

# Spring AI 智能航空助手项目

---

完整架构图

环境准备

api-key 配置

实战开始

前端：

后端

创建Controller

配置ChatClient

流式对话

解决sse长连接重复请求问题：

加入打印日志拦截器

预设角色

对话记忆

实现 退订 、 更改预定

退订业务：

难点

解决： 加入提示词

通过function-call 进行方法回调

通过function-call 进行退订

退订前确认信息是否存在

通过function-call 获取预定信息

通过function-call 修改预定信息

通过RAG(检索增强生成)， 外挂一个知识库

redis向量数据库

安装RedisStack

引入依赖

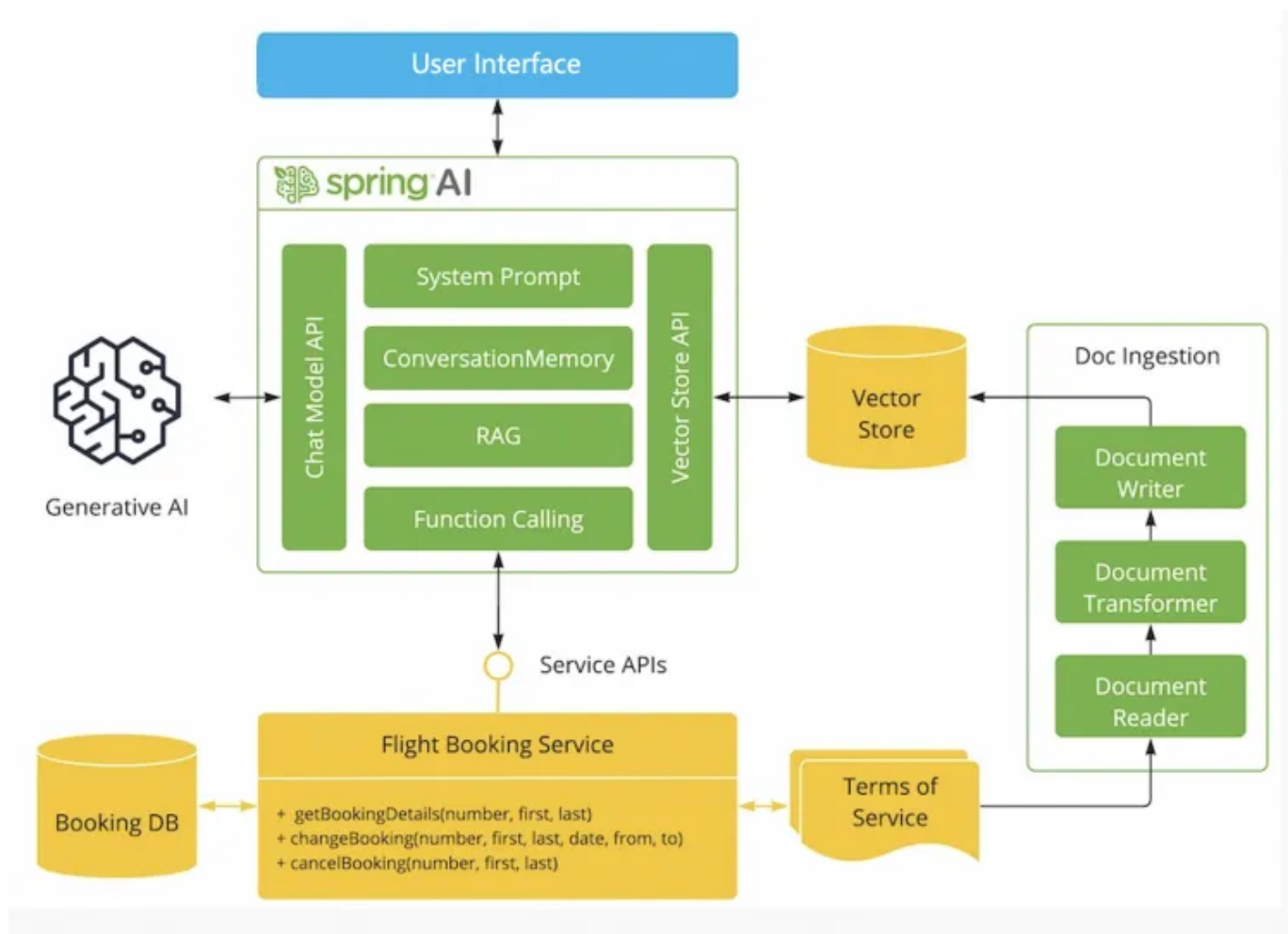
配置连接

配置向量数据库

文档嵌入

文档查询

## 完整架构图



## 环境准备

- Java 17
- Node.js 18+
- 获取API Key

api-key 配置

在正式开始体验之前，需要申请到模型的 api-key。申请地址：

<https://help.aliyun.com/zh/dashscope/developer-reference/activate-dashscope-and-create-an-api-key>

您可以通过 DashScope 提供的方式配置 api-key，SCA AI 完全兼容 DashScope 环境变量配置 key 的方式：<https://help.aliyun.com/zh/dashscope/developer-reference/api-key-settings>

## 实战开始

### 前端：

前端环境安装：<https://note.youdao.com/s/AXcWXq3v>

1. 直接下载我的代码
2. `cd app/chatgpt-demo`
3. 执行安装依赖命令：`npm i`
4. 运行：`npm run dev`

### 后端

依赖：

```
1 <dependencies>
2     <dependency>
3         <groupId>org.springframework.boot</groupId>
4         <artifactId>spring-boot-starter-web</artifactId>
5     </dependency>
6     <dependency>
7         <groupId>com.alibaba.cloud.ai</groupId>
8         <artifactId>spring-ai-alibaba-starter</artifactId>
9         <version>1.0.0-M2.1</version>
10    </dependency>
11
12    <dependency>
13        <groupId>org.springframework.boot</groupId>
14        <artifactId>spring-boot-starter-test</artifactId>
15        <scope>test</scope>
16    </dependency>
17
18
19 </dependencies>
```

配置:

```
1 spring.ai.dashscope.api-key=${TONGYI_AI_KEY}
2 spring.ai.dashscope.chat.options.model=qwen-max
```

创建Controller

配置ChatClient

```
1
2 @RestController
3 public class OpenAiController {
4     private final ChatClient chatClient;
5
6     public OpenAiController(ChatClient.Builder chatClientBuilder, VectorStore vectorStore, ChatMemory chatMemory) {
7         this.chatClient = chatClientBuilder.build();
8     }
9 }
```

## 流式对话

```
1 @CrossOrigin
2 @GetMapping(value = "/ai/generateStreamAsString", produces = MediaType.TEXT_EVENT_STREAM_VALUE)
3 public Flux<String> generateStreamAsString(@RequestParam(value = "message", defaultValue = "讲个笑话") String message) {
4     //Prompt prompt = new Prompt(new UserMessage(message));
5     //return chatClient.stream(prompt);
6     Flux<String> content = chatClient.prompt()
7         .user(message)
8         .stream()
9         .content();
10
11     return content;
12
13 }
```

解决sse长连接重复请求问题:

```
1  @CrossOrigin
2  @GetMapping(value = "/ai/generateStreamAsString", produces = MediaType
   .TEXT_EVENT_STREAM_VALUE)
3  public Flux<String> generateStreamAsString(@RequestParam(value = "mess
   age", defaultValue = "讲个笑话") String message) {
4      //Prompt prompt = new Prompt(new UserMessage(message));
5      //return chatClient.stream(prompt);
6      Flux<String> content = chatClient.prompt()
7          .user(message)
8          .stream()
9          .content();
10
11     return content
12         .doOnNext(System.out::println)
13         .concatWith(Flux.just("[complete]"));
14
15 }
```

## 加入打印日志拦截器

### 1. 加入Advisor

```
1  public class LoggingAdvisor implements RequestResponseAdvisor {
2
3      @Override
4      public AdvisedRequest adviseRequest(AdvisedRequest request, Map<String,
   Object> context) {
5          System.out.println("Request: " + request);
6          return request;
7      }
8  }
```

### 2. 使用chatClientBuilder的 `defaultAdvisors()` 方法注册Advisor

```

1 public OpenAiController(ChatClient.Builder chatClientBuilder, VectorStore v
  vectorStore, ChatMemory chatMemory) {
2     this.chatClient = chatClientBuilder.defaultAdvisors(new LoggingAdvi
  sor()).build();
3 }

```

## 预设角色

我们做的是一个智能机票助手，他主要的职责是给用户退票、改票、买票。必须职责分明

```

1 public OpenAiController(ChatClient.Builder chatClientBuilder, VectorStore
  vectorStore, ChatMemory chatMemory) {
2     this.chatClient = chatClientBuilder
3         .defaultSystem("""
4             您是“XuShu”航空公司的客户聊天支持代理。请以友好、乐于助人
  且愉快的方式来回复。
5             您正在通过在线聊天系统与客户互动。
6             请讲中文。
7             今天的日期是 {current_date}.
8             """)
9         .defaultAdvisors(new LoggingAdvisor()).build();
10 }
11

```

## 对话记忆

```

1 @Bean
2 public ChatMemory chatMemory() {
3     return new InMemoryChatMemory();
4 }

```

```

1 .defaultAdvisors(new PromptChatMemoryAdvisor(chatMemory))

```

# 实现 退订 、 更改预定

所有退订 、 更改 逻辑 都在该类中已准备好 `FlightBookingService` ， 没有操作数据库， 这个不是重点。

## 退订业务：

超过2天无法退订

退订后状态改为“取消”

#	Name	Date	From	To	Status
101	云小宝	2024-11-11	南京	南京	✓
102	李千问	2024-11-13	青岛	成都	✓
103	张百炼	2024-11-15	青岛	深圳	✓
104	王通义	2024-11-17	广州	南京	✓
105	刘魔搭	2024-11-19	成都	南京	✓

## 难点

- 1. 需要获取当前航班信息 （ 需要用户提供 姓名， 航班号）

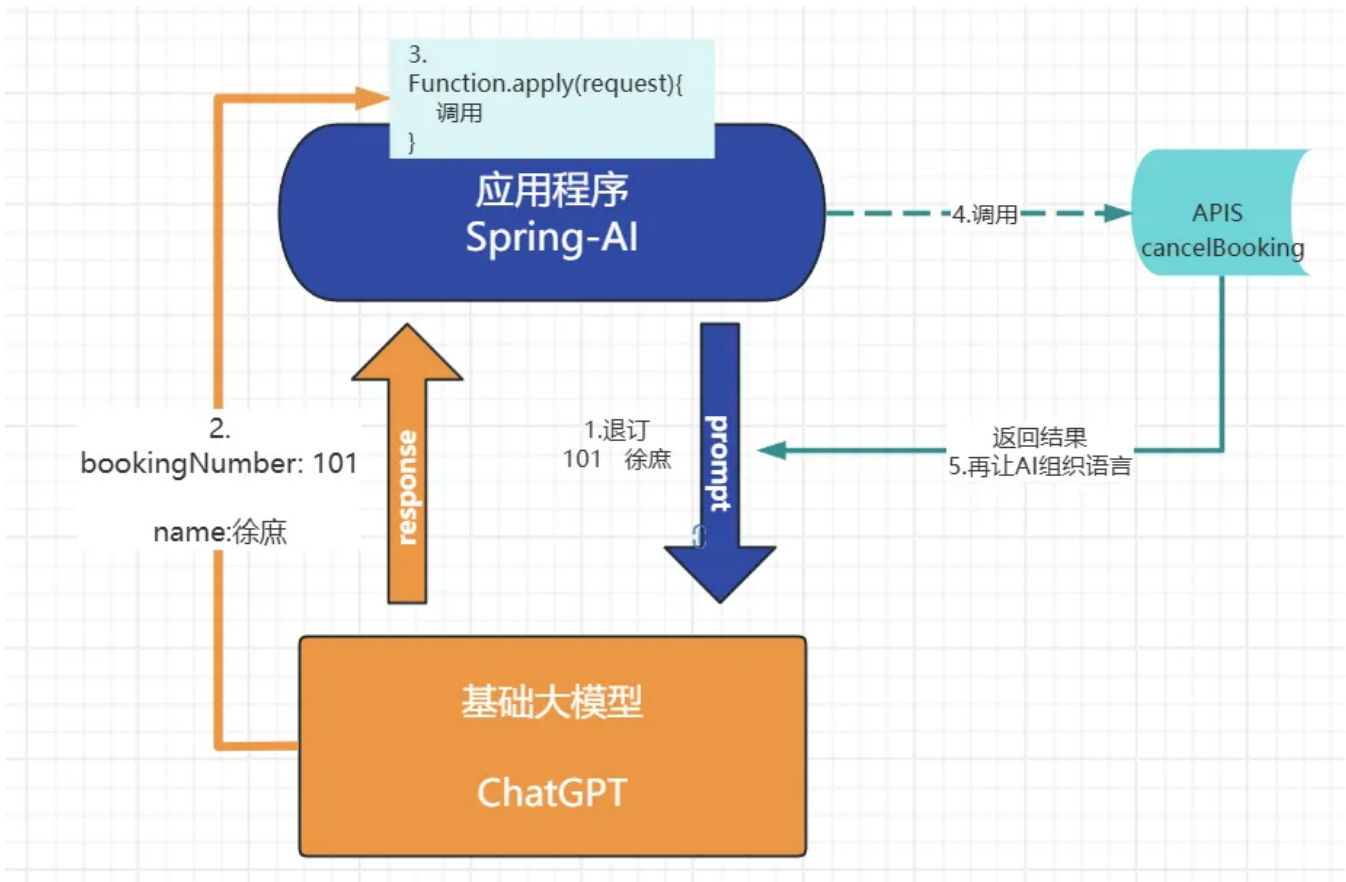
解决： 加入提示词

在进行有关预订或取消预订的信息之前， 您必须始终从用户处获取以下信息： 预订号、 客户姓名。  
在询问用户之前， 请检查消息历史记录以获取此信息。



```
1 public OpenAiController(ChatClient.Builder chatClientBuilder, VectorStore
  vectorStore, ChatMemory chatMemory) {
2     this.chatClient = chatClientBuilder
3         .defaultSystem("""
4             您是“XuShu”航空公司的客户聊天支持代理。请以友好、乐于助人
              且愉快的方式来回复。
5             您正在通过在线聊天系统与客户互动。
6             在提供有关预订或取消预订的信息之前，您必须始终
7             从用户处获取以下信息：预订号、客户姓名。
8             在询问用户之前，请检查消息历史记录以获取此信息。
9             请讲中文。
10            今天的日期是 {current_date}.
11            """)
12        .defaultAdvisors(
13            new PromptChatMemoryAdvisor(chatMemory),
14            new QuestionAnswerAdvisor(vectorStore, SearchReque
15            st.defaults()), // RAG
16            new LoggingAdvisor()).build();
17 }
```

通过function-call 进行方法回调



### 1. 需要告诉大模型：回调哪个方法

- a. 提供实现了Function接口Bean(调用apply)

### 2. 需要告诉大模型：什么对话才回调

- a. 配置Function作用（处理退订）

### 3. 需要告诉大模型：提取对话的什么关键字

- a. Function的第一个泛型去指定提取关键字的变量名

——Function 来解决

## 通过function-call 进行退订

```
1  @Bean
2  @Description("取消机票预定")
3  public Function<CancelBookingRequest, String> cancelBooking() {
4      return request -> {
5          flightBookingService.cancelBooking(request.bookingNumber(), request.name());
6          return "";
7      };
8  }
9
10 public record CancelBookingRequest(String bookingNumber, String name) {
11 }
```

## 退订前确认信息是否存在

解决： 加入提示词

如果更改需要收费，您必须在继续之前征得用户同意。

或者

在更改或退订之前，请先获取预订信息并且用户确定信息。

## 通过function-call 获取预定信息

```
1  @Bean
2      @Description("获取机票预定详细信息")
3  public Function<BookingDetailsRequest, BookingDetails> getBookingDetails() {
4      return request -> {
5          try {
6              return flightBookingService.getBookingDetails(request.bookingNumber(), request.name());
7          }
8          catch (Exception e) {
9              logger.warn("Booking details: {}", NestedExceptionUtils.getMostSpecificCause(e).getMessage());
10             return new BookingDetails(request.bookingNumber(), request.name(), null, null, null, null);
11          }
12      };
13  }
14
15
16  public record BookingDetailsRequest(String bookingNumber, String name)
17  {
18  }
19
20  public record BookingDetails(String bookingNumber, String name, LocalDate date, BookingStatus bookingStatus,
21  String from, String to, String bookingClass) {
22  }
```

## 通过function-call 修改预定信息

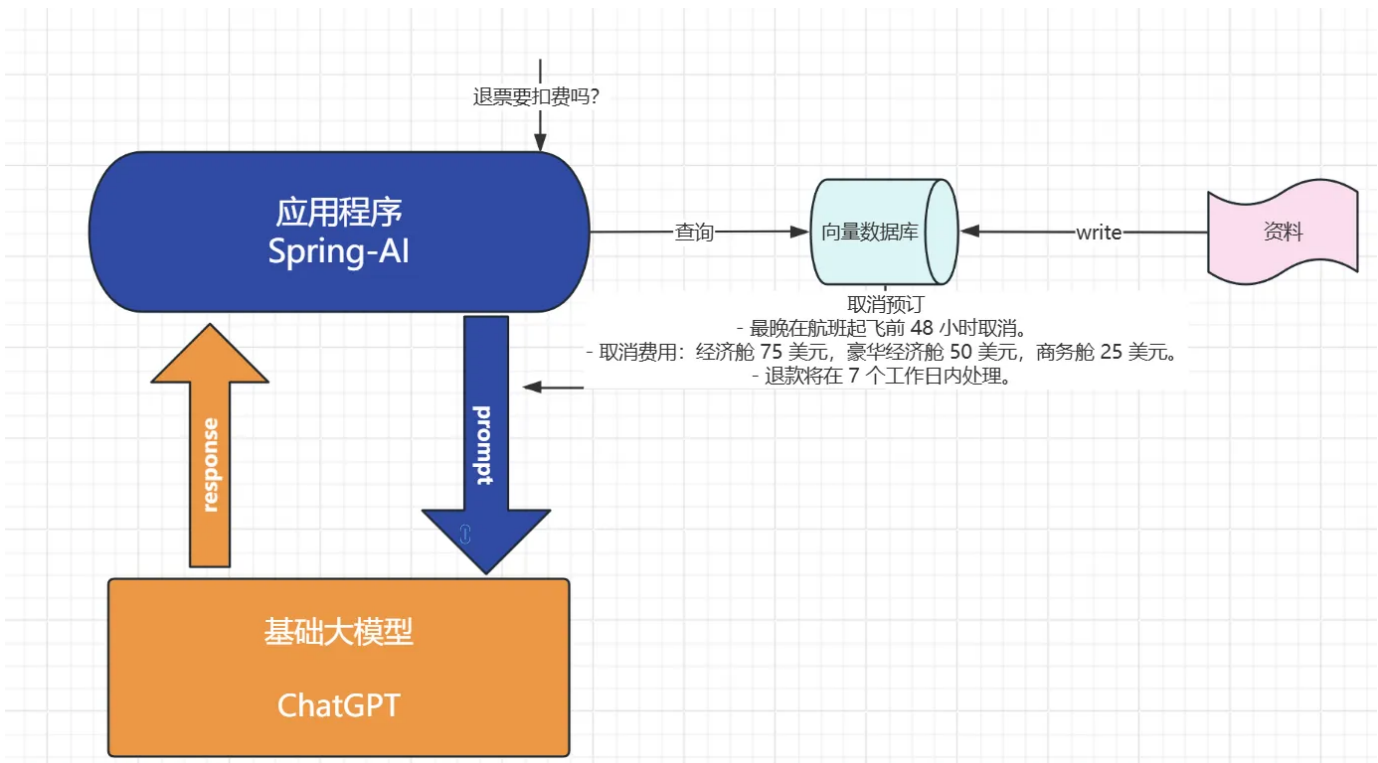
```
1      @Bean
2      @Description("修改机票预定日期")
3      public Function<ChangeBookingDatesRequest, String> changeBooking() {
4          return request -> {
5              flightBookingService.changeBooking(request.bookingNumber(), re
quest.name(), request.date(), request.from(),
6                  request.to());
7              return "";
8          };
9      }
10
```

ChangeBookingDatesRequest:

```
1
```

## 通过RAG(检索增强生成), 外挂一个知识库

向量数据库: **相似性检索**, 大数据领域 (推荐、你喜欢的)



1. 配置向量数据库
2. 写入数据 (Embedding)
3. 查询

[terms-of-service.txt](#)

Plain Text

- 1 本服务条款适用于您对 Funnair 的体验。预订航班，即表示您同意这些条款。
- 2 1. 预订航班
- 3 - 通过我们的网站或移动应用程序预订。
- 4 - 预订时需要全额付款。
- 5 - 确保个人信息（姓名、ID 等）的准确性，因为更正可能会产生 25 的费用。
- 6 2. 更改预订
- 7 - 允许在航班起飞前 24 小时更改。
- 8 - 通过在线更改或联系我们的支持人员。
- 9 - 改签费：经济舱 50，豪华经济舱 30，商务舱免费。
- 10 3. 取消预订
- 11 - 最晚在航班起飞前 48 小时取消。
- 12 - 取消费用：经济舱 75 美元，豪华经济舱 50 美元，商务舱 25 美元。
- 13 - 退款将在 7 个工作日内处理。

向量数据库

redis es ...

```

1  @Bean
2  public VectorStore vectorStore(EmbeddingModel embeddingModel) {
3      return new SimpleVectorStore(embeddingModel);
4  }
5

```

## 写入向量数据库



```

1  @Bean
2  CommandLineRunner ingestTermOfServiceToVectorStore(EmbeddingModel embeddin
3  gModel, VectorStore vectorStore,
4  @Value("classpath:rag/terms-of-service.txt") Resource termsOfServiceDocs) {
5  return args -> {
6      vectorStore.write(                                     // 3.写入
7          new TokenTextSplitter().transform(                // 2.转换
8              new TextReader(termsOfServiceDocs).read())    // 1.读取
9      );
10
11  };
12  }

```

## 配置Advisor:

```
new QuestionAnswerAdvisor(vectorStore, SearchRequest.defaults()), // RAG
```

`QuestionAnswerAdvisor` 可以在用户发起的提问时，先向数据库查询相关的文档，再把相关的文档拼接到用户的提问中，再让模型生成答案。那就是 `RAG` 的实现了。

```

1  this.chatClient = chatClientBuilder
2      .defaultSystem("""
3      您是“Funnair”航空公司的客户聊天支持代理。请以友好、乐于助人
      且愉快的方式来回复。
4      您正在通过在线聊天系统与客户互动。
5      在提供有关预订或取消预订的信息之前，您必须始终
6      从用户处获取以下信息：预订号、客户姓名。
7      在询问用户之前，请检查消息历史记录以获取此信息。
8      在更改预订之前，您必须确保条款允许这样做。
9      如果更改需要收费，您必须在继续之前征得用户同意。
10     使用提供的功能获取预订详细信息、更改预订和取消预订。
11     如果需要，可以调用相应函数调用完成辅助动作。
12     请讲中文。
13     今天的日期是 {current_date}.
14     """)
15     .defaultAdvisors(
16         new PromptChatMemoryAdvisor(chatMemory),
17         new QuestionAnswerAdvisor(vectorStore, SearchR
equest.defaults()), // RAG
18         new LoggingAdvisor())
19     .defaultFunctions("getBookingDetails", "changeBooking", "c
ancelBooking") // FUNCTION CALLING
20     .build();

```

## redis向量数据库

### 安装RedisStack

需要先禁用掉自己原本的redis，防止端口冲突。访问 `localhost:8001` 查看数据库的信息。用户名：`default`，密码：`123456`。

```

1  docker run -d --name redis-stack --restart=always -v redis-data:/data -p 6
    379:6379 -p 8001:8001 -e REDIS_ARGS="--requirepass 123456" redis/redis-stac
    k:latest

```

### 引入依赖



```

1 <dependency>
2   <groupId>org.springframework.ai</groupId>
3   <artifactId>spring-ai-redis-store</artifactId>
4 </dependency>
5 <dependency>
6   <groupId>redis.clients</groupId>
7   <artifactId>jedis</artifactId>
8 </dependency>
9 <dependency>
10  <groupId>org.springframework.ai</groupId>
11  <artifactId>spring-ai-tika-document-reader</artifactId>
12 </dependency>

```

## 配置连接

```

1  spring:
2    data:
3      redis:
4        database: 0
5        timeout: 10s
6        lettuce:
7          pool:
8            # 连接池最大连接数
9            max-active: 200
10           # 连接池最大阻塞等待时间（使用负值表示没有限制）
11           max-wait: -1ms
12           # 连接池中的最大空闲连接
13           max-idle: 10
14           # 连接池中的最小空闲连接
15           min-idle: 0
16       repositories:
17         enabled: false
18       password: 123456

```

## 配置向量数据库

如果你的项目里面有用到redis，需要先禁用 `RedisVectorStoreAutoConfiguration`。这是SpringAI自动配置RedisStack的向量数据库连接，会导致Redis的连接配置冲突。

`VectorStore` 对象需要提供 `EmbeddingModel`，这个案例提供的是阿里灵积的 `EmbeddingModel`。可以切换换成其他厂家的EmbeddingModel。

```

1  @Configuration
2  // 禁用SpringAI提供的RedisStack向量数据库的自动配置，会和Redis的配置冲突。
3  @EnableAutoConfiguration(exclude = {RedisVectorStoreAutoConfiguration.class})
4  // 读取RedisStack的配置信息
5  @EnableConfigurationProperties({RedisVectorStoreProperties.class})
6  @AllArgsConstructor
7  public class RedisVectorConfig {
8
9      /**
10       * 创建RedisStack向量数据库
11       *
12       * @param embeddingModel 嵌入模型
13       * @param properties      redis-stack的配置信息
14       * @return vectorStore 向量数据库
15       */
16     @Bean
17     public VectorStore vectorStore(EmbeddingModel embeddingModel,
18                                     RedisVectorStoreProperties properties,
19                                     RedisConnectionDetails redisConnectionD
20     tails) {
21         RedisVectorStore.RedisVectorStoreConfig config = RedisVectorStore.
22         RedisVectorStoreConfig.builder().withIndexName(properties.getIndex()).with
23         Prefix(properties.getPrefix()).build();
24         return new RedisVectorStore(config, embeddingModel,
25                                     new JedisPooled(redisConnectionDetails.getStandalone().get
26         Host(),
27                                                         redisConnectionDetails.getStandalone().getPort()
28                                                         , redisConnectionDetails.getUsername(),
29                                                         redisConnectionDetails.getPassword()),
30                                                         properties.isInitializeSchema());
31     }
32 }

```

## 文档嵌入

在上面的 `VectorStore` 配置中我们提供了 `EmbeddingModel`，调用 `vectorStore.add(splitDocuments)` 底层会把文档给 `EmbeddingModel` 把文本变成向量然后再存入向量数据库。

```

1 private final VectorStore vectorStore;
2 /**
3  * 嵌入文件
4  *
5  * @param file 待嵌入的文件
6  * @return 是否成功
7  */
8 @SneakyThrows
9 @PostMapping("embedding")
10 public Boolean embedding(@RequestParam MultipartFile file) {
11     // 从IO流中读取文件
12     TikaDocumentReader tikaDocumentReader = new TikaDocumentReader(new
InputStreamResource(file.getInputStream()));
13     // 将文本内容划分成更小的块
14     List<Document> splitDocuments = new TokenTextSplitter()
15         .apply(tikaDocumentReader.read());
16     // 存入向量数据库，这个过程会自动调用embeddingModel，将文本变成向量再存入。
17     vectorStore.add(splitDocuments);
18     return true;
19 }

```

## 文档查询

调用 `vectorStore.similaritySearch(query)` 时同样会先把用户的提问给 `EmbeddingModel`，将提问变成向量，然后与向量数据库中的文档向量进行相似度计算（cosine值）。

要注意：此时向量数据库不会回答用户的提问。要回答用户的提问需要指定advisor

```

1 /**
2  * 查询向量数据库
3  *
4  * @param query 用户的提问
5  * @return 匹配到的文档
6  */
7
8 @GetMapping("query")
9 public List<Document> query(@RequestParam String query) {
10     return vectorStore.similaritySearch(query);
11 }

```

指定advisor

```
1  return chatClient.prompt()  
2      .user(prompt)  
3      // 2. QuestionAnswerAdvisor会在运行时替换模板中的占位符`questi  
on_answer_context`, 替换成向量数据库中查询到的文档。此时的query=用户的提问+替换完的  
提示词模板;  
4      .advisors(new QuestionAnswerAdvisor(vectorStore, prompt))  
5      .stream()  
6      // 3. query发送给大模型得到答案  
7      .content()  
8      .map(chatResponse -> ServerSentEvent.builder(chatResponse)  
9          .event("message")  
10         .build());
```