool:
29     adjusted_tools = adjust_format(action_tool)
30     actions=adjusted_tools
31 else:
32     actions=[]
33 print(actions)

```

## tools模块

执行模块，依次执行actions内的动作，根据action\_name判断执行函数web\_search\_answer()还是rag()，web\_search\_answer()和rag()就相当于tools，可供agent使用

代码块

```

1  def process_actions(actions):
2      """
3      处理动作列表函数
4
5      Args:
6          actions: 动作列表，每个动作包含action_name和prompt
7
8      Returns:
9          memory: 包含每次调用结果的记忆列表
10     """
11     # 初始化记忆列表
12     memory = []
13
14     # 依次处理每个动作
15     for action in actions:
16         action_name = action['action_name']

```



```

17     prompt = action['prompt']
18
19     print(f'正在执行{action_name}: "{prompt}"')
20
21     try:
22         # 根据动作类型调用相应的函数
23         if action_name == '本地文档搜索':
24             result = rag(prompt)
25         elif action_name == '网络搜索':
26             result = web_search_answer(prompt)
27         else:
28             result = f"未知的动作类型: {action_name}"
29
30         # 将结果添加到记忆中
31         memory_item = {
32             "提问": prompt,
33             "结果": result
34         }
35         memory.append(memory_item)
36
37         # 输出结果
38         print(f"提问: {prompt}")
39         print(f"结果: {result}")
40         print("-----")
41
42     except Exception:
43         # 如果执行失败, 静默处理, 继续下一轮循环
44         print("-----")
45         continue
46
47     print("所有执行动作已完成, 结果已添加到memory中。")
48     return memory

```

## memory模块

记忆模块memory, 用来储存搜索到的信息

代码块

```

1  memory_global=[]
2  memory_global.extend(memory_new[1:])

```

## 封装成agent接口

调用deepseek r1接口, 将检索到的内容添加到用户prompt中, 生成包含思考和推理两部分

```

1  from openai import OpenAI
代码块
2  import os
3
4  # 初始化OpenAI客户端
5  def final_answer(memory_global,user_query):
6      client = OpenAI(
7          # 如果没有配置环境变量, 请用百炼API Key替换: api_key="sk-xxx"
8          api_key = 'your_bailian_api_key',
9          base_url="https://dashscope.aliyuncs.com/compatible-mode/v1"
10     )
11
12     reasoning_content = "" # 定义完整思考过程
13     answer_content = "" # 定义完整回复
14     is_answering = False # 判断是否结束思考过程并开始回复
15
16
17     final_prompt=f'''
18         你是一个星辰电动ES9的智能销售助手, 负责根据用户的问题和提供的参考内容生成回答。请
严格按照以下要求生成回答:
19         基于提供的参考内容进行回答, 如果原文没有参考内容, 根据你自己的知识进行回答
20         你需要用有打动力的销售的语言进行输出, 突出星辰电动的优势
21
22         参考内容:
23         {memory_global}
24
25         用户问题: {user_query}
26
27     '''
28     print(final_prompt)
29     print('-'*130)
30
31     # 创建聊天完成请求
32     completion = client.chat.completions.create(
33         model="deepseek-r1", # 此处以 deepseek-r1 为例, 可按需更换模型名称
34         messages=[
35             {"role": "user", "content": final_prompt}
36         ],
37         stream=True,
38         # 解除以下注释会在最后一个chunk返回Token使用量
39         # stream_options={
40             #     "include_usage": True
41         # }
42     )
43
44     print("\n" + "=" * 20 + "思考过程" + "=" * 20 + "\n")
45
46     for chunk in completion:

```

```

47         # 如果chunk.choices为空, 则打印usage
48         if not chunk.choices:
49             print("\nUsage:")
50             print(chunk.usage)
51         else:
52             delta = chunk.choices[0].delta
53             # 打印思考过程
54             if hasattr(delta, 'reasoning_content') and delta.reasoning_content
!= None:
55                 print(delta.reasoning_content, end='', flush=True)
56                 reasoning_content += delta.reasoning_content
57             else:
58                 # 开始回复
59                 if delta.content != "" and is_answering == False:
60                     print("\n" + "=" * 20 + "完整回复" + "=" * 20 + "\n")
61                     is_answering = True
62                     # 打印回复过程
63                     print(delta.content, end='', flush=True)
64                     answer_content += delta.content
65
66             # print("=" * 20 + "完整思考过程" + "=" * 20 + "\n")
67             # print(reasoning_content)
68             # print("=" * 20 + "完整回复" + "=" * 20 + "\n")
69             # print(answer_content)

```

封装成agent接口, 其中agent\_plan、process\_actions、final\_answer可以视为三个子agent, 协同工作。

代码块

```

1  def agent(user_query):
2      #记忆模块
3      memory_global=[]
4      #获得执行清单
5      action_tool=agent_plan(user_query)
6      if action_tool:
7          adjusted_tools = adjust_format(action_tool)
8          actions=adjusted_tools
9      else:
10         actions=[]
11     if actions:
12         #进行任务执行
13         memory_new=process_actions(actions)
14         memory_global.extend(memory_new[1:])
15     final_answer(memory_global,user_query)
16     user_query='比较一下和华为汽车的优劣势'

```

```
17 print(agent(user_query))
```

## 实现能自我反思的agent

在prompt中要求agent自我反思，基于已有的信息，是否还需要延伸再进行查询

代码块

```
1 def reflection(user_query,memory_global):
2     prompt=''
3     你是一个专业的汽车销售助手的规划模块。你的任务是：
4     1. 分析用户的查询:{0}
5     2. 基于已有的信息，是否还需要延伸再进行查询
6
7     ##目前已有的信息：
8     {1}
9
10
11     ## 可用工具
12     1. **本地文档搜索**：搜索本地星辰电动ES9的文档，包含以下章节：
13         - 产品概述
14         - 设计理念
15         - 技术规格
16         - 驱动系统
17         - 电池与充电
18         - 智能座舱
19         - 智能驾驶
20         - 安全系统
21         - 车身结构
22         - 舒适性与便利性
23         - 版本与配置
24         - 价格与购买信息
25         - 售后服务
26         - 环保贡献
27         - 用户评价
28         - 竞品对比
29         - 常见问题
30         - 联系方式
31
32     2. **网络搜索**：在互联网上搜索相关信息
33
34     ## 工具选择规则
35     - 当查询明确涉及星辰电动ES9的具体信息、参数、功能或服务时，优先使用**本地文档搜索**
36     - 当查询涉及以下情况时，使用**网络搜索**：
37         - 与其他品牌车型的详细对比
38         - 最新市场动态或新闻
```

```

39     - 非官方的用户体验或评测
40     - 星辰电动ES9文档中可能没有的信息
41     - 需要实时数据（如当前市场价格波动等）
42
43     ## prompt延伸的规则
44     - 本地检索的查询扩展侧重于产品信息的深度查询
45     - 网络检索的查询扩展侧重于本地无法检索到的信息
46
47     ###重要!
48     至多再扩展不超过3个查询，如果需要扩展则按照下面的输出格式输出，如果不需要则返回None
49
50
51
52
53     ## 输出格式
54     你的输出应该是一个JSON格式的列表，每个项目包含：
55     1. `action_name`：工具名称（"本地文档搜索"或"网络搜索"）
56     2. `prompts`：一个扩展的问题，如果是网络检索，prompt不包含电动ES9，如果是本地检索，
57         prompt只包含询问电动ES9，检索内容一定是一个简单问题，不包含对比
58     [
59         {{
60             "action_name": "工具名称",
61             "prompts": '查询内容'
62         }}
63     ]
64
65     ''' .format(user_query,memory_global)
66     result=(middle_json_model(prompt))
67     # print(result)
68     json_list=extract_json_content(result)
69     try:
70         structure_output=json.loads(json_list)
71     except:
72         structure_output = None
73
74     return structure_output
75
76     print(reflection(user_query,memory_global))

```

封装成agent接口，在执行完process\_actions后，agent进行反思，直到检索到的内容足够完善

代码块

```

1     def agent_reflection(user_query):
2         memory_global=[]

```

```
3     print("开始任务规划...")
4     action_tool=agent_plan(user_query)
5     if action_tool:
6         adjusted_tools = adjust_format(action_tool)
7         actions=adjusted_tools
8     else:
9         actions=[]
10    if actions:
11        print("开始任务执行...")
12        memory_new=process_actions(actions)
13        memory_global.extend(memory_new[1:])
14    action_reflect=reflection(user_query,memory_global)
15    if action_reflect:
16        print("回顾内容, 进行反思...")
17        memory_new=process_actions(actions)
18        memory_global.extend(memory_new)
19    final_answer(memory_global,user_query)
20    print(agent_reflection(user_query))
```