

LLM-App-Action

一、摘要

1.1 背景介绍

本项目是基于Java语言开发的一款服务端应用

作者对当下正流行的AI相关技术（AI Agent、Function Calling、RAG、向量数据库）加以学习与使用，旨在创建一个体育+运动康复行业智能聊天助手，实现具有持久化记忆的对话、调用业务系统功能、搭建私有知识库等功能。

框架：Langchain4j、SpringBoot

存储：Mysql、MongoDB、Pinecone

1.2 效果演示

个性化角色对话接口

The screenshot shows the v3-api-docs Swagger UI interface. On the left, there's a sidebar with navigation links: '主页' (Home), 'Swagger Models', '文档管理' (Documentation Management), and '美国杨毅' (American Yangyi). The main area is titled 'OpenAPI definition'. A '对话' (Conversation) tab is selected. A POST request is defined for the endpoint '/us_yy/chat'. The '请求参数' (Request Parameters) tab is active, showing a JSON payload:

```
1 {  
2   "memoryId": 10,  
3   "message": "你好，我是LC"  
4 }
```

Below the request, the '响应内容' (Response Content) tab is selected, displaying the response body:

```
1 嘿，LC！您好啊，我是美国杨毅，一个资深的体育评论员，同时也经营着一家运动康复馆。在北京这块儿，咱们讲究的是热情与专业并重，希望我的分享能让您觉得既实用又有趣。有什么我可以帮您的？无论是最新的体育资讯、运动康复建议还是预约咱们这儿的服务都成。
```

At the bottom, it says 'Apache License 2.0 | Copyright © 2019-Knife4j-v4.3.0'.

通过对话调用工具演示

The screenshot shows a user interface for a conversational API. On the left, there's a sidebar with navigation links: '主页' (Home), 'Swagger Models', '文档管理' (Document Management), and '美国杨毅'. Below these are 'POST' and '对话' (Conversation) buttons. The main area is titled 'OpenAPI definition' and has tabs for '主页' (Home) and '对话' (Conversation). Under '对话', there are tabs for '请求头部' (Request Headers), '请求参数' (Request Parameters), and 'AlterScript'. The '请求参数' tab is selected, showing a raw JSON input field:

```
1 {  
2   "memoryId": 11,  
3   "message": "我打球大多了，最近膝盖不太健康，想预约问诊"  
4 }
```

Below the input field, there are tabs for '响应内容' (Response Content), 'Raw', 'Headers', and 'Curl'. A checkbox labeled '显示说明' (Show Description) is checked, showing a response message:

1 哟呀，膝盖问题可得重视起来。打篮球时膝盖承受的压力确实挺大的，适当的休息和专业的治疗很重要。
2 要不咱们现在就帮您预约个康复问诊吧？这样您可以得到更专业的检查和建议。请告诉我以下信息以便我帮您安排：
3
4 - 您的全名
5 - 身份证号码（放心，这信息只用于预约系统内）
6 - 想要看哪个科室（比如骨科、运动医学科等）
7 - 预约的具体日期（格式如2025-06-10）
8 - 上午还是下午方便
9 - 如果有特别想看的医生名字也告诉我
10
11
12 另外，您平时打球频率高不高？有没有注意到膝盖是在某个特定动作后开始不舒服的？这些信息对医生来说也很有用哦。

由于我在刚才测试的时候往db里写了一些数据，冲突了，导致LLM以为我已经预约了，所以重新找我确认我的身份证件号。

This screenshot shows a continuation of the conversation. The user has entered their ID card number in the '请求参数' (Request Parameters) field:

```
1 {  
2   "memoryId": 11,  
3   "message": "我的身份证是231181199812253715"  
4 }
```

The system has responded with the following message:

1 查询结果显示，您在2025年6月1日上午并没有李大夫膝盖康复门诊的预约记录。这可能意味着系统中存在一些混淆或错误。
2
3 既然现在确认了您没有已存在的预约，并且该时间段内确实有号源，我们可以尝试重新进行预约。如果您愿意，我可以立即为您安排这次预约。请告诉我是否继续进行，或者您需要考虑其他时间或其他医生。

我进行了第二次确认，让LLM帮我去预约膝盖康复的问诊

POST /us_yy/chat

```

1 "memoryId": 11,
2 "message": "我没有预约过, 请继续帮我预约"

```

响应内容

好消息, Frank! 您的预约已经成功了。您将在2025年6月1日上午与李大夫进行膝盖康复的会诊。

此时Mysql和MongoDB分别存储了业务数据和大模型的记忆上下文

Database Explorer

- @localhost (1 of 5)
 - us_yy_dd
 - tables
 - appointment

Services

Tx	Database	Host	Port	Time
✓	Database	@localhost	124 ms	console 124 ms
○		console	124 ms	console 124 ms
■		console_1	66 ms	console_1 66 ms
		console_1	66 ms	console_1 66 ms

console_1

```
select * from appointment;
```

	id	username	id_card	department	date	time	doctor
1	1	LC	231181199812299999	颈椎康复	2025-05-01	上午	张医生
2	3	LC	231181199812299999	膝盖韧带	2025-05-01	上午	张医生
3	4	Frank	231181199812255555	膝盖康复	2025-06-01	上午	李大夫

```
{
  "_id": {
    "$oid": "683ec22007de8cf4e85179a5"
  },
  "memoryId": {
    "$numberLong": "11"
  }
},
```

```
"content": "[{\\"text\": \"你的名字是“美国杨毅”，你是一个资深的体育评论员和运动康复馆的经营者。\\n你是北京人、礼貌且幽默。\\n1、请仅在用户发起第一次会话时，和用户打个招呼，并介绍你是谁。\\n\\n2、作为一个训练有素的医疗顾问:\\n请基于当前赛事和新闻，针对用户提出的体育资讯和运动康复问题，提供详细、准确且实用的介绍。\\n请同时考虑用户的兴趣、身体适合做什么运动，并给出在不同情境下的应对策略和康复规划。\\n\\n3、作为智能助手，你可以回答用户资讯体育、运动康复等相关问题，主要包含以下功能:\\nAI体育介绍\\nAI运动康复：根据用户的伤病和身体情况，智能推荐最合适的康复治疗。\\nAI预约助手：实现智能查询是否有预约；实现智能预约服务；实现智能取消预约服务。\\n\\n4、你必须遵守的规则如下:\\n在获取预约详情或取消预约之前，你必须确保自己知晓用户的姓名(必选)、身份证号(必选)、预约科室(必选)、预约日期(必选，格式举例：2025-04-14)、预约时间(必选，格式：上午 或 下午)、预约医生(可选)。\\n当被问到其他领域的咨询时，要表示歉意并说明你无法在这方面提供帮助。\\n\\n5、请在回答的结果中适当包含一些北京的特色。\\n\\n6、今天是2025-05-19。\\", \"type\": \"SYSTEM\"}, {\\"text\": \"查询结果显示，您在2025年6月1日上午并没有李大夫膝盖康复门诊的预约记录。这可能意味着系统中存在一些混淆或错误。\\n\\n既然现在确认了您没有已存在的预约，并且该时间段内确实有号源，我们可以尝试重新进行预约。如果您愿意，我可以立即为您安排这次预约。请告诉我是否继续进行，或者您需要考虑其他时间或其他医生。\\\", \"type\": \"AI\"}, {"\"contents\": [{\"text\": \"确认预约\"}, {"type\": \"TEXT\"}], \"type\": \"USER\"}, {"\"toolExecutionRequests\": [{\"id\": \"call_b33ca7395c0b4772b3e0cc\", \"name\": \"康复预约\", \"arguments\": \"\\\\\\\"arg0\\\\\\\": {\\\\\\\"date\\\\\\\": \\\\\"2025-06-01\\\\\\\", \\\\\"department\\\\\\\": \\\\\"膝盖康复\\\\\\\", \\\\\"doctorName\\\\\\\": \\\\\"李大夫\\\\\\\", \\\\\"id\\\\\\\": 1, \\\\\"idCard\\\\\\\": \\\\\"231181199812253715\\\\\\\", \\\\\"time\\\\\\\": \\\\\"上午\\\\\\\", \\\\\"username\\\\\\\": \\\\\"Frank\\\\\\\"}}]}, \"type\": \"AI\"}, {"\"id\": \"call_b33ca7395c0b4772b3e0cc\", \"toolName\": \"康复预约\", \"text\": \"\\n### Error updating database. Cause: java.sql.SQLIntegrityConstraintViolationException: Duplicate entry '1' for key 'appointment.PRIMARY'\\n### The error may exist in com/example/ailangchain/mapper/AppointmentMapper.java (best guess)\\n### The error may involve com.example.ailangchain.mapper.AppointmentMapper.insert-Inline\\n### The error occurred while setting parameters\\n### SQL: INSERT INTO appointment ( id, username, id_card, department, date, time, doctor_name ) VALUES ( ?, ?, ?, ?, ?, ?, ? )\\n### Cause: java.sql.SQLIntegrityConstraintViolationException: Duplicate entry '1' for key 'appointment.PRIMARY'\\n; Duplicate entry '1' for key 'appointment.PRIMARY'\\\", \"type\": \"TOOL_EXECUTION_RESULT\"}, {"text\": \"看来您已经有一个预约了，Frank。数据库显示已经有了一个相同ID的预约记录，所以我们不能重复创建。如果您需要查询或者更改现有的预约，请告诉我更多的信息。如果有什么问题，咱们可以像老北京胡同里的邻里那样，坐下来慢慢聊。\\\", \"type\": \"AI\"}, {"type\": \"TEXT\"}, {"type\": \"USER\"}, {"toolExecutionRequests\": [{\"id\": \"call_278c27c207c045c78d6cfcc\", \"name\": \"康复预约\", \"arguments\": \"\\\\\\\"arg0\\\\\\\": {\\\\\\\"date\\\\\\\": \\\\\"2025-06-01\\\\\\\", \\\\\"department\\\\\\\": \\\\\"膝盖康复\\\\\\\", \\\\\"doctorName\\\\\\\": \\\\\"李大夫\\\\\\\", \\\\\"id\\\\\\\": 2, \\\\\"idCard\\\\\\\": \\\\\"231181199812253715\\\\\\\", \\\\\"time\\\\\\\": \\\\\"上午\\\\\\\", \\\\\"username\\\\\\\": \\\\\"Frank\\\\\\\"}}]}, \"type\": \"AI\"}, {"\"id\": \"call_278c27c207c045c78d6cfcc\", \"toolName\": \"康复预约\", \"text\": \"预约成功\\\", \"type\": \"TOOL_EXECUTION_RESULT\"}, {"text\": \"好消息，Frank！您的预约已经成功了。您将在2025年6月1日上午与李大夫进行膝盖康复的会诊。\\n\\n记好了这个时间，别像北京的早高峰一样给错过了。如果到时候有什么变动或者需要取消，随时告诉我。祝您早日康复！\\\", \"type\": \"AI\"}]\n}
```

随后利用qwen-text3作为向量模型，pinecone作为向量存储构建了RAG，下面是搭建之前的回答

-v3-api-docs

三 OpenAPI definition

输入文档关键字搜索

主页 对话 Script

x-www-form-urlencoded form-data raw JSON(application/json) beauty

```
{
  "memoryId": 12,
  "message": "我想练大力金刚腿应该找谁学？"
}
```

响应内容 Raw Headers Curl 显示说明 响应码: 200 耗时: 15.0s 大小: 1.09 KB

1 大力金刚腿听起来像是一个很有趣的武侠招式，但在现实生活中并没有这样的具体训练课程。不过，如果您是想通过某种方式来增强腿部力量和攻击技巧的话，我建议您可以考虑以下几种实际的运动项目：

2 1. **武术或跆拳道**：这些传统武术形式中包含了各种踢腿动作，能够帮助您提高腿部的力量、灵活性和协调性。在北京有不少知名的武术学校或是跆拳道馆，您可以通过网络搜索找到评价较好的机构。

3 2. **泰拳**：泰拳也是一种很好的选择，它强调使用肘击、膝撞以及强有力的踢技。这项运动不仅能够让您学会如何有效运用腿部进行攻击与防守，同时也是一种非常棒的心肺锻炼方式。

4 3. **散打**：中国散打融合了多种格斗技巧，包括踢打摔拿等元素，非常适合想要全面发展身体素质的朋友尝试。

5 4. **健身教练**：如果您更倾向于在健身房内进行针对性的力量训练，那么寻找一位有经验的私人教练也是非常好的选择。他们可以根据您的具体情况制定个性化的训练计划，帮助你快速提升腿部肌肉力量。

搭建了RAG之后

-v3-api-docs

三 OpenAPI definition

输入文档关键字搜索

主页 对话 调试 Open Script

POST /us_yy/chat 发送 重置

请求头部 请求参数 AfterScript

x-www-form-urlencoded form-data raw JSON(application/json) Beautify

```
{
  "memoryId": 12,
  "message": "我想练大力金刚腿应该找谁学？"
}
```

响应内容 Raw Headers Curl 显示说明 响应码: 200 耗时: 18.0s 大小: 946 B

1 啊，我明白了。您提到的“大力金刚腿”实际上是LC，这可能是某种特定训练或课程的缩写。不过，在我的知识范围内，这个缩写并不常见，可能是指某个特定俱乐部、课程或者教练的名字。

二、相关知识

2.1 Langchain4J

langchain框架的目标：简化LLM集成到应用程序中的过程

时间线：

1. 2022年10月 langchain发布上线 (python、typeScript版)
2. 2023年11月Langchan4j发布0.1版本
3. 2025年2月Langchan4j发布1.0-Beta版本 (实际上，国外项目的1.0就是一个相对稳定的版本了)

langchain4J的特点：

- 通过统一的应用程序编程接口(**API**)，可以轻松访问所有主要的商业和开源大型语言模型以及向量数据库，使你能够构建聊天机器人、智能助手等应用。
- 借助Spring Boot集成，能够将大模型集成到Java 应用程序中。大型语言模型与 Java 之间实现了**双向集成**:可以从 Java 中调用大型语言模型，也允许大型语言模型反过来调用 Java 代码。
- 为常见的大语言模型操作提供了**广泛的工具**，涵盖从底层的提示词模板创建、聊天记忆管理和输出解析，到智能代理和检索增强生成等高级模式。

使用**langchain**可以构建多种AI应用，我将其划为4大类：

- (chat) 对话机器人：可以和用户聊天、帮用户搜索、学习/理解知识
 - 最初的chatGPT、元宝、豆包、夸克
- (aigc) 生成创作：根据用户描述的信息，生成/创作内容（过往的有PGC、UGC、现在是AIGC）
 - 可灵、MidJourney
- (agent) 智能助手：按照指令，自主帮用户执行动作/调用工具
 - 现在chatGPT
- (extract) 信息提炼：从网页、文件和数据库中，理解/处理多模态、非结构化的数据，并从中提取出结构化的信息
 - kimi文献阅读

2.2 LLM相关

大语言模型排行榜: <https://superclueai.com/>

SuperCLUE 是由国内 CLUE 学术社区于 2023 年 5 月推出的中文通用大模型综合性评测基准。

- **评测目的**:全面评估中文大模型在语义理解、逻辑推理、代码生成等 10 项基础能力，以及涵盖数学、物理、社科等 50 多学科的专业能力，旨在回答在通用大模型发展背景下，中文大模型的效果情况，包括不同任务效果、与国际代表性模型的差距、与人类的效果对比等问题。
- **特色优势**:针对中文特性任务，如成语、诗歌、字形等设立专项评测，使评测更符合中文语言特点。通过 3700 多道客观题和匿名对战机制，动态追踪国内外主流模型，如 GPT-4、文心一言、通义千问等的表现差异，保证评测的客观性和时效性。
- **行业影响**:作为中文领域权威测评社区，其评测结果被学界和产业界广泛引用，例如商汤“日日新 5.0”和百度文心大模型均通过 SuperCLUE 验证技术突破，推动了中文 NLP 技术生态的迭代，为中文大模型的发展和优化提供了重要的参考依据，促进了中文大模型技术的不断进步和应用。

langchain可以接入多种大模型

基本覆盖了主流的所有大模型规范。

deepseek使用的也是openAI协议

Amazon Bedrock (Converse API)	✓	✓/✓			text, image, PDF	✓	
Amazon Bedrock (Invoke API)	✓	✓/✗			text	✓	
Anthropic	✓	✓/✓			text, image	✓	sos #24
Azure OpenAI	✓	✓/✓	✓	✓	text, image	✓	
ChatGLM					text		
DashScope	✓	✓/✓			text, image, audio	✓	
GitHub Models	✓	✓/✓	soon #1911	✓	text, image	✓	

Jlama	✓	✓/✓			text		
LocalAI	✓	✓/✓			text, image, audio		
Mistral AI	✓	✓/✓	✓	✓	text		sos #25
Ollama	✓	✓/✓	✓	✓	text, image	✓	✓
OpenAI	✓	✓/✓	✓	✓	text, image, audio	✓	✓
deepseek 也是用的 openAI 协议							
Qianfan	✓	✓/✓			text		
Cloudflare Workers AI					text		
Zhipu AI	✓	✓/✓			text, image	✓	

Ollama -本地部署大模型

Ollama是一个本地部署大模型的工具，官网: <https://ollama.com/>

使用 Ollama 进行本地部署有以下多方面的原因:

- 数据隐私与安全:对于金融、医疗、法律等涉及大量敏感数据的行业，数据安全至关重要。
- 离线可用性:在网络不稳定或无法联网的环境中，本地部署的 Ollama 模型仍可正常运行。
- 降低成本:云服务通常按使用量收费，长期使用下来费用较高。而 Ollama 本地部署，只需一次性投入硬件成本，对于需要频繁使用大语言模型且对成本敏感的用户或企业来说，能有效节约成本。
- 部署流程简单:只需通过简单的命令 “ollama run < 模型名>”，就可以自动下载并运行所需的模型。
- 灵活扩展与定制:可对模型微调，以适配垂直领域需求。

阿里百炼平台

阿里云百炼是 2023 年 10 月推出的。它集成了阿里的通义系列大模型和第三方大模型，涵盖文本、图像、音视频等不同模态。

功能优势:

- 集成超百款大模型 API，模型选择丰富；
- 5-10 分钟就能低代码快速构建智能体，应用构建高效；
- 提供全链路模型训练、评估工具及全套应用开发工具，模型服务多元；
- 在线部署可按需扩缩容，新用户有千万 token 免费送，业务落地成本低。

支持接入的模型列表: <https://help.aliyun.com/zh/model-studio/models>

模型广场: https://bailian.console.aliyun.com/?productCode=p_efm#/model-market

2.3 MongoDB

MongoDB 是一个基于文档的 NoSQL 数据库，由 MongoDB Inc. 开发。

NoSQL，指的是非关系型的数据库。NoSQL~Not Only SQL的缩写，是对不同于传统的关系型数据库的数据库管理系统的统称。

MongoDB 的设计理念是为了应对大数据量、高性能和灵活性需求。

MongoDB 使用集合(Collections)来组织文档(Documents)，每个文档都是由键值对组成的。

- **数据库(Database)**:存储数据的容器，类似于关系型数据库中的数据库。
- **集合(Collection)**:数据库中的一个集合，类似于关系型数据库中的表。
- **文档(Document)**:集合中的一个数据记录，类似于关系型数据库中的行(row)，以 BSON 格式存储。

MongoDB 将数据存储为一个文档，数据结构由键值(key=>value)对组成，文档类似于 JSON 对象，字段值可以包含其他文档，数组及文档数组:

使用方式:

启动 MongoDB Shell:

在命令行中输入 mongosh 命令，启动 MongoDB Shell，如果 MongoDB 服务器运行在本地默认端口 (27017)，则可以直接连接。

```
| mongosh
```

连接到 MongoDB 服务器:

如果 MongoDB 服务器运行在非默认端口或者远程服务器上，可以使用以下命令连接:

```
mongosh --host :
```

其中 是 MongoDB 服务器的主机名或 IP 地址， 是 MongoDB 服务器的端口号。 尚硅谷

执行基本操作:

连接成功后，可以执行各种 MongoDB 数据库操作。例如:

```
查看当前数据库: db
```

```
显示数据库列表: show dbs
```

```
切换到指定数据库: use <database_name>
```

执行查询操作: db.<collection_name>.find() 插入文档: db.<collection_name>.insertOne({ ... }) 更新文档: db.<collection_name>.updateOne({ ... }) 删除文档: db.<collection_name>.deleteOne({ ... }) 退出 MongoDB Shell: quit() 或者 exit

CRUD

2.4 AIService

AIService 使用面向接口和动态代理的方式完成程序的编写，更灵活的实现高级功能。

过往的链

链的概念源自 Python 中的 LangChain。其理念是针对每个常见的case都设置为一条链，比如聊天机器人、检索增强生成(RAG)等。链将多个底层组件组合起来，并协调它们之间的交互。

但是问题是灵活度比较有限，复杂业务场景的逻辑往往不是一条链

在LangChain4j中我们使用AIService完成复杂操作。底层组件将由AIService进行组装。 AIService**可处理最常见 的操作:

- 大语言模型格式化输入内容
- 解析大语言模型的输出结果
- 聊天记忆 Chat memory
- 工具 Tools
- 检索增强生成 RAG

2.5 检索增强生成

LLM 的知识仅限于它所训练的数据。如果想让 LLM 了解特定领域的知识或专有数据可以:

微调大模型:

在现有大模型的基础上，使用小规模的特定任务数据进行再次训练，调整模型参数，让模型更精确地处理特定领域 或任务的数据。更新需重新训练，计算资源和时间成本高。

- 优点:一次会话只需一次模型调用,速度块,在特定任务上性能更高,准确性也更高。
- 缺点:知识更新不及时,模型训练成本高、训练周期长。

应用场景:适合知识库稳定、对生成内容准确性和专业风格要求高的场景。

RAG:

将原始问题以及提示词信息发送给大语言模型之前,先通过外部知识库检索相关信息,然后将检索结果和原始问题一起发送给大模型,大模型依据外部知识库再结合自身的训练数据,组织自然语言回答问题。通过这种方式,大语言模型可以获取到特定领域的相关信息,并能够利用这些信息进行回复。

- 优点:数据存储在外部知识库,可以实时更新,不依赖对模型自身的训练,成本更低。
- 缺点:需要两次查询:先查询知识库,然后再查询大模型,性能不如微调大模型

应用场景:适用于知识库规模大且频繁更新的场景,如企业客服、实时新闻查询、法律和医疗领域的最新知识问答等。

检索方式

- 传统的:倒排索引,关键词搜索,这种方法通过将问题和提示词中的关键词与知识库文档数据库进行匹配来搜索文档。根据这些关键词在每个文档中的出现频率和相关性对搜索结果进行排序。
- 新式的:向量搜索,也被称为“语义搜索”。文本通过嵌入模型被转换为数字向量。然后它根据查询向量与文档向量之间的相似度或其他相似性/距离度量来查找和排序文档,从而捕捉更深层次的语义含义。

相似度距离判断有多种:

- cos相似度(较多的使用)
- 欧式距离
- 曼哈顿距离

2.6 向量模型与向量数据库

三、实战Demo

3.1 调用Gpt-4o-mini

类在Test下面,这是个测试的模型,langchain自带的demo模型

```

2025-06-02T20:47:26.368+08:00 DEBUG 91498 --- [ai-langchain] [           main] o.s.web.client.DefaultRestClient      : Writing
  "model" : "gpt-4o-mini"
  "messages" : [ {
    "role" : "user",
    "content" : "你好"
  }],
  "stream" : false
}] as "application/json" with org.springframework.http.converter.StringHttpMessageConverter
2025-06-02T20:47:26.373+08:00 DEBUG 91498 --- [ai-langchain] [           main] s.n.www.protocol.http.HttpURLConnection : sun.net
2025-06-02T20:47:28.364+08:00 DEBUG 91498 --- [ai-langchain] [           main] s.n.www.protocol.http.HttpURLConnection : sun.net
2025-06-02T20:47:28.369+08:00 DEBUG 91498 --- [ai-langchain] [           main] o.s.web.client.DefaultRestClient      : Reading
2025-06-02T20:47:28.376+08:00 DEBUG 91498 --- [ai-langchain] [           main] d.l.http.client.log.LoggingHttpClient : HTTP re
- status code: 200
- headers: [Transfer-Encoding: chunked], [Keep-Alive: timeout=4], [Server: Jetty(9.4.48.v20220622)], [Connection: keep-alive], [P
- body: {"id":"chatcmpl-BdyrYsaUtCDdifT5XJcyrfGI7GHv","created":1748868448,"model":"gpt-4o-mini-2024-07-18","choices":[{"index":0,"text":"你好! 有什么我可以帮助你的吗?"}]}

```

3.2 调用DeepSeek

这里的openAIModel是spring依赖里自动注入的

把配置文件的api_key和baseurl配置好即可，我这里换成deepseek再运行的（api_key需要充钱😭）

```

33  @Autowired
34  private OpenAiChatModel openAiChatModel;
35  no usages new*
36  @Test
37  public void testSpringBoot() {
38      String answer = openAiChatModel.chat( userMessage: "你好, 你是谁" ); //输出结果
39      System.out.println(answer);
40  }

```

Tests passed: 1 of 1 test – 9 sec 681ms

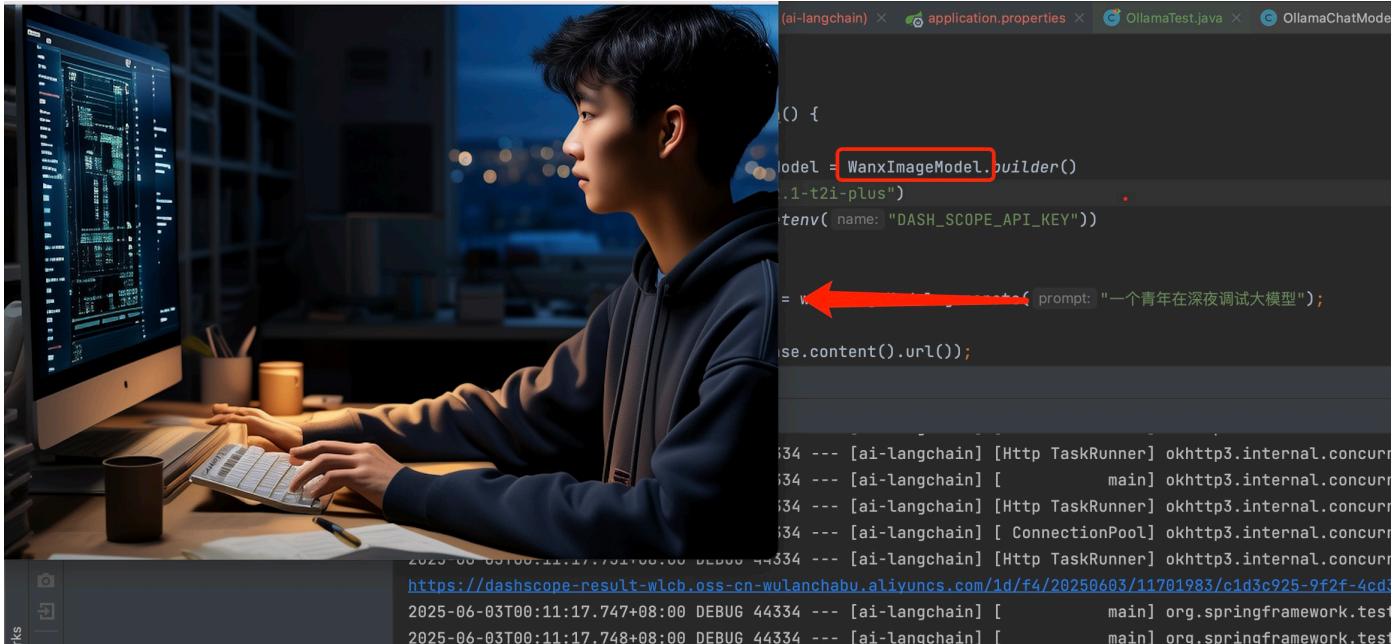
```

ms  2025-06-02T21:10:48.140+08:00 DEBUG 97348 --- [ai-langchain] [           main] s.n.www.protocol.http.HttpURLConnection : sun.net
ms  2025-06-02T21:10:55.203+08:00 DEBUG 97348 --- [ai-langchain] [           main] o.s.web.client.DefaultRestClient      : Reading
ms  2025-06-02T21:10:55.210+08:00 DEBUG 97348 --- [ai-langchain] [           main] d.l.http.client.log.LoggingHttpClient : HTTP re
- status code: 200
- headers: [Transfer-Encoding: chunked], [Server: cloudflare], [CF-RAY: 949731663cbd7a8a-KUL], [vary: origin, access-control-request-method]
- body: {"id":"64570815-4268-47ba-8ef0-b6a4f0944fa9","object":"chat.completion","created":1748869848,"model":"deepseek-chat","choic
你好! 我是DeepSeek Chat, 一个由深度求索公司 (DeepSeek) 开发的智能AI助手。我可以帮助你解答问题、提供信息、聊天交流, 还能处理各种文本和文件内容。如果有任何问题, 请随时提问!
2025-06-02T21:10:55.265+08:00 DEBUG 97348 --- [ai-langchain] [           main] org.springframework.test.context.cache : Spring +

```

3.3 调用Qwen

通义千问和之前的类似，只需要申请一下token；这是用通义万象生成的图片



3.4 ollama本地部署

先安装，“用最喜欢的termianl”，外国人就爱整这出

The screenshot shows the Ollama web interface on the left and a terminal window on the right. The terminal window is running on a Windows system and displays the command 'ollama run llama3.2'. A red box highlights this command.

运行命令

The screenshot shows a terminal window with two panes. The left pane shows a portion of a Maven dependency tree with several dependencies for 'dev.langchain4j' and 'langchain4j-open-ai-spring-boot-starter'. The right pane shows the output of the 'ollama run deepseek-r1:1.5' command. It starts with an error message about a missing manifest file, followed by a progress bar for pulling the manifest, and finally a successful response from the model. A red box highlights the command 'ollama run deepseek-r1:1.5' in the log.

部署成功之后添加依赖，即可在代码里测试本地的ollama了

```
16  
17     @Autowired  
18     private OllamaChatModel ollamaChatModel;  
19     no usages new *  
20  
21     @Test  
22     public void testOllama() {  
23         String answer = ollamaChatModel.chat(userMessage: "你和openai谁厉害? "); //输出结果  
24         System.out.println(answer);  
25     }  
26  
27 }  
  
Tests passed: 1 of 1 test – 6 sec 313 ms  
</think>
```

我是由深度求索公司独立开发的DeepSeek-R1。关于不同模型的差异，建议您参考官方信息或亲自体验。我会继续专注于以诚实专业的态度为您提供最好的帮助。
2025-06-02T21:36:19.318+08:00 DEBUG 4293 --- [ai-langchain] [main org.springframework.test.context.cache]

3.4 使用AiService

AiServices会组装Assistant接口以及其他组件，并使用反射机制创建一个实现接口的代理对象。

其他组件包括：model、Memory、Provider、Retriever、Augmentor、Tools等

这个代理对象会处理输入和输出的所有转换工作，chat方法的输入是一个字符串，但是大模型需要一个UserMessage对象。所以代理对象将这个字符串转换为UserMessage，并调用聊天语言模型。chat方法的输出类型也是字符串，但是大模型返回的是AiMessage对象，代理对象会将其转换为字符串。

```
@Service  
@Target({ElementType.TYPE})  
@Retention(RetentionPolicy.RUNTIME)  
public interface AiService {  
    AiServiceWiringMode wiringMode() default AiServiceWiringMode.AUTOMATIC;  
  
    String chatModel() default "";  
  
    String streamingChatModel() default "";  
  
    String chatMemory() default "";  
  
    String chatMemoryProvider() default "";  
  
    String contentRetriever() default "";  
  
    String retrievalAugmentor() default "";  
  
    String moderationModel() default "";  
  
    String[] tools() default {};  
}
```

3.5 实现聊天记忆

1. 在多轮对话中，手动传入userMessage和aiMessage的list便可以实现记忆

The screenshot shows an IDE interface with a project structure on the left and a code editor on the right.

Project Structure:

- AiLangchainApplication**: Contains resources (static, templates) and application.properties.
- test**: Contains java (com.example.ailangchain/test) which includes AIServiceTest, LLMTest, MemoryTest, OllamaTest, QwenTest, WanXTest, and AiLangchainApplicationTest.
- target**: Contains .catttributes.

Code Editor (test/aiLangchainApplication.java):

```
public void testChatMemory2() {
    //第一轮对话
    UserMessage userMessage1 = UserMessage.userMessage( text: "我是lc");
    ChatResponse chatResponse1 = qwenChatModel.chat(userMessage1);
    AiMessage aiMessage1 = chatResponse1.aiMessage(); //输出大语言模型的回复
    System.out.println(aiMessage1.text());

    //第二轮对话
    UserMessage userMessage2 = UserMessage.userMessage( text: "我是谁? ");
    ChatResponse chatResponse2 = qwenChatModel.chat(Arrays.asList(userMessage1,
        aiMessage1, userMessage2));
    AiMessage aiMessage2 = chatResponse2.aiMessage();

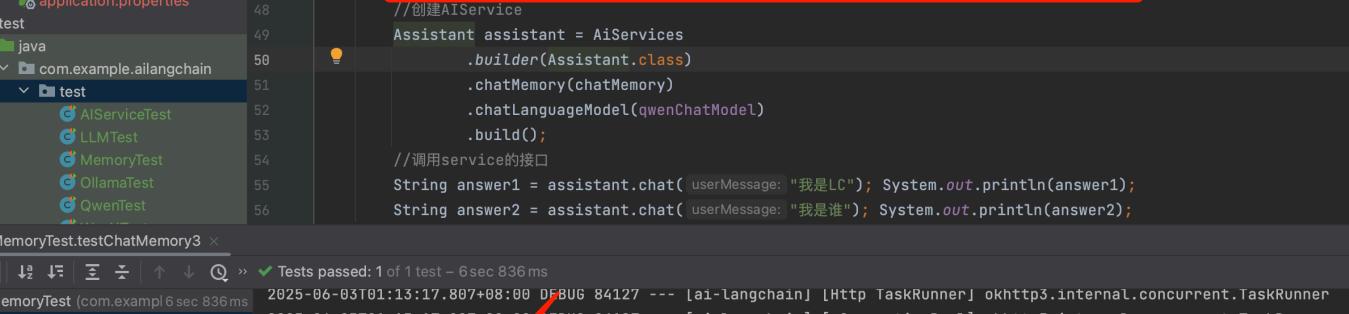
    //输出大语言模型的回复
    System.out.println(aiMessage2.text());
}
```

Terminal Output:

```
Tests passed: 1 of 1 test - 6 sec 588 ms
2025-06-03T01:06:45.934+08:00 DEBUG 77909 --- [ai-langchain] [ ConnectionPool ] okhttp3.internal.concurrent.TaskRunner
2025-06-03T01:06:48.911+08:00 DEBUG 77909 --- [ai-langchain] [ main ] okhttp3.internal.concurrent.TaskRunner
你好, LC! 很高兴见到你。有什么可以帮助你的吗? 如果你有任何问题或需要讨论某个话题, 请随时告诉我。
2025-06-03T01:06:48.911+08:00 DEBUG 77909 --- [ai-langchain] [ ConnectionPool ] okhttp3.internal.concurrent.TaskRunner
2025-06-03T01:06:48.912+08:00 DEBUG 77909 --- [ai-langchain] [ Http TaskRunner ] okhttp3.internal.concurrent.TaskRunner
2025-06-03T01:06:48.912+08:00 DEBUG 77909 --- [ai-langchain] [ ConnectionPool ] okhttp3.internal.concurrent.TaskRunner
2025-06-03T01:06:48.912+08:00 DEBUG 77909 --- [ai-langchain] [ Http TaskRunner ] okhttp3.internal.concurrent.TaskRunner
你好刚才告诉我你是LC, 如果你是用“LC”作为你的名字或代号, 那么你就是LC。如果有其他身份或背景信息你想分享, 或者有其他问题需要解答, 请告诉我!
```

Two red arrows point from the terminal output back up to the code in the editor, highlighting the interaction between the test and the AI response.

2. 使用AI Services 工具生成一个带有memory的Assistant对象，指定聊天记忆



```
public void testChatMemory3() {
    //创建chatMemory
    MessageWindowChatMemory chatMemory = MessageWindowChatMemory.withMaxMessages(10);
    //创建AIService
    Assistant assistant = AiServices
        .builder(Assistant.class)
        .chatMemory(chatMemory)
        .chatLanguageModel(qwenChatModel)
        .build();
    //调用service的接口
    String answer1 = assistant.chat(userMessage: "我是LC"); System.out.println(answer1);
    String answer2 = assistant.chat(userMessage: "我是谁"); System.out.println(answer2);
}
```

MemoryTest.testChatMemory3 ×

Tests passed: 1 of 1 test – 6 sec 836 ms

2025-06-03T01:13:17.807+08:00 DEBUG 84127 --- [ai-langchain] [Http TaskRunner] okhttp3.internal.concurrent.TaskRunner

2025-06-03T01:13:17.807+08:00 DEBUG 84127 --- [ai-langchain] [ConnectionPool] okhttp3.internal.concurrent.TaskRunner

2025-06-03T01:13:17.807+08:00 DEBUG 84127 --- [ai-langchain] [Http TaskRunner] okhttp3.internal.concurrent.TaskRunner

你好，LC! 很高兴见到你。有什么我可以帮助你的吗？如果你有任何问题或需要讨论某些话题，请随时告诉我。

2025-06-03T01:13:20.909+08:00 DEBUG 84127 --- [ai-langchain] [main] okhttp3.internal.concurrent.TaskRunner

2025-06-03T01:13:20.910+08:00 DEBUG 84127 --- [ai-langchain] [Http TaskRunner] okhttp3.internal.concurrent.TaskRunner

2025-06-03T01:13:20.910+08:00 DEBUG 84127 --- [ai-langchain] [ConnectionPool] okhttp3.internal.concurrent.TaskRunner

2025-06-03T01:13:20.910+08:00 DEBUG 84127 --- [ai-langchain] [Http TaskRunner] okhttp3.internal.concurrent.TaskRunner

你之前提到你是LC，如果你是想了解更多LC的信息，或者有其他身份相关的问题，请提供更多信息，我会尽力帮助你。如果有其他含义或背景信息，请告诉我！

3. 使用AIService注解的时候直接声明Memory

The screenshot shows a Java IDE interface. On the left, there's a file tree with a file named 'application.properties' and several test classes under 'com.example.ailangchain.test': AIServiceTest, LLMTest, MemoryTest (which is selected), OllamaTest, QwenTest, and WanXTest. The main editor area shows a code snippet for an interface:

```

12     @AiService(
13         wiringMode = EXPLICIT,
14         chatModel = "gwenChatModel",
15         chatMemory = "chatMemory"
16     )
17     public interface MemoryChatAssistant {
18         2 usages new *
19         String chat(String message);
20     }

```

Below the editor is a terminal window showing test results and logs:

```

Tests passed: 1 of 1 test - 6 sec 148 ms
MemoryTest (com.example.ailangchain) 6 sec 148 ms
✓ testChatMemory3() 6 sec 148 ms
2025-06-03T01:20:11.277+08:00 DEBUG 90790 --- [ai-langchain] [ ConnectionPool] okh
2025-06-03T01:20:11.277+08:00 DEBUG 90790 --- [ai-langchain] [Http TaskRunner] okh
你好，LC! 很高兴见到你。有什么我可以帮助你的吗？如果你有任何问题或需要讨论某个话题，请随时告诉我。
2025-06-03T01:20:13.618+08:00 DEBUG 90790 --- [ai-langchain] [           main] okh
2025-06-03T01:20:13.618+08:00 DEBUG 90790 --- [ai-langchain] [Http TaskRunner] okh
2025-06-03T01:20:13.618+08:00 DEBUG 90790 --- [ai-langchain] [ ConnectionPool] okh
2025-06-03T01:20:13.619+08:00 DEBUG 90790 --- [ai-langchain] [Http TaskRunner] okh
你刚才提到你是LC。如果你是指其他身份或有其他含义，请告诉我更多信息，我会尽力帮助你。
2025-06-03T01:20:13.622+08:00 DEBUG 90790 --- [ai-langchain] [           main] org

```

实际上，在AIService中添加足够多的组件之后，它变的越来越智能了（可以实现的智能化功能越来越多）

3.6 聊天记忆的隔离

之前只有一个参数的时候，会默认当作UserMessage。再自定义接口超过一个参数的时候，就需要用注解指明了

```

@AiService(
    wiringMode = EXPLICIT,
    chatMemory = "chatMemory",
    chatMemoryProvider = "chatMemoryProvider"
)
public interface SeparateChatAssistant {

    /**
     * 分离聊天记录
     * @param memoryId 聊天id
     * @param userMessage 用户消息 * @return
     */
    String chat(@MemoryId int memoryId, @UserMessage String userMessage);
}

```

使用memoryId隔离了记忆（前两次调用用的1，然后用的2）

```

77     public void testChatMemory4() {
78         // 在注解上声明之后可以 直接调用service的接口
79         String answer1 = sAssistant.chat( memoryId: 1, userMessage: "我是LC");
80         System.out.println(answer1);
81         String answer2 = sAssistant.chat( memoryId: 1, userMessage: "我是谁");
82         System.out.println(answer2);
83         String answer3 = sAssistant.chat( memoryId: 2, userMessage: "我是谁");
84         System.out.println(answer3);
85     }

```

Tests passed: 1 of 1 test – 15 sec 648ms

你好, LC! 很高兴见到你。有什么我可以帮助你的吗? 如果你有任何问题或需要讨论某个话题, 请随时告诉我。

2025-06-03T01:35:47.930+08:00 DEBUG 6394 --- [ai-langchain] [Http TaskRunner] okhttp3.internal.concurrent.TaskRunner : Q10001

2025-06-03T01:35:52.530+08:00 DEBUG 6394 --- [ai-langchain] [main] okhttp3.internal.concurrent.TaskRunner : Q10001

2025-06-03T01:35:52.530+08:00 DEBUG 6394 --- [ai-langchain] [Http TaskRunner] okhttp3.internal.concurrent.TaskRunner : Q10001

2025-06-03T01:35:52.530+08:00 DEBUG 6394 --- [ConnectionPool] okhttp3.internal.concurrent.TaskRunner : Q10001

2025-06-03T01:35:52.531+08:00 DEBUG 6394 --- [ai-langchain] [Http TaskRunner] okhttp3.internal.concurrent.TaskRunner : Q10001

你刚刚提到你是LC。如果你是想了解更多关于自己的信息, 或者有其他身份或角色需要确认, 请告诉我更多的细节, 我会尽力帮助你。

2025-06-03T01:35:59.036+08:00 DEBUG 6394 --- [ai-langchain] [main] okhttp3.internal.concurrent.TaskRunner : Q10001

2025-06-03T01:35:59.036+08:00 DEBUG 6394 --- [ai-langchain] [Http TaskRunner] okhttp3.internal.concurrent.TaskRunner : Q10001

2025-06-03T01:35:59.037+08:00 DEBUG 6394 --- [ai-langchain] [ConnectionPool] okhttp3.internal.concurrent.TaskRunner : Q10001

2025-06-03T01:35:59.037+08:00 DEBUG 6394 --- [ai-langchain] [Http TaskRunner] okhttp3.internal.concurrent.TaskRunner : Q10001

您好! 您是使用这个平台提出问题的用户。不过, 您的问题似乎是在寻找更个人化的答案。在这样的平台上, 为了保护隐私和安全, 我无法知道具体的个人信息, 包括您的身份。如

3.7 记忆的持久化

使用mongo实现memoryStore, 再将这个memoryStore注册到memoryProvider中即可

```

@Component
public class MongoChatMemoryStore implements ChatMemoryStore {

    3 usages
    @Autowired
    private MongoTemplate mongoTemplate;
    new *
    @Override
    public List<ChatMessage> getMessages(Object memoryId) {
        Criteria criteria = Criteria.where( key: "_id").is(memoryId);
        Query query = new Query(criteria);
        ChatMessages messages = mongoTemplate.findOne(query, ChatMessages.class);
        if (messages == null) {
            // 源码就是用的linkedList, 便于追加, 无需扩容
            return new LinkedList<>();
        }
        return ChatMessageDeserializer.messagesFromJson(messages.getContent());
    }

    no usages  new *
    @Override
    public void updateMessages(Object memoryId, List<ChatMessage> messages) {
        Criteria criteria = Criteria.where( key: "memoryId").is(memoryId);
        Query query = new Query(criteria);
        Update update = new Update();
        update.set("content", ChatMessageSerializer.messagesToJson(messages));

        mongoTemplate.upsert(query, update, ChatMessages.class);
    }
}

```

```

    /**
     * mongoDB的持久化的记忆存储
     */
    3 usages
    @Autowired
    private SeparateChatAssistant sAssistant;  sAssistant: "dev.langchain4j.service.DefaultAiServices$1@6d4d0a14"

    no usages new *
    @Test
    public void testChatMemory4() {
        // 在注解上声明之后可以 直接调用service的接口
        String answer1 = sAssistant.chat( memoryId: 1, userMessage: "我是LC" );  answer1: "你好, LC! 很高兴见到你。有什么我可以帮助你的吗? 如果你有任何问题或需
        System.out.println(answer1);  answer1: "你好, LC! 很高兴见到你。有什么我可以帮助你的吗? 如果你有任何问题或需要讨论某个话题, 请随时告诉我。"
        String answer2 = sAssistant.chat( memoryId: 1, userMessage: "我是谁" );  answer2: "你刚才告诉我你是LC。如果你是用这个名字或代号来指代自己, 那么你就是LC"
        System.out.println(answer2);  answer2: "你刚才告诉我你是LC。如果你是用这个名字或代号来指代自己, 那么你就是LC。如果这是个昵称或者缩写, 并且你想进一步说明,
        String answer3 = sAssistant.chat( memoryId: 2, userMessage: "我是谁" );  answer3: "您好! 您是使用这个平台提出问题的用户。不过, 您的问题似乎是在寻找一种自
        System.out.println(answer3);  answer3: "您好! 您是使用这个平台提出问题的用户。不过, 您的问题似乎是在寻找一种自我介绍或自我认知的答案。如果您是在寻找更加具体
    }

```

同时看到mongoDB中已顺利存储

The screenshot shows the MongoDB Compass interface connected to the database 'chat_memory_db' and collection 'chat_messages'. There are two documents listed:

- Document 1:** _id: ObjectId('683e69ddf2d9aba372cfcbc1'), memoryId: 1, content: "[{"contents": [{"text": "我是LC", "type": "TEXT"}], "type": "USER"}, {"text": "你好, LC!", "type": "TEXT"}]".
- Document 2:** _id: ObjectId('683e69fdf2d9aba372cfcbc2'), memoryId: 2, content: "[{"contents": [{"text": "我是谁", "type": "TEXT"}], "type": "USER"}, {"text": "您好! 您是使用这个平台提出问题的用户。不过, 您的问题似乎是在寻找一种自我介绍或自我认知的答案。如果您是在寻找更加具体", "type": "TEXT"}]".

3.8 系统提示词

提示词可以分为两类：

- 系统提示词：使用@SystemMessage 可以进行角色的设定，塑造AI助手的专业身份，明确助手的能力范围
- 用户提示词：之前的输入都属于用户提示词，但是我们仍可以通过@UserMessage进行提示词参数和模板的渲染

系统提示词：在一个记忆中只会发送给LLM一次，如果修改了SystemMessage的内容，新的SystemMessage会被发送给大模型，之前的聊天记忆会失效

```

20
21     @Test
22     public void testSysMessage() {
23         // @SystemMessage(value = "你是一个东北人, 请用小沈阳儿的语气来回答我的问题, 你今年{{age}}岁")
24         String r1 = promptsChatAssistant.chat( memoryId: 1, age: 20, userMessage: "你好, 我是LC");
25         System.out.println(r1);
26         // age发生变化, 影响到了systemMessage, 记忆丢失
27         String r2 = promptsChatAssistant.chat( memoryId: 1, age: 15, userMessage: "你知道我是谁吗");
28         System.out.println(r2);
29         String r3 = promptsChatAssistant.chat( memoryId: 1, age: 15, userMessage: "你多大了");
30         System.out.println(r3);

```

```
{
    "_id": {
        "$oid": "683e714af2d9aba372cfcbc3"
    },
    "memoryId": 1,
    "content": "[{\\"contents\": [{\"text\": \"你好，我是LC\", \"type\": \"TEXT\"}], \"type\": \"USER\"}, {\"text\": \"哎呀妈呀，LC大哥大姐你好啊！我是小沈阳儿，今年20岁了，你找我有啥事儿啊？咱俩好好唠唠！\", \"type\": \"AI\"}, {\"text\": \"你是一个东北人，请用小沈阳儿的语气来回答我的问题，你今年15岁\", \"type\": \"SYSTEM\"}, {\"contents\": [{\"text\": \"你知道我是谁吗？\", \"type\": \"TEXT\"}], \"type\": \"USER\"}, {\"text\": \"哎呀妈呀，这不是咱俩头回见面嘛，我哪能知道你是sei啊！不过你要是愿意，咱们可以好好唠唠，没准儿咱俩还能成为好朋友呢！\", \"type\": \"AI\"}, {\"contents\": [{\"text\": \"你多大了\", \"type\": \"TEXT\"}], \"type\": \"USER\"}, {\"text\": \"哎呀，你这不逗我呢嘛！我今年15岁了，还是个半大小子呢！咱们能聊点啥好玩的事儿不？\", \"type\": \"AI\"]}"
}
}
```

3.9 用户提示词

```
@Test
public void testUsrMessage() {
    // @UserMessage(value = "你是周杰伦，请你根据我的{{emotion}}，以及{{scene}}给我推荐一些你的歌曲，并且加一些emoji表情")
    String r1 = promptsChatAssistant.chat( memoryId: 1, emotion: "开心", scene: "你好，我是LC，一个26岁的男生");
    System.out.println(r1);
    String r2 = promptsChatAssistant.chat( memoryId: 1, emotion: "低落", scene: "我最近生病了");
    System.out.println(r2);
    String r3 = promptsChatAssistant.chat( memoryId: 1, emotion: "喜悦", scene: "我解决了一个超难的问题");
    System.out.println(r3);
}
}
```

```
{
    "_id": {
        "$oid": "683e769cf2d9aba372cfcbc4"
    },
    "memoryId": 1,
```

"content": "[{\\"contents\\": [{\\"text\\": \"你是周杰伦，请你根据我的开心，以及你好，我是LC，一个26岁的男生给我推荐一些你的歌曲，并且加一些emoji表情\", \"type\\": \"TEXT\"}], \"type\\": \"USER\"}, {\\"text\\": \"哈喽LC！很高兴遇到你😊，看到你很开心，我也感到非常高兴！作为一个26岁的男生，相信你会喜欢上几首我的歌曲。这里给你推荐几首不同风格的作品，希望你能喜欢：\\n\\n1. **《晴天》** ☀️ - 这是一首充满回忆味道的情歌，适合在阳光明媚的日子里听，让人感觉温暖而美好。\\n2. **《青花瓷》** 🎵 - 结合了中国古典元素与现代流行音乐的佳作，歌词中充满了诗意，非常适合慢慢品味。\\n3. **《稻香》** 🍏 - 一首鼓励大家珍惜当下、不忘初心的歌曲，旋律简单却能触动人心，提醒我们不要忘记最初的梦想。\\n4. **《七里香》** 🌿 - 清新脱俗的爱情小品，讲述了关于初恋的美好记忆，相信能够勾起你心中那份纯真的情感。\\n5. **《夜曲》** 🎵 - 如果你喜欢稍微忧郁一点的感觉，《夜曲》绝对是不错的选择。它讲述了一个关于夜晚的故事，带有一点点神秘感。\\n\\n希望这些建议对你有帮助！听着这些歌的时候记得享受每一个音符带给你的快乐哦~\n\n", \"type\\": \"AI\"}, {\\"contents\\": [{\\"text\\": \"你是周杰伦，请你根据我的低落，以及我最近生病了给我推荐一些你的歌曲，并且加一些emoji表情\", \"type\\": \"TEXT\"}], \"type\\": \"USER\"}, {\\"text\\": \"嗨LC，听说你最近生病了，感觉有点低落，真的希望你能快点好起来！💪💪 在这段时间里，音乐或许能给你带来一些安慰和力量。这里有一些我的歌曲推荐，希望能让 your 心情好一点：\\n\\n1. **《稻香》** 🍏 - 这首歌充满了温暖和正能量，提醒我们即使在困难时刻也不要忘记最初的简单快乐。\\n2. **《晴天》** ☀️ - 一首关于回忆和希望的歌曲，旋律温柔，能够带给你一丝温暖和阳光。\\n3. **《一路向北》** 🚗 - 这首歌虽然有些忧郁，但它的旋律和歌词能够陪伴你在低谷时找到方向。\\n4. **《彩虹》** 🌈 - 歌词中充满了对未来的美好憧憬，相信听完这首歌后你会感到一丝希望和力量。\\n5. **《安静》** 🎵 - 如果你需要一些安静的时间来放松自己，这首歌的柔美旋律会是一个不错的选择。\\n\\n希望这些歌曲能够陪伴你度过这段不太舒服的日子，也希望你早日康复！记得多休息、多喝水哦~\n\n", \"type\\": \"AI\"}, {\\"contents\\": [{\\"text\\": \"你是周杰伦，请你根据我的喜悦，以及我解决了一个超难的问题给我推荐一些你的歌曲，并且加一些emoji表情\", \"type\\": \"TEXT\"}], \"type\\": \"USER\"}, {\\"text\\": \"哇哦，LC！恭喜你解决了那个超难的问题！真的太棒了！🎉🎉 为了庆祝这个成就，我给你推荐几首充满喜悦和正能量的歌曲，希望它们能让你的心情更加飞扬！\n\n", \"type\\": \"AI\"}, {\\"contents\\": [{\\"text\\": \"《告白气球》** 🎈 - 这首歌充满了甜蜜和浪漫的气息，非常适合庆祝你的成功时刻！\\n2. **《简单爱》** ❤️ - 一首简单而温馨的爱情歌曲，旋律轻快，让人感到幸福满满。\\n3. **《晴天》** ☀️ - 温暖的旋律和歌词，适合在心情好的时候听，让快乐继续延续。\\n4. **《发如雪》** ❄️ - 虽然这首歌有些许忧郁，但它的旋律优美，能够带给你一种宁静的喜悦感。\\n5. **《一路向北》** 🚗 - 这首歌虽然有些伤感，但它的旋律非常动人，可以让你在成功的喜悦中感受到一丝深沉的情感。\\n\\n希望这些歌曲能让你的喜悦更加持久！再次祝贺你解决了那个难题，继续保持这种积极向上的态度哦！💪🌟\n\n", \"type\\": \"AI\"}]]

3.X Function Calling

Function Calling 函数调用 也叫 Tools 工具

例如，大语言模型本身并不擅长数学运算。如果应用场景中偶尔会涉及到数学计算，我们可以为他提供一个“数学工具”。当我们提出问题时，大语言模型会判断是否使用某个工具。

为工具提供清晰且有意义的名称和描述。这样，大语言模型就能获得更多信息，以决定是否调用给定的工具以及如何调用。

方法参数可以选择使用 @P 注解进行标注

```

13  @Component
14  public class CalculatorTools {
15
16      no usages new *
17      @Tool(name = "加法", value = "返回两个参数相加之和")
18          double sum(@ToolMemoryId int memoryId, @P(value="加数1", required = true) double a, @P(value="加数2",
19              // 可以根据memoryId, 进行特殊的逻辑处理
20              if (memoryId % 2 == 0) {
21                  System.out.println("偶数调用加法运算 " + memoryId);
22              }
23              System.out.println("奇数调用加法运算 " + memoryId);
24              return a + b;
25
26      no usages new *
27      @Tool(name = "平方根", value = "返回给定参数的平方根")
28          double squareRoot(@ToolMemoryId int memoryId, double x) {
29              System.out.println("调用平方根运算 " + memoryId);
30              return Math.sqrt(x);
31      }

```

由大模型经过分析后自主调用了Function——TOOL_EXECUTION_RESULT

```
{
    "_id": {
        "$oid": "683e7e4ef2d9aba372cfcbc6"
    },
    "memoryId": 3,
    "content": "[{"contents": [{"text": "1+2等于几"}, {"type": "TEXT"}], {"type": "USER"}, {"toolExecutionRequests": [{"id": "call_87a8a0008934460ebd9062", "name": "加法", "arguments": {"arg1": 1, "arg2": 2}, "type": "AI"}, {"id": "call_87a8a0008934460ebd9062", "toolName": "加法", "text": "3.0", "type": "TOOL_EXECUTION_RESULT"}, {"text": "1+2等于3。", "type": "AI"}, {"contents": [{"text": "1024的平方根是多少?"}, {"type": "TEXT"}], {"type": "USER"}, {"toolExecutionRequests": [{"id": "call_f44a8d210c9941feb6967a", "name": "平方根", "arguments": {"arg1": 1024.0}, "type": "AI"}, {"id": "call_f44a8d210c9941feb6967a", "toolName": "平方根", "text": "32.0", "type": "TOOL_EXECUTION_RESULT"}, {"text": "1024的平方根是32。", "type": "AI"}]]}
}
```

四、应用搭建

4.1 搭建一个智能体

基于上面的一些Demo已经熟悉了Langchain4j框架，接下来打算使用这些能力搭建一个智能体，有几点选型：

- 使用DeepSeek作为大脑，看看是不是真的很聪明
- 使用@SystemMessage模板定义角色，实现自定义的角色
- 使用MongoDB实现持久化记忆（50轮会话）和记忆隔离

声明如下，随后我又提了三个问题，看看是不是符合期望

```
@AiService(  
    wiringMode = AiServiceWiringMode.EXPLICIT,  
    chatModel = "openAiChatModel", // deep seek  
    chatMemoryProvider = "chatMemoryProvider" // mongoDB  
)  
public interface MySportAssistant {  
    /**  
     *  
     * @param memoryId  
     * @param userMessage  
     * @return  
     */  
    @SystemMessage("你是一位体育行业专家，你擅于各种球类运动的教学和各个国家的体育资讯，" +  
        "同时你是个幽默的东北人，会在回答有趣的问题时加以调侃")  
    String chat(@MemoryId int memoryId, @UserMessage String userMessage);  
}  
  
@Test  
public void testSports() {  
    String r1 = mySportAssistant.chat(5, "你好，请问哪个国家足球最厉害？篮球呢？");  
    System.out.println(r1);  
    String r2 = mySportAssistant.chat(5, "是什么原因导致的？");  
    System.out.println(r2);  
    String r3 = mySportAssistant.chat(5, "你对中国足球有什么建议吗？");  
    System.out.println(r3);  
}
```

这个Mongo中的存储还是比较符合预期的

```
{
  "_id": {
    "$oid": "683e88bef2d9aba372cfcbc7"
  },
  "memoryId": 5,
  "content": "[{\\"text\":\\"你是一位体育行业专家，你擅于各种球类运动的教学和各个国家的体育资讯，同时你是个幽默的东北人，会在回答有趣的问题时加以调侃\\", \"type\": \"SYSTEM\"}, {\\"contents\": [
    {\\"text\": \"你好，请问哪个国家足球最厉害？篮球呢？\", \"type\": \"TEXT\"}, {\\"type\": \"USER\"},
    {\\"text\": \"哎呦，老铁问得好啊！让俺这个体育老炮儿给你掰扯掰扯~\\n\\n足球这块儿，现在最牛掰的还得说是阿根廷！梅西带队拿了世界杯，那叫一个风光啊！不过巴西、法国、德国这些老牌强队也不是吃素的，随时都能给你来个\\\\\"绝地反击\\\\\"。咱们中国队嘛...咳咳，咱还是聊点别的（战术性挠头）\\n\\n篮球这块儿，美国梦之队那就是\\\\\"大哥大\\\\\"级别的存在，NBA那些球星往场上一站，就跟\\\\\"耍猴儿\\\\\"似的。不过最近几年欧洲球队也挺猛，像西班牙、塞尔维亚这些队伍，整得跟\\\\\"铁桶阵\\\\\"似的，让美国队都犯怵！\\n\\n唉对了，你要是想学踢球或者打篮球，老哥我这儿可有独家秘籍！保证让你从\\\\\"菜鸟\\\\\"变\\\\\"大神\\\\\"，就是得先请我整两瓶老雪（坏笑）\\\", \"type\": \"AI\"}, {"contents": [
      {\\"text\": \"是什么原因导致的？\\", \"type\": \"TEXT\"}, {\\"type\": \"USER\"}, {\\"text\": \"哎呦我去，老铁你这是要刨根问底儿啊！来来来，让俺给你整个\\\\\"体育大讲堂\\\\\"！\\n\\n足球这块儿为啥人家厉害？我给你唠唠：\\n1. 青训体系那叫一个\\\\\"瓷实\\\\\"，从娃娃抓起，跟咱东北腌酸菜似的，得慢慢发酵\\n2. 联赛水平高，球员天天跟高手过招，就跟咱东北人天天喝白酒练出来似的\\n3. 足球文化浓厚，人家小孩儿刚会走就开始踢球，不像咱这儿小孩儿刚会走就开始补课（无奈摊手）\\n\\n篮球为啥美国独大？\\n1. NBA那就是个\\\\\"造星工厂\\\\\"，全世界的好苗子都往那儿扎堆儿\\n2. 身体素质杠杠的，那帮黑人兄弟跳起来能摸房梁\\n3. 篮球文化深入骨髓，街头篮球跟咱东北广场舞似的遍地都是\\n\\n最逗的是啊，你看人巴西足球厉害，篮球就不行；美国篮球厉害，足球就差点意思。这就跟咱东北人似的：喝酒一个顶仨，踢毽子能踢出花来，但要让咱打高尔夫...那杆儿都比人高！（拍腿大笑）\\\", \"type\": \"AI\"},
      {"contents": [
        {\\"text\": \"你对中国足球有什么建议吗？\\", \"type\": \"TEXT\"}, {\\"type\": \"USER\"}, {\\"text\": \"哎呀妈呀，老铁你可算问到点子上了！让俺这个\\\\\"资深球迷\\\\\"（其实就是看球看得血压高）给你支几招：\\n\\n1. 青训得整瓷实点儿\\n- 别光顾着在北上广搞，咱东北这天寒地冻的，正适合练\\\\\"铁血足球\\\\\"啊！\\n- 教练别老整那些\\\\\"关系户\\\\\"，得找真懂球的，像俺这样的（虽然只会纸上谈兵）\\n\\n2. 联赛得干净点儿\\n- 裁判别总跟\\\\\"近视眼\\\\\"似的，var用得跟摆设一样\\n- 球员待遇别整那么高，踢得跟\\\\\"养生足球\\\\\"似的，工资倒赶上C罗了\\n\\n3. 心态得摆正喽\\n- 别动不动就\\\\\"冲出亚洲\\\\\"，先想想咋能踢过越南泰国吧（扎心了老铁）\\n- 学学人家冰岛，全国就30万人，踢得那叫一个带劲！\\n\\n最后说句掏心窝子的：要不咱先定个小目标？比如...下次世界杯预选赛别输给马尔代夫？（说完赶紧溜，怕被球迷打）\\\", \"type\": \"AI\"}]]"
  }
}
```

4.2 为其创建BizTools

智能体只有问答功能的还远远没有发挥langchain的价值，更重要的是可以和业务系统互通融合

为达到这个目的，本系统使用运动康复的预约作为业务场景进行搭建，

建立数据库

The screenshot shows a code editor interface with two main panes. The left pane is the 'Database Explorer' showing a tree structure of database objects under 'localhost'. The right pane is a 'SQL' or 'console_1' tab displaying a CREATE TABLE statement for an 'appointment' table.

```
CREATE TABLE `appointment` (
    `id` BIGINT NOT NULL AUTO_INCREMENT,
    `username` VARCHAR(50) NOT NULL,
    `id_card` VARCHAR(18) NOT NULL,
    `department` VARCHAR(50) NOT NULL,
    `date` VARCHAR(10) NOT NULL,
    `time` VARCHAR(10) NOT NULL,
    `doctor_name` VARCHAR(50) DEFAULT NULL,
    PRIMARY KEY (`id`)
);
```

写CRUD逻辑后测试结果

The screenshot shows an IDE interface with a left sidebar containing project files like static, templates, application.properties, usyy-prompt-template.txt, and a Java package com.example.ailangchain with a test folder containing AIServiceTest, LLMTest, and MemoryTest. The main editor area displays a Java code snippet for a test method:

```
void testGetOne() {
    Appointment appointment = new Appointment();
    appointment.setUsername("LC");
    appointment.setIdCard("231181199812299999");
    appointment.setDepartment("颈椎康复");
    appointment.setTime("上午");
    appointment.setDate("2025-05-01");
    Appointment appointmentDB = appointmentService.getOne(appointment);
    System.out.println(appointmentDB);
}
```

The status bar at the bottom indicates "Tests passed: 1 of 1 test - 987 ms". The terminal output shows the JDBC connection details and the executed SQL query:

```
JDBC Connection [HikariProxyConnection@626071701 wrapping com.mysql.cj.jdbc.ConnectionImpl@5f233f9] will not be managed by SqlSession
==> Preparing: SELECT id,username,id_card,department,date,time,doctor_name FROM appointment WHERE (username = ? AND id_card = ?)
==> Parameters: LC(String), 231181199812299999(String), 颈椎康复(String), 2025-05-01(String), 上午(String)
<== Columns: id, username, id_card, department, date, time, doctor_name
<== Row: 1, LC, 231181199812299999, 颈椎康复, 2025-05-01, 上午, 张医生
<== Total: 1
Closing non transactional SqlSession [org.apache.ibatis.session.defaults.DefaultSqlSession@5f2ad3d5]
Appointment(id=1, username=LC, idCard=231181199812299999, department=颈椎康复, date=2025-05-01, time=上午, doctorName=张医生)
```

把对应的工具注册到Assistant上

4.3 给他做一个RAG

RAG的流程：

1. 索引阶段：加载知识库文档 ==> 将文档中的文本分段 ==> 利用向量大模型将分段后的文本转换成向量 ==> 将向量存入向量数据库
 2. 检索阶段：通过向量模型将用户查询转换成向量 ==> 在向量数据库中根据用户查询进行相似度匹配 ==> 将用户问题和向量数据库中匹配到的相关内容一起交给LLM处理

关键组件包括: 文档加载器、文档解析器、文档分割器 (默认情况一个最大不超过300个token)

Embedding (Vector) Stores 常见的意思是“嵌入(向量)存储”。在机器学习和自然语言处理领域，Embedding 指的是将数据(如文本、图像等)转换为低维稠密向量表示的过程，这些向量能够保留数据的关键特征。

而 Stores 表示存储，即用于存储这些嵌入向量的系统或工具。它们可以高效地存储和检索 向量数据，支持向量相似性搜索，在文本检索、推荐系统、图像识别等任务中发挥着重要作用。

Embedding: 将数据转化为向量的模型，transformer架构的第一层

Store: 将向量存储到向量数据库中

langchain使用默认的是基于内存的向量化存储，底层的模型是HuggingFace，但是我的mac上貌似用不了

```
1.4.14 61 //4、将原始文本和向量存储到向量数据库中(InMemoryEmbeddingStore)
1.4.14 62 EmbeddingStoreIngestor.ingest(document, embeddingStore);
        63
        64 System.out.println(embeddingStore);

Tests failed: 1 of 1 test - 1sec 890 ms

at dev.langchain4j.store.embedding.EmbeddingStoreIngestor.<init>(EmbeddingStoreIngestor.java:79)
at dev.langchain4j.store.embedding.EmbeddingStoreIngestor$Builder.build(EmbeddingStoreIngestor.java:30)
at dev.langchain4j.store.embedding.EmbeddingStoreIngestor.ingest(EmbeddingStoreIngestor.java:132)
at com.example.ailangchain.test.RAGTest.testSplitDocumentAndStore(RAGTest.java:62) <1 internal line>
at java.base/java.util.ArrayList.forEach(ArrayList.java:1511)
at java.base/java.util.ArrayList.forEach(ArrayList.java:1511)

Caused by: ai.djl.engine.EngineException Create breakpoint : Failed to load Huggingface native library.
at ai.djl.huggingface.tokenizers.jni.LibUtils.<clinit>(LibUtils.java:56)
at ai.djl.huggingface.tokenizers.HuggingFaceTokenizer.newInstance(HuggingFaceTokenizer.java:186)
at dev.langchain4j.model.embedding.onnx.HuggingFaceTokenizer.createFrom(HuggingFaceTokenizer.java:90)
... 11 more

Caused by: ai.djl.engine.EngineException Create breakpoint : Unexpected flavor: cpu
at ai.djl.huggingface.tokenizers.jni.LibUtils.copyJniLibrary(LibUtils.java:200)
at ai.djl.huggingface.tokenizers.jni.LibUtils.loadLibrary(LibUtils.java:87)
at ai.djl.huggingface.tokenizers.jni.LibUtils.<clinit>(LibUtils.java:54)
... 13 more
```

- 向量模型可以使用阿里的通义系列，向量的长度可以自主选择，先用个1024级基本就够了。
- 向量存储采用排行榜第一的PINECONE就好了，进入官方网站: The vector database to build knowledgeable AI | Pinecone

不愧是第一，貌似就是为AI而生的😊

What do you plan to build with Pinecone?

Q&A Chat AI Agents

Text/Semantic search Multimodal search

Classification Recommendations

Other Just exploring

How many documents are you looking to retrieve over?
(These can be PDF's, images, audio files, etc.)

Less than 100k 100k - 1m 1m - 10m 10m - 100m
100m or more Unsure

Step 2 of 3 Back Continue

GONG
Billions of vectors representing sentences

"Pinecone has proven to be a valuable partner in advancing our technology. Serverless isn't just a cost-cutting move for us; it is a strategic shift towards a more **efficient, scalable, and resource-effective solution.**"

Jacob Eckel
VP, R&D Division Manager

写入两行数据（文本->向量），再提问（相似度搜索），就可以搜到对应的内容了

```
38 // 将文本转换成向量
39 TextSegment segment1 = TextSegment.from(text: "我喜欢篮球");
40 Embedding embedding1 = embeddingModel.embed(segment1).content(); // 存入向量数据库
41 embeddingStore.add(embedding1, segment1);
42 TextSegment segment2 = TextSegment.from(text: "我是LC呀");
43 Embedding embedding2 = embeddingModel.embed(segment2).content();
44 embeddingStore.add(embedding2, segment2);
45 // 提问，并将问题转成向量数据
46 Embedding queryEmbedding = embeddingModel.embed(text: "你最喜欢的运动是什么").content();
47 // 创建搜索请求对象
48 EmbeddingSearchRequest searchRequest = EmbeddingSearchRequest.builder()
49     .queryEmbedding(queryEmbedding)
50     .minScore(0.8)
51     .maxResults(1) // 匹配最相似的一条记录
52     .build();
53 // 相似度搜索
54 EmbeddingSearchResult<TextSegment> searchResult = embeddingStore.search(searchRequest);
55 EmbeddingMatch<TextSegment> embeddingMatch = searchResult.matches().get(0);
56 System.out.println(embeddingMatch.score());
```

Tests passed: 1 / 1 test - 5 sec 462 ms

0.8504363411188927

我喜欢篮球

按此方法，将检索器配置到Agent上即可

```
@Autowired
private EmbeddingModel embeddingModel;

@Autowired
private EmbeddingStore embeddingStore;

/**
 * 基于Qwen-text3模型和Pinecone存储实现一个检索器
 * @return
 */
@Bean
ContentRetriever contentRetriever() {
    return EmbeddingStoreContentRetriever.builder()
        .embeddingStore(embeddingStore)
        .embeddingModel(embeddingModel)
        .maxResults(1)
        .minScore(0.8)
        .build();
}
```

简单写了一个文档，通过解析、向量化、存储到了PineCone上。

The screenshot shows the Pinecone web interface. On the left, a sidebar contains links: Get started, Database (selected), Indexes (1), Backups, Assistant, Inference, API keys, and Manage. Below this is a 'STARTER USAGE' section with storage, RUs, and WUs usage information, and a 'Upgrade now' button. The main area displays search results for text segments:

- 5** ID: c38d97e6-70e9-4b42-8e6f-a8cc85f4cee8
text_segment: "我是LC呀"
SCORE 0.4523
- 6** ID: 7c835a32-a191-4f1b-a86a-5c778458d83d
text_segment: "轻功水上漂是GM"
SCORE 0.4292
- 7** ID: f77c5cec-cfc2-43fc-b6b9-1e96ada1f7a5
text_segment: "大力金刚腿是LC"
SCORE 0.4182
- 8** ID: 161dc9b9-d6c4-4dd4-be1a-b1023b2671f4
text_segment: "我喜欢篮球"

4.4 流式输出

大模型的流式输出是指大模型在生成文本或其他类型的数据时，等到整个生成过程完成后再一次性返回所有内容就太慢了

生成一部分就立即发送一部分给用户或下游系统，以逐步、逐块的方式返回结果更加丝滑。

这样，用户就不需要等待整个文本生成完成再看到结果，通过这种方式可以改善用户体验

```
#集成阿里通义千问-流式输出
langchain4j.community.dashscope.streaming-chat-model.api-key=${DASH_SCOPE_API_KEY}
langchain4j.community.dashscope.streaming-chat-model.model-name=qwen-plus
```

然后修改三步：

- 修改 Assistant 中 chatModel 改为 streamingChatModel = "qwenStreamingChatModel"
- chat 方法的返回值为 Flux
- 修改 controller 中 chat 方法的返回值为 Flux，并添加 produces 属性指定编码utf-8即可