



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



Note			



2/18

SPRING AI

GENERATIVE ARTIFICIAL INTELLIGENCE CON JAVA

Simone Scannapieco

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

Novembre 2025

Note			



APPLICAZIONI BUSINESS PERCHÉ APPROCCIO MULTI CONFIGURAZIONE?



- Selezione modello LLM inn base alla task
 - Ragionamento complesso richiede modelli più potenti
 - Richieste semplici affidate a modelli più contenuti
 - Creazione pipeline con LLM specializzati in sotto-task
- Strategia di fallback
 - Molteplici configurazioni permettono switch automatici a modelli secondari se il primario non risponde
- Test comparativi
 - 1 In base ad accuratezza, latenza, costi, ...
- Preferenze utente
 - Approccio user-centric all'interazione con i modelli

Note		

PROGETTO SPRING AI



- Inizializzare un progetto Spring con le seguenti caratteristiche:
 - Maven come build tool
 - Spring Boot alla versione più recente non SNAPSHOT
 - Linguaggio Java 21
 - Group it.venis.ai.spring
 - Artifact demo
 - jar packaging
 - Spring Web, OpenAI o Ollama come dipendenze
- A Riportare variabili di ambiente env in launch. json e settings. json

Note	

PROGETTO SPRING AI APPLICAZIONE E PASSAGGI



- Stub di progetto Spring Al multi-configurazione (Gemini + Ollama)
 - Creazione application.yml per applicativo multi LLM
 - 2 Creazione docker-compose.yml per servizio Docker Ollama
 - 3 Creazione file variabili di ambiente per servizio Ollama
 - 4 Creazione script per start, stop ed eliminazione servizi Docker
 - Creazione configurazione multi-LLM
 - 6 Creazione modelli per domanda e risposta
 - Creazione interfaccia ed implementazione del servizio di richiesta
 - 8 Creazione del controllore MVC
 - 9 Test delle funzionalità con Postman/Insomnia

Note		

PROGETTO SPRING AI CONFIGURAZIONE APPLICATIVO



File application.yml

```
application:
    name: demo
ai:
     chat:
         client:
              enabled: false # Spring AI auto-configures a single ChatClient.Builder bean by default.
# Disabling ChatClient.Builder auto-configuration allows to manually
# configure multiple bean and inject them where needed.
    ollama:
          base-url: http://172.19.0.2:11434
             pull-model-strategy: when_missing
              max-retries: 3
          chat:
            options:
              model: VitoF/llama-3.1-8b-italian
              temperature: 0.2
              top-k: 40
              top-p: 0.9
              repeat-penalty: 1.1
              presence-penalty: 1.0
    openai:
          api-key: ${GOOGLE_AI_API_KEY}
         base-url: https://generativelanguage.googleapis.com/v1beta/openai
         chat:
              completions-path: /chat/completions
              options:
                   model: gemini-2.0-flash-lite
                   temperature: 2.0
```

Simone Scannapieco

Spring AI - Corso avanzato

m Venis S.p.A, Venezia, IT

6/18

PROGETTO SPRING AI SERVIZIO DOCKER OLLAMA - I



7/18

```
File docker-compose.yml
              spring-ai-llm-gpu:
                 image: ollama/ollama:${OLLAMA_VERSION:-latest}
                 \verb|hostname: spring-ai-llm| \\
                 container_name: spring_ai_llm
                 environment:
                   OLLAMA_HOST: "${OLLAMA_HOST:-0.0.0.0}:${OLLAMA_PORT-11434}"
                   OLLAMA_HUS1: "${ULLAMA_DEBUG:-false}
OLLAMA_DEBUG: ${OLLAMA_DEBUG:-false}
OLLAMA_FLASH_ATTENTION: ${OLLAMA_FLASH_ATTENTION:-false}
OLLAMA_KEEP_ALIVE: ${OLLAMA_KEEP_ALIVE:-"5m"}
OLLAMA_MAX_LOADED_MODELS: ${OLLAMA_MAX_LOADED_MODELS:-1}
                   OLLAMA_NUM_PARALLEL: ${OLLAMA_NUM_PARALLEL:-1}
                expose:
- ${OLLAMA_PORT:-11434}
                 deploy:
                   resources:
                     reservations:
                        devices:
                           - driver: nvidia
                             count: all
                             capabilities: [gpu]
                   - spring_ai_llm:/root/.ollama
                 {\tt restart:\ unless-stopped}
                profiles: [llm-gpu]
Simone Scannapieco
                                                                                                                       m Venis S.p.A, Venezia, IT
```

Spring AI - Corso avanzato

PROGETTO SPRING AI SERVIZIO DOCKER OLLAMA - II

Simone Scannapieco



m Venis S.p.A, Venezia, IT €

8/18

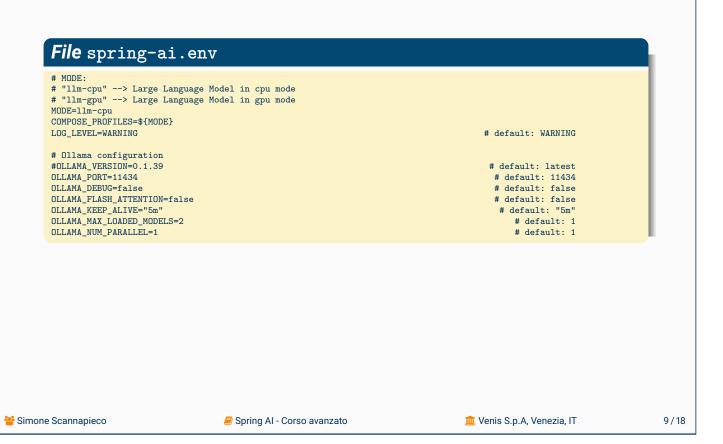
```
File docker-compose.yml
  {\tt spring-ai-llm-cpu:}
    image: ollama/ollama:${OLLAMA_VERSION:-latest}
    hostname: spring-ai-llm
    container_name: spring_ai_llm
    environment:
     OLLAMA_HOST: "${OLLAMA_HOST:-0.0.0.0}:${OLLAMA_PORT-11434}"
      OLLAMA_DEBUG: ${OLLAMA_DEBUG:-false}
      OLLAMA_FLASH_ATTENTION: ${OLLAMA_FLASH_ATTENTION:-false}
      OLLAMA_KEEP_ALIVE: ${OLLAMA_KEEP_ALIVE:-"5m"}
OLLAMA_MAX_LOADED_MODELS: ${OLLAMA_MAX_LOADED_MODELS:-1}
      OLLAMA_NUM_PARALLEL: ${OLLAMA_NUM_PARALLEL:-1}
    expose:
      - ${OLLAMA_PORT:-11434}
    volumes:
      - spring_ai_llm:/root/.ollama
    restart: unless-stopped
    profiles: [llm-cpu]
volumes:
  spring_ai_llm:
    name: spring_ai_llm
```

Spring AI - Corso avanzato

Note ______

PROGETTO SPRING AI VARIABILI DI AMBIENTE SERVIZIO DOCKER OLLAMA





Note	

PROGETTO SPRING AI SCRIPT SERVIZI DOCKER



#!/bin/bash stack=spring-ai-demo-services rmi=local file=docker-compose.yml envfile=spring-ai.env if [-z \${1+x}]; then rmi=local; else rmi=\$1; fi ## --rmi flag must be one of: all, local docker compose -f \$file -p \$stack --env-file \$envfile up --build --remove-orphans --force-recreate --detach

Simone Scannapieco

Spring AI - Corso avanzato

m Venis S.p.A, Venezia, IT

10/18

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

PROGETTO SPRING AI SCRIPT SERVIZI DOCKER



#!/bin/bash stack=spring-ai-demo-services rmi=local file=docker-compose.yml envfile=spring-ai.env if [-z \${1+x}]; then rmi=local; else rmi=\$1; fi ## --rmi flag must be one of: all, local docker compose -f \$file -p \$stack --env-file \$envfile down --rmi \$rmi

Simone Scannapieco

Spring AI - Corso avanzato

m Venis S.p.A, Venezia, IT

11/18

N	ot	e

PROGETTO SPRING AI SCRIPT SERVIZI DOCKER

Simone Scannapieco



m Venis S.p.A, Venezia, IT

12/18

File erase_spring_ai_services.sh #!/bin/bash stack=spring-ai-demo rmi=local file=docker-compose.yml envfile=spring-ai.env if [-z \{1+x\}]; then rmi=local; else rmi=\{1\}; fi ## --rmi flag must be one of: all, local docker compose -f \{file -p \{5\} stack --env-file \{5\} envfile down --rmi \{5\} rmi --volumes

ote	

Spring AI - Corso avanzato

PROGETTO SPRING AI CONFIGURAZIONE MULTI LLM



configurazione Gemini + Ollama package it.venis.ai.spring.demo.config; import org.springframework.ai.chat.client.ChatClient; import org.springframework.ai.ollama.OllamaChatModel; import org.springframework.ai.openai.OpenAiChatModel; import org.springframework.context.annotation.Bean; import org.springframework.context.annotation.Configuration; @Configuration public class ChatClientConfig { @Bean public ChatClient geminiChatClient(OpenAiChatModel geminiChatClient) { return ChatClient.create(geminiChatClient); /* * or: * ChatClient.Builder chatClientBulder = ChatClient.builder(geminiChatClient); * return chatClientBulder.build(); **/ }

Simone Scannapieco

@Rean

* or:

Spring AI - Corso avanzato

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

public ChatClient ollamaChatClient(OllamaChatModel ollamaChatModel) {

return chatClientBulder.build();

* return ChatClient.create(ollamaChatModel);

m Venis S.p.A, Venezia, IT

13 / 18

PROGETTO SPRING AI MODELLI PER DOMANDE E RISPOSTE



Modello per domanda

```
package it.venis.ai.spring.demo.model;
import java.util.UUID;
public record Question(UUID id, String question) {
    public Question(String question) {
        this(UUID.randomUUID(), question);
    }
}
```

Modello per risposta

```
package it.venis.ai.spring.demo.model;
public record Answer(String answer) {
}
```

Simone Scannapieco

Spring AI - Corso avanzato

m Venis S.p.A, Venezia, IT

14/18

PROGETTO SPRING AI SERVIZIO MULTI LLM

Note





Interfaccia servizio package it.venis.ai.spring.demo.services; import it.venis.ai.spring.demo.model.Answer; ${\tt import it.venis.ai.spring.demo.model.Question;}$ public interface QuestionService { Answer getGeminiAnswer(Question question); Answer getOllamaAnswer(Question question); 👺 Simone Scannapieco Spring AI - Corso avanzato m Venis S.p.A, Venezia, IT 15/18

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;
             this.geminiChatClient = geminiChatClient;
this.ollamaChatClient = ollamaChatClient;
             @Override
             {\tt public \ Answer \ getGeminiAnswer}({\tt Question \ question}) \ \ \{
                 return new Answer(this.geminiChatClient.prompt()
                         .user(question.question())
                         .call()
                         .content());
             }
             @Override
             public Answer getOllamaAnswer(Question question) {
                 return new Answer(this.ollamaChatClient.prompt()
                        .user(question.question())
                         .content());
Simone Scannapieco
                                            Spring AI - Corso avanzato
                                                                                             m Venis S.p.A, Venezia, IT
                                                                                                                                     16/18
```

Ν	ote
---	-----

PROGETTO SPRING AI MVC DEL SERVIZIO MULTI LLM



Implementazione controllore REST

```
package it.venis.ai.spring.demo.controllers;
{\tt import org.springframework.web.bind.annotation.PostMapping;}
import org.springframework.web.bind.annotation.RequestBody;
import org.springframework.web.bind.annotation.RestController;
import it.venis.ai.spring.demo.model.Answer;
import it.venis.ai.spring.demo.model.Question;
import it.venis.ai.spring.demo.services.QuestionService;
public class QuestionController {
    private final QuestionService service;
    public QuestionController(QuestionService service) {
        this.service = service;
    @PostMapping("/gemini/ask")
    public Answer geminiAskQuestion(@RequestBody Question question) {
        return this.service.getGeminiAnswer(question);
    @PostMapping("/ollama/ask")
    public Answer ollamaAskQuestion(@RequestBody Question question) {
        return this.service.getOllamaAnswer(question);
```

Simone Scannapieco

Spring AI - Corso avanzato

m Venis S.p.A, Venezia, IT

17 / 18



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

Branch: 1-spring-ai-gemini-ollama-configuration

🐸 Simone Scannapieco

Spring AI - Corso avanzato

m Venis S.p.A, Venezia, IT

18 / 18

Note	Ν	ot	e
------	---	----	---