

Spring AI

(Also Spring with AI without Spring AI)

Contact Info

Ken Kousen

Kousen IT, Inc.

ken.kousen@kousenit.com

<http://www.kousenit.com>

<https://kousenit.org> (blog)

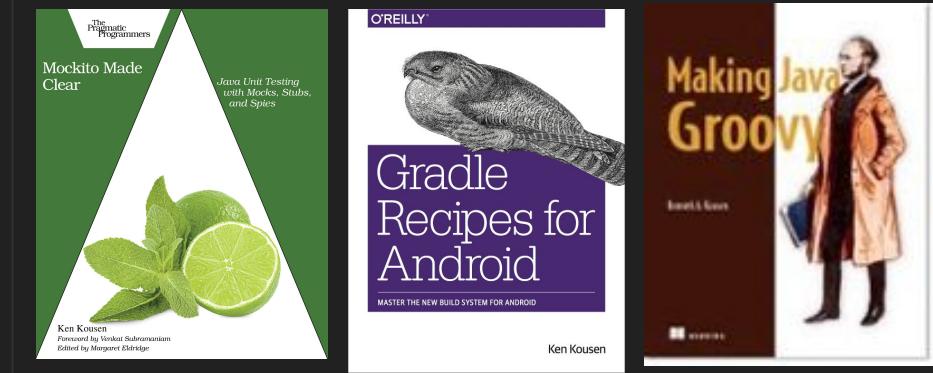
[@kenkousen](https://twitter.com/kenkousen) (twitter)

<https://foojay.social/@kenkousen> (mastodon)

Tales from the jar side (free newsletter)

<https://kenkousen.substack.com>

<https://youtube.com/@talesfromthejarside>



OpenAI API

- OpenAI
 - <https://chat.openai.com/> → regular ChatGPT
 - <https://platform.openai.com> → developer API
 - Introduction and guides: <https://platform.openai.com/docs/introduction>
 - API Reference: <https://platform.openai.com/docs/api-reference>
- Need to register and generate a key
 - Requires credit card with a small deposit (about \$10)
 - Each request costs a small amount (prices discussed later)
- Ollama → can run LLMs on your local machine
 - Also discussed later

OpenAI

- Get a key at <https://platform.openai.com/signup>
- Set as environment variable
 - SPRING_AI_OPENAI_API_KEY
- Spring reads configuration property
 - spring.ai.openai.api-key

Spring AI

- Reference page: <https://docs.spring.io/spring-ai/reference/index.html>
- Very early version: currently 0.8.0
- Based on LangChain and LlamalIndex, but not a direct port of either
- Aims to generalize to multiple LLMs

Dependencies

```
1 <repositories>
2   <repository>
3     <id>spring-milestones</id>
4     <name>Spring Milestones</name>
5     <url>https://repo.spring.io/milestone</url>
6     <snapshots>
7       <enabled>false</enabled>
8     </snapshots>
9   </repository>
10  <repository>
11    <id>spring-snapshots</id>
12    <name>Spring Snapshots</name>
13    <url>https://repo.spring.io/snapshot</url>
14    <releases>
15      <enabled>false</enabled>
16    </releases>
17  </repository>
18 </repositories>
```

```
1 repositories {
2   mavenCentral()
3   maven { url 'https://repo.spring.io/milestone' }
4   maven { url 'https://repo.spring.io/snapshot' }
5 }
```

OpenAI Dependencies

```
1 <dependency>
2   <groupId>org.springframework.ai</groupId>
3   <artifactId>spring-ai-openai-spring-boot-starter</artifactId>
4   <version>0.8.0-SNAPSHOT</version>
5 </dependency>
```

```
1 dependencies {
2   implementation 'org.springframework.ai:spring-ai-openai-spring-boot-starter:0.8.0-SNAPSHOT'
3 }
```

Spring AI

- Major features:
 - General interfaces
 - Prompt Templates
 - Output Parsing
 - Embeddings
 - Retrieval Augmented Generation (RAG)

```
1 @FunctionalInterface
2 public interface ChatClient {
3
4     default String generate(String message) {
5         Prompt prompt = new Prompt(new UserMessage(message));
6         return generate(prompt).getGeneration().getContent();
7     }
8
9     ChatResponse generate(Prompt prompt);
10 }
```

```
1 public class ChatResponse {  
2  
3     private final GenerationMetadata metadata;  
4     private final List<Generation> generations;  
5     private PromptMetadata promptMetadata;  
6  
7     // Getters for each  
8 }
```

```
1 public abstract class AbstractMessage implements Message {  
2  
3     protected String content;  
4     protected Map<String, Object> properties = new HashMap<>();  
5     protected MessageType messageType;  
6  
7     // Getters...  
8 }
```

Message Types

```
1 public enum MessageType {  
2  
3     USER("user"),  
4     ASSISTANT("assistant"),  
5     SYSTEM("system"),  
6     FUNCTION("function");  
7  
8     // ... other details ...  
9 }
```

Message Hierarchy

- ✓ * ⓘ **AbstractMessage** (org.springframework.ai.prompt.messages)
 - ⓘ **AssistantMessage** (org.springframework.ai.prompt.messages)
 - ⓘ **Generation** (org.springframework.ai.chat)
 - ⓘ **UserMessage** (org.springframework.ai.prompt.messages)
 - ⓘ **SystemMessage** (org.springframework.ai.prompt.messages)
 - ⓘ **FunctionMessage** (org.springframework.ai.prompt.messages)
 - ⓘ **ChatMessage** (org.springframework.ai.prompt.messages)

Messages

- ChatMessage: String role, String content
- UserMessage: String message or Resource
- SystemMessage: String content or Resource

OpenAI Reference

- All based on (at least) OpenAI chat model
- Request:
 - model
 - list of messages, each with role and content
- Response:
 - Complicated, but includes:
 - model
 - list of choices, each with a message
 - finish reason
 - usage, with tokens

```
1 curl https://api.openai.com/v1/chat/completions \
2   -H "Content-Type: application/json" \
3   -H "Authorization: Bearer $OPENAI_API_KEY" \
4   -d '{
5     "model": "gpt-3.5-turbo",
6     "messages": [
7       {
8         "role": "system",
9         "content": "You are a helpful assistant."
10      },
11      {
12        "role": "user",
13        "content": "Hello!"
14      }
15    ]
16  }'
```

```
1  {
2      "id": "chatcmpl-123",
3      "object": "chat.completion",
4      "created": 1677652288,
5      "model": "gpt-3.5-turbo-0613",
6      "system_fingerprint": "fp_44709d6fcb",
7      "choices": [
8          {
9              "index": 0,
10             "message": {
11                 "role": "assistant",
12                 "content": "\n\nHello there, how may I assist you today?",
13             },
14             "logprobs": null,
15             "finish_reason": "stop"
16         ],
17         "usage": {
18             "prompt_tokens": 9,
19             "completion_tokens": 12,
20             "total_tokens": 21
21     }
```

Spring AI

- Headers (like Accept and Content-Type) supplied automatically
- Authorization key read from application.properties
 - Usually set as environment variable

```
1 #logging.level.web=debug  
2 spring.ai.openai.api-key=${OPENAI_API_KEY}  
3 spring.ai.openai.chat.model=gpt-4-1106-preview
```

Prompts

- Prompts done using messages
 - System message → context
 - User message → actual request

```
1 public class Prompt {  
2  
3     private final List<Message> messages;  
4  
5     // constructors and utility methods omitted  
6 }
```

Prompt Templates

- Boilerplate with template variables
- Based on StringTemplate open source library, <https://www.stringtemplate.org/>
- Uses { ... } to specify variables

Tell me a {adjective} joke about {topic}

- Developer supplies a Map of keys and values

Prompt Templates

```
1 PromptTemplate promptTemplate = new PromptTemplate(jokeResource);  
2 Prompt prompt = promptTemplate.create(  
3     Map.of("adjective", adjective, "topic", topic));  
4 chatClient.generate(prompt).getGeneration();
```

```
1 String userText = """  
2     Tell me about three famous pirates from the Golden Age of Piracy and why they did.  
3     Write at least a sentence for each pirate.  
4     """;  
5  
6 Message userMessage = new UserMessage(userText);  
7  
8 String systemText = """  
9     You are a helpful AI assistant that helps people find information.  
10    Your name is {name}  
11    You should reply to the user's request with your name and also in the style of a {voice}.  
12    """;  
13  
14 SystemPromptTemplate systemPromptTemplate = new SystemPromptTemplate(systemText);  
15 Message systemMessage = systemPromptTemplate.createMessage(Map.of("name", name, "voice", voice));  
16  
17 Prompt prompt = new Prompt(List.of(userMessage, systemMessage));  
18  
19 List<Generation> response = aiClient.generate(prompt).getGenerations();
```

Spring Resources

- Spring supports the **Resource** abstraction
- Can put prompts in files and inject them with **@Value**

```
1 @Value("classpath:/prompts/system-message.st")
2 private Resource systemResource;
```

Tokens

- Cost is measured in tokens
- Prices are per 1000 tokens
- Rough estimate is 1000 tokens is about 750 words
- Each model has a cost (see <https://openai.com/pricing>)
 - GPT-4 Turbo: \$0.01 / 1K input, \$0.03 / 1K output
 - GPT-4 → multiply by 3 (seriously, **turbo is cheaper**)
 - GPT-3.5 Turbo: \$0.001 / 1K input, \$0.002 / 1K output (**default**)
 - Embedding (more about that later) is about \$0.0001 / 1K tokens
 - DALL-E 3, standard resolution, 1024x1024: \$0.040 / image
 - Whisper (transcribe audio): \$0.006 / minute
 - Text-to-speech: \$0.015 / 1K characters, or double that for HD

Output Parsing

- Output is a String
- Spring AI tries to format the string as requested

Output Parsers

- Similar to RowMapper in Spring JDBC
- **OutputParser** interface, extends Parser<T> and FormatProvider
- Available implementations:
 - BeanOutputParser
 - MapOutputParser
 - ListOutputParser

```
1 public interface OutputParser<T> extends Parser<T>,
2                                         FormatProvider {
3 }
4
5 @FunctionalInterface
6 public interface Parser<T> {
7     T parse(String text);
8 }
9
10 public interface FormatProvider {
11     String getFormat();
12 }
```

Adding your own data

- Read in data → Spring has document readers
 - JsonReader
 - TextReader
 - PagePdfDocumentReader
 - ParagraphPdfDocumentReader
 - TokaDocumentReader
- Divide documents to fit into AI model's context window
 - TextSplitter, several others
- Write out results for storage

Embedding Client

- Convert text to a series of vectors, known as embeddings
- Embedding is a numerical vector of similarity scores

Vector Databases

- Specialized database to handle similarity searches
- Do a similarity search to find similar documents
- Add the similar documents as the context for a query
- This is called Retrieval Augmented Generation (RAG)

Available Vector Stores

- Azure Vector Search
- Chroma
- Milvus
- Neo4j
- PgVector (PostgreSQL)
- PineCone
- Redis
- SimpleVectorStore
- Weaviate

Ollama

- Ollama (<https://ollama.ai/>) lets you run LLMs on your local machine
 - Many models available
 - Default is llama2 (from Meta/Facebook)
- Installations available for macOS and Linux
 - Windows coming soon, but you can use WSL
- Download the model you want
- Ollama exposes it at **localhost:11434**

Ollama

- Download models (can be quite large)
 - Uncensored (!) models available
- Tends to be slow, but at least it's free
 - Also avoids sending data offsite
- Use additional Spring Boot Starter
 - `spring-ai-ollama-spring-boot-starter`
- Base URL and model are set as properties
 - `spring.ai.ollama.base-url` (defaults to `localhost:11434`)
 - `spring.ai.ollama.chat.model` (defaults to `llama2`)
- Same classes as OpenAI

Spring with AI without Spring AI

Spring with AI without Spring AI

- Always good to know one layer below your abstraction
 - Also useful for functionality not yet in Spring AI
- Spring includes **REST clients**
 - RestTemplate legacy
 - WebClient async in webflux module
 - RestClient sync version of WebClient, available in Spring 6.1
- Spring provides a **JSON parser** (Jackson)
- Spring Boot has a **validation** starter
 - Can use jakarta.validation.constraints annotations
- Spring Boot has a **GraalVM native image compiler**

Spring with AI without Spring AI

- What we can do:
 - List the models
 - Generate audio from text
 - Image generation using DALL-E 3
 - Work with LLMs Spring AI does not yet support (Claude, Gemini, ...)

Process

- Map request and response JSON to Java records
- Define HTTP Exchange interface
 - Add abstract methods
 - `@GetExchange`, `@PostExchange`, etc
- Customize rest clients to submit headers
- Use boilerplate to tell Spring to generate implementation classes
- Inject interfaces into services to do pre- and post-processing
- Profit!

List the OpenAI Models

- All the models are available in a GET request to
<https://api.openai.com/v1/models>

Map JSON to Records

JSON Model
and Model List:

```
1 {
2   "id": "davinci",
3   "object": "model",
4   "created": 1686935002,
5   "owned_by": "openai"
6 }
```

```
1 {
2   "object": "list",
3   "data": [
4     {
5       "id": "model-id-0",
6       "object": "model",
7       "created": 1686935002,
8       "owned_by": "organization-owner"
9     },
10    {
11      "id": "model-id-1",
12      "object": "model",
13      "created": 1686935002,
14      "owned_by": "organization-owner",
15    }
16  ]
17 }
```

Java Records

- Can nest records to match JSON structure
- Camel-case to snake-case conversion:
 - Either `@JsonProperty` specifying names, or
 - `@JsonNaming` annotation on class
 - Spring has a property for this, but it doesn't appear to work on records

Model and ModelList

```
1 public record ModelList(List<Model> data) {  
2     @JsonNaming(PropertyNamingStrategies.SnakeCaseStrategy.class)  
3     public record Model(String id, long created, String ownedBy) {  
4 }
```

HTTP Exchange Interface

```
1 @HttpExchange("/v1")
2 public interface OpenAIInterface {
3
4     // Models
5     @GetExchange(value = "/models", accept = "application/json")
6     ModelList listModels();
7 }
```

Configure Beans

- Need a RestClient (or WebClient or RestTemplate)
 - Set Authorization header
 - Set base URL
- Need bean to tell Spring to implement the exchange interface
 - Boilerplate -- copy and paste
 - Just change interface class
- NOTE: Only has to be done once → add extra methods to interface later

```
1 @Configuration
2 public class AppConfig {
3     @Bean
4     public RestClient openAIRestClient(@Value("${openai.baseUrl}") String baseUrl,
5                                         @Value("${OPENAI_API_KEY}") String apiKey) {
6         return RestClient.builder()
7             .defaultHeader("Authorization", "Bearer %s".formatted(apiKey))
8             .baseUrl(baseUrl)
9             .build();
10    }
11
12    @Bean
13    public OpenAIInterface openAIInterface(RestClient client) {
14        RestClientAdapter adapter = RestClientAdapter.create(client);
15        HttpServiceProxyFactory factory = HttpServiceProxyFactory.builderFor(adapter).build();
16        return factory.createClient(OpenAIInterface.class);
17    }
18 }
```

Service classes

- This is actually enough, but:
 - Might want to pre- or post-process requests and responses
- Model list is a simple GET request with no inputs
- But output can be extracted from result and organized
- Therefore, autowire interface into a service

```
1 @Service
2 public class OpenAIService {
3     public static final String GPT35 = "gpt-3.5-turbo";
4     public static final String GPT4 = "gpt-4-1106-preview";
5
6     private final OpenAIInterface openAIInterface;
7
8     @Autowired
9     public OpenAIService(OpenAIInterface openAIInterface) {
10         this.openAIInterface = openAIInterface;
11     }
12
13     public List<String> getModelNames() {
14         return openAIInterface.listModels().data().stream()
15             .map(ModelList.Model::id)
16             .sorted()
17             .toList();
18     }
19 }
```

```
1 @SpringBootTest
2 class OpenAIServiceTest {
3     @Autowired
4     private OpenAIService openAIService;
5
6     @Test
7     void getGPTModels() {
8         List<String> modelNames = openAIService.getModelNames();
9         assertThat(modelNames).anyMatch(name → name.contains(OpenAIService.GPT35));
10        assertThat(modelNames).anyMatch(name → name.contains(OpenAIService.GPT4));
11        modelNames.stream()
12            .filter(name → name.contains("gpt"))
13            .forEach(System.out::println);
14    }
15
16    @Test
17    void getAllModels() {
18        openAIService.getModelNames()
19            .forEach(System.out::println);
20    }
21 }
```

Generate audio from text

- Now that we have the interface, the service, and the necessary beans:
 - Add methods to exchange interface
 - Add methods to service to pre-process requests and post-process responses
- Text to speech
 - Use speech endpoint
 - Send POST request with text
 - Choose built-in voice (six available)
 - Save returned byte[] to file
 - (Optional) Add an mp3 player
- <https://platform.openai.com/docs/guides/text-to-speech>

Sample curl request

Generate spoken audio from input text

curl ▾  Copy

```
1 curl https://api.openai.com/v1/audio/speech \
2   -H "Authorization: Bearer $OPENAI_API_KEY" \
3   -H "Content-Type: application/json" \
4   -d '{
5     "model": "tts-1",
6     "input": "Today is a wonderful day to build something people love!",
7     "voice": "alloy"
8   }' \
9   --output speech.mp3
```

Create speech

- Five total parameters; three required
 - **model** → tts-1 or tts-1-hd
 - **input** → text for audio
 - **voice** → one of: *alloy*, *echo*, *fable*, *onyx*, *nova*, *shimmer*
 - **response_format** → one of: mp3, opus, aac, flac (mp3 is default)
 - **speed** → value between 0.25 and 4.0 (1.0 is default)

Map request JSON to record

```
1 public record TTSRequest(  
2     String model,  
3     String input,  
4     Voice voice,  
5     ResponseFormat responseFormat,  
6     double speed  
7 ) {  
8     public TTSRequest(String model, String input, Voice voice) {  
9         this(model, input, voice, ResponseFormat.MP3, 1.0);  
10    }  
11 }
```

Useful enums

```
1 public enum Voice {  
2     @JsonProperty("alloy") ALLOY,  
3     @JsonProperty("echo") ECHO,  
4     @JsonProperty("fable") FABLE,  
5     @JsonProperty("onyx") ONYX,  
6     @JsonProperty("nova") NOVA,  
7     @JsonProperty("shimmer") SHIMMER;  
8 }
```

```
1 public enum ResponseFormat {  
2     @JsonProperty("mp3") MP3,  
3     @JsonProperty("opus") OPUS,  
4     @JsonProperty("aac") AAC,  
5     @JsonProperty("flac") FLAC  
6 }
```

Add method to exchange interface

```
1 @HttpExchange("/v1")
2 public interface OpenAIInterface {
3
4     // Models
5     @GetExchange(value = "/models", accept = "application/json")
6     ModelList listModels();
7
8     // Text-to-Speech
9     @PostExchange(value = "/audio/speech",
10                  accept = "audio/mpeg", contentType = "application/json")
11     byte[] getTextToSpeechResponse(@RequestBody TTSRequest ttsRequest);
12 }
```

Add methods to service

```
 1 public byte[] getAudioResponse(TTSRequest ttsRequest) {  
 2     byte[] bytes = openAIInterface.getTextToSpeechResponse(ttsRequest);  
 3     String fileName = FileUtils.writeSoundBytesToFile(bytes);  
 4     logger.info("Saved {} to {}", fileName, FileUtils.AUDIO_DIRECTORY);  
 5     return bytes;  
 6 }  
 7  
 8 public byte[] getAudioResponse(String prompt) {  
 9     return getAudioResponse(new TTSRequest(DEFAULT_MODEL, prompt, DEFAULT_VOICE));  
10 }  
11  
12 // .. add convenient overloads ...
```

Try it out

```
 1 @SpringBootTest
 2 class OpenAIServiceTest {
 3     @Autowired
 4     private OpenAIService service;
 5
 6     private final String audioPrompt = """
 7         The YouTube channel, "Tales from the jar side" is your best
 8         source for learning about Java, Spring, and other open source
 9         technologies, especially when combined with AI tools.
10         The companion newsletter of the same name is also a lot of fun.
11         """ .replaceAll("\\s+", " ") .trim();
12
13     @Test
14     void getAudioResponse() {
15         byte[] audioResponse = service.getAudioResponse(audioPrompt);
16         assertThat(audioResponse.length).isPositive();
17     }
18 }
```

(Optional) Add MP3 Player

- JLayer: <https://github.com/umjammer/jlayer>
- Available through jitpack.io
- Instantiate **Player** class
 - Requires `InputStream`

JLayer dependency

```
1 repositories {  
2     mavenCentral()  
3     maven { setUrl("https://jitpack.io") }  
4 }  
5  
6 dependencies {  
7     // JLayer for MP3 support  
8     implementation("com.github.umjammer:jlayer:1.0.2")  
9 }
```

Add method to service

```
1 public void createAndPlay(String text, Voice voice) {  
2     TTSRequest ttsRequest = new TTSRequest(TTS_1_HD, text, voice);  
3     byte[] bytes = getAudioResponse(ttsRequest);  
4     try {  
5         Player player = new Player(buffer);  
6         player.play();  
7     } catch (JavaLayerException e) {  
8         throw new RuntimeException(e);  
9     }  
10 }
```

(Optional) Validate Request

- Add `spring-boot-starter-validation` dependency
- Can now use JPA validation annotations

```
1 import jakarta.validation.constraints.*;
2
3 public record TTSRequest(
4     @Pattern(regexp = "tts-1(-hd)?(-1106)?") String model,
5     @NotBlank @Size(max = 4096) String input,
6     Voice voice,
7     ResponseFormat responseFormat,
8     @DecimalMin("0.25") @DecimalMax("4.0") double speed
9 ) {
10     public TTSRequest(String model, String input, Voice voice) {
11         this(model, input, voice, ResponseFormat.MP3, 1.0);
12     }
13 }
```

Validation

- Inject **Validator** instance into service
 - Careful! There are several Validator classes in classpath
 - Be sure to pick the one from jakarta.validation
- Validate request

```
1 private void validateRequest(TTSRequest ttsRequest) {  
2     Set<ConstraintViolation<TTSRequest>> violations = validator.validate(ttsRequest);  
3     if (!violations.isEmpty()) {  
4         throw new IllegalArgumentException(violations.toString());  
5     }  
6 }
```

(Optional) Add a controller

```
1 @RestController
2 public class TextToSpeechController {
3     // ... autowire OpenAIService ...
4
5     @PostMapping("/tts")
6     public ResponseEntity<String> convertTextToSpeech(@RequestBody String text) {
7         service.getAudioResponse(TextToSpeechService.TTS_1_HD, text, Voice.randomVoice());
8         return ResponseEntity.ok("%s... converted to mp3".formatted(text.substring(0, 24)));
9     }
10
11    @PostMapping("/tts/play")
12    public ResponseEntity<String> convertTextAndPlay(@RequestBody String text) {
13        service.createAndPlay(text, Voice.randomVoice());
14        return ResponseEntity.ok("%s... converted to mp3".formatted(text.substring(0, 24)));
15    }
16 }
```

(Optional) GraalVM native image compiler

- Spring supports GraalVM native image compiler
- For stateless services, this can be convenient
- Use with Function As A Service cloud providers
- Rapid startup and shutdown of stateless services

Process

- This process can be used for all the services
 - Map inputs and outputs to records
 - Add methods to HTTP exchange interface
 - Autowire interface into service for preprocessing requests and postprocessing responses
- Optionally:
 - Add enums or other constants
 - Add helper functions to extract desired data from responses
 - Validate requests
 - Add controllers

DALL-E Image Generation

- **Image generation is included in GitHub repository**
 - Request has only a few parameters
 - Response can either be a URL for the image (default), or
 - Response can be the 64-bit binary encoded image itself
- **Valid image sizes are all square:**
 - 1024x1024 (default for DALL-E 3)
 - Other sizes cost more
 - High definition costs more
- **Can request only 1 image at a time for DALL-E 3**
- **Endpoint:**
 - POST to <https://api.openai.com/v1/images/generations>

Image Request

```
1 {
2   "prompt": "A cute baby sea otter",
3   "n": 2,
4   "size": "1024x1024",
5   "response_format": "b64_json"
6 }
```

```
1 public record ImageRequest(String prompt,
2                             Integer n,
3                             String size,
4                             String responseFormat) {
```

```
1 public record ImageResponse(Long created,
2                             List<Image> data) {
```

- **response_format:**
 - url (default)
 - b64_json → returns the binary encoded string directly

Saving an image to a file

- May return multiple images, so need unique file names
- My approach is to use a combination of a timestamp and a counter
- From FileUtils:
 - int counter is an attribute of this class

```
1 public static boolean writeImageToFile(String imageData) {  
2     String timestamp = LocalDateTime.now()  
3         .format(DateTimeFormatter.ofPattern("yyyyMMddHHmmss"));  
4     String fileName = String.format("image_%s_%d.png", timestamp, counter++);
```

Saving an image to a file

- Java library has `java.util.Base64` class
 - `Base64.getDecoder().decode(imageData) → byte[]`

```
1 try {
2     Files.createDirectories(directory);
3     byte[] bytes = Base64.getDecoder().decode(imageData);
4     Files.write(filePath, bytes, StandardOpenOption.CREATE_NEW);
5     System.out.printf("Saved %s to src/main/resources/images%n", fileName);
6     return true;
7 } catch (IOException e) {
8     throw new UncheckedIOException("Error writing image to file", e);
9 }
```

Test

```
1  @Test
2  public void testDallE() {
3      DallE dallE = new DallE();
4      long num = dallE.getImages("""
5          A realistic photo of a
6          robot leaping into the air
7          in joy after accomplishing a
8          difficult task successfully
9          """, 2);
10     System.out.printf("Downloaded %d images%n", num);
11     assertThat(num).isEqualTo(2);
12 }
```

Spring and AI without Spring AI

- See [my AI playlist](#) on the *Tales from the jar side* [YouTube channel](#) :)

Conclusions

- Spring AI works, but very early in development process
 - Expect frequent changes
- Best features:
 - Simple chat requests and responses
 - Prompt templates for input
 - Output parsers for output
 - RAG with embedding clients and vector databases
 - Works with multiple LLMs (though OpenAI most mature)

Conclusions

- Otherwise, Spring provides helpful features
 - HTTP Exchange interfaces
 - JavaConfig files to configure rest clients and implement exchange interfaces
 - Autowiring into services
 - Jackson JSON parsing
 - Validation using the Bean Validation specification
 - Rest controllers
 - GraalVM native image compilation
- Expect Spring AI project to improve