

# LangChain 模块解析 (上)

LangChain Modules Hands-on

n:

Press Space for next page →



# 《AGI 应用实践》系列教程

## AGI OVERVIEW

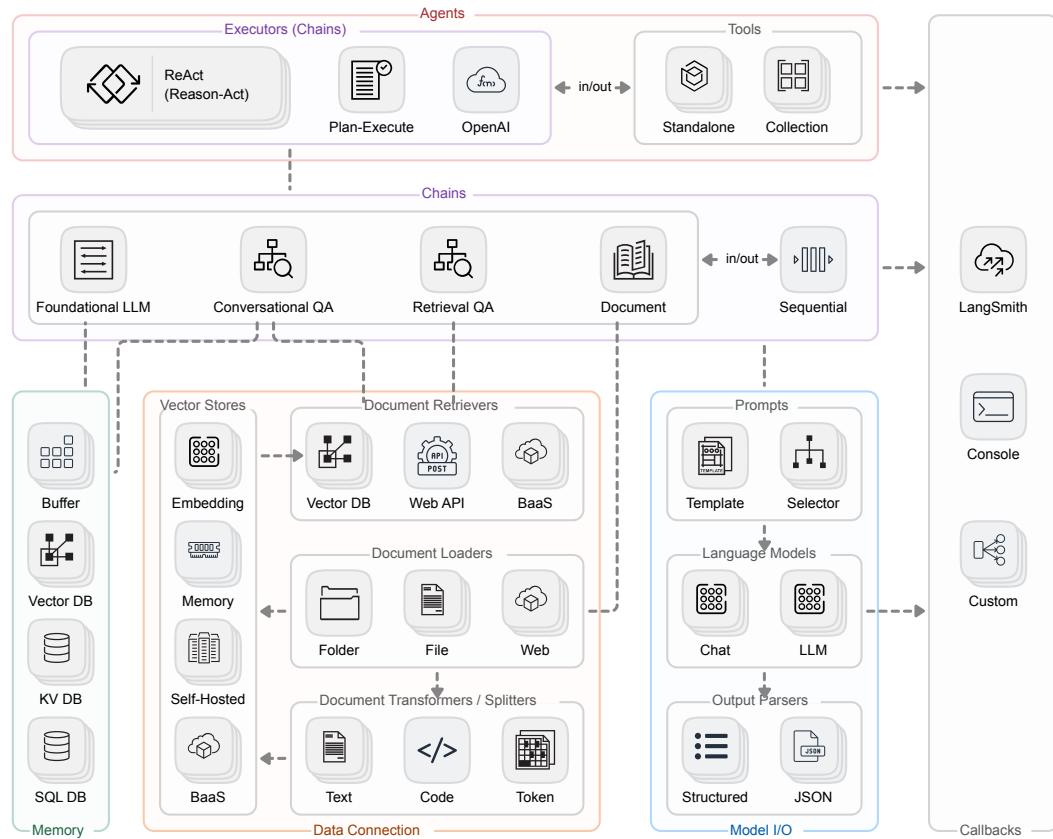
- 101: AGI 提示工程指北
- 102: OpenAI API 浅出

## LANGCHAIN IN ACTION

- 201: LangChain 功能模块解析（上篇）
- 202: LangChain 功能模块解析（下篇）
- 203: LangChain ChatBot 示例

# 教程大纲

- LangChain 核心模块概览
- Model I/O · 选择你的模型
- Data Connection · 准备你的数据
- Memory · 记录你的对话



LangChain (JS) Modules Overview v0.2.0 for npm:langchain@0.0.114 (Last Updated: July 22, 2023) by [@zhanqhaili0610](#)

# LangChain 核心模块概览

基于 [LangChain Modules Overview](#) 模块架构图一起来认识一下六大模块，参见 [官方文档](#)

## 201 教程的关注模块

- **Model I/O**: 管理大语言模型（Models），及其输入（Prompts）和格式化输出（Output Parsers）
  - 数据流：（Parser 生成的输出格式提示 → ）Prompt → Model → Output Parser
- **Data Connection**: 管理主要用于建设私域知识（库）的向量数据存储（Vector Stores）、内容数据获取（Document Loaders）和转化（Transformers），以及向量数据查询（Retrievers）
  - 数据流：Loaders → （Transformers → ）Embeddings Model → Vector Stores → Retrievers
- **Memory**: 用于存储和获取对话历史记录的功能模块

## 202 教程的关注模块

- **Chains**: 用于串联 Memory ↔ Model I/O ↔ Data Connection，以实现 串行化 的连续对话、推测流程
- **Agents**: 基于 Chains 进一步串联工具（Tools），从而将大语言模型的能力和本地、云服务能力结合
- **Callbacks**: 提供了一个回调系统，可连接到 LLM 申请的各个阶段，便于进行日志记录、追踪等数据导流

# 演示环境准备

基于 Val Town 平台，调用 [LangChain](#) 来实现 LLM 应用流程



@webup.getModelBuilder v6

```
1  async function getModelBuilder(spec: {
2    type?: "llm" | "chat" | "embedding";
3    provider?: "openai" | "huggingface";
4  } = { type: "llm", provider: "openai" }, options?: any) {
5    const { extend, cond, matches, invoke } = await import("npm:lodash-es");
6    const logger = {
7      handleLLMStart: (llm, prompts: string[]) =>
8        console.log("🕒 LLM Started:", llm),
9      handleLLMEnd: (llm) => console.log("🕒 LLM Ended: ", llm),
10    };
11    const callbacks = options?.verbose !== false ? [logger] : [];
12    const args = extend({ callbacks }, options);
13    if (spec?.provider === "openai")
```

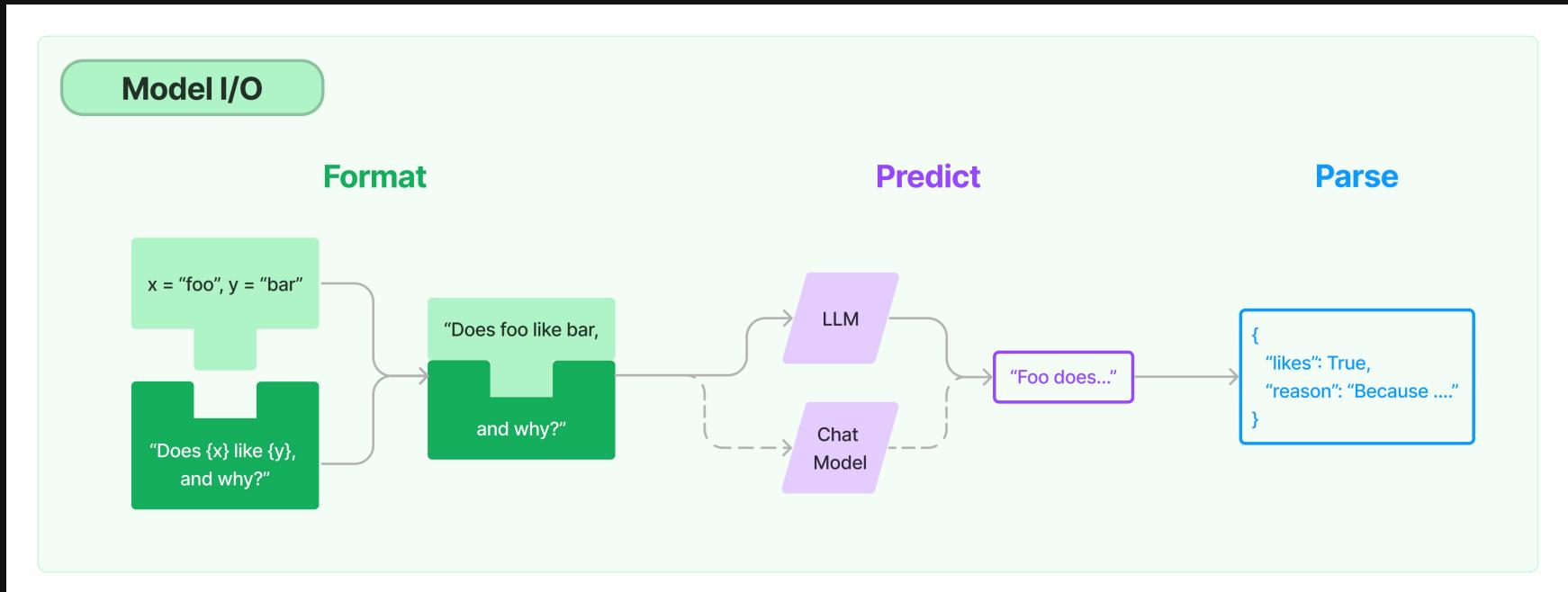
💡 官方也提供 Python 类库，以及庞大的社区 Integrations (实名羡慕啊 🍋 )

# Model I/O · 选择你的模型

同时也准备好你的输入、设定好输出的格式（结构）

# Model I/O 三元组

- Prompts: 模板化、动态选择和管理模型输入
- Language Models: 通过通用接口调用语言模型
- Output Parsers: 从模型输出中提取信息



# Models：一切的缘起

在 LangChain 中，把模型分成三类：LLM 大语言模型，Chat 对话模型，Embeddings 嵌入模型

- LLM：将文本字符串作为输入，并返回续写文本的应用模型



@webup.modelSampleLLMCall v0

```
1 const modelSampleLLMCall = (async () => {
2   const builder = await @webup.getModelBuilder();
3   const model = await builder();
4   return await model.call("Tell me a famous saying");
```

- Chat：基于 LLM 支持并将聊天消息序列作为输入，并返回聊天消息的应用模型



@webup.modelSampleChatCall v1

```
1 const modelSampleChatCall = (async () => {
2   const builder = await @webup.getModelBuilder({
3     type: "chat",
4     provider: "openai",
5   });
6   const model = await builder();
7   const { SystemMessage, HumanMessage } = await import("npm:langchain/schema");
```

# Models 提高多个通用接口

通用接口也支持批处理调用，并返回更丰富的响应内容



@webup.modelSampleLLMGenerate v0

```
1 const modelSampleLLMGenerate = (async () => {
2   const builder = await @webup.getModelBuilder();
3   const model = await builder();
4   return await model.generate(["Tell me a joke", "Tell me a poem"]);
5 })();
```



@webup.modelSampleChatGenerate v3

```
1 const modelSampleChatGenerate = (async () => {
2   const builder = await @webup.getModelBuilder({
3     type: "chat",
4     provider: "openai",
5   });
6   const model = await builder();
7   const { SystemMessage, HumanMessage } = await import("npm:langchain/schema");
8   return await model.generate([
```

# Prompts 的模板能力

Prompt Template 是一个带标记的文本字符串，可以接收来自最终用户的一组参数并生成提示



@webup.promptSampleTemplates v6

```
1 const promptSampleTemplates = (async () => {
2   const { PromptTemplate } = await import("npm:langchain/prompts");
3   const multipleInputPrompt = new PromptTemplate({
4     inputVariables: ["adjective", "content"],
5     template: "Tell me a {adjective} joke about {content}.",
6   });
7   const formattedMultipleInputPrompt = await multipleInputPrompt.format({}
```

模板也提供 Partial 的形式，用于分阶段的注入参数（变量内容）



@webup.promptSampleTemplatesPartial v2

```
1 const promptSampleTemplatesPartial = (async () => {
2   const { PromptTemplate } = await import("npm:langchain/prompts");
3   const date = () => new Date().toISOString();
4   const prompt = new PromptTemplate({
5     template: "Tell me a {adjective} joke about the day {date}",
6     inputVariables: ["adjective", "date"],
```

# Prompts 模板的组合

可以通过 Pipeline 管道来完成模板的组合

Pipeline 是一组提示模板，其中的每个模板将格式化，并按最终模板格式进行组合



@webup.promptSampleTemplatesPipeline v0

```
1 const promptSampleTemplatesPipeline = (async () => {
2   const { PromptTemplate, PipelinePromptTemplate } = await import(
3     "npm:langchain/prompts"
4   );
5   const fullPrompt = PromptTemplate.fromTemplate(`{introduction}
6
7 {example}
8
9 {start}`);
10  const introductionPrompt = PromptTemplate.fromTemplate(
11    `You are impersonating {person}.`,
12  );
13  const examplePrompt = PromptTemplate.fromTemplate(
14    `Here's an example of an interaction:
15 Q: {example_q}
16 A: {example_a}`)
```

# Prompts 的示例选择器

当您有大量示例时，可能需要选择哪些要包含在提示中，示例选择器是负责执行此操作的工具

如下示例选择器 根据长度选择 要使用的示例，当我们担心构建的提示会超过上下文窗口的长度时，这非常有用



@webup.promptSampleSelectorLength v1

```
1 const promptSampleSelectorLength = (async () => {
2   const { LengthBasedExampleSelector, PromptTemplate, FewShotPromptTemplate } =
3     await import("npm:langchain/prompts");
4   // Create a prompt template that will be used to format the examples.
5   const examplePrompt = new PromptTemplate({
6     inputVariables: ["input", "output"],
7     template: "Input: {input}\nOutput: {output}",
8   });
9   // Create a LengthBasedExampleSelector that will be used to select the examples.
10  const exampleSelector = await LengthBasedExampleSelector.fromExamples([
11    { input: "happy", output: "sad" },
12    { input: "tall", output: "short" },
13    { input: "energetic", output: "lethargic" },
14    { input: "sunny", output: "gloomy" },
15    { input: "windy", output: "calm" },
16  ]);
```

# Output Parsers 结构化输出：列表

如下示例中，我们将模型的输出解析为具有特定长度和分隔符的列表



@webup.parserSampleListCustom v1

```
1 const parserSampleListCustom = (async () => {
2   const { PromptTemplate } = await import("npm:langchain/prompts");
3   const { CustomListOutputParser } = await import(
4     "npm:langchain/output_parsers"
5   );
6   // With a `CustomListOutputParser`, we can parse a list with a specific length and separator.
7   const parser = new CustomListOutputParser({ length: 3, separator: "\n" });
8   const formatInstructions = parser.getFormatInstructions();
9   const prompt = new PromptTemplate({
10     template: "Provide a list of {subject}.\n{format_instructions}",
11     inputVariables: ["subject"],
12     partialVariables: { format_instructions: formatInstructions },
13   );
14   const builder = await @webup.getModelBuilder();
15   const model = await builder();
```

# Output Parsers 结构化输出：JSON（粗粒度）

如下示例中，我们将模型的输出解析为一个 包含多个字段的 JSON 对象



@webup.parserSampleJSON v0

```
1 const parserSampleJSON = (async () => {
2   const { PromptTemplate } = await import("npm:langchain/prompts");
3   const { StructuredOutputParser } = await import(
4     "npm:langchain/output_parsers"
5   );
6   // With a `StructuredOutputParser` we can define a schema for the output.
7   const parser = StructuredOutputParser.fromNamesAndDescriptions({
8     answer: "answer to the user's question",
9     source: "source used to answer the user's question, should be a website.",
10   });
11   const formatInstructions = parser.getFormatInstructions();
12   const prompt = new PromptTemplate({
13     template:
14       "Answer the users question as best as possible.\n{format_instructions}\n{question}",
15     inputVariables: ["question"],
16   });
17   const response = await llm.call(prompt);
18   const parsedResponse = parser.parse(response);
19   return parsedResponse;
20 });
21
22 module.exports = parserSampleJSON;
```

# Output Parsers 结构化输出：JSON（细粒度）

如下示例中，我们将模型的输出解析为一个结构被严格要求的 JSON 对象（[Zod 格式](#)）



@webup.parserSampleJSONZod v1

```
1 const parserSampleJSONZod = (async () => {
2   const { z } = await import("npm:zod");
3   const { PromptTemplate } = await import("npm:langchain/prompts");
4   const { StructuredOutputParser } = await import(
5     "npm:langchain/output_parsers"
6   );
7   // We can use zod to define a schema for the output using the `fromZodSchema` method of `StructuredOutput
8   const parser = StructuredOutputParser.fromZodSchema(z.object({
9     answer: z.string().describe("answer to the user's question"),
10    sources: z.array(z.string()).describe(
11      "sources used to answer the question, should be websites.",
12    ),
13  }));
14  const formatInstructions = parser.getFormatInstructions();
15  const prompt = new PromptTemplate({
```

# Expression Language : LLM Chain

0.0.121

现在可以用管道直接把 Prompt ➡ Model ➡ Output Parser 串联起来



@webup.pipeSampleLLM v1

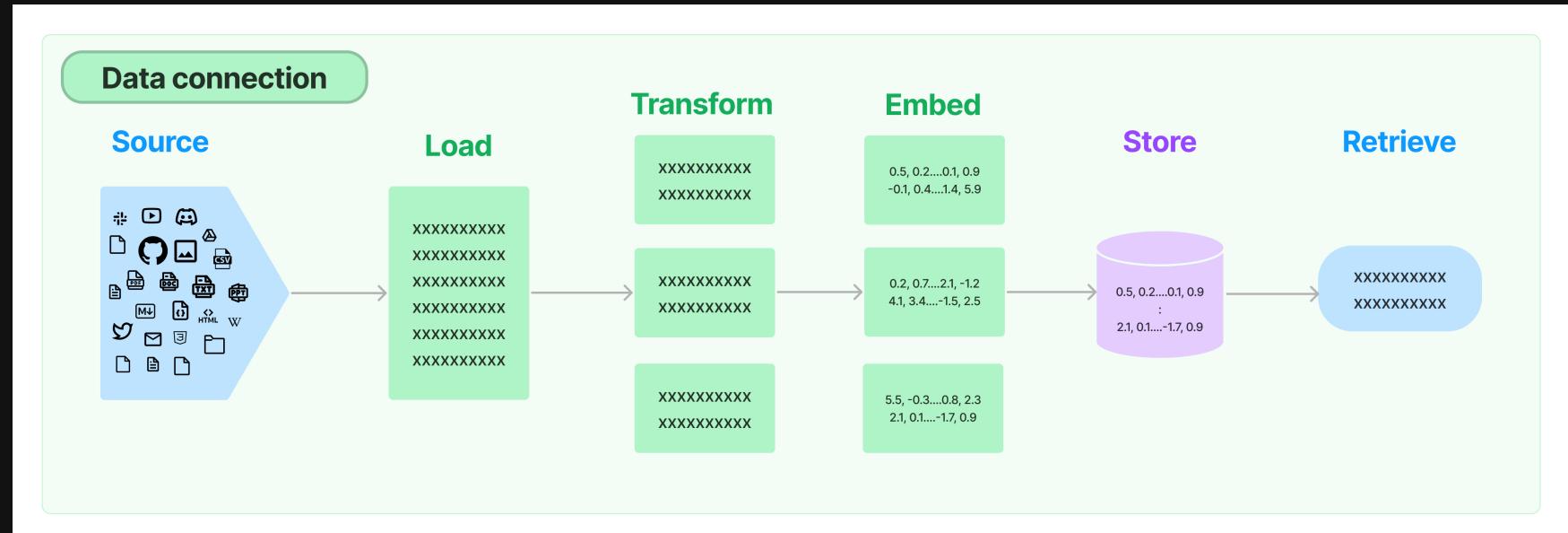
```
1 const pipeSampleLLM = (async () => {
2   const { PromptTemplate } = await import("npm:langchain/prompts");
3   const { RunnableSequence } = await import("npm:langchain/schema/runnable");
4   const { StringOutputParser } = await import(
5     "npm:langchain/schema/output_parser"
6   );
7   const builder = await @webup.getModelBuilder();
8   const model = await builder();
9   const promptTemplate = PromptTemplate.fromTemplate(
10     "Tell me a joke about {topic}",
11   );
12   const outputParser = new StringOutputParser();
13   const chain = RunnableSequence.from([promptTemplate, model, outputParser]);
14   const result = await chain.invoke({ topic: "bears" });
15   return result;
16 }
```

# Data Connection · 准备你的数据

获取数据，整理数据，存储数据，查询数据

# Data Connection 的核心组件

- Document Loaders & Transformers: 从许多不同的源加载文档；然后拆分文档、删除冗余文档等
- Embedding Models: 获取非结构化文本并将其转换为向量数据
- Vector Stores: 存储和搜索向量数据
- Retrievers: 从 Vector Stores 和其它数据源查询您的数据



# Document Loaders 的丰富数据源

Document Loaders 从文档源加载数据，并将数据转化成“文档”格式（一段文本和关联的元数据）

数据源可以是 本地文件，也可以是 在线文件；在如下示例中，我们分别从 GitHub 和网页加载我们教程的内容



@webup.loaderSampleWeb v0

```
1 const loaderSampleWeb = (async () => {
2   let builder, loader, docs;
3   // Load markdown files from GitHub
4   builder = await @webup.getWebLoaderBuilder(
5     "https://github.com/webup/agi-talks/tree/master/pages/common/",
6     "github",
7     { branch: "master", recursive: false, unknown: "warn" },
8   );
9   loader = await builder();
10  docs = await loader.load();
11  console.log(docs[0].pageContent);
12  // Load the same file as raw webpage
13  builder = await @webup.getWebLoaderBuilder(
14    "https://raw.githubusercontent.com/webup/agi-talks/master/pages/common/end.md",
15  );
```

# Document Transforms 的文本切分

受限于大语言模型的文本处理量，对于特别长的文本有必要将其分割成块，并尽量将语义相关的片段保留在一起

文本切分的主要方式大体包括“按字符”、“按 Token”、“按代码”；在如下示例中，我们主要演示前两种切分方式



@webup.splitterSampleCharToken v2

```
1 const splitterSampleCharToken = (async () => {
2   const { RecursiveCharacterTextSplitter, TokenTextSplitter } = await import(
3     "npm:langchain/text_splitter"
4   );
5   const text = `Hi.\n\nI'm Harrison.\n\nHow? Are? You?\nOkay then f f f f.
6     This is a weird text to write, but gotta test the splittingggg some how.\n\n
7     Bye!\n\n-H.`;
8   const options = {
9     chunkSize: 30,
10    chunkOverlap: 5,
11  };
12  const charSplitter = new RecursiveCharacterTextSplitter(options);
13  const chars = await charSplitter.createDocuments([text]);
14  const tokenSplitter = new TokenTextSplitter(options);
15  const tokens = await tokenSplitter.createDocuments([text]);
16 . . . . .
```

# Vector Stores 向量存储

存储和搜索非结构化数据的最常见方法之一是嵌入它并存储生成的嵌入向量

向量存储 负责存储嵌入式数据并为您执行向搜索 —— 查询“最相似”的嵌入向量；向量由 Embeddings 模型生成



@webup.getVectorStoreBuilder v0

```
1  async function getVectorStoreBuilder(docs: object[], spec: {
2    type: "memory" | "baas";
3    provider?: "pinecone";
4  } = { type: "memory" }, embed: "openai" | "huggingface" = "openai") {
5    const { cond, matches } = await import("npm:lodash-es");
6    const builder = await @webup.getModelBuilder({
7      type: "embedding",
8      provider: embed,
9    });
10   const model = await builder();
11   const setup = cond([
12     [
13       matches({ type: "memory" }),
14       async () => {
15         const { MemoryVectorStore } = await import(
16           "
```

# 向量存储的相似度搜索

在如下示例中，我们分别展示 内存 和 Pinecone 云服务 两种向量存储的使用（写入和搜索）



@webup.vectorStoreSampleMemory v5

```
1 const vectorStoreSampleMemory = (async () => {
2   const docs = await @webup.getSampleDocuments();
3   const builder = await @webup.getVectorStoreBuilder(docs);
4   const store = await builder();
5   return store.similaritySearch("pinecone", 1);
6 })();
```



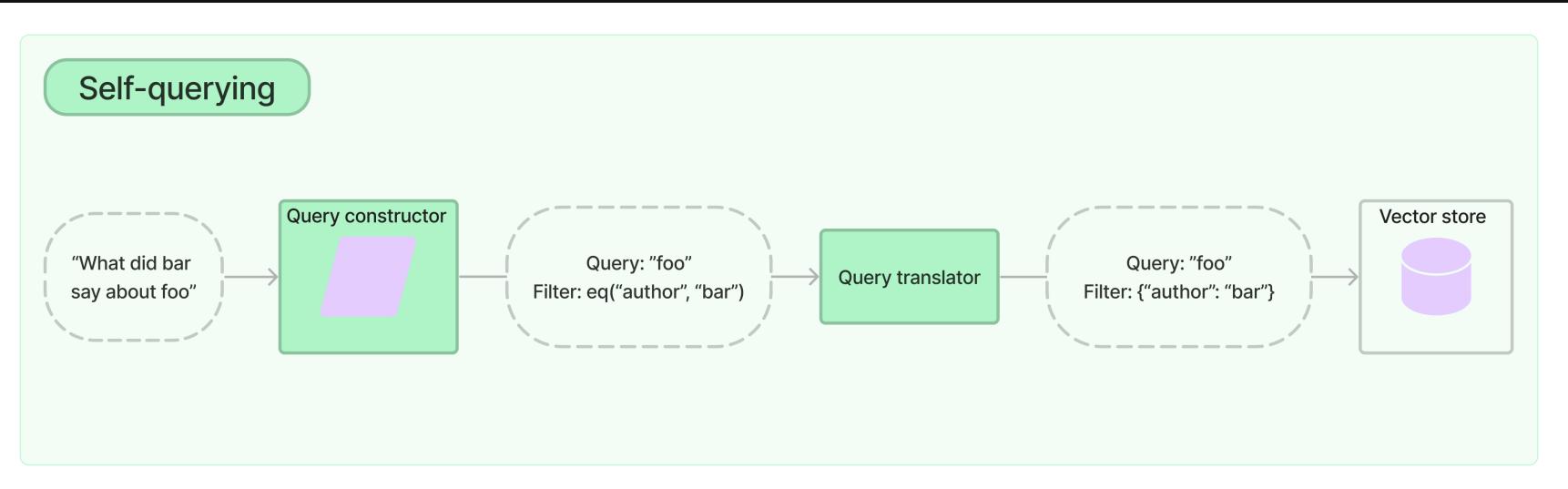
@webup.vectorStoreSamplePinecone v0

```
1 const vectorStoreSamplePinecone = (async () => {
2   const builder = await @webup.getVectorStoreBuilder([], {
3     type: "baas",
4     provider: "pinecone",
5   });
6   const store = await builder();
7   return store.similaritySearch("pinecone", 1);
```

# Retrievers 从各个文档源检索

Retrievers 比向量存储更通用，它不需要能够存储文档，只需返回（或检索）文档

- 向量存储是最常用的文档源，但也可以从其它本地、云上的文档源进行数据的检索
- Self-Query Retriever: 顾名思义，这是一种具有查询自身能力的检索器
  - 即给定任何自然语言查询，Retriever 通过构造 LLM Chain 编写结构化查询，并将其应用于底层向量存储



# Retrievers 从各个文档源检索（续）

检索的一项挑战是，通常不知道将数据引入系统时，文档存储系统将面临哪些特定查询

- 与查询最相关的信息可能被隐藏在包含大量不相关文本中，传递完整文档成本更高且响应较差
- Contextual Compression Retriever: 上下文压缩旨在解决此问题
  - 使用给定查询的上下文来压缩它们，以便只返回相关信息，而不是立即按原样返回检索到的文档

## Contextual compression



# Expression Language : LLM Chain + Retriever

0.0.121

可以用管道直接把 Prompt ➡ Model ➡ Output Parser ➡ Retriever 串联起来



@webup.pipeSampleLLMRetriever v0

```
1 const pipeSampleLLMRetriever = (async () => {
2   const { PromptTemplate } = await import("npm:langchain/prompts");
3   const { RunnableSequence, RunnablePassthrough } = await import(
4     "npm:langchain/schema/runnable"
5   );
6   const { StringOutputParser } = await import(
7     "npm:langchain/schema/output_parser"
8   );
9   const { Document } = await import("npm:langchain/document");
10  const modelBuilder = await @webup.getModelBuilder();
11  const model = await modelBuilder();
12  const docs = await @webup.getSampleDocuments();
13  const vectorBuilder = await @webup.getVectorStoreBuilder(docs);
14  const vector = await vectorBuilder();
15  const retriever = vector.asRetriever();
```

# Expression Language 0.0.121 : LLM Chain + Retriever

在上个示例中，我们将字符串输入直接传递到链路中；如希望链路接受多个输入，可以传递函数映射来解析输入



@webup.pipeSampleLLMRetrieverInputs v3

```
1 const pipeSampleLLMRetrieverInputs = (async () => {
2   const { PromptTemplate } = await import("npm:langchain/prompts");
3   const { RunnableSequence, RunnablePassthrough } = await import(
4     "npm:langchain/schema/runnable"
5   );
6   const { StringOutputParser } = await import(
7     "npm:langchain/schema/output_parser"
8   );
9   const { Document } = await import("npm:langchain/document");
10  const modelBuilder = await @webup.getModelBuilder();
11  const model = await modelBuilder();
12  const docs = await @webup.getSampleDocuments();
13  const vectorBuilder = await @webup.getVectorStoreBuilder(docs);
14  const vector = await vectorBuilder();
15  const retriever = vector.asRetriever();
```

# Expression Language 0.0.121 : LLM Chain + Retriever

在上个示例中，我们使链路接受了多个输入；我们还可以通过格式化函数添加对话历史记录



@webup.pipeSampleLLMRetrieverConversation v1

```
1 const pipeSampleLLMRetrieverConversation = (async () => {
2   const { PromptTemplate } = await import("npm:langchain/prompts");
3   const { RunnableSequence, RunnablePassthrough } = await import(
4     "npm:langchain/schema/runnable"
5   );
6   const { StringOutputParser } = await import(
7     "npm:langchain/schema/output_parser"
8   );
9   const { Document } = await import("npm:langchain/document");
10  const modelBuilder = await @webup.getModelBuilder();
11  const model = await modelBuilder();
12  const docs = await @webup.getSampleDocuments();
13  const vectorBuilder = await @webup.getVectorStoreBuilder(docs);
14  const vector = await vectorBuilder();
15  const retriever = vector.asRetriever();
```

# Memory · 记录你的对话

对话记录是保持链路状态的基石

# Memory 对话状态存储

Memory 负责在用户与大语言模型的交互过程中保留状态

- Memory 的作用可以归结为从一系列聊天消息中摄取、捕获、转换和提取知识
- Memory 可以返回多条信息（例如最近的 N 条消息和摘要）或者所有以前的消息



@webup.getMemoryBuilder v3

```
1  async function getMemoryBuilder(spec: {
2    type: "buffer" | "summary" | "vector";
3    provider?: "openai";
4  } = { type: "buffer" }, options = {}) {
5    const { cond, matches } = await import("npm:lodash-es");
6    const setup = cond([
7      [
8        matches({ type: "buffer" }),
9        async () => {
10          const { BufferMemory } = await import("npm:langchain/memory");
11          return new BufferMemory();
12        },
13      ],
14    ]);
  
```

# Memory 的多样化存取方式

Memory 不仅可以整存整取，还可以利用大语言模型加工并返回处理后的内容

- 通过 LLM 模型对记录内容进行总结



@webup.memorySampleSummary v3

```
1 const memorySampleSummary = (async () => {
2   const builder = await @webup.getMemoryBuilder({
3     type: "summary",
4     provider: "openai",
5   });
6   const memory = await builder();
7   await memory.saveContext({ input: "My favorite sport is soccer" });
8 }
```

- 通过 Vector Store 对记录内容进行相似度匹配后返回



@webup.memorySampleVector v0

```
1 const memorySampleVector = (async () => {
2   const builder = await @webup.getMemoryBuilder({
3     type: "vector",
4     provider: "openai",
5   });
6   const memory = await builder();
```

# 参考资料

本教程在制作过程中参考和引用了以下资料（排名不分先后）的内容，特此鸣谢！

## ■ 视频资料

- [Short Courses | Learn Generative AI from DeepLearning.AI](#)

## ≡ 图文资料

- [Core Concepts | !\[\]\(08c7b22d50947e00457b407a54a4ee1a\_img.jpg\) LangChain](#)
  - [JS/TS Docs, Python Docs, LangSmith Docs](#)
- [入门：Prompts（提示词） | 通往 AGI 之路](#)

## </> 代码资料

- [openai/openai-cookbook: Examples and guides for using the OpenAI API](#)
- [datawhalechina/prompt-engineering-for-developers: 吴恩达大模型系列课程中文版](#)
- [slidevjs/slidev: Presentation Slides for Developers](#)

感谢聆听 ❤

⌚ webup | ⚡ serviceup