

# RAG Spring AI

Java can also handle LLMs

# Przemek Skwiercz

- Programista i cyklista.
- Jeden z założycieli i organizatorów bydgoskiego JUG'a.
- Od ponad 20 lat zarabiam na kromkę chleba i kufelek piwa tworząc oprogramowanie więc jestem już chyba seniorem (wiekowo na pewno) :)
- W wolnym czasie miłośnik: wypraw rowerowych, kraftowego piwa, starego punk rocka i teatru (niekoniecznie w tej kolejności).



# RAG Spring AI Presentation

**RAG - What, Why and How**

How **SPRING AI** support creation of RAG App

**RAG Spring AI - Demo app**

# RAG

# What, Why & How

# LLM

## Main Problems

- “Static knowledge” of LLMs
- No real time updates (dynamic state)
- No access to highly specific data
- Security of customers data

# LLM

## Providing custom data to LLMs

- **Fine tune model** - rather for people with knowledge in AI areas, not for common users of LLMs
- **Function Calling** - permits the model to request the execution of client-side functions, thereby accessing necessary information or performing tasks dynamically as required
- **“Stuff the prompt”** - adds custom data to prompts -> **RAG**

# RAG

## Use cases

- Customer feedback analysis
- Speeding up onboarding new employees
- Boosting customer support
- Search engines
- Medical analysis

# RAG

## What is RAG?

**RAG - Retrieval-Augmented Generation**

### Retrieval Phase

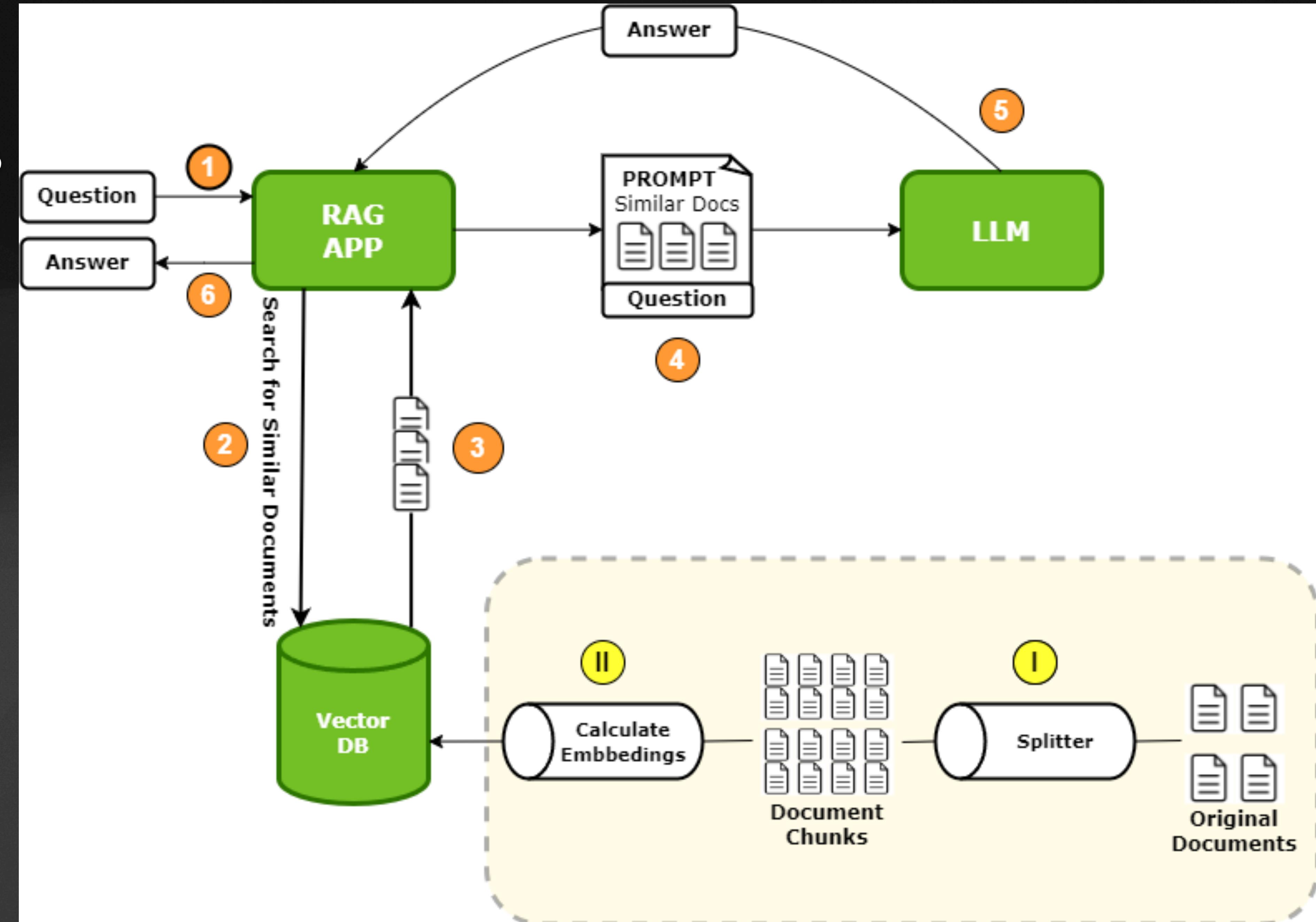
Search local data (stored in vector DB) to find relevant information to query

### Generation Phase

LLM which works with custom data and generates answer

# RAG

## How it works?

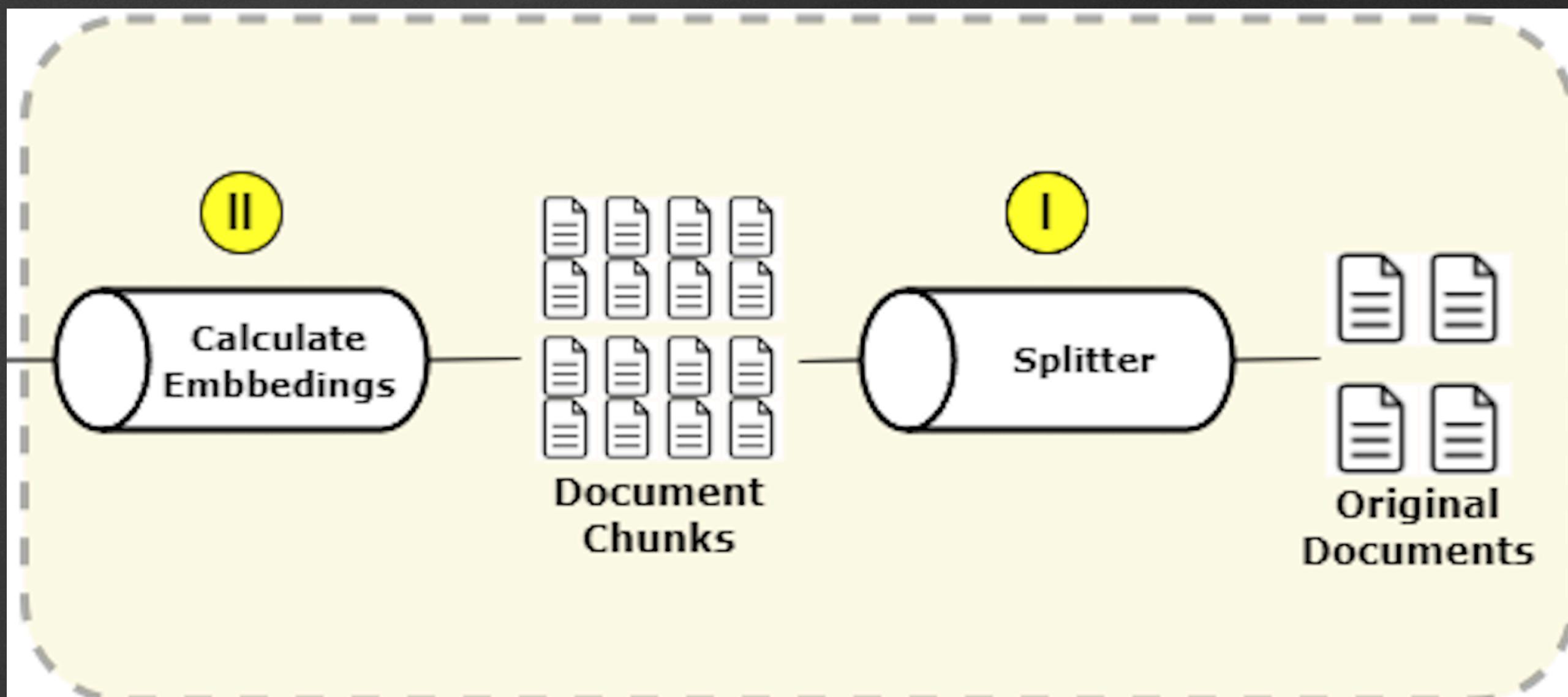


# RAG

## Embeddings

**Embeddings** - convert data into vector representation. It is performed in two steps:

1. Tokenisation - (algorithm **BPE** - Byte Pair Encoding)
2. Transformation - tokens to vectors (libraries provided by LLMs)



# RAG Embeddings

experiment with  
tokenisation:



GPT-4o & GPT-4o mini    GPT-3.5 & GPT-4    GPT-3 (Legacy)

Many words map to one token, but some don't: indivisible.

Unicode characters like emojis may be split into many tokens containing the underlying bytes: 🖐

Sequences of characters commonly found next to each other may be grouped together: 1234567890

[Clear](#) [Show example](#)

| Tokens | Characters |
|--------|------------|
| 53     | 252        |

Many words map to one token, but some don't: indivisible.

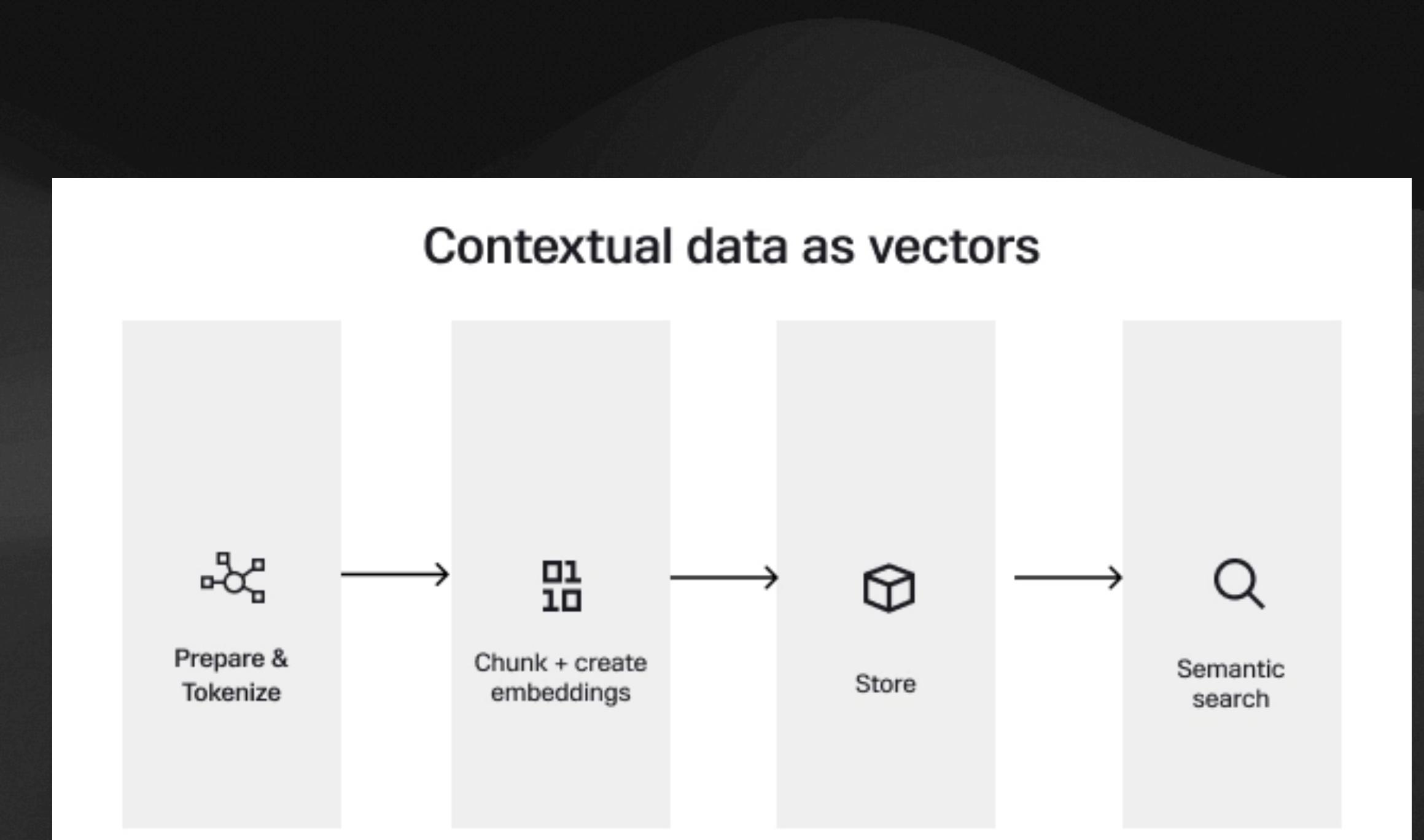
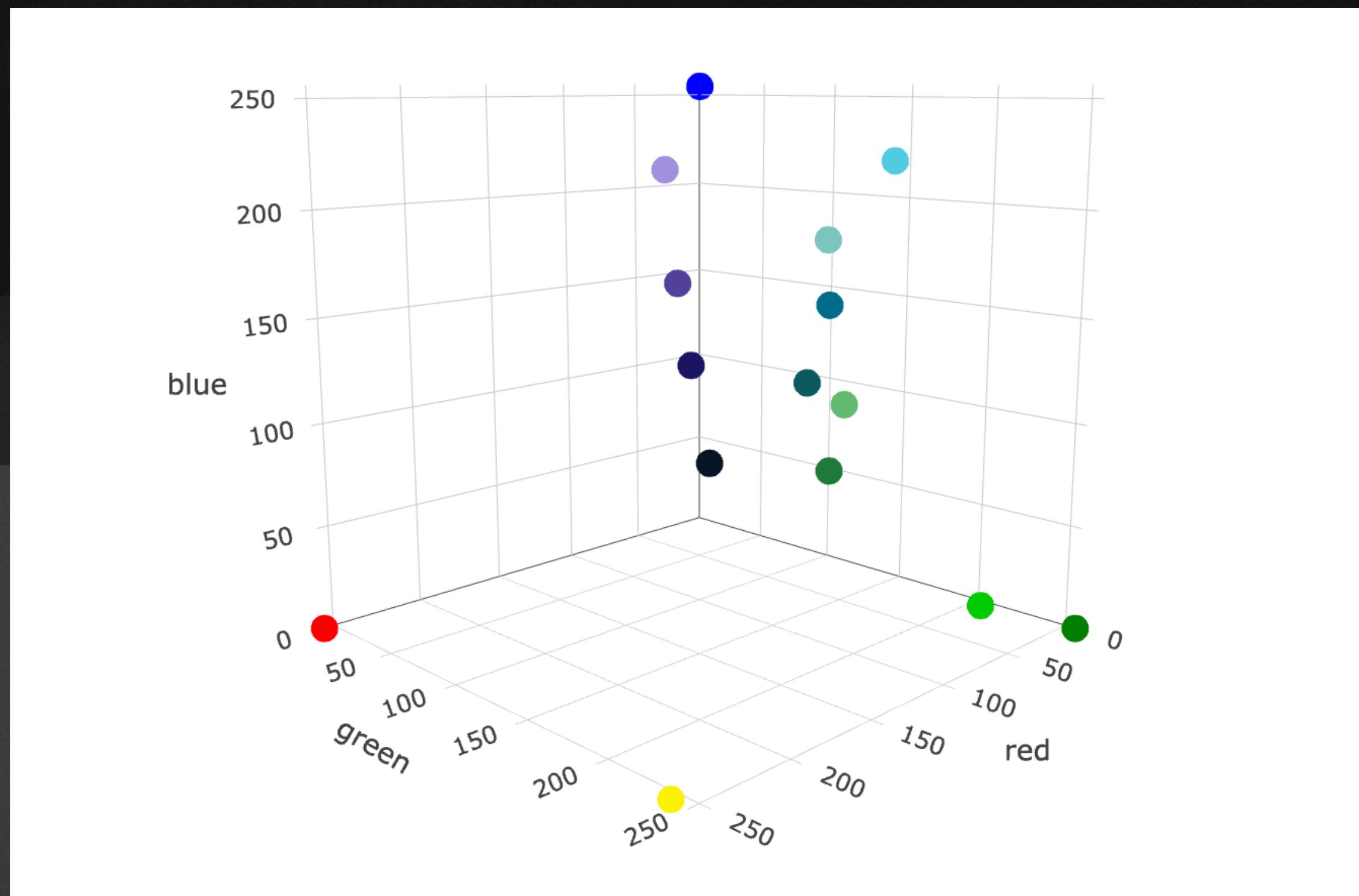
Unicode characters like emojis may be split into many tokens containing the underlying bytes: 🖐

Sequences of characters commonly found next to each other may be grouped together: 1234567890

[Text](#) [Token IDs](#)

# RAG

## Data as Vectors



Source: <https://weaviate.io/blog/what-is-a-vector-database>

Source: <https://www.singlestore.com/blog/a-complete-guide-to-vector-databases/>

# RAG

## Similarity Search

Vector embeddings allow us to find and retrieve similar objects from the vector database by searching for objects that are close to each other in the vector space, which is called vector search, **similarity search**, or semantic search.

***Paris - France + Poland = Warsaw***

***King - Man + Woman = Queen***

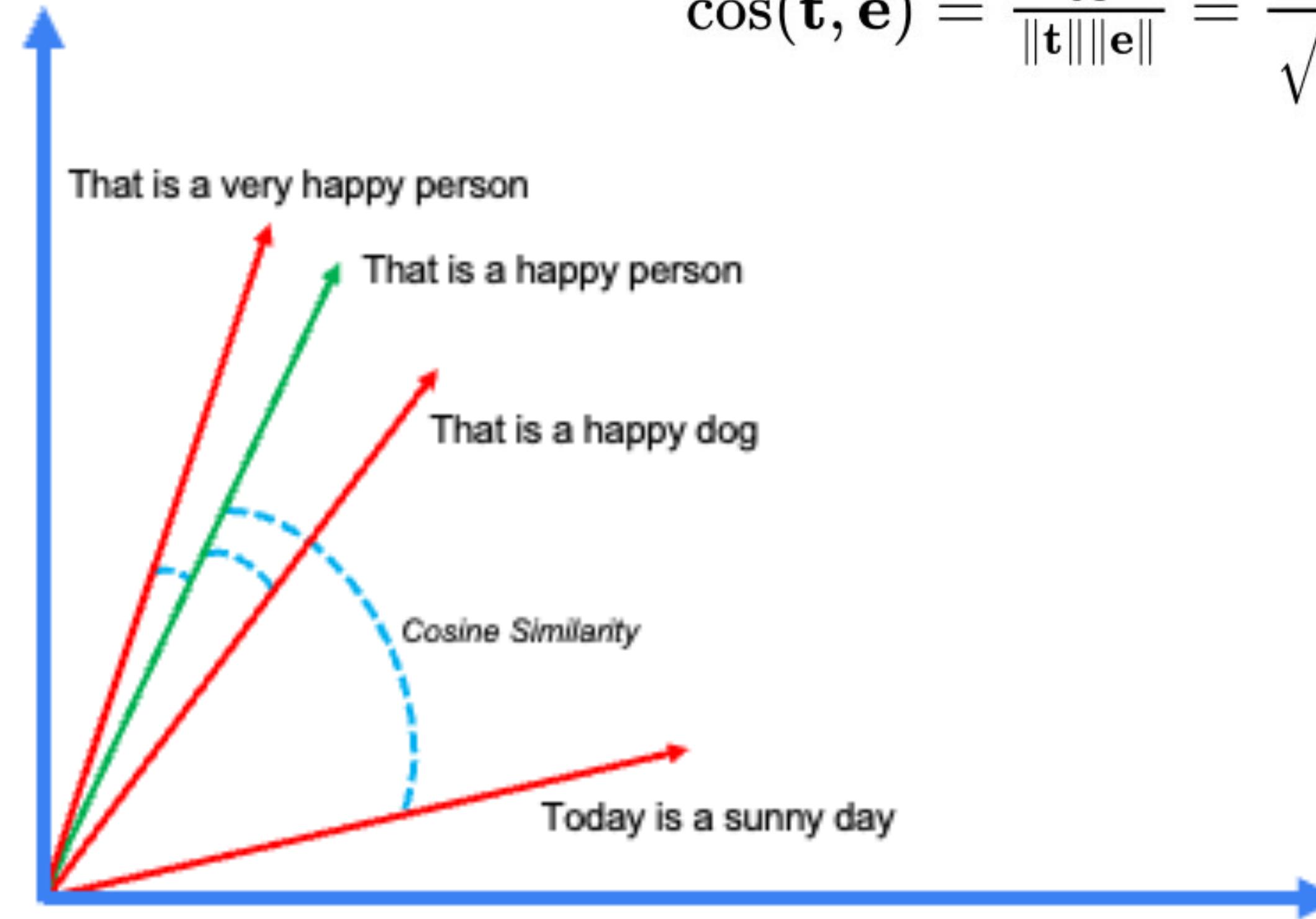


# RAG

## Cosine

### Similarity

$$\cos(\mathbf{t}, \mathbf{e}) = \frac{\mathbf{t} \cdot \mathbf{e}}{\|\mathbf{t}\| \|\mathbf{e}\|} = \frac{\sum_{i=1}^n t_i e_i}{\sqrt{\sum_{i=1}^n (t_i)^2} \sqrt{\sum_{i=1}^n (e_i)^2}}$$



# RAG

## Cosine Similarity

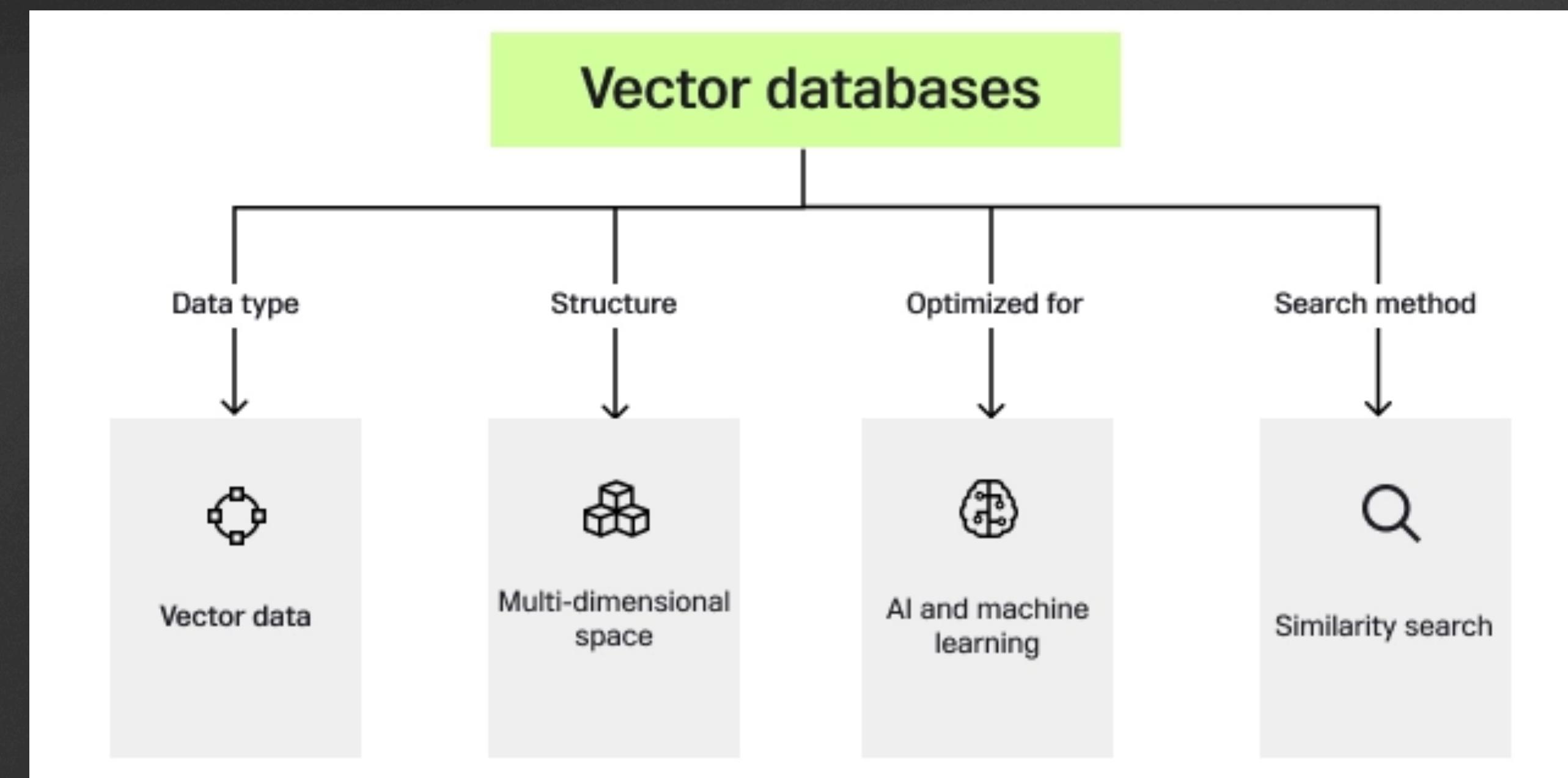
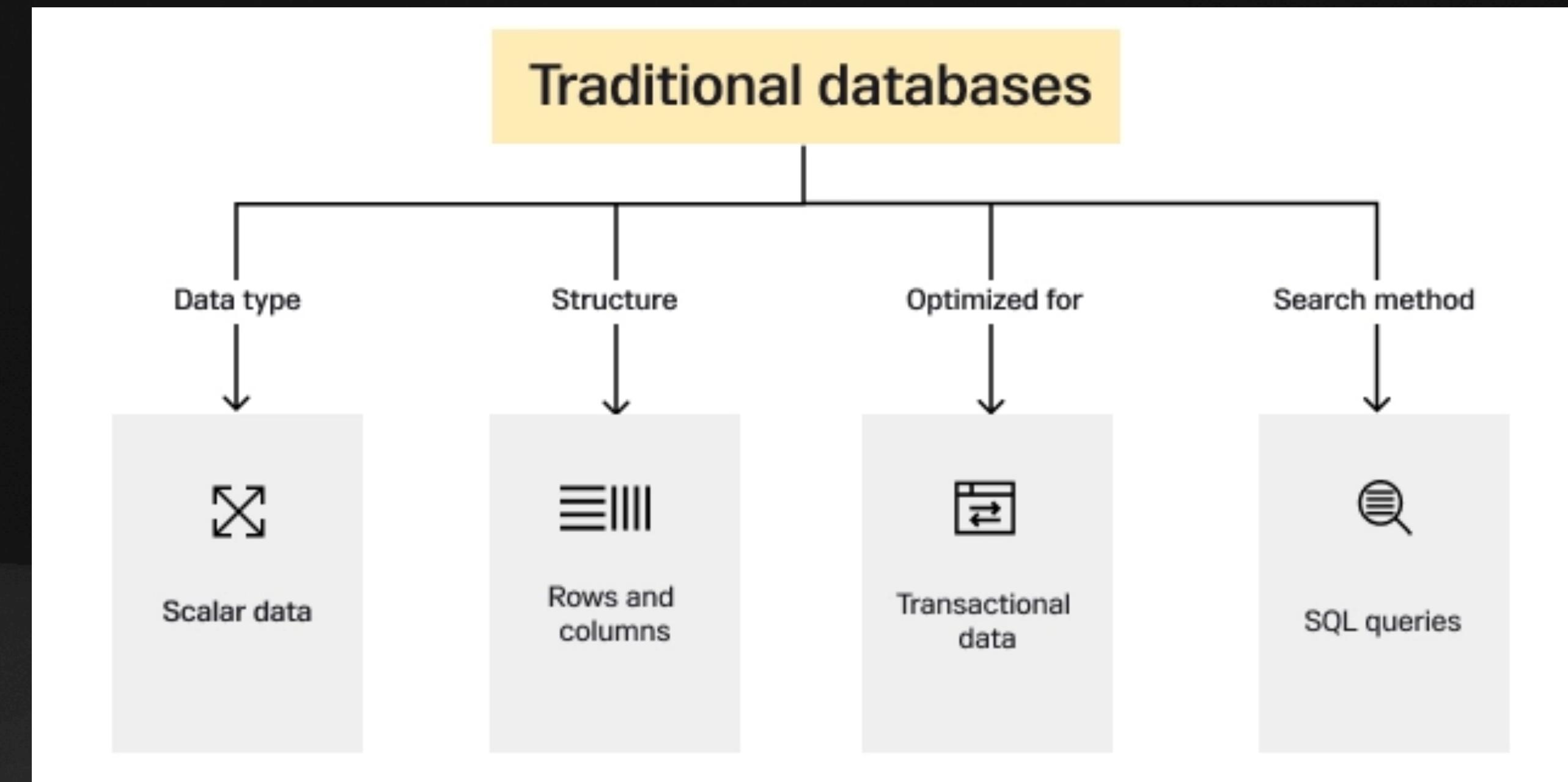
$$\cos(\mathbf{t}, \mathbf{e}) = \frac{\mathbf{t}\mathbf{e}}{\|\mathbf{t}\| \|\mathbf{e}\|} = \frac{\sum_{i=1}^n \mathbf{t}_i \mathbf{e}_i}{\sqrt{\sum_{i=1}^n (\mathbf{t}_i)^2} \sqrt{\sum_{i=1}^n (\mathbf{e}_i)^2}}$$

# RAG

## Traditional DB

### vs

## Vector DB



# RAG

## Prompts

**Prompts** - are the inputs that guide an AI model to generate specific outputs

**Prompt Types:**

SYSTEM

USER

ASSISTANT

TOOL

# RAG

## System prompt vs User prompt

### System Prompt

- System prompt is a fixed prompt providing context and instructions to the model.
- The system prompt is always included in the provided input to the LLM, regardless of the user prompt.
- The idea is that the system prompt should not be editable by any party other than the developers.

### User Prompt

- User prompts are the input queries from the user.
- This is how the customer interacts with LLM application, and it enables the user to shape this application's input.

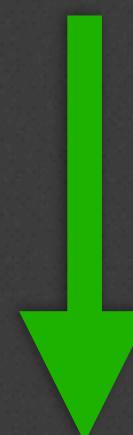
# Spring AI

# Spring AI

## When was it made?

In **mid-2023**

Spring revealed that it was working on a module that makes it easier to create AI applications using the Spring framework



**12.2023**

Spring AI left the experimental phase and was officially released



**01.2024**

Version 0.8.0 was released (most often used in courses, YT videos, etc.)



**04.2024**

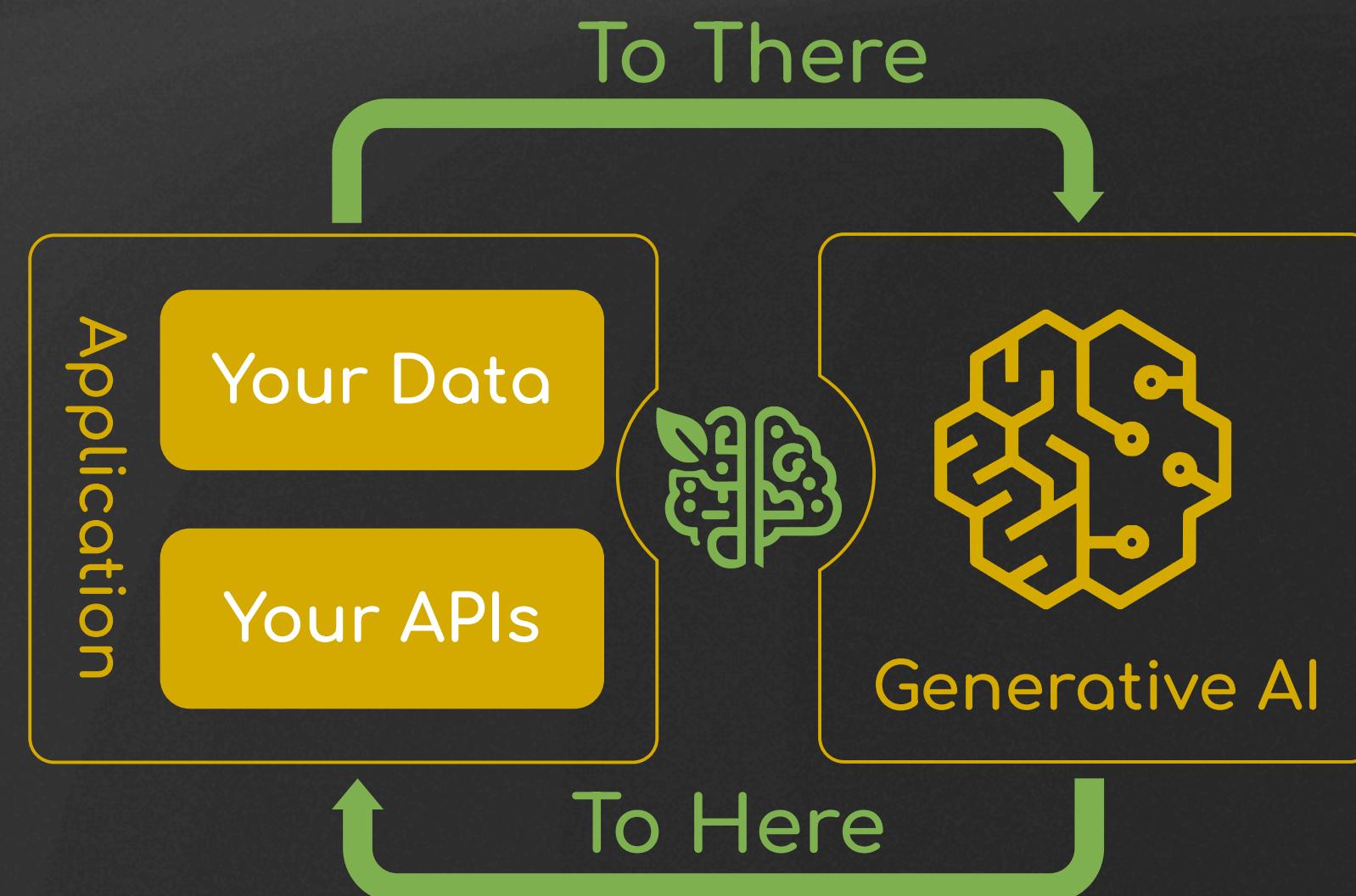
Version 1.0.0 (1.0.0RC) was released, introducing changes causing backward incompatibility (mainly ChatClient/ChatModel)



# Spring AI

## What is it used for?

- Spring AI is designed to integrate AI models with applications written in Spring
- It allows us to connect our data and API with AI models
- Spring AI is an abstraction on LLMs



# Spring AI Supported LLMs

- OpenAI Chat Completion (streaming, multi-modality & function-calling support)
- Microsoft Azure Open AI Chat Completion (streaming & function-calling support)
- Ollama Chat Completion (streaming, multi-modality & function-calling support)
- Hugging Face Chat Completion (no streaming support)
- Google Vertex AI PaLM2 Chat Completion (no streaming support)
- Google Vertex AI Gemini Chat Completion (streaming, multi-modality & function-calling support)
- Amazon Bedrock (Cohere, Llama, Titan, Anthropic, Jurassic2)
- Mistral AI Chat Completion (streaming & function-calling support)
- Anthropic Chat Completion (streaming & function-calling support)
- watsonx.AI

# Spring AI Supported Vector DB

- **SimpleVectorStore** - A simple implementation of persistent vector storage, good for educational purposes.
- **PgVector Store**
- **Azure Vector Search**
- **Apache Cassandra**
- **Elasticsearch Vector Store**
- **GemFire Vector Store**
- **Milvus Vector Store**
- **MongoDB Atlas Vector Store**
- **Neo4j Vector Store**
- **OpenSearch Vector Store**
- **Oracle Vector Store**
- **Pinecone Vector Store**
- **Qdrant Vector Store**
- **Redis Vector Store**
- **SAP Hana Vector Store**

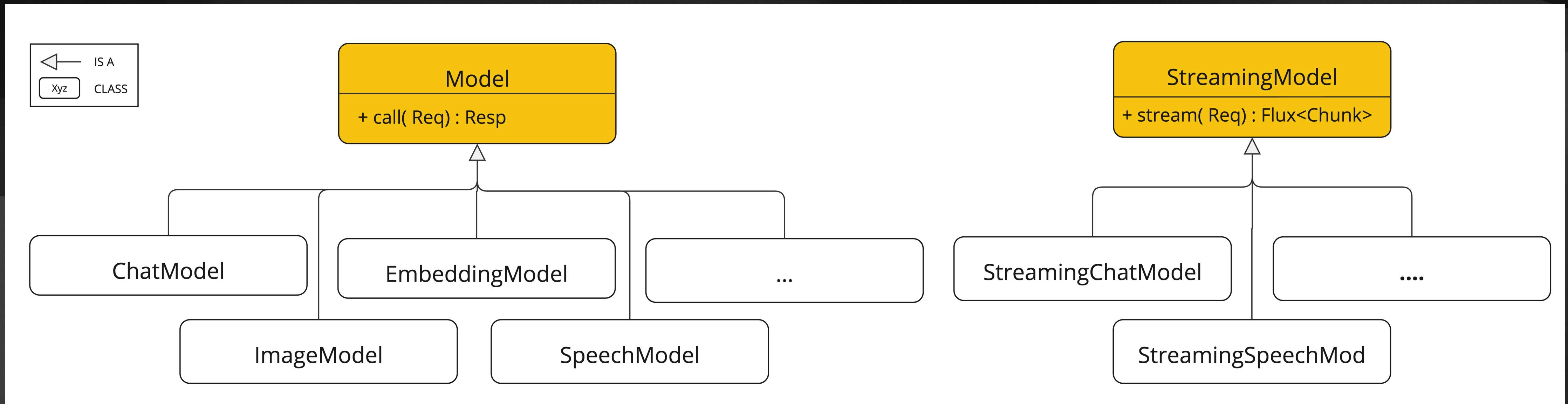
# Spring AI

## Alternatives for Spring AI

|   | JDK              | GraalVM   | Release              | RAG support   | Doc   | Dev Activity  |
|---|------------------|---|----------------------|---|---|---|
| <br><b>LangChain</b>       | 8, 11,<br>17, 21 |  | General Availability |  |  |  |
| <br><b>Spring AI</b>       | 17               |  | Experimental         |  |  |  |
| <br><b>Semantic Kernel</b> | 17, 21           |  | Alpha                |  |  |  |

February 2024

# Spring AI API



# Spring AI Simple app

```
@Service no usages 1 inheritor - pskwiercz *
public class OpenAiServiceExample {

    private final ChatClient chatClient; 2 usages

    public OpenAiServiceExample(ChatClient.Builder chatClient) { new *
        this.chatClient = chatClient.build();
    }

    ⚡ Rename usages
    public String ask(String question) { new *

        return chatClient
            .prompt() ChatClientRequest
            .user(question)
            .call() CallResponseSpec
            .content();
    }
}
```

```
spring.application.name=spring-ai-ud
server.port=8080
spring.ai.openai.api-key=${OPEN_AI_KEY}
spring.ai.openai.chat.options.model=gpt-4o
```

```
<parent>
    <groupId>org.springframework.boot</groupId>
    <artifactId>spring-boot-starter-parent</artifactId>
    <version>3.3.0</version>
    <relativePath/> <!-- lookup parent from repository -->
</parent>
<groupId>com.pskwiercz</groupId>
<artifactId>spring-ai-ud</artifactId>
<version>0.0.1-SNAPSHOT</version>
<packaging>war</packaging>
<name>spring-ai-ud</name>
<description>spring-ai-ud</description>
<properties>
    <java.version>21</java.version>
    <spring-ai.version>1.0.0-M1</spring-ai.version>
</properties>
<dependencies> ⚡ Add Starters...
    <dependency>
        <groupId>org.springframework.boot</groupId>
        <artifactId>spring-boot-starter-web</artifactId>
    </dependency>
    <dependency>
        <groupId>org.springframework.ai</groupId>
        <artifactId>spring-ai-openai-spring-boot-starter</artifactId>
    </dependency>
</dependencies>
```

# SPRING AI

## Vector DB Configuration

- Vector DB are used to integrate customer data with AI models.
- When a user query is to be sent to the AI model, a set of similar documents is first retrieved. These documents then serve as the context for the user's question and are sent to the AI model, along with the user's query.

```
public interface VectorStore {  
  
    void add(List<Document> documents);  
  
    Optional<Boolean> delete(List<String> idList);  
  
    List<Document> similaritySearch(String query);  
  
    List<Document> similaritySearch(SearchRequest request);  
}
```

JAVA

# SPRING AI

## Vector DB configuration

```
ai:  
  openai:  
    api-key: ${OPENAI_API_KEY}  
    embedding:  
      options:  
        model: text-embedding-ada-002  
    chat:  
      options:  
        model: gpt-4-turbo  
  
vectorstore:  
  milvus:  
    initialize-schema: true  
    client:  
      host: "localhost"  
      port: 19530  
      username: "root"  
      password: "milvus"  
      databaseName: "default"  
      collectionName: "vector_store"  
      embeddingDimension: 1536  
      indexType: IVF_FLAT  
      metricType: COSINE
```

# SPRING AI

## Prompts

```
public class Prompt implements ModelRequest<List<Message>> {  
  
    private final List<Message> messages;  
  
    private ChatOptions chatOptions;  
}
```

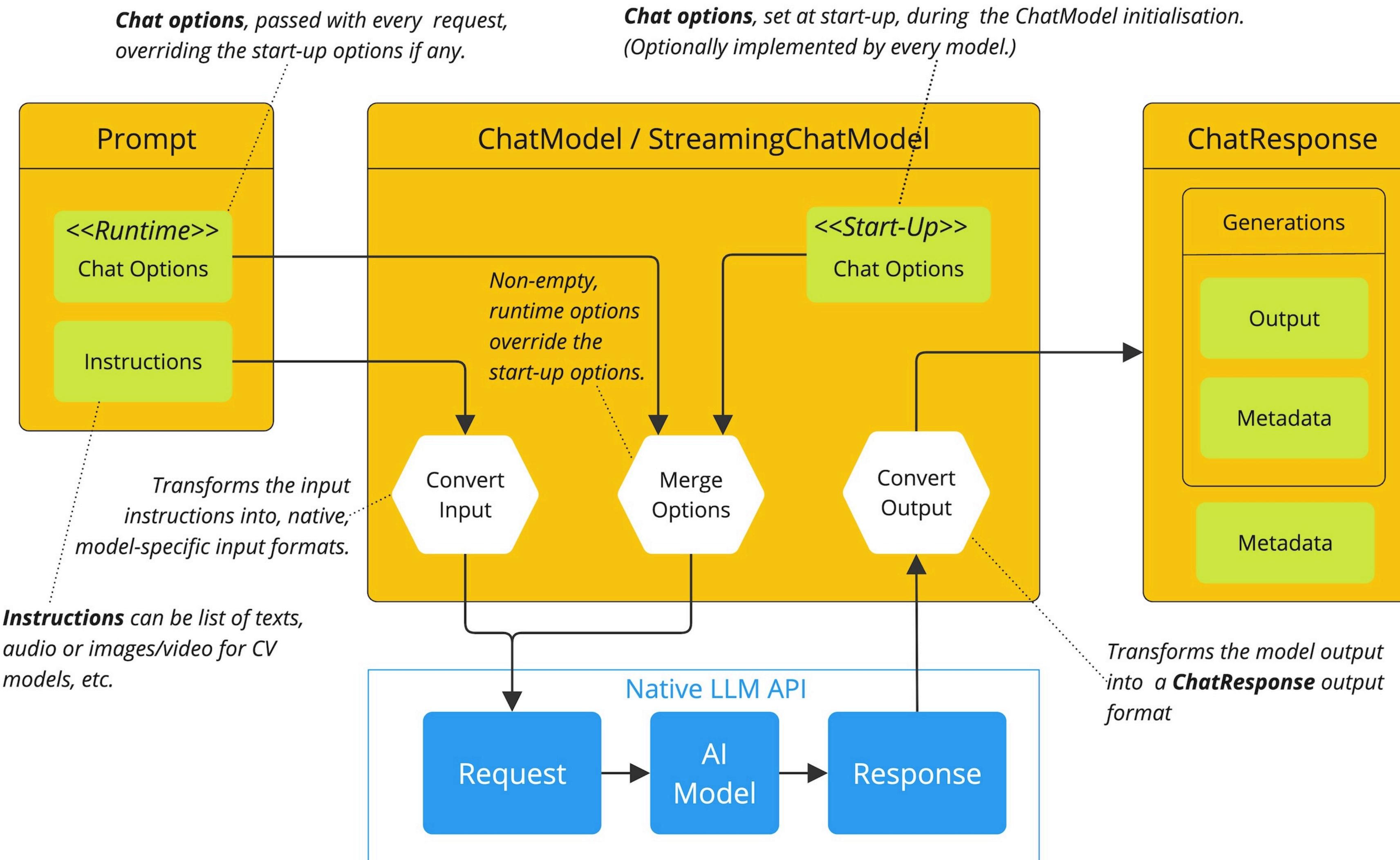
JAVA

```
public interface ChatOptions extends ModelOptions {  
  
    String getModel();  
    Float getFrequencyPenalty();  
    Integer getMaxTokens();  
    Float getPresencePenalty();  
    List<String> getStopSequences();  
    Float getTemperature();  
    Integer getTopK();  
    Float getTopP();  
    ChatOptions copy();  
}
```

JAVA

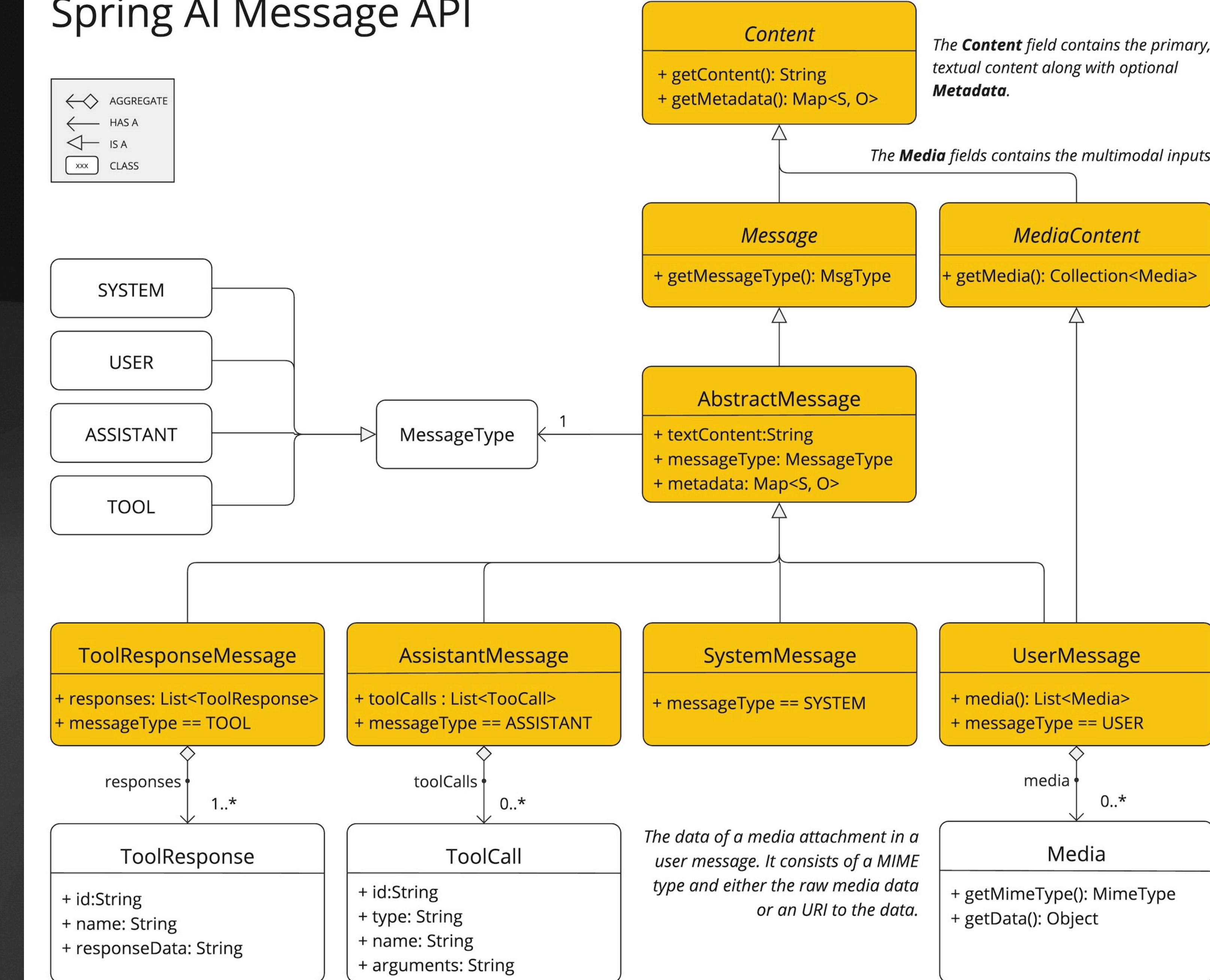


# Spring AI Chat Model



# SPRING AI Prompts

# Spring AI Message AP



*The **Content** field contains the primary, textual content along with optional **Metadata**.*

*The **Media** fields contains the multimodal inputs.*

## *Message*

## *MediaContent*

AbstractMessage

---

- + textContent:String
- + messageType: MessageType
- + metadata: Map<S, O>

| ToolResponseMessage                                      | AssistantMessage   | SystemMessage           |
|--|--|-------------------------|
| + responses: List<ToolResponse><br>+ messageType == TOOL | + toolCalls : List<ToolCall><br>+ messageType == ASSISTANT | + messageType == SYSTEM |

UserMessage

---

- + media(): List<Media>
- + messageType == USER

*The data of a media attachment in a user message. It consists of a MIME type and either the raw media data or an URI to the data.*

- + getMimeType(): MimeType
- + getData(): Object

# SPRING AI

## String Templates

```
1  What is the longest river in {country}? Response in JSON with the name of the river in a property named
2  Do not wrap the json codes in JSON

© OpenAiService.java ×

12
13  @Service 2 usages  ↗ pskwiercz *
14  public class OpenAiService {
15
16      private final ChatClient chatClient; 2 usages
17
18      @Value("classpath:templates/river-prompt.st")
19      private Resource riverPrompt;
20
21  ➜ >     public OpenAiService(ChatClient.Builder chatClient) { this.chatClient = chatClient.build(); }
24
25  @
26
27      public Answer getRiver(GetRiverRequest getRiverRequest) { 1 usage new *
28
29          Prompt prompt = new PromptTemplate(riverPrompt)
30              .create(Map.of( k1: "country", getRiverRequest.country()));
31
32          String answer = chatClient
33              .prompt(prompt) ChatClientPromptRequest
34              .call() CallPromptResponseSpec
35              .content();
36
37          return new Answer(answer);
38      }
39  }
```

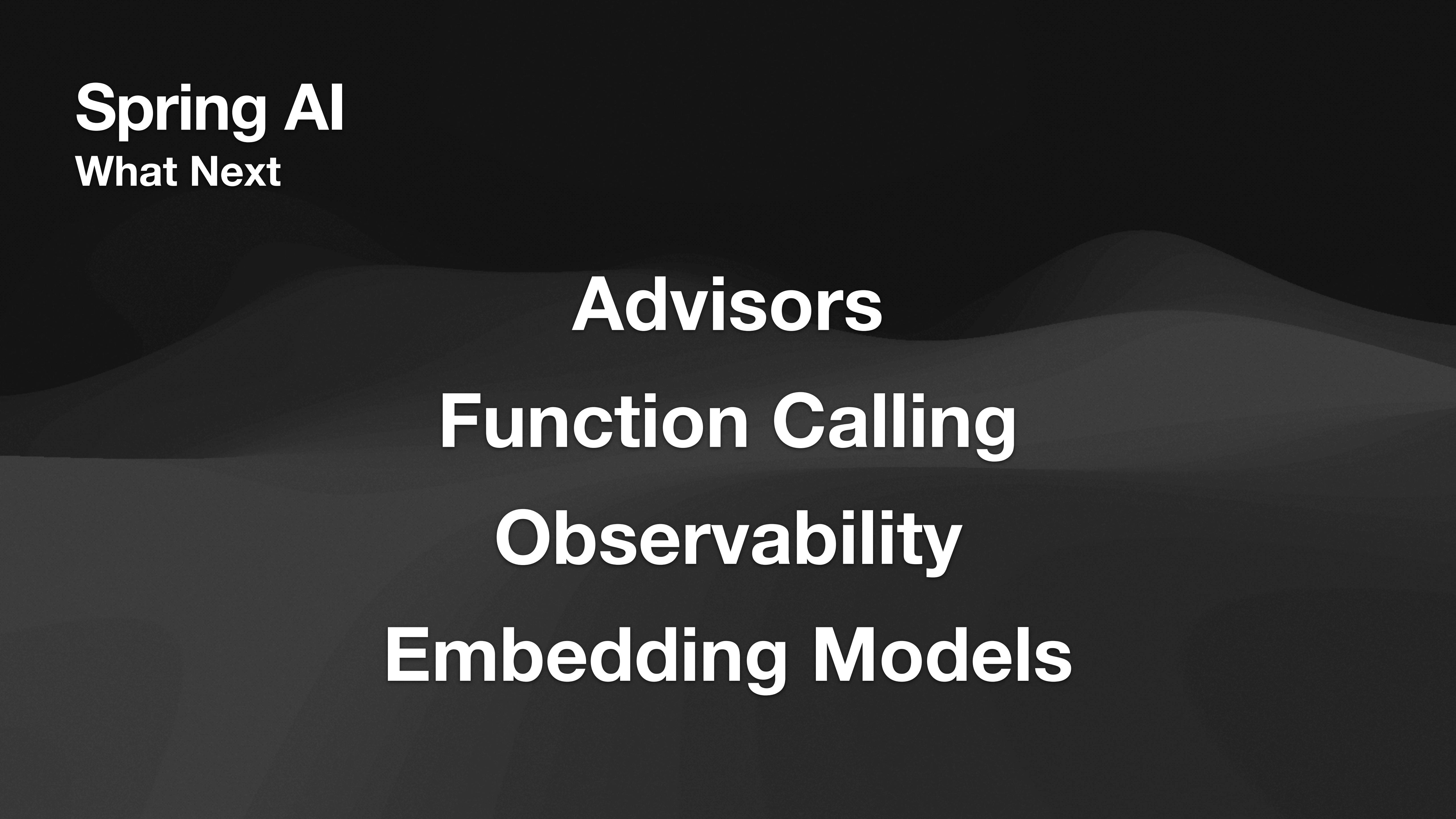
# Demo



[https://github.com/pskwiercz/  
jug-rag-spring-ai](https://github.com/pskwiercz/jug-rag-spring-ai)

# Spring AI

## What Next



Advisors  
Function Calling  
Observability  
Embedding Models

The background features a minimalist design with abstract, overlapping circular shapes in various shades of green. The colors transition from dark green on the left to light teal on the right.

# Q&A