



Il tuo partner per la Formazione e la Trasformazione digitale della tua azienda



Note			



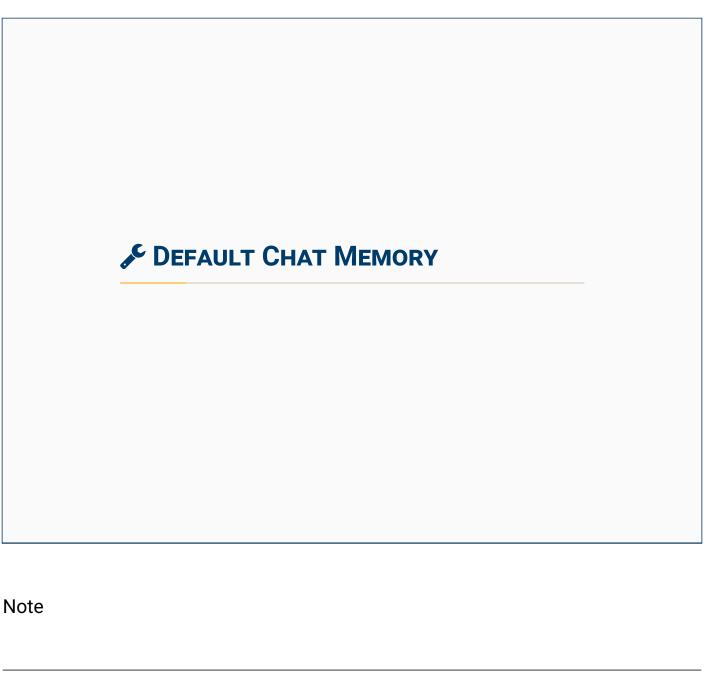
SPRING AI

GENERATIVE ARTIFICIAL INTELLIGENCE CON JAVA

Simone Scannapieco

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

Novembre 2025



Note	

SPRING AI CHAT MEMORY PRINCIPI FONDAMENTALI



- ▲ LLM sono stateless
 - Non sono concepiti per ricordare conversazioni precedenti
 - Ogni interazione con LLM avviene senza supporto di una memoria di sistema
 - A Passaggio da interazione domanda/risposta a una conversazione???
- Approccio astratto Spring Al
 - ChatMemory decide la strategia di gestione della memoria
 - Quali messaggi salvare
 - Quando rimuovere un messaggio
 - ChatMemoryRepository si occupa del CRUD dei messaggi
- Strategie di gestione memoria (WIP)
 - Salvataggio ultimi n messaggi
 - Salvataggio messaggi entro una finestra temporale
 - 3 Salvataggio entro un numero prestabilito di token

Note	

SPRING AI CHAT MEMORY/MEMORY REPOSITORY IMPLEMENTAZIONI DI DEFAULT



- ChatMemory
 - MessageWindowChatMemory implementa la strategia 1
 - Finestra di default impostata a 20
 - Auto-configurazione di Spring Al per bean ChatMemory
 - Modificabile invocando maxMessages()

Configurazione statica

MessageWindowChatMemory memory = MessageWindowChatMemory.builder()
 .maxMessages(10)
 .build();

- ChatMemoryRepository
 - ☼ InMemoryChatMemoryRepository recupera e salva messaggi in memoria
 - Utilizzo di ConcurrentHashMap per gestione multi-sessione

Simone Scannapieco

Spring AI - Corso avanzato

m Venis S.p.A, Venezia, IT

4/14

ote	

SPRING AI CHAT MEMORY ADVISORS MEMORIA ATTRAVERSO CHATCLIENT



- Come usare tutto questo con ChatClient?!
- Chat memory advisor
 - Anello di congiunzione fra
 - Sistema di gestione/CRUD memoria
 - Sistema di aggiunta di contesto al prompt
 - Implementazioni di BaseChatMemoryAdvisor
 - Precuperano lo storico della conversazione dalla memoria e la includono
 - nel prompt come lista di messaggi (MessageChatMemoryAdvisor)
 - nel prompt di sistema in formato testuale (PromptChatMemoryAdvisor)

Advisor	Formato storico	Caso utilizzo
MessageChatMemoryAdvisor	Lista di messaggi strutturati	Chatbot più performanti
${\tt PromptChatMemoryAdvisor}$	Testo puro	Limitazioni economiche

Simone Scannapieco
 Spring Al - Corso avanzato
 ✓ Venis S.p.A, Venezia, IT
 5/14

Note		

SPRING AI CHAT MEMORY ADVISORS CASO D'USO





6/14

Setup minimale

```
*\ \textit{By automatically configuring a Chat \textit{Memory bean, Spring AI instantiates a InMemory Chat Repository and a}
           * MessageWindowChatMemory under-the-hood. The only thing to decide is to use either a MessageChatMemoryAdvisor * or a PromptChatMemoryAdvisor on top of the ChatMemory bean.
          public ChatClient chatClient(ChatModel chatModel, ChatMemory chatMemory) {
              ChatClient.Builder chatClientBuilder = ChatClient.builder(chatModel);
              return chatClientBuilder
                   .defaultAdvisors(MessageChatMemoryAdvisor
                                         .builder(chatMemory)
                   .build();
          }
          /*
 * Use the ChatClient as usual.
          @PostMapping("/client/ask/memory")
          public Answer getMemoryAwareAnswer(@RequestBody Question question) {
              return new Answer(this.ChatClient
                                     .prompt()
                                     .user(question.question())
                                     .call()
                                     .content()
Simone Scannapieco
                                                Spring AI - Corso avanzato
                                                                                                      m Venis S.p.A, Venezia, IT
```


PROGETTO SPRING AI APPLICAZIONE E PASSAGGI

Simone Scannapieco



m Venis S.p.A, Venezia, IT

7/14

- Onfigurazioni bean per default chat memory per ChatClient Ollama
 - Modifica configurazioni di ChatClient per Ollama
 - Modifica interfaccia e implementazione del servizio di risposta
 - 3 Modifica controllore MVC
 - Test delle funzionalità con Postman/Insomnia

Spring AI - Corso avanzato

PROGETTO SPRING AI ADVISORS





Configurazione Gemini + Ollama - I

```
package it.venis.ai.spring.demo.config;
@Configuration
public class MemoryChatClientConfig {
     * Done under-the-hood via Spring auto-configuration.
     public ChatMemoryRepository chatMemoryRepository() {
        return new InMemoryChatMemoryRepository();
     * Custom bean for chat memory, based on the (auto-)configured chat memory repository.

* Done under-the-hood via Spring auto-configuration with maxMessages = 20.
     {\tt public \ Chat Memory \ Chat Memory \ Repository \ chat Memory \ Repository) \ \{}
         return MessageWindowChatMemory.builder()
                   .chatMemoryRepository(chatMemoryRepository)
.maxMessages(10)
                   .build();
     }
```

Simone Scannapieco

Spring AI - Corso avanzato

m Venis S.p.A, Venezia, IT

8/14

PROGETTO SPRING AI ADVISORS



Configurazione Gemini + Ollama - II

Simone Scannapieco

Spring AI - Corso avanzato

m Venis S.p.A, Venezia, IT

9/14

PROGETTO SPRING AI ADVISORS



Configurazione Gemini + Ollama - III

```
* Custom chat client based on the configured chat advisor.
@Bean
public ChatClient ollamaMemoryChatClient(OllamaChatModel ollamaChatModel,
             ChatMemory chatMemory,

@Qualifier("promptChatMemoryAdvisor") BaseChatMemoryAdvisor chatMemoryAdvisor) {
    ChatClient.Builder chatClientBuilder = ChatClient.builder(ollamaChatModel);
    {\tt return} \ {\tt chatClientBuilder}
             .defaultAdvisors(List.of(
                     new SimpleLoggerAdvisor(),
                     new OllamaCostSavingsAdvisor(),
                     chatMemoryAdvisor))
             .defaultSystem(
                          Sei un assistente AI di nome LLamaMemoryBot, addestrato per intrattenere una
                          conversazione con un umano.
             .defaultOptions(ChatOptions.builder()
    .temperature(0.1)
                      .build())
             .build();
```

Simone Scannapieco

Spring AI - Corso avanzato

m Venis S.p.A, Venezia, IT

10 / 14

PROGETTO SPRING AI SERVIZIO MULTI LLM



```
Interfaccia servizio

package it.venis.ai.spring.demo.services;
import it.venis.ai.spring.demo.model.Ansver;
import it.venis.ai.spring.demo.model.Question;
public interface QuestionService {

Answer getGeminiAnsver(Question question);

Answer getOllamaAnsver(Question question);

Answer getOllamaDefaultAnsver(Question question);

Answer getOllamaMemoryAwareAnsver(Question question);
}

Simone Scannapieco

Spring Al - Corso avanzato
```

Note	

PROGETTO SPRING AI SERVIZIO MULTI LLM





```
Implementazione servizio
package it.venis.ai.spring.demo.services;
@Service
@Configuration
public class QuestionServiceImpl implements QuestionService {
     private final ChatClient geminiChatClient;
     private final ChatClient ollamaChatClient;
     private final ChatClient ollamaMemoryChatClient;
      \begin{array}{ll} \textbf{public QuestionServiceImpl(@Qualifier("geminiChatClient") ChatClient geminiChatClient,} \\ & @Qualifier("ollamaChatClient") \ ChatClient \ ollamaChatClient, \\ \end{array} 
               {\tt QQualifier("ollamaMemoryChatClient")} \ \ {\tt ChatClient \ ollamaMemoryChatClient)} \ \ \{
          this.geminiChatClient = geminiChatClient;
this.ollamaChatClient = ollamaChatClient;
          this.ollamaMemoryChatClient = ollamaMemoryChatClient;
```

Simone Scannapieco

}

@Override

Spring AI - Corso avanzato

public Answer getOllamaMemoryAwareAnswer(Question question) { return new Answer(this.ollamaMemoryChatClient.prompt()

.user(question.question())

.content());

m Venis S.p.A, Venezia, IT

12/14

PROGETTO SPRING AI MVC DEL SERVIZIO MULTI LLM



m Venis S.p.A, Venezia, IT

13 / 14

Implementazione controllore REST

```
package it.venis.ai.spring.demo.controllers;
...

@RestController
public class QuestionController {
    private final QuestionService service;
    public QuestionController(QuestionService service) {
        this.service = service;
    }
    ...

@PostMapping("/ollama/ask/default")
    public Answer ollamaAskDefaultQuestion(@RequestBody Question question) {
        return this.service.getOllamaDefaultAnswer(question);
    }

@PostMapping("/ollama/ask/memory")
    public Answer getOllamaMemoryAwareAnswer(@RequestBody Question question) {
        return this.service.getOllamaMemoryAwareAnswer(question);
    }
}
```

Spring AI - Corso avanzato

N	ot	e
---	----	---

Simone Scannapieco



https://github.com/simonescannapieco/spring-ai-advanced-dgroove-venis-code.git

Branch: 4-spring-ai-gemini-ollama-default-chat-memory

Simone Scannapieco

Spring AI - Corso avanzato

m Venis S.p.A, Venezia, IT

14/14

Ν	ot	e
---	----	---