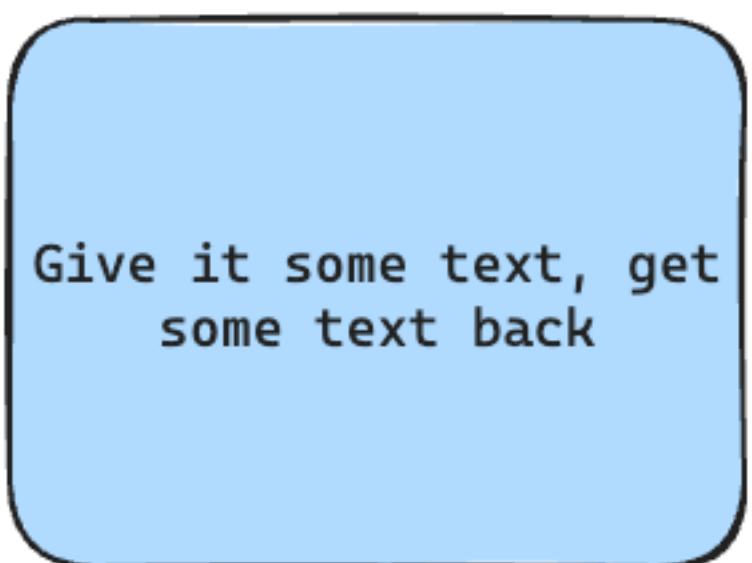


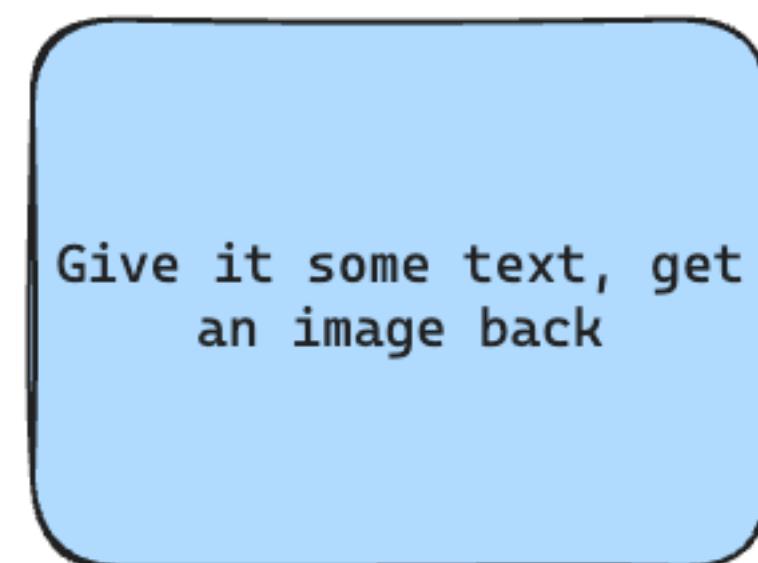
# **SPRING AI API**

**The Spring AI API covers a wide range of functionalities**

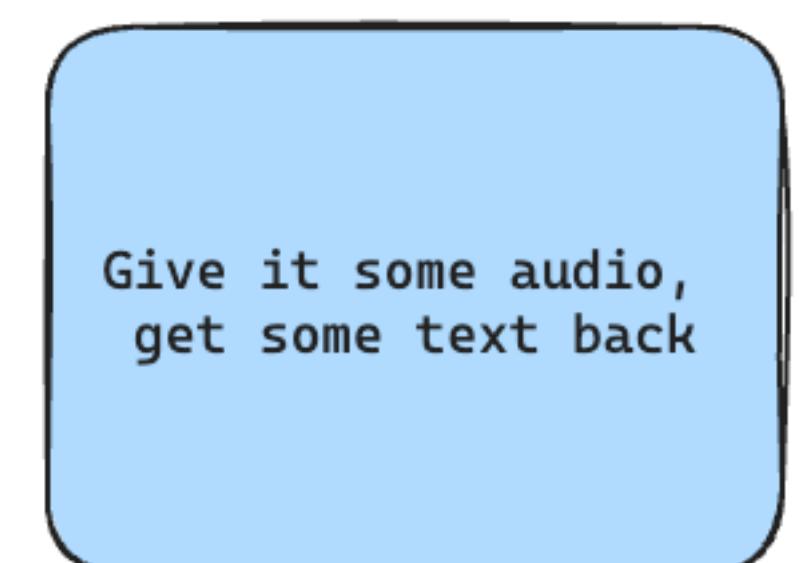
Chat Model



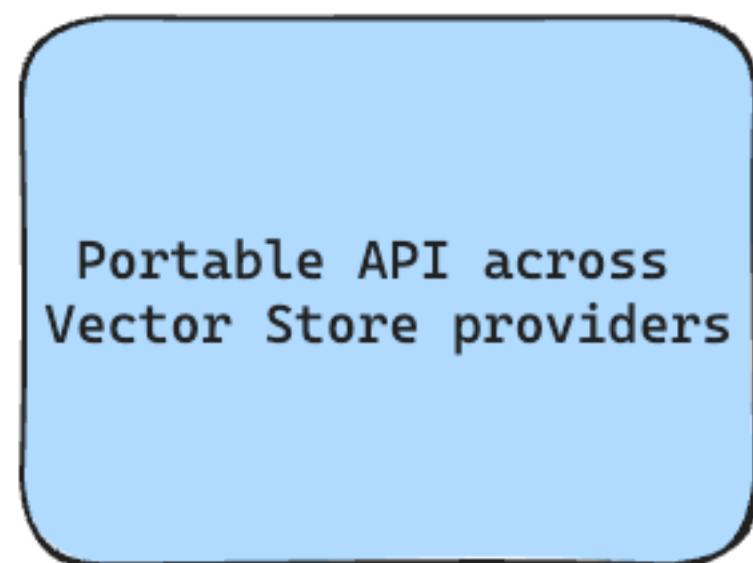
Text to Image



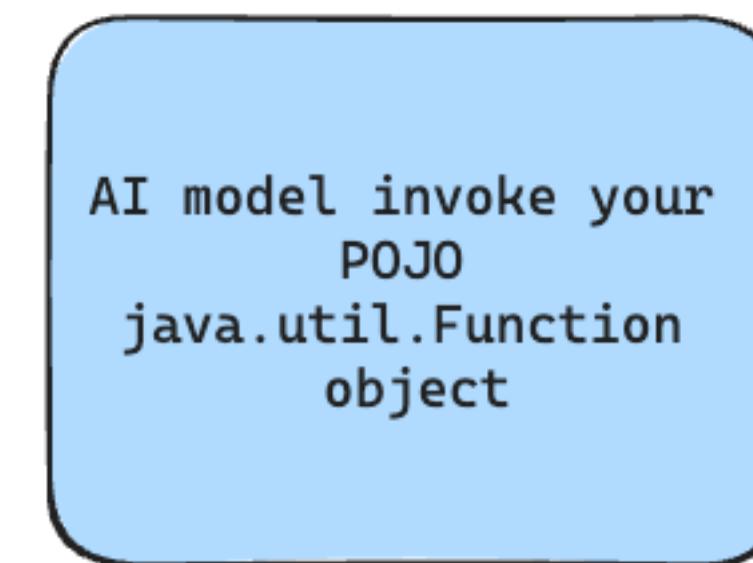
Transcription



Embedding



Functions



# **CHAT MODEL**

**Give it some text, get some text back**

- Open AI
- Azure Open AI
- Amazon Bedrock
- Google Vertex AI Palm
- Google Gemini
- HuggingFace - Access to thousands of models, including those from Meta such as Llama 2
- Ollama - Run AI Models on your local machine
- MistralAI

# TEXT- TO- IMAGE MODELS

Give it some text, get an image back

- OpenAI with DALL-E
- StabilityAI

# **TRANSCRIPTION**

**Give it some audio, get some text back**

- Open AI

# Embedding Models

## Portable API across Vector Store providers

- Open AI
- Azure Open AI
- Ollama
- ONNX
- PostgresML
- Bedrock Cohere
- Bedrock Titan
- Google VertexAI
- MistalAI



# GETTING STARTED



**Project**

- Gradle - Groovy     Gradle - Kotlin  
 Maven

**Language**

- Java     Kotlin     Groovy

**Spring Boot**

- 3.3.0 (SNAPSHOT)     3.3.0 (M3)     3.2.5 (SNAPSHOT)     3.2.4  
 3.1.11 (SNAPSHOT)     3.1.10

**Project Metadata**

Group dev.danvega

Artifact hello-ai

Name hello-ai

Description Demo project for Spring Boot

Package name dev.danvega

Packaging  Jar     War

Java  22     21     17

**Dependencies**

**ADD DEPENDENCIES...** ⌘ + B

**Spring Web** WEB

Build web, including RESTful, applications using Spring MVC. Uses Apache Tomcat as the default embedded container.

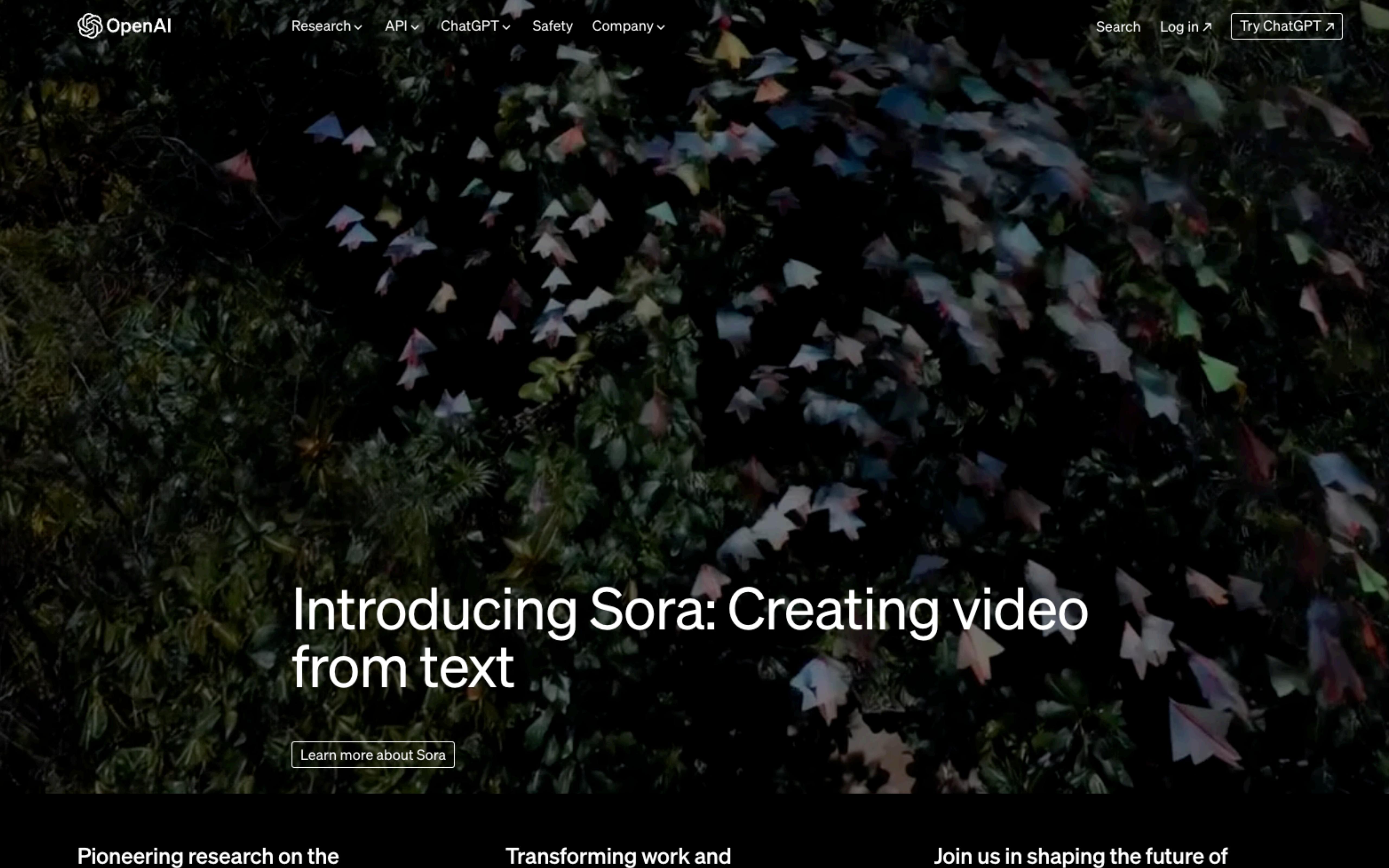
**OpenAI** AI

Spring AI support for ChatGPT, the AI language model and DALL-E, the Image generation model from OpenAI.

**GENERATE** ⌘ + ↵

**EXPLORE** CTRL + SPACE

**SHARE...**



# Introducing Sora: Creating video from text

[Learn more about Sora](#)



## API keys

Your secret API keys are listed below. Please note that we do not display your secret API keys again after you generate them.

Do not share your API key with others, or expose it in the browser or other client-side code. In order to protect the security of your account, OpenAI may also automatically disable any API key that we've found has leaked publicly.

Enable tracking to see usage per API key on the [Usage page](#).

NAME	SECRET KEY	TRACKING ⓘ	CREATED	LAST USED ⓘ	PERMISSIONS	⋮
Secret key	sk-...aWYE	+ Enable	Feb 6, 2023	Never	All	
spring-ai	sk-...NKve	Enabled	Mar 12, 2024	Mar 13, 2024	All	

[+ Create new secret key](#)

## Default organization

If you belong to multiple organizations, this setting controls which organization is used by default when making requests with the API keys above.

Personal



Note: You can also specify which organization to use for each API request. See [Authentication](#) to learn more.

# TOKENS



Show prices per 1K tokens

# Language models

Multiple models, each with different capabilities and price points. Prices can be viewed in units of either per 1M or 1K tokens. You can think of tokens as pieces of words, where 1,000 tokens is about 750 words. This paragraph is 35 tokens.

## GPT-4

With broad general knowledge and domain expertise, GPT-4 can follow complex instructions in natural language and solve difficult problems with accuracy.

[Learn about GPT-4](#)

Model	Input	Output
gpt-4	\$30.00 / 1M tokens	\$60.00 / 1M tokens
gpt-4-32k	\$60.00 / 1M tokens	\$120.00 / 1M tokens

## GPT-3.5 Turbo

GPT-3.5 Turbo models are capable and cost-effective.

`gpt-3.5-turbo-0125` is the flagship model of this family, supports a 16K context window and is optimized for dialog.

`gpt-3.5-turbo-instruct` is an Instruct model and only supports a 4K context window.

[Learn about GPT-3.5 Turbo ↗](#)

Model	Input	Output
gpt-3.5-turbo-0125	\$0.50 / 1M tokens	\$1.50 / 1M tokens
gpt-3.5-turbo-instruct	\$1.50 / 1M tokens	\$2.00 / 1M tokens

<https://platform.openai.com/tokenizer>

## Tokenizer

### Learn about language model tokenization

OpenAI's large language models (sometimes referred to as GPT's) process text using **tokens**, which are common sequences of characters found in a set of text. The models learn to understand the statistical relationships between these tokens, and excel at producing the next token in a sequence of tokens.

You can use the tool below to understand how a piece of text might be tokenized by a language model, and the total count of tokens in that piece of text.

It's important to note that the exact tokenization process varies between models. Newer models like GPT-3.5 and GPT-4 use a different tokenizer than previous models, and will produce different tokens for the same input text.

GPT-3.5 & GPT-4   GPT-3 (Legacy)

Hello, My name is Dan Vega, Java Champion, Spring Developer Advocate, Husband and #GirlDad based outside of Cleveland OH. I created this website as a place to document my journey as I learn new things and share them with you. I have a real passion for teaching and I hope that one of blog posts, videos or courses helps you solve a problem or learn something new.

[Clear](#)   [Show example](#)

Tokens   Characters

78   363

Hello, My name is Dan Vega, Java Champion, Spring Developer Advocate, Husband and #GirlDad based outside of Cleveland OH. I created this website as a place to document my journey as I learn new things and share them with you. I have a real passion for teaching and I hope that one of blog posts, videos or courses helps you solve a problem or learn something new.

[Text](#)   [Token IDs](#)

A helpful rule of thumb is that one token generally corresponds to ~4 characters of text for common English text. This translates to roughly  $\frac{3}{4}$  of a word (so 100 tokens  $\approx$  75 words).

**spring.application.name=he**  
**spring.ai.openai.api-key=Y**  
**spring.ai.openai.chat.opti**

## Configuration Properties

The prefix `spring.ai.openai.chat` is the property prefix that lets you configure the chat client implementation for OpenAI.

Property	Description	Default
<code>spring.ai.openai.chat.enabled</code>	Enable OpenAI chat client.	true
<code>spring.ai.openai.chat.base-url</code>	Optional overrides the <code>spring.ai.openai.base-url</code> to provide chat specific url	-
<code>spring.ai.openai.chat.api-key</code>	Optional overrides the <code>spring.ai.openai.api-key</code> to provide chat specific api-key	-
<code>spring.ai.openai.chat.options.model</code>	This is the OpenAI Chat model to use  <small>gpt-3.5-turbo (the gpt-3.5-turbo, gpt-4, and gpt-4-32k point to the latest model versions)</small>	<code>gpt-3.5-turbo</code>
<code>spring.ai.openai.chat.options.temperature</code>	The sampling temperature to use that controls the apparent creativity of generated completions. Higher values will make output more random while lower values will make results more focused and deterministic. It is not recommended to modify temperature and top_p for the same completions request as the interaction of these two settings is difficult to predict.	0.8
<code>spring.ai.openai.chat.options.frequencyPenalty</code>	Number between -2.0 and 2.0. Positive values penalize new tokens based on their existing frequency in the text so far, decreasing the model's likelihood to repeat at the same line verbatim.	0.0f
<code>spring.ai.openai.chat.options.logitBias</code>	Modify the likelihood of specified tokens appearing in the completion.	-
<code>spring.ai.openai.chat.options.maxTokens</code>	The maximum number of tokens to generate in the chat completion. The total length of input tokens and generated tokens is limited by the model's context length.	-
<code>spring.ai.openai.chat.options.n</code>	How many chat completion choices to generate for each input message. Note that you will be charged based on the number of generated tokens across all of the choices. Keep n as 1 to minimize costs.	1
<code>spring.ai.openai.chat.options.presencePenalty</code>	Number between -2.0 and 2.0. Positive values penalize new tokens based on whether they appear in the text so far, increasing the model's likelihood to talk about new topics.	-
<code>spring.ai.openai.chat.options.responseFormat</code>	An object specifying the format that the model must output. Setting to <code>{ "type": "json_object" }</code> enables JSON mode, which guarantees the message the model generates is valid JSON.	-



**DEMO - GETTING STARTED**

# PROMPTS

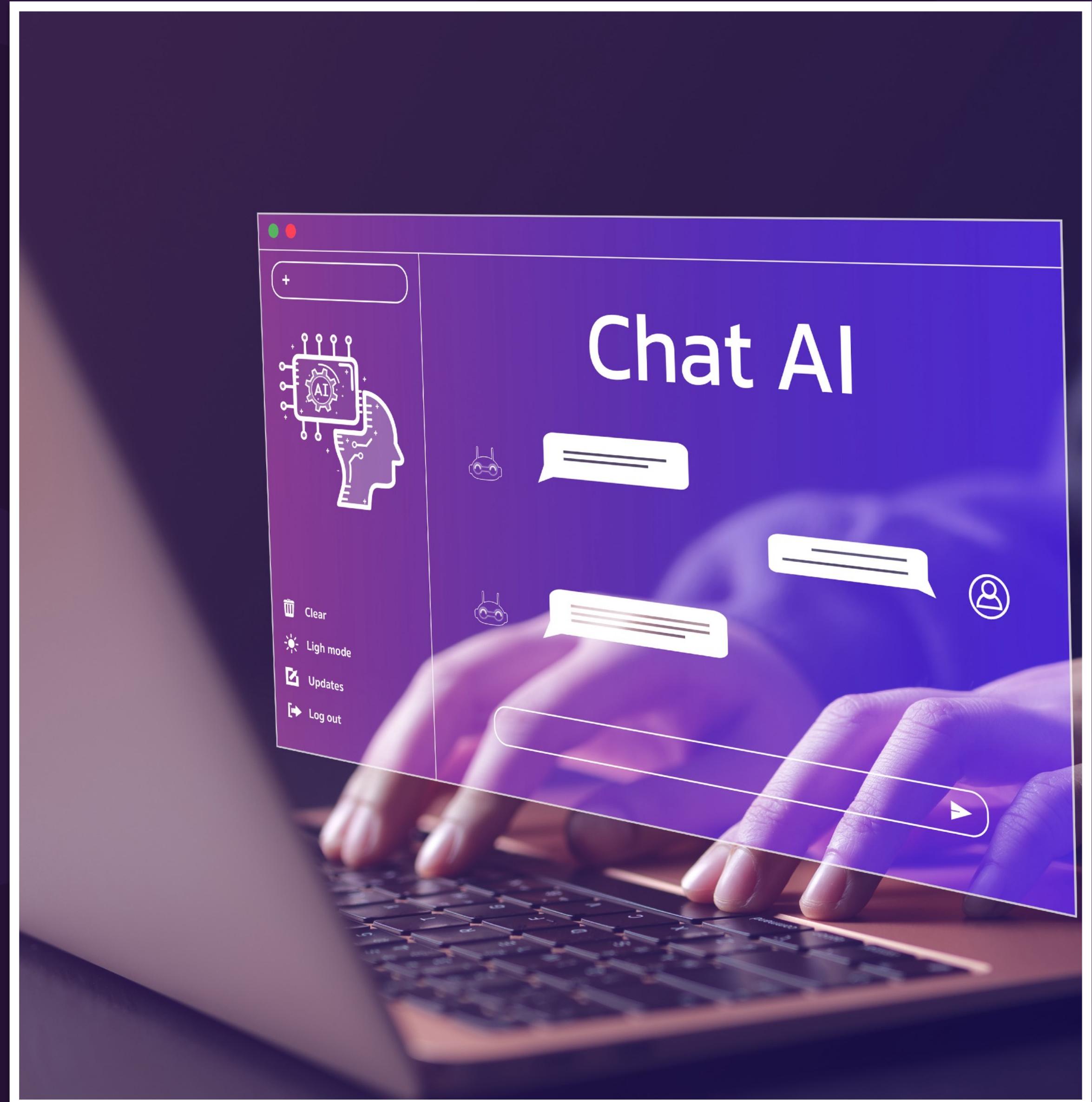


Prompts serve as the foundation for the language-based inputs that guide an AI model to produce specific outputs. For those familiar with ChatGPT, a prompt might seem like merely the text entered into a dialog box that is sent to the API. However, it encompasses much more than that. In many AI Models, the text for the prompt is not just a simple string.

# Prompt Engineering

## Effective Communication

- Input directing an AI model to produce specific outputs
- Model outputs are greatly inference by prompt style and wording
- Prompt Engineering
  - Prompt techniques and effective prompts are share in the community
  - [OpenAI Guidelines](#)
  - Mini Course: [ChatGPT Prompt Engineering for developers](#)
- Prompts and Spring
  - Prompt management relies on Text Template Engines
  - Analogous to the view in Spring MVC



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

    public Prompt(String contents) {
        this((Message)(new UserMessage(contents)));
    }

    public Prompt(Message message) {
        this(Collections.singletonList(message));
    }

    public Prompt(List<Message> messages) {
        this.messages = messages;
    }

    public Prompt(String contents, ChatOptions modelOptions) {
        this((Message)(new UserMessage(contents)), modelOptions);
    }

    public Prompt(Message message, ChatOptions modelOptions) {
        this(Collections.singletonList(message), modelOptions);
    }

    public Prompt(List<Message> messages, ChatOptions modelOptions) {
        this.messages = messages;
        this.modelOptions = modelOptions;
    }
}
```

```
public class Prompt implements ModelRequest<List<Message>> {  
    private final List<Message> messages;  
    private ChatOptions modelOptions;  
  
    public Prompt(String contents) {  
        this((Message)(new UserMessage(contents)));  
    }  
  
    public Prompt(Message message) {  
        this(Collections.singletonList(message));  
    }
```

#### Choose Implementation of Message (6 found)

- © AbstractMessage (org.springframework.ai.chat.messages) Maven: org.springframework.ai:spring-ai-core:0.8.1 (spring-ai-core-0.8.1.jar)
- © AssistantMessage (org.springframework.ai.chat.messages) Maven: org.springframework.ai:spring-ai-core:0.8.1 (spring-ai-core-0.8.1.jar)
- © ChatMessage (org.springframework.ai.chat.messages) Maven: org.springframework.ai:spring-ai-core:0.8.1 (spring-ai-core-0.8.1.jar)
- © FunctionMessage (org.springframework.ai.chat.messages) Maven: org.springframework.ai:spring-ai-core:0.8.1 (spring-ai-core-0.8.1.jar)
- © SystemMessage (org.springframework.ai.chat.messages) Maven: org.springframework.ai:spring-ai-core:0.8.1 (spring-ai-core-0.8.1.jar)
- © UserMessage (org.springframework.ai.chat.messages) Maven: org.springframework.ai:spring-ai-core:0.8.1 (spring-ai-core-0.8.1.jar)

```
}
```

```
public Prompt(Message message, ChatOptions modelOptions) {  
    this(Collections.singletonList(message), modelOptions);  
}  
  
public Prompt(List<Message> messages, ChatOptions modelOptions) {  
    this.messages = messages;  
    this.modelOptions = modelOptions;  
}
```

# ROLES

- **System Role:** Guides the AI's behavior and response style, setting parameters or rules for how the AI interprets and replies to the input. It's akin to providing instructions to the AI before initiating a conversation.
- **User Role:** Represents the user's input – their questions, commands, or statements to the AI. This role is fundamental as it forms the basis of the AI's response.
- **Assistant Role:** The AI's response to the user's input. More than just an answer or reaction, it's crucial for maintaining the flow of the conversation. By tracking the AI's previous responses (its 'Assistant Role' messages), the system ensures coherent and contextually relevant interactions.
- **Function Role:** This role deals with specific tasks or operations during the conversation. While the System Role sets the AI's overall behavior, the Function Role focuses on carrying out certain actions or commands the user asks for. It's like a special feature in the AI, used when needed to perform specific functions such as calculations, fetching data, or other tasks beyond just talking. This role allows the AI to offer practical help in addition to conversational responses.



**DEMO - PROMPTS**

# STRUCTURED OUTPUT



# **OUTPUT PARSING**

## **Challenges with handling the response**

- The Challenge
  - Output of Generative LLM is a `java.util.String`
  - Even if you ask for JSON, you get a JSON String
  - ChatGPT wants to chat, not reply in JSON
- OpenAI has introduced a new feature to help with this
- Spring AI's OutputParser uses refined prompts to desired results

# **STRUCTURED OUTPUT API**

## **Challenges with handling the response**

- The Challenge
  - Output of Generative LLM is a `java.util.String`
  - Even if you ask for JSON, you get a JSON String
  - ChatGPT wants to chat, not reply in JSON
- OpenAI has introduced a new feature to help with this
- Spring AI's OutputParser uses refined prompts to desired results

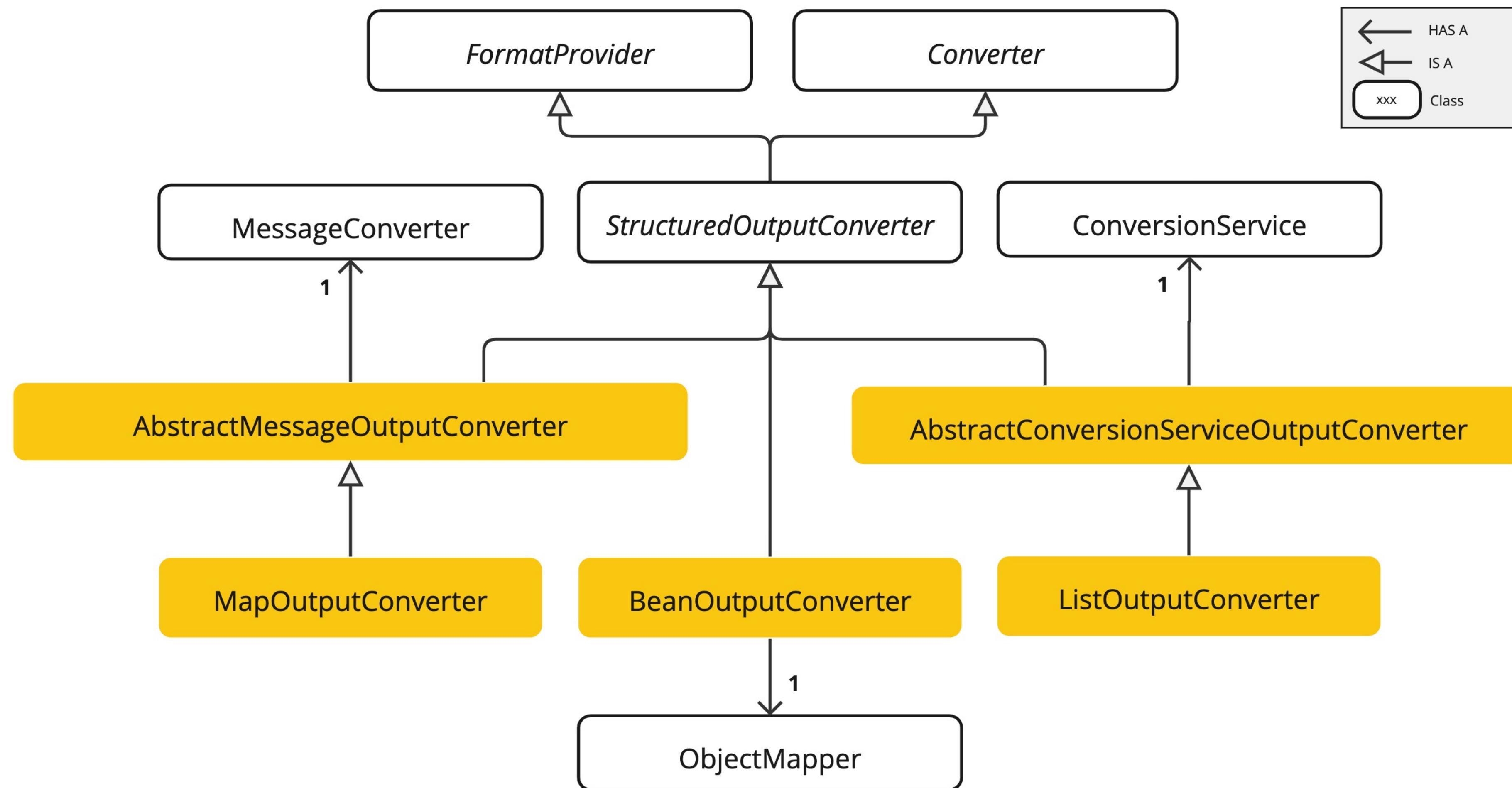
# **STRUCTURED OUTPUT API**

## **Structured Output Converter Interface**

```
public interface StructuredOutputConverter<T> extends Converter<String, T>, FormatProvider {  
  
    /**  
     * @deprecated Use the {@link #convert(Object)} instead.  
     */  
    default T parse(@NonNull String source) {  
        return this.convert(source);  
    }  
  
}  
  
public interface FormatProvider {  
  
    /**  
     * @return Returns a string containing instructions for how the output of a language  
     * generative should be formatted.  
     */  
    String getFormat();  
  
}
```

# STRUCTURED OUTPUT API

## Available Converters





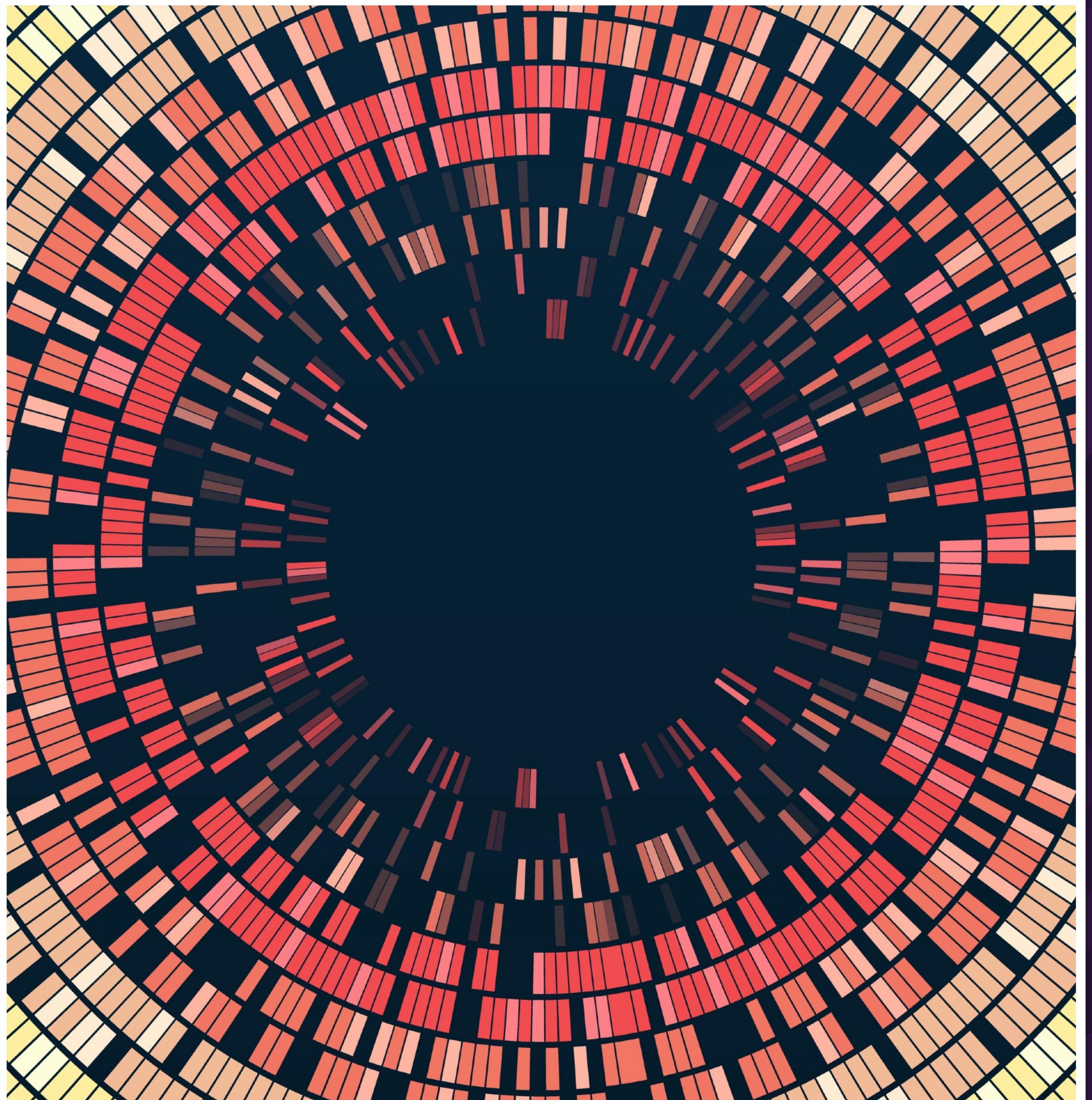
**DEMO - OUTPUT PARSER**

# **BRING YOUR OWN DATA**

# **BRING YOUR OWN DATA**

## How to use your own data in AI Applications

- AI Models have limitations
  - They are trained with public knowledge up to a certain date.
  - They don't know about your private / corporate data.
- What can we do about this problem?
  - Fine Tune the Model
  - "Stuff the prompt" - add your data into the prompt
  - Function Calling
- Retrieval Augmented Generation (RAG)
  - How to retrieve the relevant data for the user input and add it to your prompt
- There are many strategies



# STUFFING THE PROMPT



# What sports are being included in the 2024 Summer Olympics?



Use the following pieces of context to answer the question at the end. If you don't know the answer just say "I'm sorry but I don't know the answer to that".

{context} ←

Question: {question}

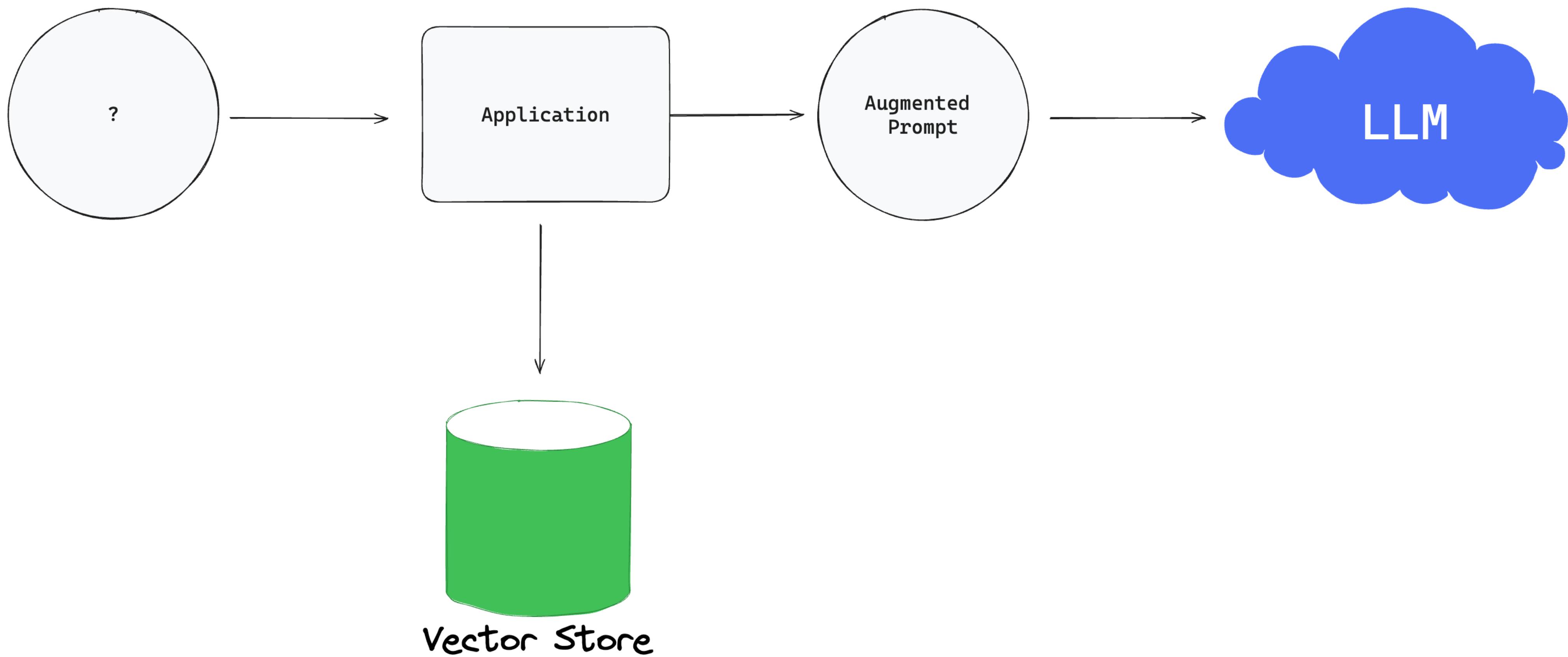
Archery, athletics, badminton, basketball , basketball 3x3, boxing, canoe slalom, canoe sprint, road cycling, cycling track, mountain biking, BMX freestyle, BMX racing, equestrian, fencing, football, golf, artistic gymnastics, rhythmic gymnastics, trampoline, handball, hockey, judo, modern pentathlon, rowing, rugby, sailing, shooting, table tennis, taekwondo, tennis, triathlon, volleyball, beach volleyball, diving, marathon swimming, artistic swimming, swimming, water polo, weightlifting,wrestling,breaking, sport climbing, skateboarding, and surfing.

I don't know the answer to that

Answers Correctly

# RAG (RETRIEVAL AUGMENTED GENERATION)

# AI APPLICATION ARCHITECTURE - RAG



# VECTOR STORES

Not just for text search

- Azure Vector Search
- ChromaVectorStore
- MilvusVectorStore
- Neo4JVectorStore
- PgVectorStore
- QdrantVectorStore
- RedisVectorStore
- WeaviateVecvtorStore
- SimpleVectorStore



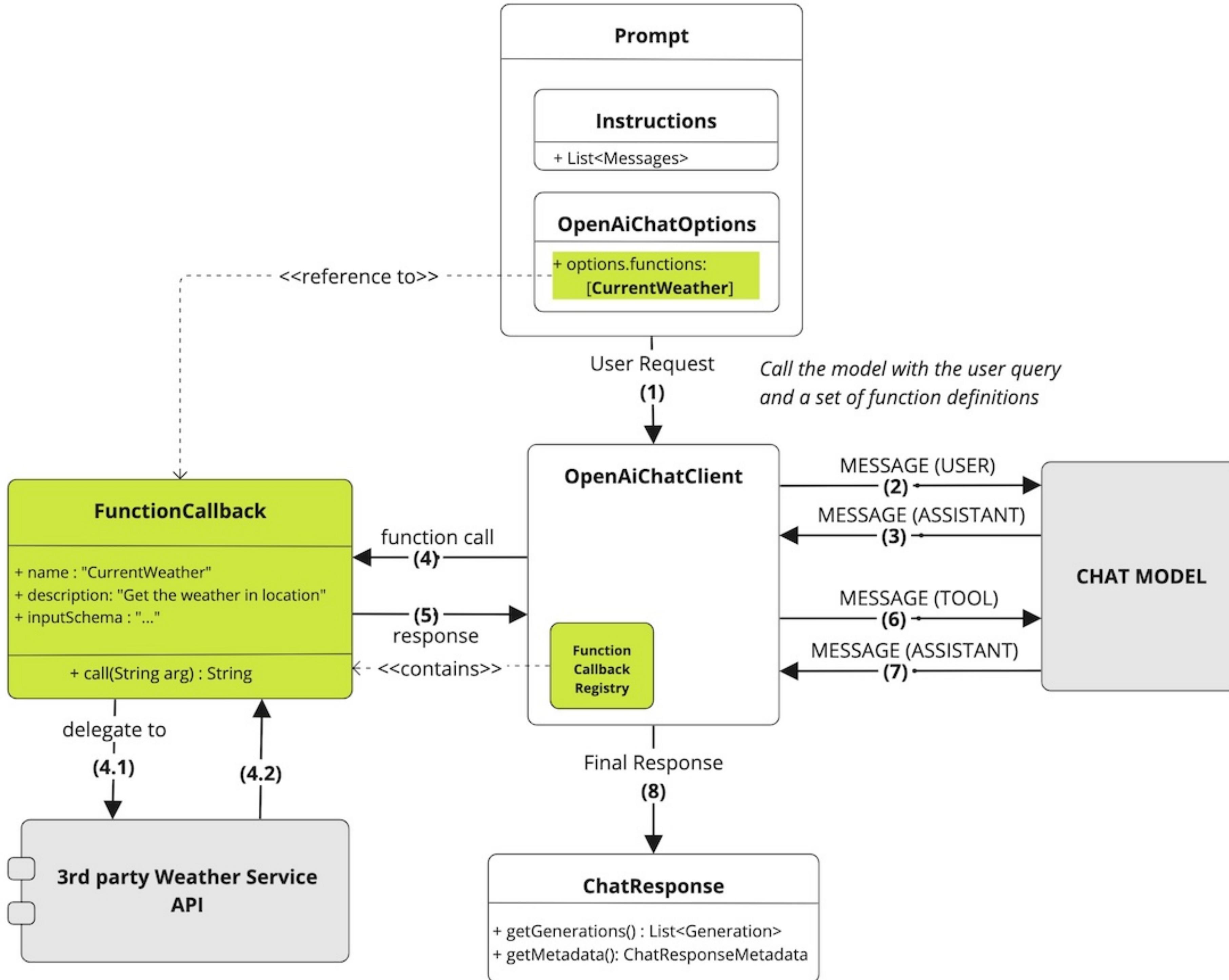
# FUNCTION CALLING



# **FUNCTION CALLING**

## Registering Custom Functions

- You can register custom Java functions with the **OpenAiChatClient** and have the OpenAI model intelligently choose to output a JSON object containing arguments to call one or many of the registered functions.
- The OpenAI API does not call the function directly; instead, the model generates JSON that you can use to call the function in your code and return the result back to the model to complete the conversation.





**DEMO - YOUR OWN DATA**

# CHAT MEMORY



# THE WEB IS STATELESS



ChatGPT 4o ▾



Hello, my name is Dan



Memory updated

Hello, Dan! How can I assist you today?

< 2/2 >

What is my name?



Your name is Dan. How can I help you today?

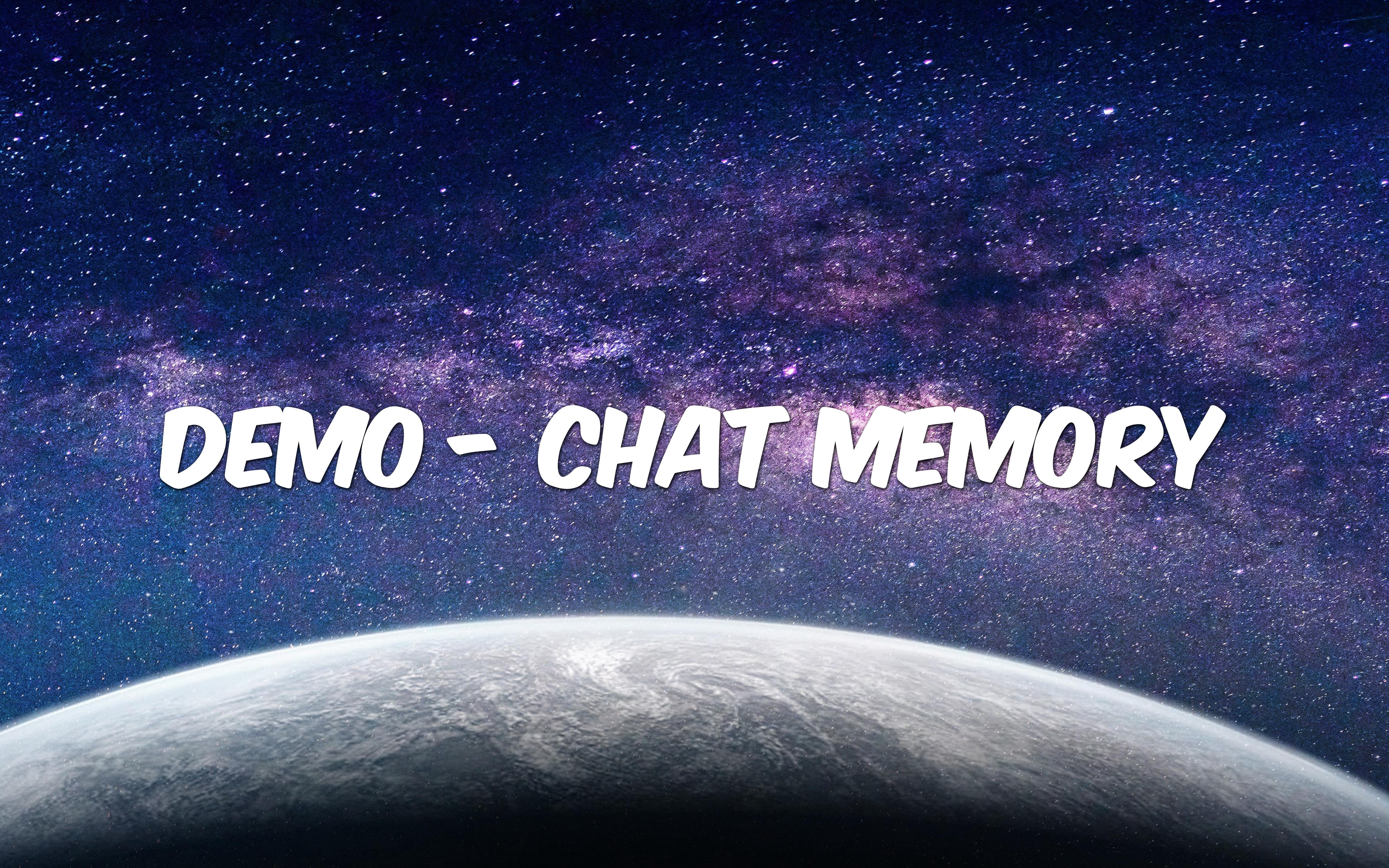
➡️ 🗑️ 🔍 ⌂ ✨

Message ChatGPT



ChatGPT can make mistakes. Check important info.





**DEMO - CHAT MEMORY**

# ★REFERENCES

## Links & Citations

- Spring AI Reference Documentation
  - <https://docs.spring.io/spring-ai/reference/>
- Google Course
  - [https://www.cloudskillsboost.google/course\\_templates/536](https://www.cloudskillsboost.google/course_templates/536)
- Spring Team
  - Mark Pollack
  - Christian Tzolov
  - Craig Walls
  - Josh Long

# THANK YOU

[dan.vega@broadcom.com](mailto:dan.vega@broadcom.com)

@therealdanvega

<https://www.danvega.dev>

