

[illegible]

SPRING AI

GENERATIVE ARTIFICIAL INTELLIGENCE CON JAVA

Simone Scannapieco

Corso avanzato per Venis S.p.A, Venezia, Italia

Novembre 2025

Note

This image shows a blank sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

This image shows a blank sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

- Simone Scannapieco

Note

This image shows a blank sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

 Simone Scannapieco
 Spring AI - Corso avanzato
 Venis S.p.A, Venezia, IT
 4 / 19

[illegible]

-
- The diagram illustrates the Spring AI Tool Dispatching Architecture, showing the flow of data and control between various components:
- Chat Request** (Left): Contains a **Tool Definition** (Green box). It lists attributes: `name`, `description`, and `input schema`.
 - Spring AI** (Center): A large container housing the **Dispatch Tool Call Requests** component (White box with yellow border).
 - Tool** (Top): A green box representing the tool being dispatched.
 - AI Model** (Bottom): A yellow box representing the AI model that processes requests.
 - Chat Response** (Right): The final output of the process.
- The flow is numbered 1 through 6:
- 1**: A solid arrow from **Tool Definition** to **Dispatch Tool Call Requests**.
 - 2**: A solid arrow from **Dispatch Tool Call Requests** to **AI Model**.
 - 3**: A solid arrow from **AI Model** to **Tool**.
 - 4**: A solid arrow from **Tool** to **Dispatch Tool Call Requests**.
 - 5**: A solid arrow from **Dispatch Tool Call Requests** to **AI Model**.
 - 6**: A solid arrow from **Dispatch Tool Call Requests** to **Chat Response**.
- Dashed arrows indicate feedback loops from the **AI Model** back to the **Tool** (via step 3) and back to the **Dispatch Tool Call Requests** (via step 5).

 Simone Scannapieco

Note

This image shows a blank sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

-
- The diagram illustrates the Spring AI Tool Dispatching Architecture, showing the flow of data and control between various components:
- Chat Request** (Yellow box) contains a **Tool Definition** (Green box). The Tool Definition includes **name**, **description**, and **input schema**.
 - Spring AI** (Large Yellow box) contains the **Dispatch Tool Call Requests** component (White box with yellow border).
 - Tool** (Green box) represents the external tool being used.
 - AI Model** (Orange box) represents the AI model that processes requests.
- The flow is numbered 1 through 6:
- 1**: A solid arrow from **Tool Definition** to **Dispatch Tool Call Requests**.
 - 2**: A solid arrow from **Dispatch Tool Call Requests** to **AI Model**.
 - 3**: A solid arrow from **AI Model** to **Tool**.
 - 4**: A solid arrow from **Tool** to **Dispatch Tool Call Requests**.
 - 5**: A solid arrow from **Dispatch Tool Call Requests** to **AI Model**.
 - 6**: A solid arrow from **Dispatch Tool Call Requests** to **Chat Response** (Yellow box).
- Dashed arrows indicate return paths from **Tool** to **Dispatch Tool Call Requests** and from **AI Model** to **Dispatch Tool Call Requests**.

 Simone Scannapieco

Note

[illegible]

-
- The diagram illustrates the Spring AI Tool Dispatching Flow, showing the interaction between a Chat Request, Tool Definition, Spring AI, Tool, AI Model, and Chat Response.
- Components:**
- Chat Request:** Contains a **Tool Definition** (green box).
 - Tool Definition:** Includes **name**, **description**, and **input schema**.
 - Spring AI:** Contains the **Dispatch Tool Call Requests** component.
 - Tool:** A green box representing the tool interface.
 - AI Model:** A yellow box representing the AI model.
 - Chat Response:** The final output of the process.
- Flow Steps:**
- 1:** Chat Request (Tool Definition) sends data to the AI Model.
 - 2:** AI Model sends data to Dispatch Tool Call Requests.
 - 3:** Dispatch Tool Call Requests sends data to Tool.
 - 4:** Tool sends data to Dispatch Tool Call Requests.
 - 5:** Dispatch Tool Call Requests sends data to AI Model.
 - 6:** Dispatch Tool Call Requests sends data to Chat Response.

 Simone Scannapieco

Spring AI - Corso avanzato

 Venis S.p.A, Venezia, IT

5 / 19

Note

This image shows a blank sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

-
- The diagram illustrates the Spring AI Tool Dispatching Architecture, showing the flow of data and control between various components:
- Chat Request** (Left): Contains a **Tool Definition** (green box) with attributes: `name`, `description`, and `input schema`.
 - Spring AI** (Center): A large container housing the **Dispatch Tool Call Requests** component (white box with yellow border).
 - Tool** (Top): A green box representing the tool implementation.
 - AI Model** (Bottom): A yellow box representing the AI model.
 - Chat Response** (Right): The final output of the process.
- The flow is numbered 1 through 6:
- 1**: A solid arrow from **Tool Definition** to **Dispatch Tool Call Requests**.
 - 2**: A solid arrow from **Dispatch Tool Call Requests** to **AI Model**.
 - 3**: A solid arrow from **AI Model** to **Tool**.
 - 4**: A solid arrow from **Tool** to **Dispatch Tool Call Requests**.
 - 5**: A solid arrow from **Dispatch Tool Call Requests** to **AI Model**.
 - 6**: A solid arrow from **Dispatch Tool Call Requests** to **Chat Response**.
- Dashed arrows indicate feedback loops from the **AI Model** back to the **Tool** (via step 3) and back to the **Dispatch Tool Call Requests** (via step 5).

 Simone Scannapieco

Spring AI - Corso avanzato

 Venis S.p.A, Venezia, IT

5 / 19

Note

This image shows a blank sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

-
- The diagram illustrates the Spring AI Tool Dispatching Architecture, showing the flow of data and control between various components:
- Chat Request** (Yellow box) contains a **Tool Definition** (Green box).
 - Tool Definition** (Green box) lists attributes: `name`, `description`, and `input schema`.
 - Spring AI** (Large Yellow box) contains the **Dispatch Tool Call Requests** (White box).
 - Tool** (Green box) is an external component.
 - AI Model** (Large Orange box) is an external component.
 - Chat Response** (Yellow box) is the final output.
- The flow is numbered 1 through 6:
- 1**: A solid arrow from **Tool Definition** to **Dispatch Tool Call Requests**.
 - 2**: A solid arrow from **AI Model** to **Dispatch Tool Call Requests**.
 - 3**: A solid arrow from **Dispatch Tool Call Requests** to **Tool**.
 - 4**: A solid arrow from **Tool** to **Dispatch Tool Call Requests**.
 - 5**: A solid arrow from **Dispatch Tool Call Requests** to **AI Model**.
 - 6**: A solid arrow from **Dispatch Tool Call Requests** to **Chat Response**.
- Feedback loops are indicated by dashed arrows:
- A dashed arrow from **Tool** back to **Tool Definition**.
 - A dashed arrow from **AI Model** back to **Tool Definition**.

 Simone Scannapieco

Note

[illegible]

-
- The diagram illustrates the Spring AI Tool Call Flow, showing the interaction between a Chat Request, Spring AI, an AI Model, and a Tool.
- Components:**
- Chat Request:** Contains a **Tool Definition** (green box) with attributes: `name`, `description`, and `input schema`.
 - Spring AI:** Contains a **Dispatch Tool Call Requests** component (white box).
 - AI Model:** (Yellow box) Receives the request and returns a response.
 - Tool:** (Green box) Receives the request and returns a response.
- Flow Steps:**
- 1:** The **Tool Definition** is passed from the **Chat Request** to the **Dispatch Tool Call Requests** component.
 - 2:** The **Dispatch Tool Call Requests** component sends the request to the **AI Model**.
 - 3:** The **AI Model** sends the response back to the **Dispatch Tool Call Requests** component.
 - 4:** The **Dispatch Tool Call Requests** component sends the request to the **Tool**.
 - 5:** The **Tool** sends the response back to the **Dispatch Tool Call Requests** component.
 - 6:** The **Dispatch Tool Call Requests** component sends the final **Chat Response**.

 Simone Scannapieco

Spring AI - Corso avanzato

 Venis S.p.A, Venezia, IT

5 / 19

Note

[illegible]

-
- The diagram illustrates the flow of a tool call request in Spring AI. It shows the interaction between a Chat Request, Spring AI, an AI Model, a Tool, and a Chat Response.
- Chat Request:** Contains a **Tool Definition** (green box) and a **returnDirect=true** flag (blue box).
 - Spring AI:** The central component that manages the tool call process.
 - AI Model:** The model that generates the initial request.
 - Tool:** The external tool that performs the requested action.
 - Chat Response:** The final output of the process.
- The flow is numbered 1 through 5:
- 1:** The Chat Request is sent to Spring AI.
 - 2:** Spring AI sends the request to the AI Model.
 - 3:** The AI Model sends the request to the Tool.
 - 4:** The Tool sends the response back to Spring AI.
 - 5:** Spring AI sends the final Chat Response.

 Simone Scannapieco

Spring AI - Corso avanzato

 Venis S.p.A, Venezia, IT

6 / 19

Note

This image shows a blank sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

- 7 / 19

Note

[illegible]

TOOL CALLING

METHOD AS TOOL

Note

This image shows a blank sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

```
@Component
public class TimeTools {

    @Tool(name="getCurrentLocalTime",
        description="Ottieni l'ora corrente nel fuso orario dell'utente.")
    String getCurrentLocalTime() {
        ...
    }

    @Tool(name="getCurrentTime",
        description="Ottieni l'ora corrente nel fuso orario specificato.",
        returnDirect=true)
    public String getCurrentTime(@ToolParam(description = "Valore che rappresenta il fuso orario.") String timeZone) {
        ...
    }
}
```

- ➔ Istanziati come `@Component`
- ➔ Ogni metodo a disposizione del LLM come *tool* deve essere annotato come `@Tool`
 - ➔ se `name` non specificato, Spring AI popola il parametro con il nome del metodo
 - ⚠ `name` come **identificativo univoco** per Spring AI
 - ➔ `description` fornisce il **contesto** al LLM per determinare se utilizzare il *tool* in base alla richiesta utente!
 - ➔ `@ToolParam` fornisce ulteriore contesto al LLM per iniettare parametri al metodo
 - ➔ `returnDirect` per sovrascrivere il comportamento di *default* se necessario

Note

[illegible]

```
@Component
public class TimeTools {

    @Tool(name="getCurrentLocalTime",
        description="Ottieni l'ora corrente nel fuso orario dell'utente.")
    String getCurrentLocalTime() {
        ...
    }

    @Tool(name="getCurrentTime",
        description="Ottieni l'ora corrente nel fuso orario specificato.",
        returnDirect=true)
    public String getCurrentTime(@ToolParam(description = "Valore che rappresenta il fuso orario.") String timeZone) {
        ...
    }
}
```

- >_ (thinking) – Scansiono i tool a disposizione... –
- >_ (thinking) – Vedo che ho dei tool che gestiscono il recupero dell'ora... –
- >_ (thinking) – Il primo tool è relativo alla posizione dell'utente... –
- >_ (thinking) – Il secondo tool è più astratto e determina la posizione attraverso parametrizzazione del fuso orario... –
- >_ (thinking) – ...Il fuso orario di New York è America/New York... –
- >_ "R: Utilizza il tool `getCurrentTime('America/New York')`."

 Simone Scannapieco

Note

[illegible]

Tool-aware ChatClient - approccio statico

```
@Bean
public ChatClient ollamaTimeToolsChatClient(OllamaChatModel ollamaChatModel, TimeTools timeTools) {

    ChatClient.Builder chatClientBuilder = ChatClient.builder(ollamaChatModel);

    return chatClientBuilder
        .defaultTools(timeTools)
        ...
        .build();

}
```

Tool-aware ChatClient - approccio dinamico

```
@Override
public Answer getOllamaTimeToolsLocalTimeAnswer(QuestionRequest request) {

    return new Answer(this.ollamaTimeToolsChatClient
        .prompt()
        .tools(timeTools)
        ...
        .call()
        .content());
}
```

Note

[illegible]

```
public class TimeTools {

    String getCurrentLocalTime() {
        ...
    }

    public String getCurrentTime(
        @ToolParam(description = "Valore che rappresenta il fuso orario.") String timeZone) {
        ...
    }

}
```

➡ **@ToolParam** permesso in approccio sia dichiarativo che programmatico

Note

This image shows a blank sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

```
import java.lang.reflect.Method;

Method method = MethodToolCallbackReflectionUtils.findMethod(TimeTools.class, "getCurrentLocalTime");
ToolCallback toolCallback = MethodToolCallback.builder()
    .toolDefinition(ToolDefinitions.builder(method)
        .name("getCurrentLocalTime")
        .description("Ottieni l'ora corrente nel fuso orario dell'utente.")
        .build())
    .toolMethod(method)
    .toolObject(new TimeTools())
    .toolMetadata(ToolMetadata.builder()
        .returnDirect(true)
        .build())
    .build();
```

- ➡ Il sistema usa *reflection* per linkare il metodo
- ➡ Costruito un `ToolCallback` che definisce la logica di utilizzo del `tool`
 - ➡ `ToolDefinitions.Builder` permette di definire nome, descrizione e schema del `tool`
 - ➡ `ToolMetadata.Builder` responsabile della definizione strategia *default* vs *direct*
- ⚠ Definizione `ToolCallback` *under-the-hood* anche con approccio dichiarativo

Note

This image shows a blank sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

```
@Bean
public ChatClient ollamaTimeToolsChatClient(OllamaChatModel ollamaChatModel) {

    ChatClient.Builder chatClientBuilder = ChatClient.builder(ollamaChatModel);

    return chatClientBuilder
        .defaultToolCallbacks(toolCallback)
        ...
        .build();
}
```

```
@Override
public Answer getOllamaTimeToolsLocalTimeAnswer(QuestionRequest request) {

    return new Answer(this.ollamaTimeToolsChatClient
        .prompt()
        .toolCallbacks(toolCallback)
        ...
        .call()
        .content());
}
```

Note

[illegible]

-  Simone Scannapieco

Note

This image shows a blank sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

TOOL CALLING

FUNCTION AS TOOL

Note

This image shows a blank sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

```
@Configuration(proxyBeanMethods = false)
class WeatherTools {

    public static final String GET_TEMP_IN_LOCATION_FUNCTION_NAME = "getTemperatureInLocation";

    @Bean(GET_TEMP_IN_LOCATION_FUNCTION_NAME)
    @Description("Ottieni la temperatura corrente nella località specificata.")
    Function<WeatherRequest, TemperatureResponse> currentTemperature() {
        return new TemperatureService();
    }
}

public enum Unit { C, F }

public record WeatherRequest(
    @ToolParam(description = "Il nome di una città o di una nazione.") String location, Unit unit) {}

public record TemperatureResponse(double temp, Unit unit) {}
```

- ⚠️ **Tipi funzionali supportati:** `Function<I, O>`, `Supplier<O>`, `Consumer<I>`, `BiFunction<I1, I2, O>`
- ➡️ Ciascun funzionale disponibile come *tool* attraverso definizione `@Bean`
 - ⚠️ Definire nome come **costante esportabile** per mancanza di garanzia *type safety*
- ➡️ `@Description` come *description* di MAT
- ➡️ `@ToolParam` con medesima logica di MAT

Note

[illegible]

```
@Bean
public ChatClient ollamaWeatherToolsChatClient(OllamaChatModel ollamaChatModel) {

    ChatClient.Builder chatClientBuilder = ChatClient.builder(ollamaChatModel);

    return chatClientBuilder
        .defaultToolNames(WeatherTools.GET_TEMP_IN_LOCATION_FUNCTION_NAME)
        ...
        .build();
}
```

```
@Override
public TemperatureResponse getOllamaTemperatureToolAnswer(@RequestBody WeatherRequest request) {

    return this.ollamaWeatherToolsChatClient
        .toolNames(WeatherTools.GET_TEMP_IN_LOCATION_FUNCTION_NAME)
        ...
        .call()
        .content();
}
```

Note

This image shows a single sheet of white paper with horizontal blue or grey ruling lines. The lines are evenly spaced and run across the width of the page. There are approximately 20 lines visible. The paper has a slight shadow on the right side, suggesting it's resting on a surface. There is no handwriting or other markings on the paper.


```
public class TemperatureService implements Function<WeatherRequest, TemperatureResponse> {

    public WeatherResponse apply(WeatherRequest request) {
        ...
    }

}

public enum Unit { C, F }

public record WeatherRequest(
    @ToolParam(description = "Il nome di una città o di una nazione.") String location, Unit unit) {}

public record TemperatureResponse(double temp, Unit unit) {}
```

- ➡ Richiesta sola la logica pertinente al funzionale scelto (es. `apply` per `Function`)

Note

This image shows a blank sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

```
public static final String GET_TEMP_IN_LOCATION_FUNCTION_NAME = "getTemperatureInLocation";

ToolCallback toolCallback = FunctionToolCallback
    .builder(GET_TEMP_IN_LOCATION_FUNCTION_NAME, new TemperatureService())
    .description("Ottieni la temperatura corrente nella località specificata.")
    .inputType(WeatherRequest.class)
    .toolMetadata(ToolMetadata.builder()
        .returnDirect(true)
        .build())
    .build();
```

- Simone Scannapieco

Note

[illegible]

```
@Bean
public ChatClient ollamaWeatherToolsChatClient(OllamaChatModel ollamaChatModel) {

    ChatClient.Builder chatClientBuilder = ChatClient.builder(ollamaChatModel);

    return chatClientBuilder
        .defaultToolCallbacks(toolCallBack)
        ...
        .build();
}
```

```
@Override
public TemperatureResponse getOllamaTemperatureToolAnswer(@RequestBody WeatherRequest request) {

    return this.ollamaWeatherToolsChatClient
        .toolCallBacks(toolCallback)
        ...
        .call()
        .content();
}
```

Note

[illegible]

- ➔ Tipi primitivi
- ➔ Optional
- ➔ Collezioni (List, Map, Array, Set)
- ➔ Tipi asincroni e reattivi

This image shows a blank sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.