



Spring6.2最新特性

1. 支持loc 容器 bean 异步并行的创建

这可以显着减少启动时间。

<https://github.com/spring-projects/spring-framework/issues/13410>

1.1. 使用方式

1. 配置需要异步创建的bean

```
1 @Bean(bootstrap = Bean.Bootstrap.BACKGROUND)
```

1. 配置线程池

```
1 @Bean
2 public Executor bootstrapExecutor(){
3     return Executors.newCachedThreadPool();
4 }
```

1.2. 测试:

1.2.1. 批量注册异步bean

```
1 @Component
2 public class BatchBeanDefinitionRegister implements
3     BeanDefinitionRegistryPostProcessor, BeanFactoryAware {
4     @Override
5     public void postProcessBeanDefinitionRegistry(BeansDefinitionRegistry registry)
6         throws BeansException {
7         for (int i=0;i<10;i++){
8             RootBeanDefinition rootBeanDefinition = new
9             RootBeanDefinition(AService.class);
10             rootBeanDefinition.setBackgroundInit(true);
11             registry.registerBeanDefinition("AService"+i,rootBeanDefinition);
12         }
13     }
14 }
```

1.2.2. 为了更好演示让bean睡几秒

```
1 public class AService {
2     @Autowired
3     BService bService;
4
5     public AService() throws InterruptedException {
6         Thread.sleep(5000);
7     }
8 }
```

1.2.3. 可以分别测试用bootstrapExecutor和不用的用时时间

```
1 public static void main(String[] args) {
2     long beginTime = System.currentTimeMillis();
3     AnnotationConfigApplicationContext ioc = new
    AnnotationConfigApplicationContext(Main2.class);
4     long endTime = System.currentTimeMillis();
5     System.out.println(endTime-beginTime);
6     AService aService = (AService) ioc.getBean("AService9");
7     System.out.println(aService);
8 }
```

可以发现如果bean很多，确实速度提升了~!!!❀❀、(°▽°)ノ❀

源码：

org.springframework.beans.factory.support.DefaultListableBeanFactory#preInstantiateSingleton

涉及到CompletableFuture不会可以参考：<https://www.processon.com/view/link/6638e5d996857d67d2e3e998?cid=6618dbe5c04c1e02ff4f0754>

1.3. 问题

1.3.1. 问题1： 如果异步beanA正在通过ioc容器加载创建， 此时另外一个线程2获取beanA.会怎么样？

线程2会阻塞吗？

在异步beanA创建时， 会往三级缓存中加入future.join();,

```

1 addSingletonFactory(beanName, () -> {
2
3         try {
4
5             future.join();
6
7         }
8         catch (CompletionException ex) {
9
10            ReflectionUtils.rethrowRuntimeException(ex.getCause());
11
12        }
13
14        return future; // not to be exposed, just to
15        lead to ClassCastException in case of mismatch
16    });

```

线程2在获取bean时，会等待异步beanA线程执行完毕，再去一级缓存中捞一遍

```

1 ObjectFactory<?> singletonFactory = this.singletonFactories.get(beanName);
2 if (singletonFactory != null) {
3     // 如果backgroundInit beanA正在创建，此时有其他线程获取beanA或回调future.join 等待
4     singletonObject = singletonFactory.getObject();
5     // Singleton could have been added or removed in the meantime.
6     if (this.singletonFactories.remove(beanName) != null) {
7         this.earlySingletonObjects.put(beanName, singletonObject);
8     }
9     else {
10         singletonObject = this.singletonObjects.get(beanName);
11     }
12 }

```



spring5的时候，源码在创建的时候会锁，按道理第2个线程获取不到锁啊。其实 ...

在spring6为了支持并行createBean，不再进行上锁了，所以getSingleton中2个线程都可以trylock

tryLock()返回值表示的是用来尝试获取锁：成功获取则返回true；获取失败则返回false，**这个方法无论如何都会立即返回**。不会像synchronized一样，一个线程获取锁之后，其他锁只能等待那个线程释放之后才能有获取锁。

1.3.2. 问题2：如果同一个Bean懒加载一起创建会怎么样？

这个时候不会往三级缓存存入 join 。而且也不会阻塞

为了支持并行createBean，在创建bean时不再使用互斥锁

```
1 // 创建前标记当前bean正在创建
2 // 由于现版本并发创建允许并行，所以可能存在同一个bean在正在创建，其他线程报错，以前版本是不存在的因为以前是互斥锁
3 // 但是这种情况极少 因为很少实际情况级别不会一起创建一个bean,除非程序员人为
4 // 相当于悲观锁（之前互斥等待）改乐观锁（现在不互斥比较替换，异常跳出）了。
5 beforeSingletonCreation(beanName);
```

在标记正在创建的bean时会报错：

```
1 if (!this.inCreationCheckExclusions.contains(beanName) &&
    !this.singletonsCurrentlyInCreation.add(beanName)) {
2
3     throw new BeanCurrentlyInCreationException(beanName);
4 }
```

相当于悲观锁（之前互斥等待）改乐观锁（现在不互斥比较替换，异常跳出）了。

1.3.3. 问题3： 如果BeanA\BeanB 都是异步lazy bean，并且循环依赖， 2个线程，线程1 getBean(A) 线程2 getBean(B) ，会互相join（死锁）？？？？？

不会，因为在创建bean的时候，实例化后会把三级缓存给覆盖掉！

2. defaultCandidate

<https://github.com/spring-projects/spring-framework/issues/26528>

当声明了bean, 不想被依赖注入，可以设置该属性

解决问题：

当声明了bean，有些时候不想被DI，以前做不到，现在可以：

```
1 @Bean(defaultCandidate = true)
```

Spring-AI

完整示例代码

<https://gitee.com/xscodeit/ai-openai-examples.git>

官网：

<https://docs.spring.io/spring-ai/reference/index.html>

整合chatgpt

前置准备

1. open-ai-key:

<https://api.xty.app/register?aff=PuZD>

<https://eylink.cn/>

淘宝搜：open ai key

2. 魔法软件，由于国家相关法律规定，建议大家自行准备。

实现

创建完后会发现加入了依赖：

```
1 <dependencyManagement>
2   <dependencies>
3     <dependency>
4       <groupId>org.springframework.ai</groupId>
5       <artifactId>spring-ai-bom</artifactId>
6       <version>1.0.0-SNAPSHOT</version>
7       <type>pom</type>
8       <scope>import</scope>
9     </dependency>
10  </dependencies>
11 </dependencyManagement>
12
13 <dependencies>
14   <dependency>
15     <groupId>org.springframework.boot</groupId>
16     <artifactId>spring-boot-starter-web</artifactId>
17   </dependency>
18
19
20   <dependency>
21     <groupId>org.springframework.ai</groupId>
22     <artifactId>spring-ai-openai-spring-boot-starter</artifactId>
23   </dependency>
24
25   <dependency>
26     <groupId>org.springframework.boot</groupId>
27     <artifactId>spring-boot-starter-test</artifactId>
28     <scope>test</scope>
29   </dependency>
30
31   <dependency>
32     <groupId>org.projectlombok</groupId>
33     <artifactId>lombok</artifactId>
34   </dependency>
35
36 </dependencies>
```


设置代理, 如果你请求的大模型的api接口不是国内的。需要将程序设置代理:

```
1 public static void main(String[] args) {
2     // 设置代理
3     String proxy = "127.0.0.1"; // 如果代理在你本机就127.0.0.1 , 如果代理是其
4     他服务器相应设置
5     int port = 7890; //设置科学上网代理的端口,
6     System.setProperty("proxyType", "4");
7     System.setProperty("proxyPort", Integer.toString(port));
8     System.setProperty("proxyHost", proxy);
9     System.setProperty("proxySet", "true");
10
11     SpringApplication.run(Application.class, args);
12
13 }
```

设置key:

```
1 spring:
2     ai:
3         openai:
4             api-key: ${OPEN_AI_KEY}
5             base-url: ${OPEN_AI_URL}
```

示例代码:

```
1
2 private final ChatClient chatClient;
3 private final OpenAiChatModel chatClient2;
4 private final OpenAiImageModel imageClient;
5 private final OpenAiAudioTranscriptionModel audioClient;
6 private final OpenAiAudioApi openAiAudioApi;
7
8
9 @Value("${OPEN_AI_KEY}")
10 private String openAiKey;
11
12
13 @GetMapping("/ai/simple")
14 public Map<String, String> completion(@RequestParam(value = "message", defaultValue =
    "给我讲个笑话") String message) {
15     System.out.println(openAiKey);
16     var value=chatClient.prompt()
17         .user(message).call().content();
18     return Map.of("generation",value );
19 }
20
21 @GetMapping(value="/ai/stream",produces="text/sse;charset=UTF-8")
22 public Flux<String> stream(@RequestParam(value = "message", defaultValue = "给我讲个笑
    话") String message ) {
23     System.out.println(openAiKey);
24     return chatClient.prompt()
25         .user(message)
26         .stream()
27         .content();
28 }
29
30
31 @GetMapping(value="/ai/img",produces="text/html")
32 public String image(@RequestParam(value = "message", defaultValue = "猫") String
    message ) {
33
34     ImageResponse response = imageClient.call(
35         new ImagePrompt(message,
```

```
36         OpenAiImageOptions.builder()
37             .withQuality("hd")
38             .withN(1)
39             .withModel(OpenAiImageApi.ImageModel.DALL_E_2.getValue())
40             // dall-e-2 256
41             .withHeight(256)
42             .withWidth(256).build()));
43
44     String url = response.getResult().getOutput().getUrl();
45     System.out.println(url);
46     return "<img src='"+url+"'>";
47
48 }
49
50
51 @GetMapping(value="/ai/audit2text")
52 public String audit2text() {
53
54     var transcriptionOptions = OpenAiAudioTranscriptionOptions.builder()
55         .withResponseFormat(OpenAiAudioApi.TranscriptResponseFormat.TEXT)
56         .withTemperature(0f)
57         .build();
58
59     // flac、mp3、mp4、mpeg、mpga、m4a、ogg、wav 或 webm。
60     var audioFile = new ClassPathResource("/hello.mp3");
61
62     AudioTranscriptionPrompt transcriptionRequest = new
63     AudioTranscriptionPrompt(audioFile, transcriptionOptions);
64
65     AudioTranscriptionResponse response = audioClient.call(transcriptionRequest);
66
67     //openAiAudioApi.createTranscription()
68     return response.getResult().getOutput();
69
70 }
71
72 @GetMapping(value="/ai/text2audit")
73 public String text2audit() {
```

```
74     ResponseEntity<byte[]> speech = openAiAudioApi.createSpeech(  
75         OpenAiAudioApi.SpeechRequest.builder().  
76             withVoice(OpenAiAudioApi.SpeechRequest.Voice.ONYX).  
77             withInput("你好，我是徐庶").build());  
78     byte[] body = speech.getBody();  
79  
80  
81     // 将byte[]存为 mp3文件  
82     try {  
83         writeByteArrayToMp3(body, System.getProperty("user.dir"));  
84     } catch (IOException e) {  
85         throw new RuntimeException(e);  
86     }  
87  
88     return "ok";  
89 }  
90  
91 public static void writeByteArrayToMp3(byte[] audioBytes, String outputFilePath) throws  
92     IOException {  
93     // 创建FileOutputStream实例  
94     FileOutputStream fos = new FileOutputStream(outputFilePath+"/xushu.mp3");  
95  
96     // 将字节数组写入文件  
97     fos.write(audioBytes);  
98  
99     // 关闭文件输出流  
100    fos.close();  
101 }  
102  
103 @GetMapping(value="/ai/mutil")  
104 public String mutilModel(String message,String imgUrl) throws IOException {  
105  
106     byte[] imageData = new ClassPathResource("/test.png").getContentAsByteArray();  
107  
108     var userMessage = new UserMessage(  
109         "这个图片你看出什么?", // content  
110         List.of(new Media(MimeTypeUtils.IMAGE_PNG, imageData))); // media  
111 }
```

```

112     ChatResponse response = chatClient.call(new Prompt(userMessage,
113         OpenAiChatOptions.builder()
114             .withModel(OpenAiApi.ChatModel.GPT_4_TURBO_PREVIEW.getValue())
115             .build()));
116
117     return response.getResult().getOutput().getContent();
118 }
119

```

function-call

AI本身是不具备实时消息能力的，比如问“现在北京的天气是什么”，AI是不知道的，这个时候我们需要通过接口来帮助AI完成,大致流程：

function-call实现的代码：

1. 在发送文本时，同时设置Function；关键代码：`.withFunction("getWaitTime")`

```

1 ChatResponse response = chatClient.call(new Prompt(
2     List.of(userMessage),
3     OpenAiChatOptions.builder()
4         .withFunction("getWaitTime")
5         .build()));

```

2. 当代码执行完call时（AI响应之后），会再调用 `getWaitTime` 对应的bean的apply方法。

```

1 @Bean
2 public Function<WaitTimeService.Request, WaitTimeService.Response> getWaitTime() {
3     return new WaitTimeService();
4 }

```

3. 并且会把 `getWaitTime` 该bean的Request作为function-call的返回参数，即可在 `apply` 方法中获取Request对应的参数

```
1 @Override
2 public Response apply(Request request) {
3     String name = request.name();
4     String location = request.location();
5     // todo...
6     return new Response(location+"有10个! ");
7 }
8
9 public record Request(String name,String location) {}
10 public record Response(String weather) {}
```

随后Response会再次丢给AI组织语言，进行响应，最终

源码：

看源码之前，了解下该接口需要哪些参数：<https://platform.openai.com/docs/api-reference/chat/create>

请求：userMessage=“长沙有多少叫徐庶的”

调用call方法，执行openai的远程api请求

```

1  @Override
2  public ChatResponse call(Prompt prompt) {
3
4      ChatCompletionRequest request = createRequest(prompt, false);
5
6      return this.retryTemplate.execute(ctx -> {
7
8          ResponseEntity<ChatCompletion> completionEntity =
9              this.callWithFunctionSupport(request);
10
11              var chatCompletion = completionEntity.getBody();
12              if (chatCompletion == null) {
13                  logger.warn("No chat completion returned for prompt: {}", prompt);
14                  return new ChatResponse(List.of());
15              }
16
17              RateLimit rateLimits =
18                  OpenAiResponseHeaderExtractor.extractAiResponseHeaders(completionEntity);
19
20              List<Generation> generations = chatCompletion.choices().stream().map(choice -> {
21                  return new Generation(choice.message().content(), toMap(chatCompletion.id(),
22                      choice))
23                  .withGenerationMetadata(ChatGenerationMetadata.from(choice.finishReason().name(),
24                      null));
25              }).toList();
26
27              return new ChatResponse(generations,
28                  OpenAiChatResponseMetadata.from(completionEntity.getBody().withRateLimit(rateLimits));
29              });

```

2. 通过 `createRequest` 将withFunction参数解析到request.tools

```

1  ChatCompletionRequest request = createRequest(prompt, false);

```

tools属性: 包含了

- i. 函数名, 也是需要调用 的bean的name: getWaitTime

ii. properties是告诉AI，要返回的向量的名字

iii. @Description

iv. 函数描述，函数功能的描述，模型使用它来选择何时以及如何调用函数。

关于open-ai对tools参数的说明: <https://platform.openai.com/docs/api-reference/chat/create#chat-create-tools>

3. 来到 `callWithFunctionSupport` 真正调用远程api接口

```
1 protected Resp callWithFunctionSupport(Req request) {  
2     Resp response = this.doChatCompletion(request);  
3     return this.handleFunctionCallOrReturn(request, response);  
4 }
```

这里2句关键代码：

1. `this.doChatCompletion(request);` 方法会正常请求**chat completion** 接口并且会带上function-call参数并携带**tools**属性，并且返回对话中的function-call所需参数（即 `WaitTimeService.Request` 的参数）
2. `handleFunctionCallOrReturn` 执行Function-callback方法, 此时会调用 `WaitTimeService.apply` 方法
 - a. 拿到之前解析的functionCallback即
 - b. 将arguments从Json转换为对象调用 `WaitTimeService.apply`
 - c. 将返回的数据再次请求大模型

流程图：

<https://www.processon.com/view/link/6654257b40e88034a42b837c?cid=6652ddc65a676876a64350e5>

SpringAI社区非常活跃

在后续的版本都会更新国内常用的大模型

<https://github.com/spring-projects/spring-ai/issues?q=Moonshot>

Add QianFan model client ✓

#751 opened 4 days ago by mxsl-gr ➡ 1.0.0-M2

Add DeepSeek model client ✓ **model client**

#702 opened 2 weeks ago by mxsl-gr ➡ 1.0.0-M2

add MiniMax model client **model client**

#628 by mxsl-gr was closed last week ➡ 1.0.0-M1

Add ZhiPu AI model client **model client**

#623 by mxsl-gr was closed 5 days ago ➡ 1.0.0-M1

Add Moonshot model client **model client**

#596 opened on Apr 17 by mxsl-gr ➡ 1.0.0-M1

有道云链接: <https://note.youdao.com/s/9t8IUhiR>