



CHATBOT ASSISTANT + RAG

TECHNICAL TEST 2025

Prepared By:

NICOLAS FARIZANO

nicolasfarizano10@gmail.com

+54 9 379 493-3651
[linkedin.com/in/nicolás-farizano-8534a0233/](https://www.linkedin.com/in/nicolás-farizano-8534a0233/)

Index

Promptior Chatbot Assistant Documentation	3
Introduction	3
Project Objectives	3
Architecture Overview	3
Modular Architecture	4
Approach	4
1. Data Sources	4
2. Preprocessing	4
3. Vector Storage	4
4. Query Processing	4
5. Response Generation	5
Implementation	5
Development Process	5
1. Environment Setup	5
2. Preloading the LLaMA2 Model	5
3. Integration of RAG Architecture with LangServe	6
4. Dockerization	6
5. Automated Build with GitHub Actions	6
Deployment with Docker on Local Machines	6
Prerequisites	6
Steps to Deploy	6
Test and Verify	7
Challenges Encountered	7
Component Diagram	11
Description	11
Diagram	12
Conclusion	12
References	13

Promtior Chatbot Assistant Documentation

Introduction

Promtior's technical test evaluates a candidate's ability to design, develop, and deploy a chatbot assistant leveraging the Retrieval-Augmented Generation (RAG) architecture. This document provides a comprehensive overview of the development process for the chatbot, built with LangChain and LangServe libraries. The chatbot is enhanced using Ollama's pre-trained LLaMA2 model to ensure accurate and contextual responses.

Due to resource constraints encountered during deployment on cloud platforms such as Azure and Railway, the solution has been containerized for seamless local deployment using Docker. This approach ensures modularity, scalability, and ease of evaluation while providing a ready-to-run environment for future development or deployment efforts.

Project Objectives

1. **Develop a functional chatbot** capable of answering questions based on *Promtior's* website and PDF content, ensuring accuracy and relevance.
2. **Implement Retrieval-Augmented Generation (RAG) architecture** using *LangChain* and *Ollama* to enhance the chatbot's contextual understanding.
3. **Prepare a containerized solution** that can be deployed locally or adapted for other hosting platforms.

Architecture Overview

1. **Chatbot Service:**
 - Built with *FastAPI* and *LangChain*.
 - Handles user queries and integrates *Ollama's LLaMA2* model.
2. **Ollama Service:**
 - Provides the *LLaMA2* model for inference.
 - Preloaded in a *Docker* container for optimized performance.
3. **Local Deployment:**
 - Dockerized chatbot and Ollama services for local execution.
 - Ready-to-use Docker Compose configuration for streamlined setup.

Modular Architecture

- **main.py**: Orchestrates the application lifecycle and manages API routes using *FastAPI*.
 - **model.py**: Encapsulates all logic for query embedding, retrieval, and language model interaction.
 - **loaders.py**: Centralizes preprocessing logic for loading web and PDF data, ensuring extensibility for additional data sources.
 - **config.py**: Manages environment-specific settings to simplify deployment and configuration changes.
-

Approach

To address this challenge, we followed a modular and scalable approach:

1. Data Sources

- **Web Content**: Extracted relevant information from *Promtior's* website.
- **PDF Content**: Parsed the "AI Engineer.pdf" document provided during the technical test.

2. Preprocessing

loaders.py module:

- Extracted data from both web and PDF sources.
- Split the extracted text into manageable chunks using **RecursiveCharacterTextSplitter** for better embedding and retrieval performance.

3. Vector Storage

- Embedded text chunks using **GPT4AllEmbeddings** from the LangChain library.
- Leveraged **FAISS (Facebook AI Similarity Search)** to store these embeddings for efficient similarity searches.

4. Query Processing

model.py module:

- Handled incoming user queries by retrieving the most relevant text chunks from the *FAISS*-based vector store.
- Integrated **Ollama LLM with LLaMA2** for language model inference.
- Implemented a robust pipeline to ensure query relevance, avoiding unnecessary details.

5. Response Generation

Developed a custom **summarization method** in `model.py` module:

- Ensured responses remained concise, focused on user intent, and provided accurate information.
 - Designed contextual prompts to guide the *LLaMA2* model toward generating precise answers.
-

Implementation

Development Process

1. Environment Setup

- **Technologies Used:**
 - **Python:** Core programming language for building and orchestrating the chatbot's logic.
 - **LangChain:** Implemented the Retrieval-Augmented Generation (RAG) pipeline.
 - **LangServe:** Simplified deployment of *LangChain* pipelines by exposing them as REST APIs.
 - **Docker:** Ensured containerization for both the chatbot and the *LLaMA2* model services.
 - **Railway:** Explored as an alternative cloud platform for hosting the Ollama service due to its flexible resource options.
 - **Azure App Services:** Initial deployment platform for the cloud-based chatbot and *LLaMA2* services, later deprioritized due to resource constraints.
- **Setup Details:**
 - Defined and installed dependencies in `requirements.txt`.
 - Established a modular codebase (`main.py`, `model.py`, `loaders.py`, `config.py`) for better maintainability and scaling.
 - Updated `config.py` to handle containerized deployments and service-based URLs for inter-container communication.

2. Preloading the LLaMA2 Model

- Created a Docker image to preload *LLaMA2* into an Ollama container.
- Steps:
 1. Pulled the base Ollama image.
 2. Used `ollama pull llama2` to download and configure the model.
 3. Tagged the preloaded image for deployment.

3. Integration of RAG Architecture with LangServe

- **Integrated LangServe:**
 - Used LangServe to expose the LangChain workflow (`/invoke`) as an API.
 - Configured LangServe to process user queries by chaining FAISS-based similarity search with LLaMA2 inference.
- **LangChain's RAG Pipeline:**
 - Employed `RecursiveCharacterTextSplitter` for preprocessing.
 - Used `GPT4AllEmbeddings` for generating vector embeddings.
 - Applied custom prompt engineering to improve response relevance.

4. Dockerization

- Built Docker images for both services:
 - **Ollama:** Image preloaded with LLaMA2.
 - **Chatbot:** Handles API requests and interacts with Ollama.
- Created a `docker-compose.yml` to orchestrate multi-container deployment, enabling developers to test the chatbot and Ollama integration seamlessly.

5. Automated Build with GitHub Actions

- Configured a GitHub Actions workflow to automate the creation and publishing of Docker images.
- Workflow triggers on every commit to the repository, ensuring the latest changes are built and pushed to Docker Hub.

Deployment with Docker on Local Machines

Prerequisites

- Ensure Docker and Docker Compose are installed on your machine.
- Verify that you have sufficient memory and CPU resources to run both the chatbot and the LLaMA2 model.

Steps to Deploy

1. **Clone the Repository:** Clone the project repository to your local machine:
 - `git clone https://github.com/nicofarizano/chatbot-RAG.git`
 - `cd <repository_directory>`
2. **Start the Services:** Run the following command to start the containers using the provided `docker-compose.yml` file:
 - `docker-compose up`

3. **Access the Chatbot:** Once the containers are running, access the chatbot service at:
 - `http://localhost:11435`

Test and Verify

- Open a browser and navigate to the chatbot endpoint.
Ask sample questions, such as:
 - `What services does Promtior offer?`
 - `When was the company founded?`

Promtior Chatbot

What services does Promtior offer?

Promtior offers technological and organizational consulting services to help businesses maximize opportunities presented by transversal disruption. They provide customized GenAI solutions to optimize processes, enhance decision-making, and boost operational efficiency. Their services include process optimization, advanced personalization, continuous innovation, and GenAI product delivery. Additionally, they offer a GenAI department as a service, providing expert talent from LATAM, and a GenAI adoption consulting service to connect AI with business goals. Promtior has delivered successful cases in various industries, transforming complex data into actionable insights and automating repetitive tasks.

When was the company founded?

Promtior was founded in May 2023, shortly after the release of ChatGPT, which had a significant impact on the business world. The company helps companies navigate this period of disruption by providing technological and organizational consulting services to create new business models and stay at the forefront of their industry. With over \$50 billion in assets and 2,000 employees, Promtior

- **Shut Down the Services:** When you're done testing, stop and remove the containers using:
 - `docker-compose down`

Challenges Encountered

1. No OpenAI Subscription:

- Initially, we considered using OpenAI's API; however, the lack of a subscription required us to pivot to alternative solutions.

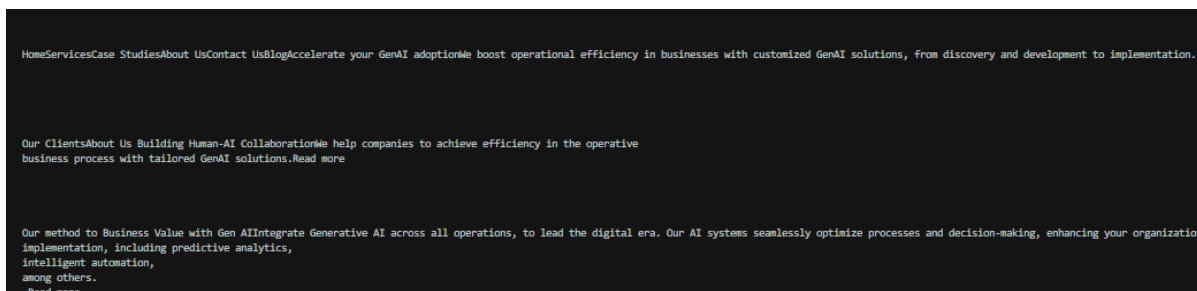
- **Solution:** Adopted Ollama and Llama2, which provided robust language model capabilities without requiring a subscription.

2. Python Version Compatibility and Dependency Conflicts:

- Several dependencies were incompatible with the latest version of Python, like `chromadb` and `sentence-transformers`, which caused significant installation issues.
- **Solution:** Instead of downgrading Python, we opted for libraries and tools that supported the latest version, ensuring long-term maintainability. Switched to `GPT4AllEmbeddings` and ensured all packages were compatible.

3. Limitation of Retrieved Documents:

- Retrieving too many documents caused irrelevant or excessive context to be included in the response.



- **Solution:** Limited the number of documents retrieved (`k`) to ensure only the most relevant information is processed.

```
75     # Recuperar los documentos más relevantes (limitar a 2)
76     docs = vectorstore.similarity_search(question, k=2)
77     print(f"\nDocumentos recuperados: {len(docs)}")
78
```

```
Escribe tu pregunta (o escribe 'salir' para terminar): What services does Promtior offer?

Documentos recuperados: 2

Respuesta generada:
In November 2022, ChatGPT was released, causing a significant impact and challenging previously unquestionable principles: Creativity is an exclusively human trait. The speed of these technological advancements has put unprecedented pressure on leaders to incorporate AI into their businesses. In May 2023, Promtior was founded facing this context, where the key question is: how to approach a scenario of transversal disruption and maximize the opportunities it presents? Through its technological and organizational consulting, Promtior offers a way to generate new business models, answering this question and bringing companies at the forefront of their sector.
About Promtior
https://promtior.ai/
Contact
Ignacio.acuna@promtior.ai
TECHNICAL TEST
```


4. Summarization of Retrieved Text:

- Initially, the chatbot returned paragraphs verbatim from the sources, which were often too verbose or irrelevant like the response generated in the previous solution.
- New solution:** Implemented a summarization method to condense the retrieved text while retaining key information.

```

90 # Resumir los documentos recuperados
91 print("\nRespuesta generada: ")
92 summarize_documents(llm, docs, question, max_words=200)
93
15
16 def summarize_documents(llm, docs, question, max_words=200): # Resume los documentos recuperados utilizando el modelo de lenguaje, enfocado en la pregunta inicial.
17     combined_text = " ".join([doc.page_content for doc in docs])
18     summary_prompt = (
19         f"A continuación, tienes un texto. Responde a la siguiente pregunta basándote en el texto. "
20         f"Limita tu respuesta a no más de {max_words} palabras.\n\n"
21         f"Pregunta: {question}\n\n"
22         f"Texto:\n{combined_text}"
23     )
24     summary = llm.invoke(summary_prompt)
25     return summary

```

Escribe tu pregunta (o escribe 'salir' para terminar): What services does Promtior offer?

Documentos recuperados: 3

Respuesta generada:

In November 2022, ChatGPT was released, revolutionizing the tech industry. Promtior was founded in May 2023 to address the challenge of integrating AI into businesses. Through technological and organizational consulting, Promtior offers customized GenAI solutions to maximize opportunities and boost operational efficiency. They help clients achieve business value through seamless optimization of processes and decision-making, enhancing organization capabilities in various areas. Their services include GenAI product delivery, department as a service, adoption consulting, and more.

5. Maintaining Focus on the Question:

- The chatbot sometimes lost focus and provided general summaries instead of directly addressing the user's question.
- Solution:** Tailored the language model's prompts to explicitly prioritize the question, ensuring answers were concise and directly relevant.

```

def summarize_documents(llm, docs, question, max_words=200): # Resume los documentos recuperados utilizando el modelo de lenguaje, enfocado en la pregunta inicial.
    combined_text = " ".join([doc.page_content for doc in docs])
    summary_prompt = (
        f"A continuación, tienes un texto. Responde a la siguiente pregunta basándote en el texto. "
        f"Prioriza información relevante y no incluyas contenido que no esté relacionado con la pregunta. "
        f"Limita tu respuesta a no más de {max_words} palabras.\n\n"
        f"Pregunta: {question}\n\n"
        f"Texto:\n{combined_text}"
    )
    summary = llm.invoke(summary_prompt)
    return summary

```

6. Misidentification of Data:

- The chatbot initially confused the founding date of Promtior (May 2023) with the release date of ChatGPT (November 2022) due to overlapping context in the retrieved documents.

```

Escribe tu pregunta (o escribe 'salir' para terminar): When was the company founded?

Documentos recuperados: 4

Respuesta generada:
According to the text, Promtior was founded in November 2022.
Escribe tu pregunta (o escribe 'salir' para terminar): When was Promtior founded?

Documentos recuperados: 4

Respuesta generada:
According to the text, Promtior was founded in May 2023.

```

7. Insufficient Resources on Azure and Railway:

-

- courageous-radianceproduction

ArchitectureObservabilityLogsSettingsTRIAL\$5.00|Share

Jan 27, 2025, 6:13 PM - Jan 27, 2025, 8

Filter logs using "I, AND, OR..."

Date (GMT+3)ServiceMessage

You reached the start of the range = Jan 27, 2025 8:11 PM

Jan 27 20:13:17cslisp-sptime=2025-01-27T21:13:16.840Z level=WFO source-server-gpu104 msg="system memory" total="104.2 GiB" free="171.8 GiB" res_usage="0 B"

Jan 27 20:13:17cslisp-spllms_model_loader: loaded meta data with 23 key-value pairs and 201 tensors from /root/.allms/models/hubs/sha256-893d646d7f6803e92207a2c26af73d7cb8de6dc7bc6bf49246 (version GGUF V3 (latest))

Jan 27 20:13:17cslisp-spllms_model_loader: warning: device memory usage may be too high. Memory usage will not scale in this context.

Jan 27 20:13:17cslisp-sptime=2025-01-27T21:13:16.842Z level=WFO source-server-gpu104 msg="default to gpu layers, required=1 layers, offload=1 layers, available="171.8 GiB" memory_gpu_overhead="0 B" memory_required_full="8.3 GiB" memory_required_partial="0 B" memory_required_kv="4.6 GiB" memory_required_allocations="(8.3 GiB)" memory_weights_total="7.4 GiB" memory_weights_repeating="(7.3 GiB)" memory_weights_nonrepeating="102.4 MiB" memory_graph_full="100.8 MiB" memory_graph_partial="681.8 MiB"

Jan 27 20:13:17cslisp-spllms_model_loader: -w 8:general.architecture str = llama

Jan 27 20:13:17cslisp-spllms_model_loader: -w 11:general.name str = LLAMA v2

Jan 27 20:13:17cslisp-spllms_model_loader: -w 5:llm.context_length int = 4096

Jan 27 20:13:17cslisp-sptime=2025-01-27T21:13:16.842Z level=WFO source-server-gpu104 msg="starting llms server" cmd="/usr/lib/llama/runners/gpu_xrt/ollama_llama_server_runner --model /root/.allms/models/hubs/sha256-893d646d7f6803e92207a2c26af73d7cb8de6dc7bc6bf49246 --ctx-size 8192 --batch-size 512 --threads 48 --no-numpy-partial4 --port 41897"

Jan 27 20:13:17cslisp-spllms_model_loader: -w 3:llms.embedding_length x32 = 4096

Jan 27 20:13:17cslisp-spllms_model_loader: -w 4:llms.block_count x32 = 32

Jan 27 20:13:17cslisp-sptime=2025-01-27T21:13:16.842Z level=WFO source-server-gpu104 msg="loaded runners' count".

Jan 27 20:13:17cslisp-spllms_model_loader: -w 5:llms.feed_forward_length x32 = 12800

Jan 27 20:13:17cslisp-sptime=2025-01-27T21:13:16.842Z level=WFO source-server-gpu104 msg="waiting for llms runner to start responding"

Jan 27 20:13:17cslisp-spllms_model_loader: -w 6:llms.rope.dimension_count x32 = 128

Jan 27 20:13:17cslisp-spllms_model_loader: -w 7:llms.attention_head_count x32 = 32

Jan 27 20:13:17cslisp-sptime=2025-01-27T21:13:16.842Z level=WFO source-server-gpu104 msg="waiting for server to become available" status="lls server error"

Jan 27 20:13:17cslisp-spllms_model_loader: -w 8:llms.attention_head_group_w x32 = 32

Jan 27 20:13:17cslisp-sptime=2025-01-27T21:13:16.852Z level=WFO source-server-gpu104 msg="starting gpus runner"

Jan 27 20:13:17cslisp-spllms_model_loader: -w 9:llms.attention_layer_norm_eps_outside F32 = 0.000019

- **Solution:** Given these constraints, I concluded that completing the challenge without investing additional resources in alternative technologies or services was not feasible. Instead, I ensured the entire project is fully containerized and prepared for local deployment using Docker. This approach maintains the modularity and scalability of the architecture, allowing future developers or teams to easily run and evaluate the solution locally or deploy it to a more suitable platform with sufficient resources.
 - The Azure App Services and Railway deployments remain accessible for review. Below are the links to the services:
 - Azure App Services - Chatbot Service:
<https://chatbot-service-faf8h3dtdna8qcey.brazilsouth-01.azurewebsites.net>
 - Azure App Services - Ollama Service:
<https://ollama-service-erehazckfbdshhgz.brazilsouth-01.azurewebsites.net>
 - Railway - Promtior Chatbot:
<https://promtior-chatbot-production.up.railway.app>
 - Railway - Ollama Preloaded:
<https://ollama-preloaded-production.up.railway.app>
-

Component Diagram

Below is a detailed description and diagram illustrating the interactions between components in the solution:

Description

User Query:

The user sends a query to the chatbot API endpoint via HTTP POST requests.

Query Processing:

- The query is processed by the chatbot service.
- The FAISS vector store retrieves the most relevant text chunks based on query similarity.
- The retrieved chunks are passed to the LLaMA2 language model for further processing.

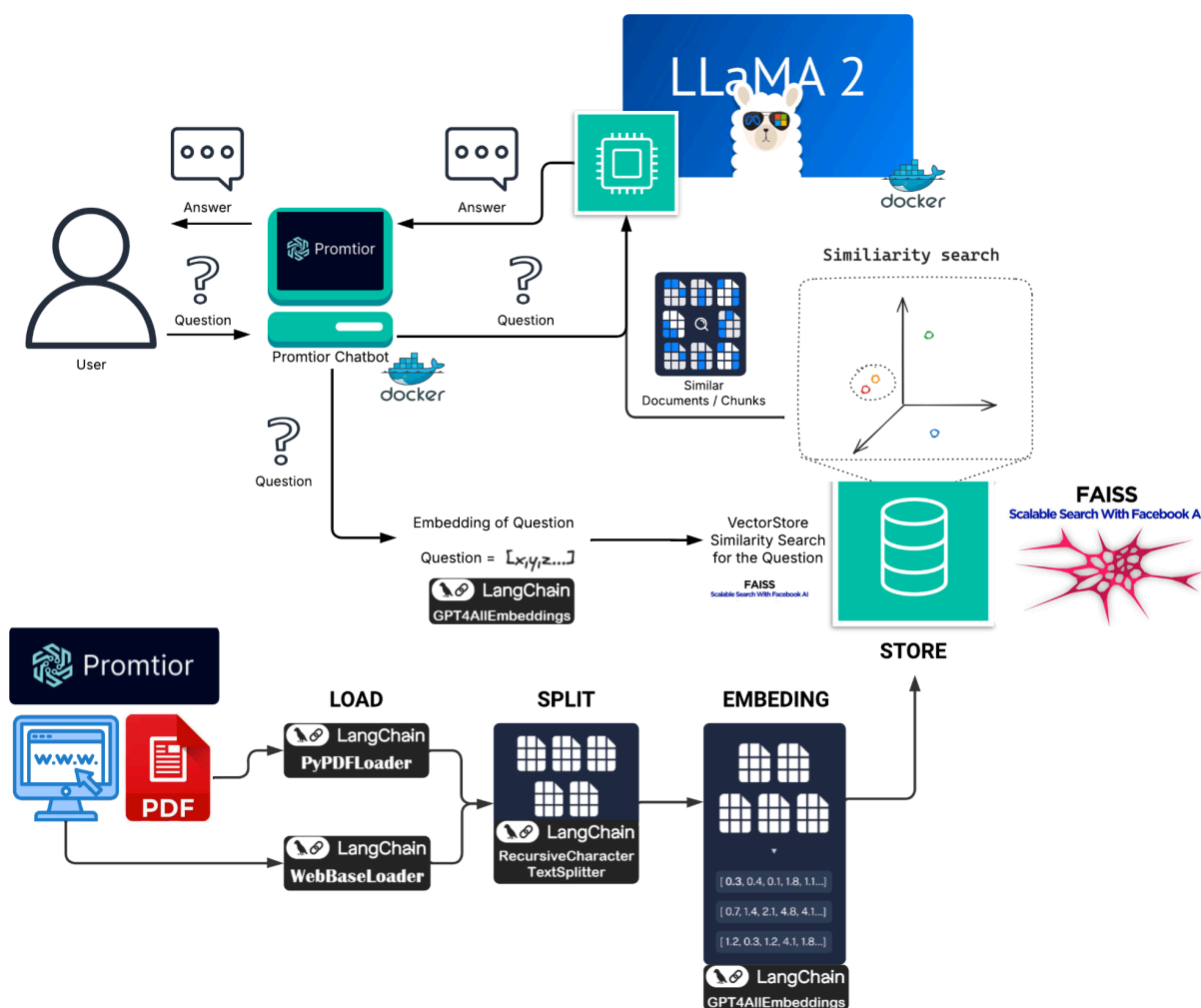
Response Generation:

The language model uses the query and retrieved context to generate a concise and accurate response.

Response Delivery:

The chatbot service formats the generated response and sends it back to the user via the API.

Diagram



Conclusion

This project showcases the development of a modular, containerized chatbot using the Retrieval-Augmented Generation (RAG) architecture. By leveraging Docker, we ensured that the solution can be deployed easily on local machines or cloud platforms. While Azure and Railway were explored for cloud deployment, technical limitations and resource constraints made local Docker deployment the most practical option for this project.

Despite not achieving deployment on Azure or Railway, the process of seeking alternatives allowed me to deepen and refine my knowledge of containerization, cloud infrastructure, and deployment strategies. These learnings have enhanced my technical expertise and will undoubtedly contribute to the success of future projects.

References

- [LangChain Documentation: For implementing the RAG pipeline and vector store integration.](#)
- [FAISS Documentation: For efficient similarity searches.](#)
- [Ollama Documentation: For leveraging the LLaMA2 language model.](#)
- [Docker Documentation: For containerizing and orchestrating services.](#)