

从LangChain+LLM的本地知识库问答到LLM与知识图谱、数据库的结合

前言

过去半年，随着ChatGPT的火爆，直接带火了整个LLM这个方向，然LLM毕竟更多是基于过去的经验数据 **预训练** 而来，没法获取最新的知识，以及各企业私有的知识

- 为了获取最新的知识，ChatGPT plus版集成了bing搜索的功能，有的模型则会调用一个定位于“链接各种AI模型、工具的langchain”的bing功能
- 为了处理企业私有的知识，要么基于开源模型微调，要么也可以基于langchain的思想调取一个外挂的向量知识库(类似存在本地的数据库一样)

所以越来越多的人开始关注langchain并把它与LLM结合起来应用，更直接推动了 **数据库** 、知识图谱与LLM的结合应用

本文侧重讲解

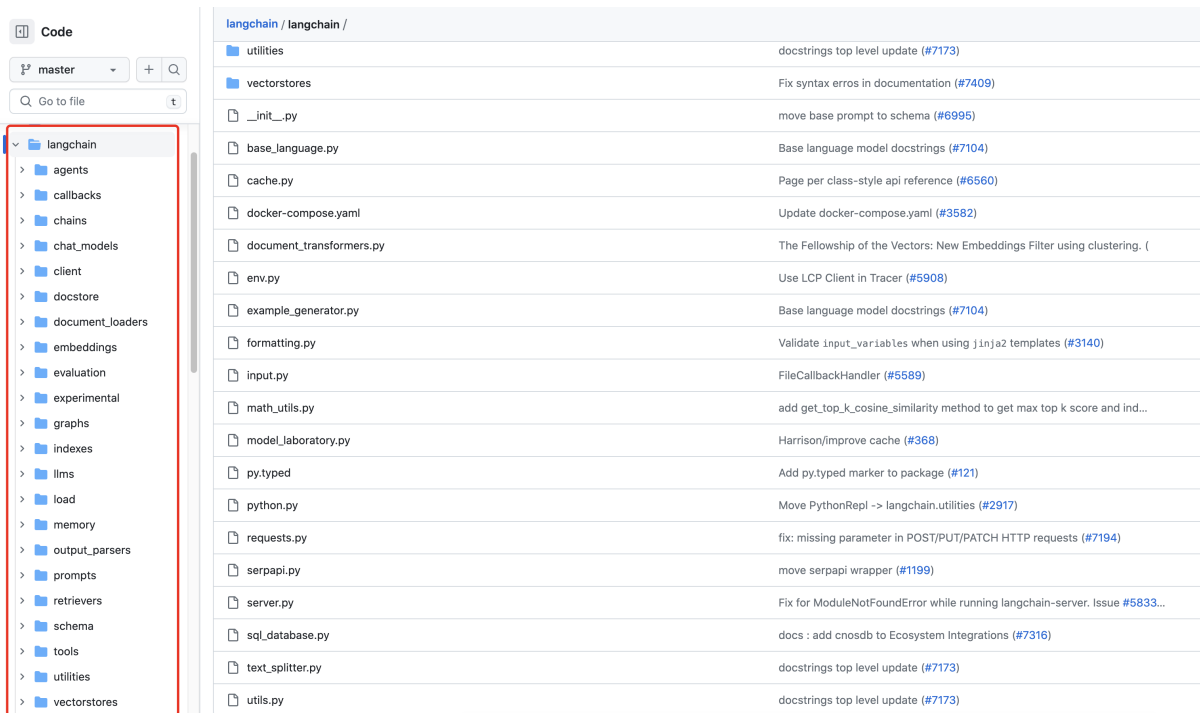
- LLM与langchain/数据库/知识图谱的结合应用
比如，虽说基于知识图谱的问答 早在2019年之前就有很多研究了，但谁会想到今年KBQA因为LLM如此突飞猛进呢
再比如，还会解读langchain-ChatGLM项目的关键源码，不只是把它当做一个工具使用，因为对工具的原理更了解，则对工具的使用更顺畅
- 其中，解读langchain-ChatGLM项目源码其实不易，因为涉及的项目、技术点不少，所以一开始容易绕晕，好在根据该项目的流程一步步抽丝剥茧之后，给大家呈现了清晰的代码架构
过程中，我从**接触该langchain-ChatGLM项目到整体源码梳理清晰并写清楚**历时了近一周，而大家有了本文之后，可能不到一天便可以理清了(提升近7倍效率)，这便是本文的价值和意义之一

阅读过程中若有任何问题，欢迎随时留言，会一一及时回复/解答，共同探讨、共同深挖

第一部分 什么是LangChain：LLM的外挂/功能库

1.1 langchain的整体组成架构

通俗讲，所谓langchain ([官网地址](#)、[GitHub地址](#))，即把AI中常用的很多功能都封装成库，且有调用各种商用模型API、开源模型的接口，支持以下各种组件



初次接触的朋友一看这么多组件可能直接晕了(封装的东西非常多，感觉它想把LLM所需要用到的功能/工具都封装起来)，为方便理解，我们可以先从大的层面把整个langchain库划分为三个大层：基础层、能力层、应用层

1.1.1 基础层：models、LLMs、index

Models：模型

各种类型的模型和模型集成，比如OpenAI的各个API/GPT-4等等，为各种不同基础模型提供统一接口
比如通过API完成一次问答

```
1 import os
2 os.environ["OPENAI_API_KEY"] = '你的api key'
3 from langchain.llms import OpenAI
4
5 llm = OpenAI(model_name="text-davinci-003",max_tokens=1024)
6 llm("怎么评价人工智能")
```

得到的回答如下图所示

```
>>> import os
>>> os.environ["OPENAI_API_KEY"] = 'sk-PFwMvMFcU7VXaL5cYM9vT'
>>> from langchain.llms import OpenAI
>>>
>>> llm = OpenAI(model_name="text-davinci-003",max_tokens=1024)
>>> llm("怎么评价人工智能")
'\n\n人工智能是一门极具潜力的学科，可以帮助人类解决许多复杂的问题。近年来，人工智能的发展取得了长足的进步，已经成功应用于诸多领域，为人们的生活带来了极大的便利。未来，人工智能将发挥更大的作用，为人类的发展做出重大贡献。'
```

LLMS层

这一层主要强调对models层能力的封装以及服务化输出能力，主要有：

- 各类LLM模型管理平台：强调的模型的种类丰富度以及易用性
- 一体化服务能力产品：强调开箱即用
- 差异化能力：比如聚焦于Prompt管理(包括提示管理、提示优化和提示序列化)、基于共享资源的模型运行模式等等

比如Google's PaLM Text APIs，再比如 `llms/openai.py` 文件下

```
1 |         model_token_mapping = {
2 |             "gpt-4": 8192,
3 |             "gpt-4-0314": 8192,
4 |             "gpt-4-0613": 8192,
5 |             "gpt-4-32k": 32768,
6 |             "gpt-4-32k-0314": 32768,
7 |             "gpt-4-32k-0613": 32768,
8 |             "gpt-3.5-turbo": 4096,
9 |             "gpt-3.5-turbo-0301": 4096,
10 |            "gpt-3.5-turbo-0613": 4096,
11 |            "gpt-3.5-turbo-16k": 16385,
12 |            "gpt-3.5-turbo-16k-0613": 16385,
13 |            "text-ada-001": 2049,
14 |            "ada": 2049,
15 |            "text-babbage-001": 2040,
16 |            "babbage": 2049,
17 |            "text-curie-001": 2049,
18 |            "curie": 2049,
19 |            "davinci": 2049,
20 |            "text-davinci-003": 4097,
21 |            "text-davinci-002": 4097,
22 |            "code-davinci-002": 8001,
23 |            "code-davinci-001": 8001,
24 |            "code-cushman-002": 2048,
25 |            "code-cushman-001": 2048,
26 |        }
```

Index(索引): Vector方案、KG方案

对用户私域文本、图片、PDF等各类文档进行存储和检索(相当于结构化文档，以便让外部数据和模型交互)，具体实现上有两个方案：一个Vector方案、一个KG方案

对于Vector方案：即对文件先切分为Chunks，在按Chunks分别编码存储并检索，可参考此代码文件：[langchain/libs/langchain/langchain/indexes/vectorstore.py](#) 该代码文件依次实现

模块导入：导入了各种类型检查、数据结构、预定义类和函数

接下来，实现了一个函数 `_get_default_text_splitter`，两个类 `VectorStoreIndexWrapper`、`VectorstoreIndexCreator`

`_get_default_text_splitter` 函数：

这是一个私有函数，返回一个默认的文本分割器，它可以将文本递归地分割成大小为1000的块，且块与块之间有重叠

```
1 | # 默认的文本分割器函数
2 | def _get_default_text_splitter() -> TextSplitter:
3 |     return RecursiveCharacterTextSplitter(chunk_size=1000, chunk_overlap=0)
```

接下来是，`VectorStoreIndexWrapper` 类：

这是一个包装类，主要是为了方便地访问和查询向量存储（Vector Store）

1. vectorstore: 一个向量存储对象的属性

```
1 |         vectorstore: VectorStore # 向量存储对象
2 |
3 |     class Config:
4 |         """Configuration for this pydantic object."""
5 |
6 |         extra = Extra.forbid # 额外配置项
```

```
7 | arbitrary_types_allowed = True # 允许任意类型
```

2. query: 一个方法, 它接受一个问题字符串并查询向量存储来获取答案

```
1 # 查询向量存储的函数
2 def query(
3     self,
4     question: str,                                # 输入的问题字符串
5     llm: Optional[BaseLanguageModel] = None,       # 可选的语言模型参数, 默认为None
6     retriever_kwargs: Optional[Dict[str, Any]] = None, # 提取器的可选参数, 默认为None
7     **kwargs: Any                                # 其他关键字参数
8 ) -> str:
9     """Query the vectorstore."""                # 函数的文档字符串, 描述函数的功能
10
11     # 如果没有提供语言模型参数, 则使用OpenAI作为默认语言模型, 并设定温度参数为0
12     llm = llm or OpenAI(temperature=0)
13
14     # 如果没有提供提取器的参数, 则初始化为空字典
15     retriever_kwargs = retriever_kwargs or {}
16
17     # 创建一个基于语言模型和向量存储提取器的检索QA链
18     chain = RetrievalQA.from_chain_type(
19         llm, retriever=self.vectorstore.as_retriever(**retriever_kwargs), **kwargs
20     )
21
22     # 使用创建的QA链运行提供的问题, 并返回结果
23     return chain.run(question)
```



解释一下上面出现的提取器

提取器首先从大型语料库中检索与问题相关的文档或片段, 然后生成器根据这些检索到的文档生成答案。

提取器可以基于许多不同的技术, 包括:

- 基于关键字的检索: 使用关键字匹配来查找相关文档
- 向量空间模型: 将文档和查询都表示为向量, 并通过计算它们之间的相似度来检索相关文档
- 基于深度学习的方法: 使用预训练的神经网络模型 (如BERT、RoBERTa等) 将文档和查询编码为向量, 并进行相似度计算
- 索引方法: 例如倒排索引, 这是搜索引擎常用的技术, 可以快速找到包含特定词或短语的文档

这些方法可以独立使用, 也可以结合使用, 以提高检索的准确性和速度

3. query_with_sources: 类似于query, 但它还返回与查询结果相关的数据源

```
1 # 查询向量存储并返回数据源的函数
2 def query_with_sources(
3     self,
4     question: str,
5     llm: Optional[BaseLanguageModel] = None,
6     retriever_kwargs: Optional[Dict[str, Any]] = None,
7     **kwargs: Any
8 ) -> dict:
9     """Query the vectorstore and get back sources."""
10     llm = llm or OpenAI(temperature=0) # 默认使用OpenAI作为语言模型
11     retriever_kwargs = retriever_kwargs or {} # 提取器参数
12     chain = RetrievalQAWithSourcesChain.from_chain_type(
13         llm, retriever=self.vectorstore.as_retriever(**retriever_kwargs), **kwargs
14     )
15     return chain({chain.question_key: question})
```



最后是VectorstoreIndexCreator 类:

这是一个创建向量存储索引的类

1. vectorstore_cls: 使用的向量存储类, 默认为Chroma

```
vectorstore_cls: Type[VectorStore] = Chroma # 默认使用Chroma作为向量存储类
```

一个简化的向量存储可以看作是一个大型的表格或数据库, 其中每行代表一个项目 (如文档、图像、句子等), 而每个项目则有一个与之关联的高维向量。向量的维度可以从几十到几千, 取决于所使用的嵌入模型

例如:

2. embedding: 使用的嵌入类, 默认为OpenAIEmbeddings

```
embedding: Embeddings = Field(default_factory=OpenAIEmbeddings) # 默认使用OpenAIEmbeddings作为嵌入类
```

3. text_splitter: 用于分割文本的文本分割器

```
text_splitter: TextSplitter = Field(default_factory=_get_default_text_splitter) # 默认文本分割器
```

4. from_loaders: 从给定的加载器列表中创建一个向量存储索引

```
1 | # 从加载器创建向量存储索引的函数
2 | def from_loaders(self, loaders: List[BaseLoader]) -> VectorStoreIndexWrapper:
3 |     """Create a vectorstore index from loaders."""
4 |     docs = []
5 |     for loader in loaders: # 遍历加载器
6 |         docs.extend(loader.load()) # 加载文档
7 |     return self.from_documents(docs)
```

5. from_documents: 从给定的文档列表中创建一个向量存储索引

```
1 | # 从文档创建向量存储索引的函数
2 | def from_documents(self, documents: List[Document]) -> VectorStoreIndexWrapper:
3 |     """Create a vectorstore index from documents."""
4 |     sub_docs = self.text_splitter.split_documents(documents) # 分割文档
5 |     vectorstore = self.vectorstore_cls.from_documents(
6 |         sub_docs, self.embedding, **self.vectorstore_kwargs # 从文档创建向量存储
7 |     )
8 |     return VectorStoreIndexWrapper(vectorstore=vectorstore) # 返回向量存储的包装对象
```

对于KG方案：这部分利用LLM抽取文件中的三元组，将其存储为KG供后续检索，可参考此代码文件：[langchain/libs/langchain/langchain/indexes/graph.py](#)

```
1 | """Graph Index Creator.""" # 定义"图索引创建器"的描述
2 |
3 | # 导入相关的模块和类型定义
4 | from typing import Optional, Type # 导入可选类型和类型的基础类型
5 | from langchain import BasePromptTemplate # 导入基础提示模板
6 | from langchain.chains.llm import LLMChain # 导入LLM链
7 | from langchain.graphs.networkx_graph import NetworkxEntityGraph, parse_triples # 导入Networkx实体图和解析三元组的功能
8 | from langchain.indexes.prompts.knowledge_triplet_extraction import ( # 从知识三元组提取模块导入对应的提示
9 |     KNOWLEDGE_TRIPLE_EXTRACTION_PROMPT,
10 | )
11 | from langchain.pydantic_v1 import BaseModel # 导入基础模型
12 | from langchain.schema.language_model import BaseLanguageModel # 导入基础语言模型的定义
13 |
14 | class GraphIndexCreator(BaseModel): # 定义图索引创建器类，继承自BaseModel
15 |     """Functionality to create graph index.""" # 描述该类的功能为"创建图索引"
16 |
17 |     llm: Optional[BaseLanguageModel] = None # 定义可选的语言模型属性，默认为None
18 |     graph_type: Type[NetworkxEntityGraph] = NetworkxEntityGraph # 定义图的类型，默认为NetworkxEntityGraph
19 |
20 |     def from_text(
21 |         self, text: str, prompt: BasePromptTemplate = KNOWLEDGE_TRIPLE_EXTRACTION_PROMPT
22 |     ) -> NetworkxEntityGraph: # 定义一个方法，从文本中创建图索引
23 |         """Create graph index from text.""" # 描述该方法的功能
24 |         if self.llm is None: # 如果语言模型为None，则抛出异常
25 |             raise ValueError("llm should not be None")
26 |         graph = self.graph_type() # 创建一个新的图
27 |         chain = LLMChain(llm=self.llm, prompt=prompt) # 使用当前的语言模型和提示创建一个LLM链
28 |         output = chain.predict(text=text) # 使用LLM链对文本进行预测
29 |         knowledge = parse_triples(output) # 解析预测输出得到的三元组
30 |         for triple in knowledge: # 遍历所有的三元组
31 |             graph.add_triple(triple) # 将三元组添加到图中
32 |         return graph # 返回创建的图
33 |
34 |     async def afrom_text( # 定义一个异步版本的from_text方法
35 |         self, text: str, prompt: BasePromptTemplate = KNOWLEDGE_TRIPLE_EXTRACTION_PROMPT
36 |     ) -> NetworkxEntityGraph:
37 |         """Create graph index from text asynchronously.""" # 描述该异步方法的功能
38 |         if self.llm is None: # 如果语言模型为None，则抛出异常
39 |             raise ValueError("llm should not be None")
40 |         graph = self.graph_type() # 创建一个新的图
41 |         chain = LLMChain(llm=self.llm, prompt=prompt) # 使用当前的语言模型和提示创建一个LLM链
42 |         output = await chain.apredict(text=text) # 异步使用LLM链对文本进行预测
43 |         knowledge = parse_triples(output) # 解析预测输出得到的三元组
44 |         for triple in knowledge: # 遍历所有的三元组
45 |             graph.add_triple(triple) # 将三元组添加到图中
46 |         return graph # 返回创建的图
```

另外，为了索引，便不得不牵涉以下这些能力

- **Document Loaders**，文档加载的标准接口

与各种格式的文档及数据源集成，比如Arxiv、Email、Excel、Markdown、PDF(所以可以做类似ChatPDF这样的应用)、Youtube ...

相近的还有

docstore，其中包含wikipedia.py等

document_transformers

- **embeddings**([langchain/libs/langchain/langchain/embeddings](#))，则涉及到各种embeddings算法，分别体现在各种代码文件中：

[elasticsearch.py](#)、[google_palm.py](#)、[gpt4all.py](#)、[huggingface.py](#)、[huggingface_hub.py](#)

[llamacpp.py](#)、[minimax.py](#)、[modelscope_hub.py](#)、[mosaicml.py](#)

[openai.py](#)

[sentence_transformer.py](#)、[spacy_embeddings.py](#)、[tensorflow_hub.py](#)、[vertexai.py](#)

1.1.2 能力层：Chains、Memory、Tools

如果基础层提供了最核心的能力，能力层则给这些能力安装上手、脚、脑，让其具有记忆和触发万物的能力，包括：Chains、Memory、Tool三部分

- **Chains：链接**

简言之，相当于包括一系列对各种组件的调用，可能是一个 Prompt 模板，一个语言模型，一个输出解析器，一起工作处理用户的输入，生成响应，并处理输出

具体而言，则相当于按照不同的需求抽象并定制化不同的执行逻辑，Chain可以相互嵌套并串行执行，通过这一层，让LLM的能力链接到各行各业

比如与Elasticsearch数据库交互的：[elasticsearch_database](#)

比如基于知识图谱问答的：[graph_qa](#)

其中的代码文件：[chains/graph_qa/base.py](#) 便实现了一个基于知识图谱实现的问答系统，具体步骤为

首先，根据提取到的实体在知识图谱中查找相关的信息「这是通过 `self.graph.get_entity_knowledge(entity)` 实现的，它返回的是与实体相关的所有信息，形式为三元组」

然后，将所有的三元组组合起来，形成上下文

最后，将问题和上下文一起输入到qa_chain，得到最后的答案

```
1  entities = get_entities(entity_string) # 获取实体列表。
2  context = "" # 初始化上下文。
3  all_triplets = [] # 初始化三元组列表。
4  for entity in entities: # 遍历每个实体
5      all_triplets.extend(self.graph.get_entity_knowledge(entity)) # 获取实体的所有知识并加入到三元组列表中。
6  context = "\n".join(all_triplets) # 用换行符连接所有的三元组作为上下文。
7
8  # 打印完整的上下文。
9  _run_manager.on_text("Full Context:", end="\n", verbose=self.verbose)
10 _run_manager.on_text(context, color="green", end="\n", verbose=self.verbose)
11
12 # 使用上下文和问题获取答案。
13 result = self.qa_chain(
14     {"question": question, "context": context},
15     callbacks=_run_manager.get_child(),
16 )
17 return {self.output_key: result[self.qa_chain.output_key]} # 返回答案
```



比如能自动生成代码并执行的：[llm_math](#)等等

比如面向私域数据的：[qa_with_sources](#)，其中的这份代码文件 [chains/qa_with_sources/vector_db.py](#) 则是使用向量数据库的问题回答，核心在于以下两个函数 `reduce_tokens_below_limit`

```
1 # 定义基于向量数据库的问题回答类
2 class VectorDBQAWithSourcesChain(BaseQAWithSourcesChain):
3     """Question-answering with sources over a vector database."""
4
5     # 定义向量数据库的字段
6     vectorstore: VectorStore = Field(exclude=True)
7
8     """Vector Database to connect to."""
9     # 定义返回结果的数量
10    k: int = 4
11
12    # 是否基于令牌限制来减少返回结果的数量
13    reduce_k_below_max_tokens: bool = False
14
15    # 定义返回的文档基于令牌的最大限制
16    max_tokens_limit: int = 3375
17
18    # 定义额外的搜索参数
19    search_kwargs: Dict[str, Any] = Field(default_factory=dict)
20
21    # 定义函数来根据最大令牌限制来减少文档
22    def _reduce_tokens_below_limit(self, docs: List[Document]) -> List[Document]:
```

```

23 |         num_docs = len(docs)
24 |         # 检查是否需要根据令牌减少文档数量
25 |         if self.reduce_k_below_max_tokens and isinstance(
26 |             self.combine_documents_chain, StuffDocumentsChain
27 |         ):
28 |             tokens = [
29 |                 self.combine_documents_chain.llm_chain.llm.get_num_tokens(
30 |                     doc.page_content
31 |                 )
32 |                 for doc in docs
33 |             ]
34 |             token_count = sum(tokens[:num_docs])
35 |
36 |             # 减少文档数量直到满足令牌限制
37 |             while token_count > self.max_tokens_limit:
38 |                 num_docs -= 1
39 |                 token_count -= tokens[num_docs]
40 |
41 |             return docs[:num_docs]
42 |

```



_get_docs

```

1 | # 获取相关文档的函数
2 | def _get_docs(
3 |     self, inputs: Dict[str, Any], *, run_manager: CallbackManagerForChainRun
4 | ) -> List[Document]:
5 |     question = inputs[self.question_key]
6 |
7 |     # 从向量存储中搜索相似的文档
8 |     docs = self.vectorstore.similarity_search(
9 |         question, k=self.k, **self.search_kwargs
10 |     )
11 |     return self._reduce_tokens_below_limit(docs)

```

比如面向SQL数据源的：[sql_database](#)，可以重点关注这份代码文件：[chains/sql_database/query.py](#)

比如面向模型对话的：[chat_models](#)，包括这些代码文件：[__init__.py](#)、[anthropic.py](#)、[azure_openai.py](#)、[base.py](#)、[fake.py](#)、[google_palm.py](#)、[human.py](#)、[jinachat.py](#)、[openai.py](#)、[promptlayer_openai.py](#)、[vertexai.py](#)

另外，还有比较让人眼前一亮的：

[constitutional_ai](#)：对最终结果进行偏见、合规问题处理的逻辑，保证最终的结果符合价值观

[llm_checker](#)：能让LLM自动检测自己的输出是否有问题的逻辑

• Memory：记忆

简言之，用来保存和模型交互时的上下文状态，处理长期记忆

具体而言，这层主要有两个核心点：

→ 对Chains的执行过程中的输入、输出进行记忆并结构化存储，为下一步的交互提供上下文，这部分简单存储在Redis即可

→ 根据交互历史构建知识图谱，根据关联信息给出准确结果，对应的代码文件为：[memory/kg.py](#)

```

1 | # 定义知识图谱对话记忆类
2 | class ConversationKGMemory(BaseChatMemory):
3 |     """知识图谱对话记忆类
4 |
5 |     在对话中与外部知识图谱集成，存储和检索对话中的知识三元组信息。
6 |     """
7 |
8 |     k: int = 2 # 考虑的上下文对话数量
9 |     human_prefix: str = "Human" # 人类前缀
10 |    ai_prefix: str = "AI" # AI前缀
11 |    kg: NetworkxEntityGraph = Field(default_factory=NetworkxEntityGraph) # 知识图谱实例
12 |    knowledge_extraction_prompt: BasePromptTemplate = KNOWLEDGE_TRIPLE_EXTRACTION_PROMPT # 知识提取提示
13 |    entity_extraction_prompt: BasePromptTemplate = ENTITY_EXTRACTION_PROMPT # 实体提取提示
14 |    llm: BaseLanguageModel # 基础语言模型
15 |    summary_message_cls: Type[BaseMessage] = SystemMessage # 总结消息类
16 |    memory_key: str = "history" # 历史记忆键
17 |
18 |    def load_memory_variables(self, inputs: Dict[str, Any]) -> Dict[str, Any]:
19 |        """返回历史缓冲区。"""
20 |        entities = self._get_current_entities(inputs) # 获取当前实体
21 |
22 |        summary_strings = []
23 |        for entity in entities: # 对于每个实体
24 |            knowledge = self.kg.get_entity_knowledge(entity) # 获取与实体相关的知识

```

```

25         if knowledge:
26             summary = f"0n {entity}: {' '.join(knowledge)}." # 构建总结字符串
27             summary_strings.append(summary)
28     context: Union[str, List]
29     if not summary_strings:
30         context = [] if self.return_messages else ""
31     elif self.return_messages:
32         context = [
33             self.summary_message_cls(content=text) for text in summary_strings
34         ]
35     else:
36         context = "\n".join(summary_strings)
37
38     return {self.memory_key: context}
39
40 @property
41 def memory_variables(self) -> List[str]:
42     """始终返回记忆变量列表。"""
43     return [self.memory_key]
44
45 def _get_prompt_input_key(self, inputs: Dict[str, Any]) -> str:
46     """获取提示的输入键。"""
47     if self.input_key is None:
48         return get_prompt_input_key(inputs, self.memory_variables)
49     return self.input_key
50
51 def _get_prompt_output_key(self, outputs: Dict[str, Any]) -> str:
52     """获取提示的输出键。"""
53     if self.output_key is None:
54         if len(outputs) != 1:
55             raise ValueError(f"One output key expected, got {outputs.keys()}")
56         return list(outputs.keys())[0]
57     return self.output_key
58
59 def get_current_entities(self, input_string: str) -> List[str]:
60     """从输入字符串中获取当前实体。"""
61     chain = LLMChain(llm=self.llm, prompt=self.entity_extraction_prompt)
62     buffer_string = get_buffer_string(
63         self.chat_memory.messages[-self.k * 2 :],
64         human_prefix=self.human_prefix,
65         ai_prefix=self.ai_prefix,
66     )
67     output = chain.predict(
68         history=buffer_string,
69         input=input_string,
70     )
71     return get_entities(output)
72
73 def _get_current_entities(self, inputs: Dict[str, Any]) -> List[str]:
74     """获取对话中的当前实体。"""
75     prompt_input_key = self._get_prompt_input_key(inputs)
76     return self.get_current_entities(inputs[prompt_input_key])
77
78 def get_knowledge_triplets(self, input_string: str) -> List[KnowledgeTriple]:
79     """从输入字符串中获取知识三元组。"""
80     chain = LLMChain(llm=self.llm, prompt=self.knowledge_extraction_prompt)
81     buffer_string = get_buffer_string(
82         self.chat_memory.messages[-self.k * 2 :],
83         human_prefix=self.human_prefix,
84         ai_prefix=self.ai_prefix,
85     )
86     output = chain.predict(
87         history=buffer_string,
88         input=input_string,
89         verbose=True,
90     )
91     knowledge = parse_triples(output) # 解析三元组
92     return knowledge
93
94 def _get_and_update_kg(self, inputs: Dict[str, Any]) -> None:
95     """从对话历史中获取并更新新知识图谱。"""
96     prompt_input_key = self._get_prompt_input_key(inputs)
97     knowledge = self.get_knowledge_triplets(inputs[prompt_input_key])
98     for triple in knowledge:
99         self.kg.add_triple(triple) # 向知识图谱中添加三元组
100
101 def save_context(self, inputs: Dict[str, Any], outputs: Dict[str, str]) -> None:
102     """将此对话的上下文保存到缓冲区。"""
103     super().save_context(inputs, outputs)

```



```

104         self._get_and_update_kg(inputs)105 |
106     def clear(self) -> None:
107         """清除记忆内容。"""
108         super().clear()
109         self.kg.clear() # 清除知识图谱内容

```



- **Tools层，工具**

其实Chains层可以根据LLM + Prompt执行一些特定的逻辑，但是如果要用Chain实现所有的逻辑不现实，可以通过Tools层也可以实现，Tools层理解为技能比较合理，典型的比如搜索、Wikipedia、天气预报、ChatGPT服务等等

1.1.3 应用层：Agents

- **Agents：代理**

简言之，有了基础层和能力层，我们可以构建各种各样好玩的，有价值的服务，这里就是Agent

具体而言，Agent 作为代理人去向 LLM 发出请求，然后采取行动，且检查结果直到工作完成，包括LLM无法处理的任务的代理 (例如搜索或计算，类似ChatGPT plus 的插件有调用bing和计算器的功能)

比如，Agent 可以使用维基百科查找 Barack Obama 的出生日期，然后使用计算器计算他在 2023 年的年龄

```

1 | # pip install wikipedia
2 | from langchain.agents import load_tools
3 | from langchain.agents import initialize_agent
4 | from langchain.agents import AgentType
5 |
6 | tools = load_tools(["wikipedia", "llm-math"], llm=llm)
7 | agent = initialize_agent(tools,
8 |                         llm,
9 |                         agent=AgentType.ZERO_SHOT_REACT_DESCRIPTION,
10 |                        verbose=True)
11 |
12 |
13 | agent.run("奥巴马的生日是哪天? 到2023年他多少岁了?")

```

此外，关于Wikipedia可以关注下这个代码文件：[langchain/docstore/wikipedia.py](#) ...

最终langchain的整体技术架构可以如下图所示 (查看高清图，此外，这里还有另一个架构图)

