



2023 - 2024

GRADUATION PROJECT

NATIONAL ENGINEERING DEGREE

SPECIALTY : INFORMATION TECHNOLOGY

TITLE : AI Search Engine

By: *Mohamed Amine TAIEB*

Academic supervisor: *Chifaa GHANMI*

Corporate Internship Supervisor: *Hazem RAOUAFI*



Abstract

This report presents the development of a retrieval-augmented generation system using LangChain, Large Language Models and a vector store. The goal is to enhance enterprise information retrieval by generating more accurate and relevant responses to user questions in an enterprise environment.

With the recent advancements in Generative Artificial Intelligence and the emergence of open-source Large Language Models, a Generative AI-powered chatbot became an interesting approach for question-answering tasks and faster access to relevant information. Because Large Language Models would require intensive computing power and time to train (in addition to a highly-competitive, multi-disciplinary large teams), using already pre-trained models is the go-to solution for most enterprises. This comes with the drawback of LLMs being confined to the data they were trained on and the general-purpose aspect of their nature. Thus, the need to further develop and improve their functionality with retrieval-augmented generation and other techniques. These would allow to augment the data accessible by large language models and customize how to retrieve it from specific databases (vector stores).

Keywords: *Generative AI, Large Language Models (LLMs), Retrieval-augmented Generation (RAG), LangChain, Prompt Engineering, Embeddings, Vector Store*

Acknowledgements

I would like to thank...

Contents

Acknowledgements	1
Contents	2
List of Figures	4
List of Tables	7
Keywords and Abbreviations	8
1 Introduction	10
1.1 Background	11
1.2 Hosting Company	12
1.3 Problem Statement	13
1.4 Objectives	15
1.5 Significance of the study	16
2 Literature Review	17
2.1 Text Generative AI	18
2.2 Retrieval-augmented Generation	20
2.2.1 Retrieval-augmented Generation paradigms	21
2.2.2 Significance of Retrieval-augmented Generation	24
2.3 Components of a typical RAG pipeline	27
2.3.1 Vector Stores	28
2.3.2 Large Language Model	31
2.3.3 Prompts and Prompt Engineering	33
3 Methodology	34
3.1 System Architecture	34
3.2 LangChain Overview	34
3.3 Vector Store Implementation	34

4	Implementation	35
4.1	Workstation	36
4.1.1	Hardware	36
4.1.2	OS and Software	36
4.2	Components	37
4.2.1	Libraries and Frameworks	37
4.2.2	LLMs	40
4.2.3	Vector Stores	45
4.2.4	Embedding Model	46
4.3	Incorporation into a functional RAG system	49
4.3.1	Data Ingestion Methods	50
4.3.2	Vector Store	51
4.3.3	LLMs and Prompts	53
4.4	Testing and Validation	57
5	Conclusion and Future Work	58
5.1	Summary of Findings	58
5.2	Future Research Directions	58
	References	59

List of Figures

1.1	Standard Sharing Software logo	12
2.1	Relationship between AI, ML, DL, NLP, and Conversational AI terms. [miarec]	18
2.2	A representative instance of the RAG process applied to question answering. It mainly consists of 3 steps. 1. Indexing. Documents are split into chunks, encoded into vectors, and stored in a vector database. 2. Retrieval. Retrieve the Top k chunks most relevant to the question based on semantic similarity. 3. Generation. Input the original question and the retrieved chunks together into LLM to generate the final answer. ([Retrieval-Augmented Generation for Large Language Models: A Survey] , March 2024)	20
2.3	Comparison between the three paradigms of RAG. (Left) Naive RAG mainly consists of three parts: indexing, retrieval and generation. (Middle) Advanced RAG proposes multiple optimization strategies around pre-retrieval and post-retrieval, with a process similar to the Naive RAG, still following a chain-like structure. (Right) Modular RAG inherits and develops from the previous paradigm, showcasing greater flexibility overall. This is evident in the introduction of multiple specific functional modules and the replacement of existing modules. The overall process is not limited to sequential retrieval and generation; it includes methods such as iterative and adaptive retrieval ([Retrieval-Augmented Generation for Large Language Models: A Survey] , March 2024)	22

2.4	RAG compared with other model optimization methods in the aspects of “External Knowledge Required” and “Model Adaption Required”. Prompt Engineering requires low modifications to the model and external knowledge, focusing on harnessing the capabilities of LLMs themselves. Fine-tuning, on the other hand, involves further training the model. In the early stages of RAG (Naive RAG), there is a low demand for model modifications. As research progresses, Modular RAG has become more integrated with fine-tuning techniques. ([Retrieval-Augmented Generation for Large Language Models: A Survey] , March 2024)	25
2.5	The workflow of the RLHF algorithm. ([A Survey of Large Language Models] , November 2023)	26
2.6	Vector Store Process Diagram. [LangChain Documentation] . . .	29
2.7	Embeddings Generation Diagram. (Microsoft, 2023)	29
2.8	Similarity Metrics for Vector Search. [Zilliz Blog]	30
2.9	An illustration of main components of the transformer model from the original paper, where layer normalization was performed after multiheaded attention. [On Layer Normalization in the Transformer Architecture (2020)]	31
2.10	Types of Prompts: [X]: context, [Y]: Abbreviation, [Z]: Expanded Form ([A Study on the Implementation of Generative AI Services Using an Enterprise Data-Based LLM Application Architecture] , 2023)	33
4.1	Streamlit logo.	37
4.2	Streamlit App Gallery.	38
4.3	LangChain logo	38
4.4	LangChain Components. (LangChain Documentation)	39
4.5	Hugging Face logo. [Hugging Face models Hub]	40
4.6	A short list of popular text-generation models available on HF Hub. HF models filtered by ‘text-generation’ task	40
4.7	An illustration of different files of a locally downloaded model (Phi-3-mini-4k-instruct) from Hugging Face Hub	41
4.8	TogetherAI playground web interface, an example of a cloud-based environment allowing multiple LLM’s customization options and parameters. TogetherAI Playground service	42
4.9	Artificial Analysis LLM Performance Leaderboard [ArtificialAnalysis/LLM-Performance-Leaderboard space on Hugging Face]	43
4.10	FAISS logo.	46
4.11	A short list of popular embedding models available on HF Hub. HF models filtered by ‘sentence-similarity’ task	47

4.12	Sentence-Transformers model performance comparison [Pre-Trained Sentence Transformers models' performance]	48
4.13	The overall pipeline of the system to be implemented, showcasing the interaction between the different components, from knowledge augmentation methods (document/data loading and chunking etc...), vector database (index), RAG pipeline (User-Index-LLM interaction), re-ranking and evaluation processes	49
4.14	Various data ingestion methods	50
4.15	Vector Store implementation classes and methods	52
4.16	Local Directory for storing vector stores.	52
4.17	Embedding Model Implementation	53
4.18	Implementation of LLMs and Prompts, focusing on extensibility .	54
4.19	Incorporating prompts into model definition.	55
4.20	LangChain Hub web interface.	55
4.21	List of implemented LLMs	56

List of Tables

1.1	Comparison of existing solutions	13
4.1	Comparison of popular vector databases	45

Keywords and Abbreviations

- RAG : Retrieval-Augmented Generation
- LLM : Large Language Model
- LL : Large Language
- FM : Foundational Model
- AI : Artificial Intelligence
- GAI : Generative Artificial Intelligence
- PE : Prompt Engineering
- XoT : Everything of Thought
- CoT : Chain of Thought
- ToT : Tree of Thoughts
- RL : Reinforcement Learning
- RLHF : Reinforcement Learning from Human Feedback
- NLP : Natural Language Processing
- ML : Machine Learning
- DL : Deep Learning
- GPT : Generative Pre-trained Transformer
- URL : Uniform Resource Locator
- API : Application Programming Interface
- HF : Hugging Face

- SBERT : Sentence Bidirectional Encoder Representations from Transformers (or Sentence Transformers for short)
- MTEB : Massive Text Embedding Benchmark
- OS : Operating System
- CPU : Central Processing Unit
- GPU : Graphics Processing Unit
- RAM : Random-access memory
- CUDA : Compute Unified Device Architecture
- VM : Virtual Machine
- WSL : Windows Subsystem for Linux
- JS : JavaScript
- FOSS : Free and Open-Source Software
- a.k.a. : also known as
- e.g. : for example
- i.e. : that is
- TTM : Time To Market
- 3S : Standard Sharing Software (The hosting company)

Chapter 1

Introduction

This chapter is a prologue to the conceptualization of this end-of-studies project. It lays out the context and influences which made it relevant in current settings. First, it describes the background and sequence of developments that brought much of the employed technologies into existence. Next, it introduces the hosting company and discusses its necessity for a such solution. And finally, it walks through the pursued objectives and innovations which address the insufficiencies of existing solutions.

1.1 Background

Business Knowledge is both integral and proprietary for any enterprise. The contents of private documents are useful for internal employees who are constantly consuming it to accomplish their workstreams. It becomes evident when considering that modern software development and network infrastructure deployment (among many other fields) are often based on exhaustive documentation and lengthy research papers. Paradoxically, this upfront research and learning introduce a sort of bottleneck requiring much time before a new project is initiated. In most cases, a further looking up for previously reviewed information is often required. Accelerating delivery, however, remains a key objective for organizations seeking a competitive edge. This together emphasizes the need to consider and propose innovative solutions which would help reduce the preliminary requirements and achieve a faster time-to-market (TTM).

In line with this discussion, it is worth highlighting the evolving landscape of web search in the contemporary era. We used to input a query into a search engine (e.g. Google) and it would look through its indexed webpages and then return a list of the most relevant pages which we would read until we have a satisfying answer. This paradigm has veered towards generative AI chatbot solutions, like ChatGPT, Copilot and Gemini, which leverage Machine Learning algorithms to mimic the natural language understanding capability of humans to generate short and accurate answers.

Auspiciously, the underlying technology which empower this kind of chatbots is discrete from their corresponding platforms, i.e., it can be integrated in other projects as a library or software component rather than being exclusive to their native applications. This presents an interesting topic for an internship project which will both provide a much needed solution for an enterprise with many activities and projects to accomplish and for a data science student aspiring to continuously learn the newest trends in AI and Machine Learning.

1.2 Hosting Company

The hosting of this project was managed by Standard Sharing Software (3S).



Figure 1.1: Standard Sharing Software logo

3S is a Tunisian IT services company specializing in the integration of IT infrastructures. Since its foundation in 1988, it has supported its clients in their digital transformation, management and improvement of their IT infrastructure, from the study to the deployment of the most innovative and high value-added technological solutions. It is based (headquarters) in Montplaisir, Tunis, Tunisia, and in Charguia 1, Tunis (the office where this internship was conducted), with different teams specializing in Network and Telecommunication infrastructures, Cyber security, and Cloud Computing. This project was organized in AI and development department, but close collaboration with these different units was required since they constitute the target users of this project.

1.3 Problem Statement

Traditional enterprise search engines often have many limitations due to the complexity of user query context comprehension and human language understanding. This is true because traditional methods of content searching is based on exact text matching without semantic-based retrieval. Even with the emergence of AI tools like ChatGPT, which address Natural Language Processing/Analysis, access to large enterprise data is still a major hurdle for such systems. We have witnessed in recent years the amount of pressure governments have put on AI products in order to evade potential privacy and copyright infringements. This is good overall as such regulations ensure fair competition and publishers' rights. However, this also means that companies need to develop their own solutions to effectively leverage the recent breakthroughs of this field.

In the table below, an identification and comparison between a few existing solutions has been undertaken to understand their limitations and how to improve on them.

Solution	LLM Chatbots	NotebookLM	Verba
Developer	Google, OpenAI and Microsoft	Google	Weaviate
RAG-based (persistent storage)	✗	✓	✓ / ✗
Enterprise focused	✗	✗	✓ / ✗
Extensible	✓ / ✗	✗	✓ / ✗
Multiple LLMs	✗	✗	✓
User-LLM feedback	✓	✗ / ✓	✗
Stable (Not experimental)	✓	✗	✗

Table 1.1: Comparison of existing solutions

In general, one can use any ready-to-use LLM-powered chatbot (e.g. ChatGPT, Microsoft Copilot, Google Gemini, and many others), but these solutions, even though allowing document-answering, don't come with document persistence,

which means by leaving the chat, the documents are gone unless re-uploaded. Google is experimenting with a new solution (NotebookLM) that allows that by connecting to cloud-hosted documents, in addition to methods allowing to upload local files and webpage content from URLs, but this comes short for enterprise use as it is intended for casual use cases like note taking, document editing suggestions and personal project research. On the other hand, Weaviate, a company behind a vector store development, maintains [a repository on GitHub](#) which addresses the same objectives as this project. Their solution, called Verba, is a well-organized open-source project that can be easily deployed on-premises or in a cloud environment. However, several challenges have to be faced before a company can adapt it to their specific needs: customization is difficult given that the project uses many dependencies from packages they developed (goldenverba, weaviate), requiring many API keys from various providers (vector database, LLMs, cloud providers, etc...), in addition to the inability to create and manage different databases for different teams.

This scarcity and inadequacy of available projects raises the need to design and develop a custom solution suitable for the needs and requirements set up by individual enterprises (which is articulated in the next section).

1.4 Objectives

This project aims to address general-purpose large language models' limitations (data confinement to training phase and private data access) by developing a novel search functionality that leverages retrieval-augmented generation to deliver insightful results for enterprise data. The system will address the following objectives:

- **Enhance AI-generated responses with External Knowledge:** Retrieve relevant passages from a curated enterprise knowledge base each time an LLM is prompted to improve its factual grounding and reduce hallucination.
- **Flexible Data Ingestion and Organization:** Implement diverse methods for multi-source, offline and online, format-agnostic file content and document indexing into manageable, domain-specific vector databases to tailor to enterprise needs and potential evolution.
- **Diverse Answering Options:** Allow for the selection of Large Language Models for answering and just-in-time data scraping and fetching to diversify and improve results each time.
- **User-Driven Feedback:** Integrate a mechanism for providing feedback by users for LLM-generated responses which can be used to rank future results and further train and improve LLMs (when supported).
- **Leverage RAG-based metrics for evaluation:** Employ retrieval-augmented generation evaluation metrics to assess the quality of pipeline stages (retrieval quality/accuracy, context-generation consistency, answer relevancy)

1.5 Significance of the study

The proposed Generative AI solution would help employees to find answers based on enterprise knowledge without having to browse vast amount of documents. This is enabled by LLM knowledge augmentation using techniques of retrieval-augmented generation and prompt engineering for context precision which would also reduce LLM hallucination. It would also allow them to interfere in the steps of answer-generation by providing ways to customize the external knowledge base and conduct online researching, powered by traditional search engines and generative AI researching tools, through a unified interface.

Chapter 2

Literature Review

For the purpose of this project, an extensive investigation and comparative study into recently released Large Language Models has been conducted. This chapter delves into a comprehensive exploration of Text-Generative AI in general: it examines the concept and diverse applications of generative AI and a Retrieval-augmented Generation pipeline system, its elements and components alongside a comparison between their technical aspects, and its eventual limitations.

2.1 Text Generative AI

The term Text Generative AI refers to the artificial-generation of textual content, it can allude to the tasks of next word suggestion, summarization, rewriting in different tones, cross-language translation, question answering, text or code generation etc... It is a type of Artificial Intelligence that utilizes models trained on massive natural language data (Large Language Models). Prominent examples of such models include GPT-3, Gemini, Llama. Even though these models are multi-modal (not only text-generative AI, can also analyse media content...), they are in essence language models whose functionality provide a good foundation for natural language processing and generation tasks.

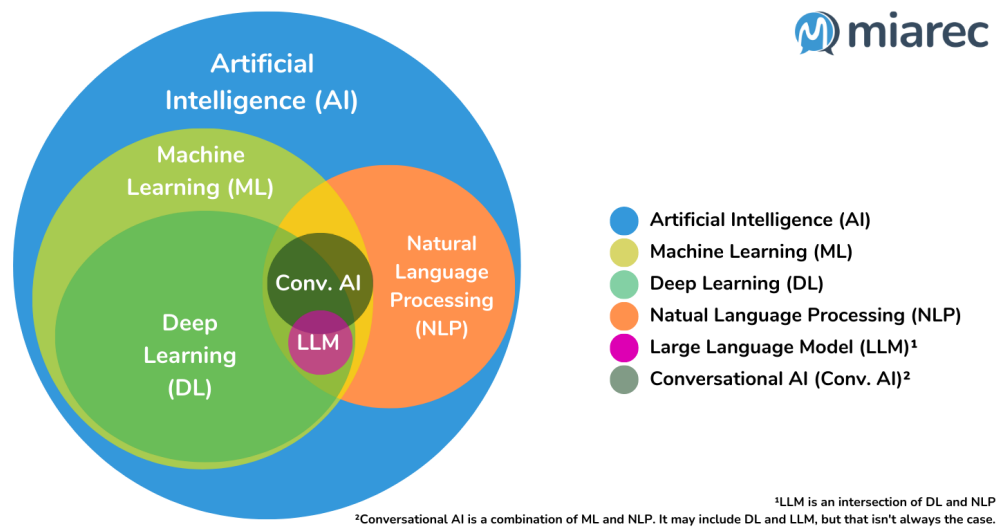


Figure 2.1: Relationship between AI, ML, DL, NLP, and Conversational AI terms. [miarec]

Recent trends in Large Language Models have shown how well they can perform as searching tools or assisting agents. They mimic humanistic thinking and actions by providing chat-like responses to messages. This technology (LLM) can achieve this by its ability to understand the context and flow of conversations; It can recognize whether the prompt is a question, a request to perform an action, or casual natural language tasks like translation and other text modification. Many cloud-based solutions and open-source initiatives have been materialized which has unlocked many possibilities for these models by customizing how they function and augmenting their knowledge.

These customizations can be summarized in two main points:

- **Prompt Engineering:** This represents how the model is prompted, which has a high impact on the quality of generated responses. There are many techniques related to this; we can precise the context and possibly the history of a chat thread on which the LLM base their answers, specify the tone or length of a response, instruct the model on how it should behave, or provide it with examples on how questions should be answered.
- **Reinforcement learning from human feedback:** For open-source Language Models, one can further train and fine-tune the model. This technique refers to the training of AI model by providing it with direct feedback from humans.

2.2 Retrieval-augmented Generation

RAG refers to the process of augmenting LLM knowledge by automatically fetching and providing relevant context for the question. It is a technique used to address hallucination and data confinement of generative AI models. This is achieved without re-training the models, which provides a cost-effective approach to improving LLM output so it remains useful and accurate in various scenarios.

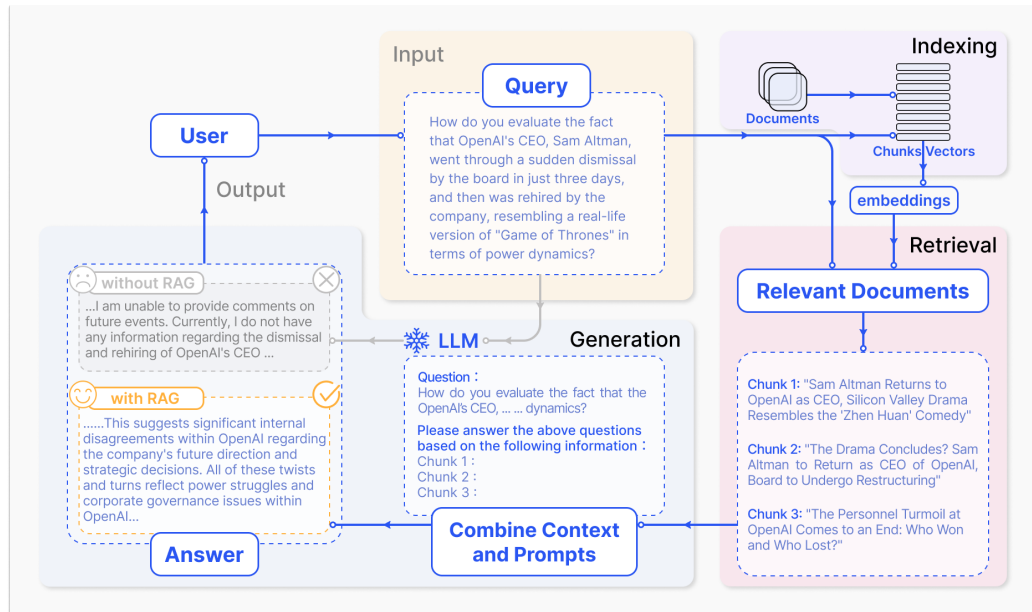


Figure 2.2: A representative instance of the RAG process applied to question answering. It mainly consists of 3 steps. 1. Indexing. Documents are split into chunks, encoded into vectors, and stored in a vector database. 2. Retrieval. Retrieve the Top k chunks most relevant to the question based on semantic similarity. 3. Generation. Input the original question and the retrieved chunks together into LLM to generate the final answer. ([Retrieval-Augmented Generation for Large Language Models: A Survey], March 2024)

2.2.1 Retrieval-augmented Generation paradigms

Retrieval-augmented Generation process is constantly evolving since the emergence of transformer-based LLMs (GPT models). It has significantly improved on the limitations of LLMs by developing their performance while being cost-effective. Even so, a simple RAG pipeline still exhibit various limitations such as ‘retrieval-echoing generation’ which means outputting the retrieved documents rather than synthesizing answers based on these (which adds more insights to generated text). Moreover, this naive RAG paradigm still introduce hallucinations into answers by confining the generative process to the retrieved documents, which in many cases may not be pertinent to the query (such as in case relevant information not present in the database). In this respect, as more research and engineering is continually being conducted, advanced spinoffs of RAG has emerged to reduce the mis-reckoning and limitations of naive RAG.

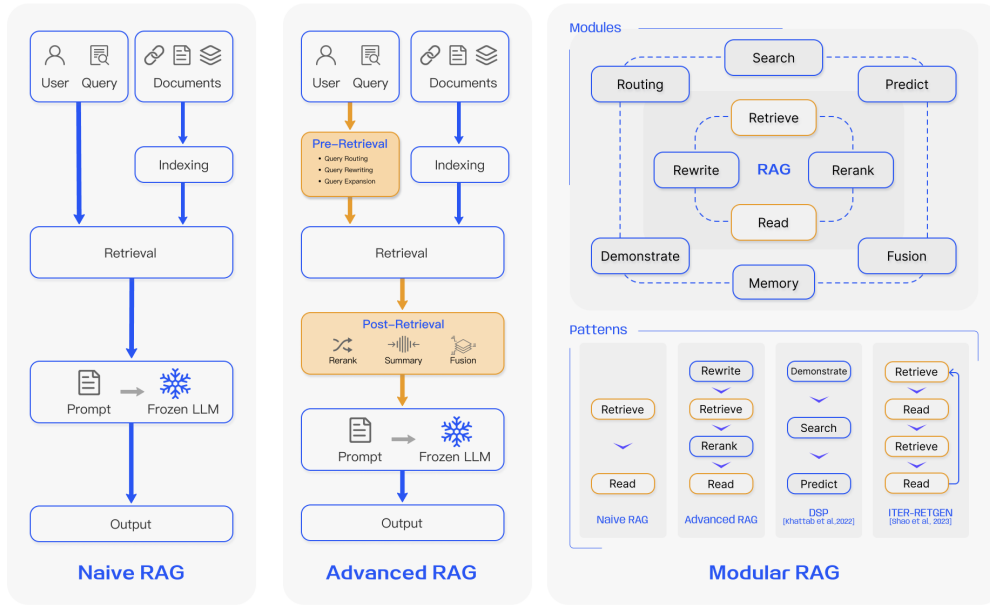


Figure 2.3: Comparison between the three paradigms of RAG. (Left) Naive RAG mainly consists of three parts: indexing, retrieval and generation. (Middle) Advanced RAG proposes multiple optimization strategies around pre-retrieval and post-retrieval, with a process similar to the Naive RAG, still following a chain-like structure. (Right) Modular RAG inherits and develops from the previous paradigm, showcasing greater flexibility overall. This is evident in the introduction of multiple specific functional modules and the replacement of existing modules. The overall process is not limited to sequential retrieval and generation; it includes methods such as iterative and adaptive retrieval ([Retrieval-Augmented Generation for Large Language Models: A Survey], March 2024)

In the light of this, we focus on the advanced RAG pipeline structure, which addresses many of the shortcomings of a naive RAG chain without introducing too much complexity to the overall pipeline, (which results in very long question-to-generation delay). This paradigm involves knowledge base augmentation by routing the input query to other processes and external knowledge sources such as web content, summarization, translation, generative research conducting and LLM instructing through prompt engineering in addition to the re-ranking of retrieval and generation phases' outputs. This is achieved through introducing two additional steps or processes into the overall pipeline: pre-retrieval and post-retrieval. The pre-retrieval phase focuses on obtaining better retrieved results by transforming and interpolating more relevant context into the prompt. The post-retrieval phase focuses on better ranking and incorporating the analyzed context into the prompt. This step (post-retrieval) also addresses content overload by filtering recurrent and repetitive information, compressing or summarizing context, and ranking post-generation output (in our case of multiple LLM selections).

These extra steps would result in better generative AI experience in case knowledge base augmentation is required just-in-time of answer-generation or a re-ranking and compression of either retrieval or generation phase output is required.

2.2.2 Significance of Retrieval-augmented Generation

RAG technology can bring many benefits for businesses:

- **Private and up-to-date data:** Even with the high-quality LLM products, it is challenging to maintain relevancy for an enterprise environment. By allowing AI models to connect to external knowledge bases and personal data, much more relevant results can be achieved.
- **Cost-effective:** Even with the availability of open-source LLMs and cloud-based FMs (Foundational Models), the computational and financial costs of re-training such models on domain-specific or private data can be high. By circumventing the re-training phase of LLMs, a huge gain can be brought off using this technique.
- **More control and customization:** It is easy to customize RAG pipelines by developers: Adapting to changing requirements and sources of information, troubleshooting model performance and results, restricting access to some information for authorized users, etc...
- **User trust:** By enabling user choices and interference in the overall process of RAG, users can look up how their choices in information sources and LLM selections affect directly the system results and performance.

RAG technique, in many ways, provides an alternative framework to fine-tuning (even though these two techniques are not mutually exclusive) with better adaptability and performance on knowledge-intensive tasks. This is because LLMs struggle to capture new factual information through unsupervised fine-tuning. RAG technique, on the other hand, is designed specifically to control which data is prioritized over another, in addition to instructing the model behavior. By combining these two techniques, we can enhance model capabilities at different levels. The following figure demonstrates the different optimization methods that affect LLM performance. One can leverage RAG technique for external knowledge augmentation, Fine-Tuning for adapting the model to different domains of expertise, and Prompt Engineering for leveraging LL model's capabilities and functions (e.g. summarization, translation, answering tasks).

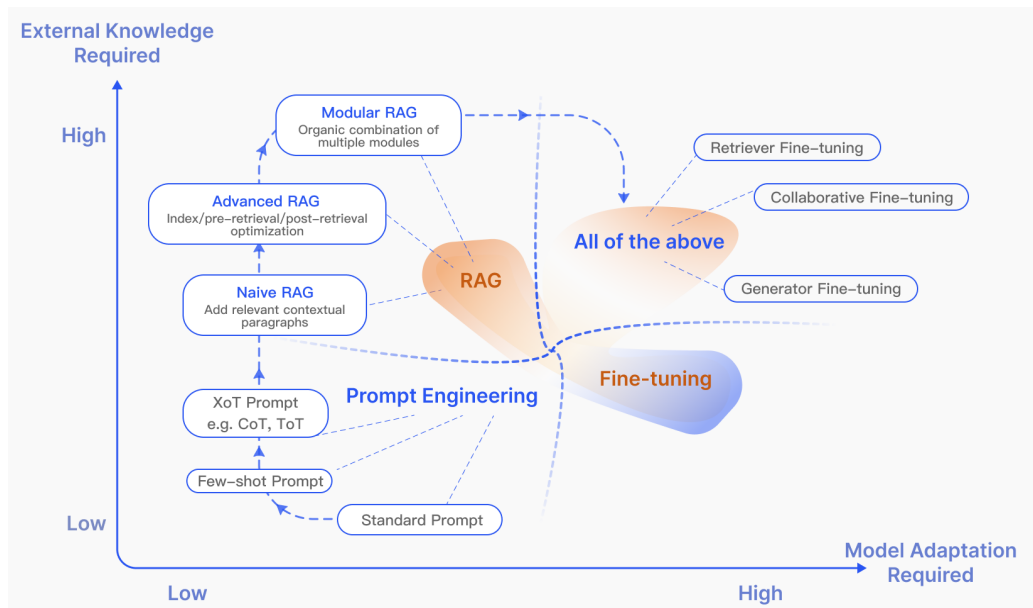


Figure 2.4: RAG compared with other model optimization methods in the aspects of “External Knowledge Required” and “Model Adaption Required”. Prompt Engineering requires low modifications to the model and external knowledge, focusing on harnessing the capabilities of LLMs themselves. Fine-tuning, on the other hand, involves further training the model. In the early stages of RAG (Naive RAG), there is a low demand for model modifications. As research progresses, Modular RAG has become more integrated with fine-tuning techniques. ([Retrieval-Augmented Generation for Large Language Models: A Survey], March 2024)

In general, the RAG technique require low model adaptation, rendering it suitable for the interaction with cloud-based LLMs, making the best blend between performance and knowledge augmentation when applied with Prompt Engineering techniques.

Fine-tuning (or in our case RLHF), on the other hand, can provide LLMs with a reward signal from human feedback to guide them on how to improve their performance in accordance to users' specifications, but this still requires more model adaptability which make it more difficult to implement in different models.

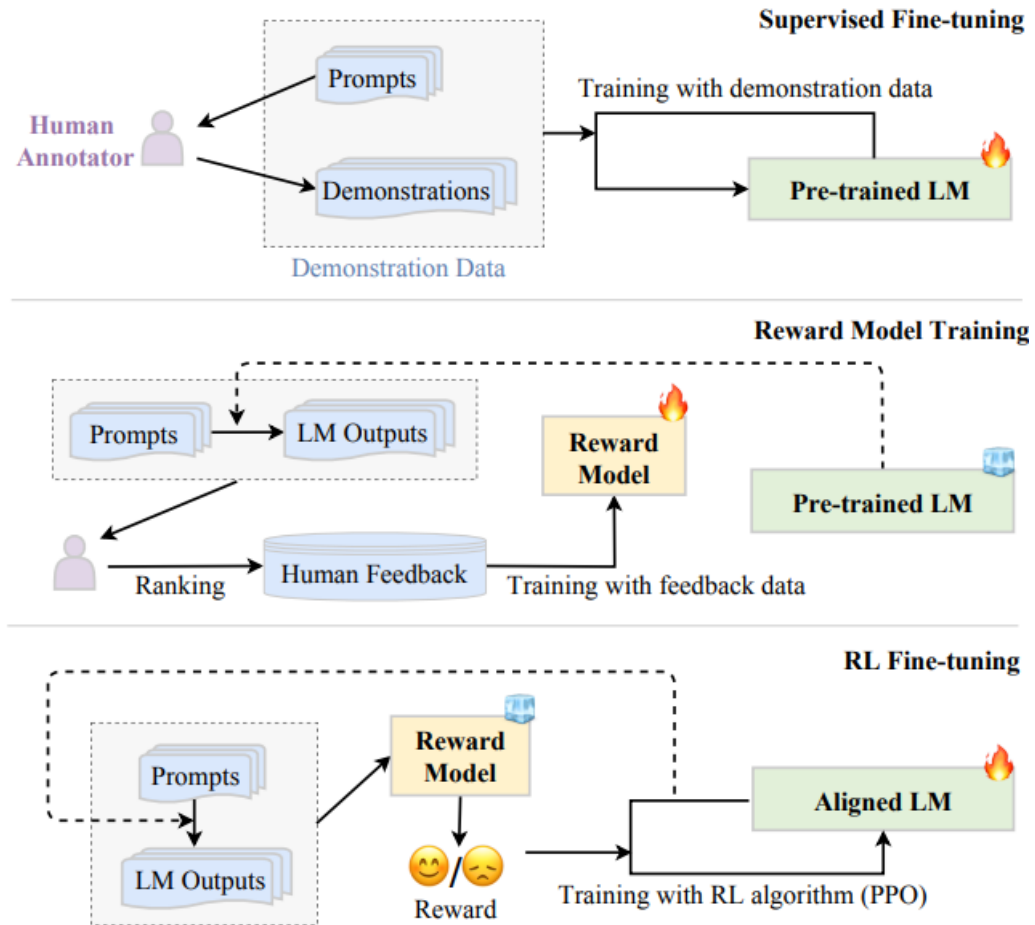


Figure 2.5: The workflow of the RLHF algorithm. ([A Survey of Large Language Models], November 2023)

2.3 Components of a typical RAG pipeline

The RAG techniques often involves many stages and processes which take place between prompting and answer-generation, as seen in previous sections. This section introduce high-level (easier to understand) definitions of the discrete components' concepts and their orchestrated functionality.

2.3.1 Vector Stores

A vector store refers to a type of database that allows the storage of documents and unstructured data in numerical representations suitable for large volume indexing and retrieval through similarity search algorithms.

”Traditional databases are made up of structured tables containing symbolic information. For example, an image collection would be represented as a table with one row per indexed photo. Each row contains information such as an image identifier and descriptive text. Rows can be linked to entries from other tables as well, such as an image with people in it being linked to a table of names.

AI tools, like text embedding (word2vec) or convolutional neural network (CNN) descriptors trained with deep learning, generate high-dimensional vectors. These representations are much more powerful and flexible than a fixed symbolic representation, as we’ll explain in this post. Yet traditional databases that can be queried with SQL are not adapted to these new representations. First, the huge inflow of new multimedia items creates billions of vectors. Second, and more importantly, finding similar entries means finding similar high-dimensional vectors, which is inefficient if not impossible with standard query languages.” ([\[Faiss: A library for efficient similarity search\]](#), 2017)

Vector, as in vector store, refers to the embedding vectors of data (arrays of floating-point numbers) which are produced by AI algorithms called embedding models. The output of these models captures the semantics of what is being embedded, which is quite suitable for LLM use cases.

As in LLMs’ case, embedding models can be designed specifically to handle textual data or multimedia files. For our case, it is obvious that a suitable embedding model would be the one that best captures the semantics with as little as possible loss of these insightful information.

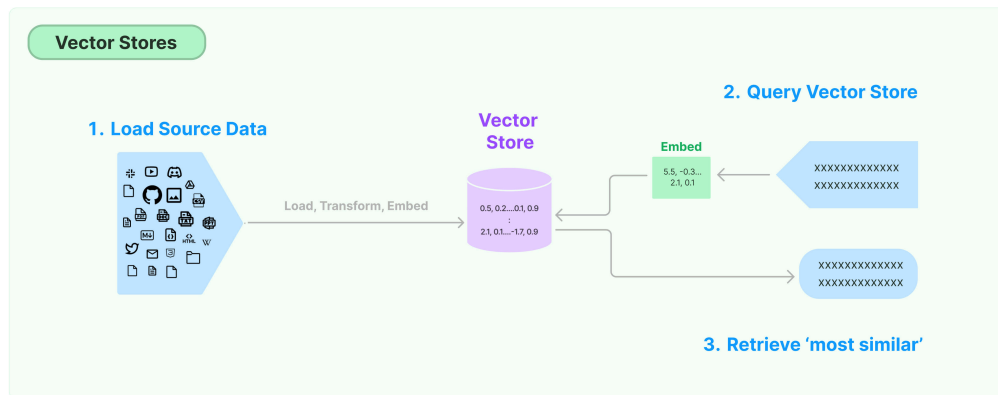


Figure 2.6: Vector Store Process Diagram. [LangChain Documentation]

In addition to embedding documents, one should consider the size of these data and the context window limitation of LLMs (each model has some limit on the amount of information it can receive as context). In these conditions, a chunking of documents should be implemented before the embedding and storing phases. This technique allows to divide documents into smaller passages, and marking these with metadata (document id, order etc...) which allow the vector store to structure and store these chunks as they were a single monolithic record. The numerical representations of document sections is adequate to be stored in a vector index in the final phase of indexation.

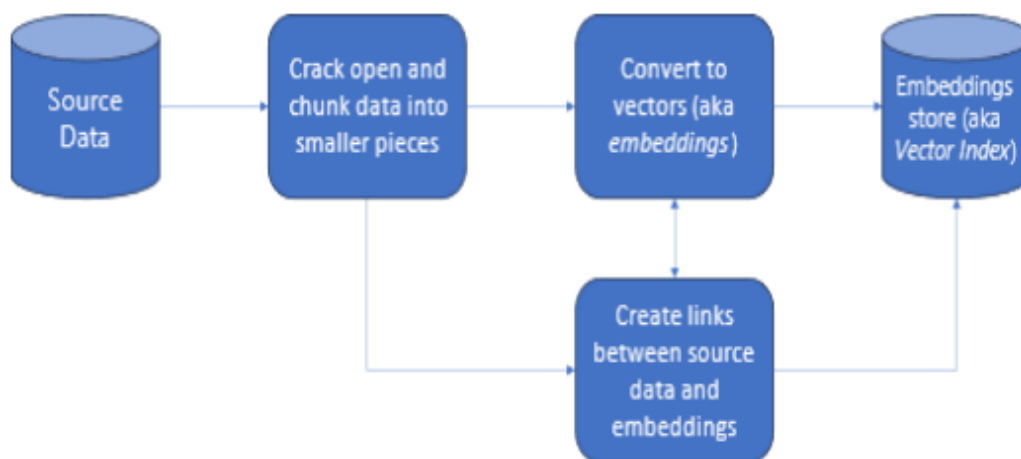


Figure 2.7: Embeddings Generation Diagram. (Microsoft, 2023)

In addition to the ability to store documents' sections in a suitable form, the vector store index should also be able to search these elements and decode them back

into their original textual form. This is achieved through similarity searching algorithms often implemented as constituent functionalities with the vector store. These functions return a number of passages semantically similar to the search term. This is achieved due to the concept of distance calculation between vectors in mathematics.

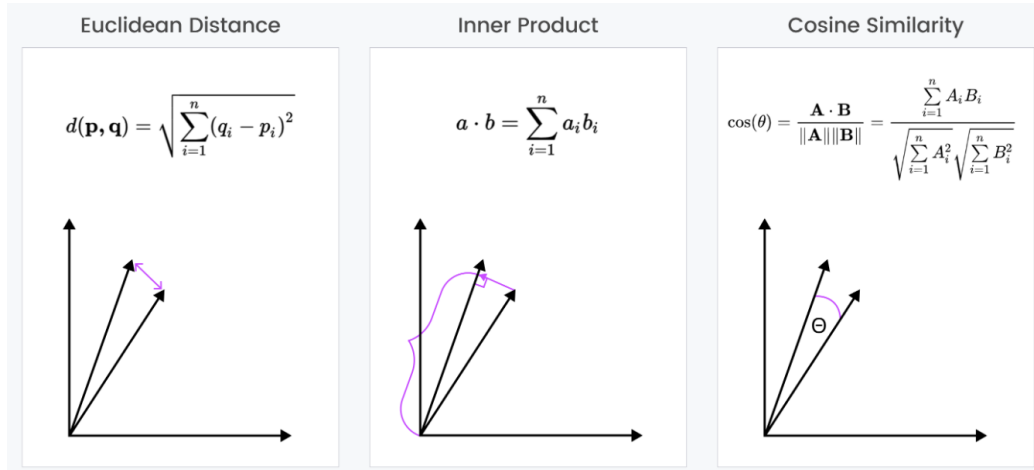


Figure 2.8: Similarity Metrics for Vector Search. [\[Zilliz Blog\]](#)

2.3.2 Large Language Model

The principal component of a RAG system is a textual GAI model which acts as tool that understands user prompts' context and generate textual responses accordingly. Essentially, these models are transformer-based, which allows them to learn complex relationships between words in sentences. For example, a prompt can provide the LLM with context about what kind of text it should generate. The LLM can then use this context to inform its attention mechanism which parts of the prompt it should focus on to generate the most relevant response.

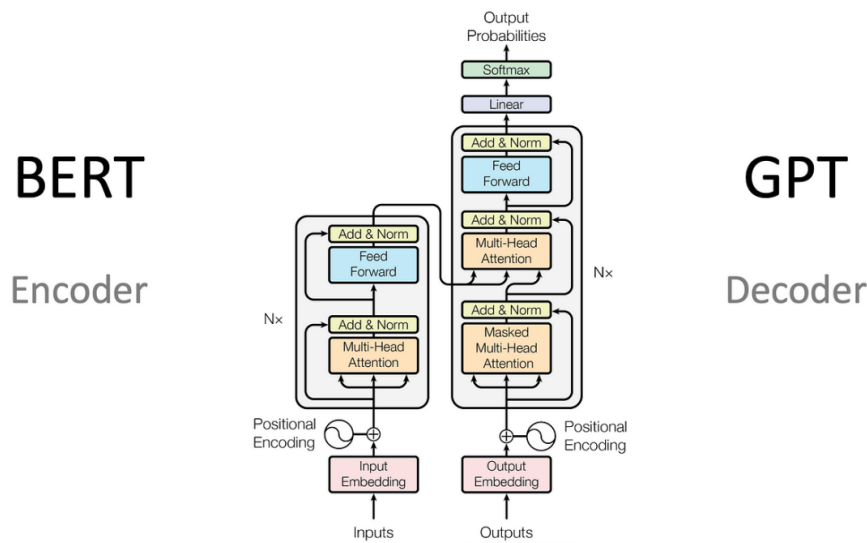


Figure 2.9: An illustration of main components of the transformer model from the original paper, where layer normalization was performed after multiheaded attention. [\[On Layer Normalization in the Transformer Architecture \(2020\)\]](#)

The main steps that a transformer-based model pass through are the following:

- **Input Embedding:** The transformer first converts the input text into a series of numerical representations, which are then passed through a positional encoding layer. This layer adds information about the word's position in the sentence, which is important because word order matters in a language like English.
- **Encoder Layers:** The core building block of the transformer encoder is the “encoder layer”. An encoder layer typically consists of two sub-layers: a multi-head attention layer and a feed-forward layer.
- **Multi-Head Attention Layer:** The multi-head attention layer allows the model to attend to different parts of the input sentence simultaneously. This

is important for understanding the relationships between words in a sentence.

- **Feed Forward Layer:** The feed-forward layer is a simple neural network that further processes the information from the attention layer.
- **Decoder Layers:** After the encoder has processed the input text, the decoder generates the output text. The decoder also uses encoder layers, but with an additional masked multi-head attention layer. This layer prevents the decoder from attending to future words in the output sentence, which would allow it to “cheat” by looking ahead.
- **Softmax Layer:** The softmax layer converts the decoder’s output into a probability distribution over all the words in its vocabulary. This allows the model to predict the next word in the sentence for instance.

2.3.3 Prompts and Prompt Engineering

Prompt Engineering (PE) plays a crucial rule in RAG context and LLMs in general. It refers to how the model is prompted, i.e. what does it receive as input. It may include several instructions to the LLM that guides it on how to perform the generation stage. It also can include different parts or steps, such as passing the retrieved context or chat history directly into the prompt, provide it with examples which it should consider when providing answers, or instruct it on how long the answer should be or in which tone it should respond.

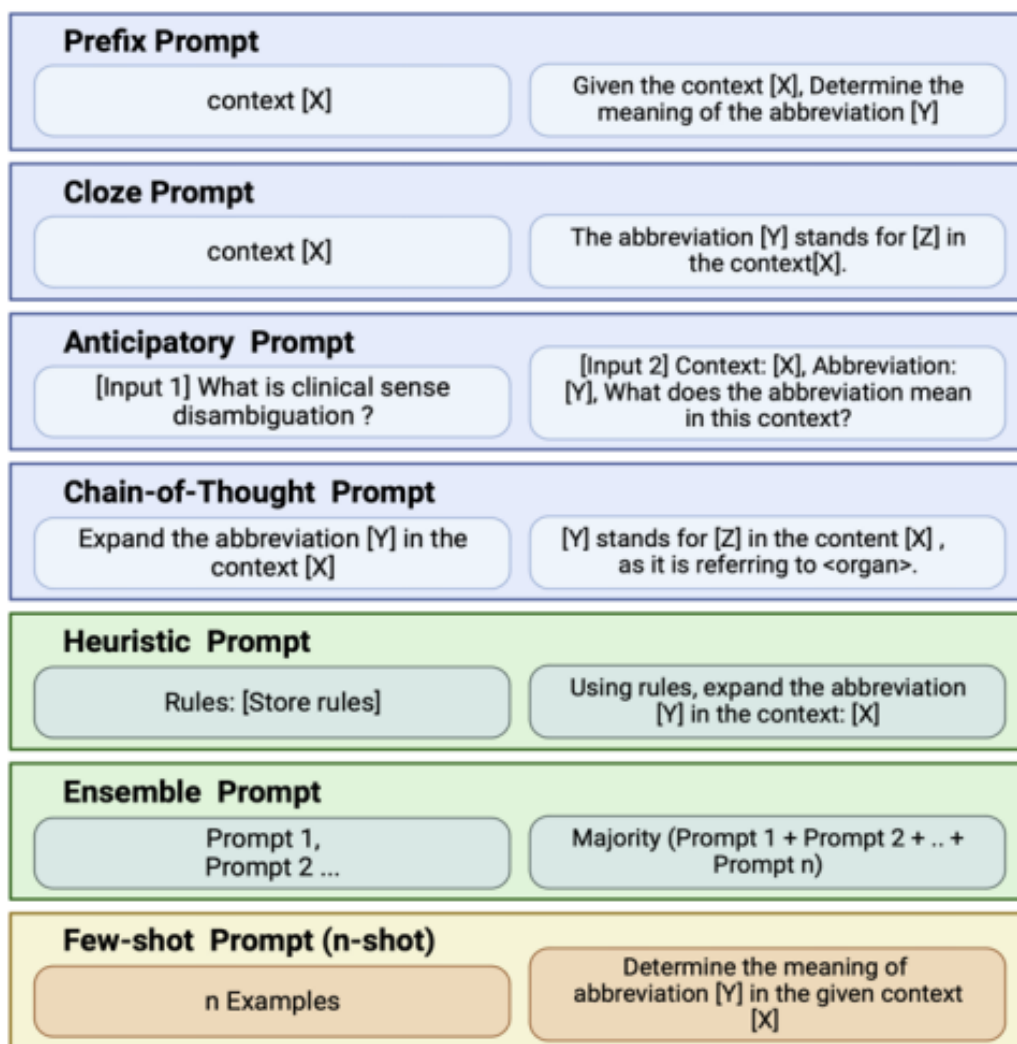


Figure 2.10: Types of Prompts: [X]: context, [Y]: Abbreviation, [Z]: Expanded Form ([A Study on the Implementation of Generative AI Services Using an Enterprise Data-Based LLM Application Architecture], 2023)

Chapter 3

Methodology

3.1 System Architecture

3.2 LangChain Overview

3.3 Vector Store Implementation

Chapter 4

Implementation

This chapter details the implementation process of the solution discussed in the previous chapter: The work environment, tools and libraries employed to accomplish the tasks, the technical aspects of design choices execution, in addition to some of the major results accomplished and how to build upon them.

4.1 Workstation

This section identifies the characteristics, both hardware and software, of the computer system on which this project was implemented.

4.1.1 Hardware

- Device : DELL Inspiron 3593
- Processor : Intel(R) Core(TM) i7-1065G7 CPU @ 1.30GHz 1.50 GHz
- RAM : 24 GB
- GPU : NVIDIA GeForce MX230

4.1.2 OS and Software

- OS : Debian trixie inside a Windows Subsystem for Linux (WSL2) VM
- Source-code Editor : Visual Studio Code (with extensions to enable interacting with WSL, Jupyter Notebooks and Python virtual environments) + NeoVim
- Web Browser : Microsoft Edge (for web application testing and troubleshooting)
- Drivers and Toolkits : NVIDIA GPU driver + CUDA Toolkit among others

4.2 Components

In this section, we research available tools making it possible to develop the different functionalities, provide comparison between alternatives when choices were made, go through how they were implemented and end up with testing the effectiveness of the developed solution in addition to laying out its unfortunate limitations.

4.2.1 Libraries and Frameworks

Web interface and services

There is a plethora of frameworks of this kind, from customary JavaScript frameworks (React, Angular, Vue) to others designed to enable faster delivery of interactive web apps (such as Streamlit and Chainlit).

The choice was made to use Streamlit as it provides a much faster way to develop LLM chat interfaces than JS frameworks, while also being more advanced than Chainlit in terms of flexibility and building custom interfaces by supplying developers with many customizable and ready-to-use interface components like chat and message containers.



Figure 4.1: Streamlit logo.

We can look through its app gallery to find an abundance of templates that provide many examples of built apps which interact with LLMs, LangChain and other frameworks.

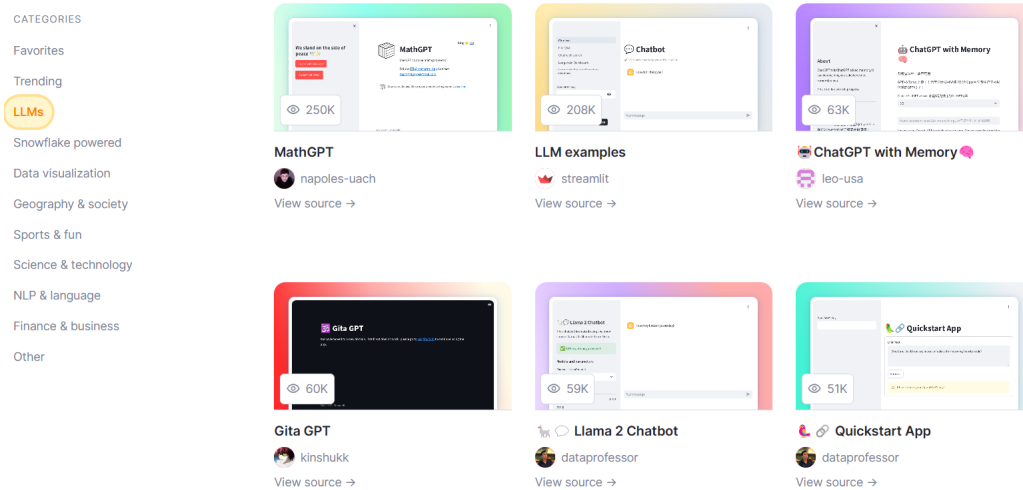


Figure 4.2: Streamlit App Gallery.

RAG pipelines development

There are a few frameworks enabling developers to interact with LLMs and RAG pipelines, each of which provide different integrations with external APIs and tools: LangChain, LlamaIndex, Haystack, Langroid.

The choice was made by 3S project coordinators to use LangChain for their solution. It is a great choice given that this platform provides all the tools to build complex RAG pipelines and personalize the different steps of these pipelines. Also, since its emergence in October 2022, this library has gained remarkable prominence with courses available on DeepLearning.ai, tutorials from NVIDIA, OpenAI, Google and others in addition to the exhaustive documentation available online.



Figure 4.3: LangChain logo

This framework has all the required components to build the most advanced RAG pipelines: Data loading from various sources, LLM and vector store integrations

from different providers (both locally and on the cloud), prompt templates for different LLMs and tasks, etc...

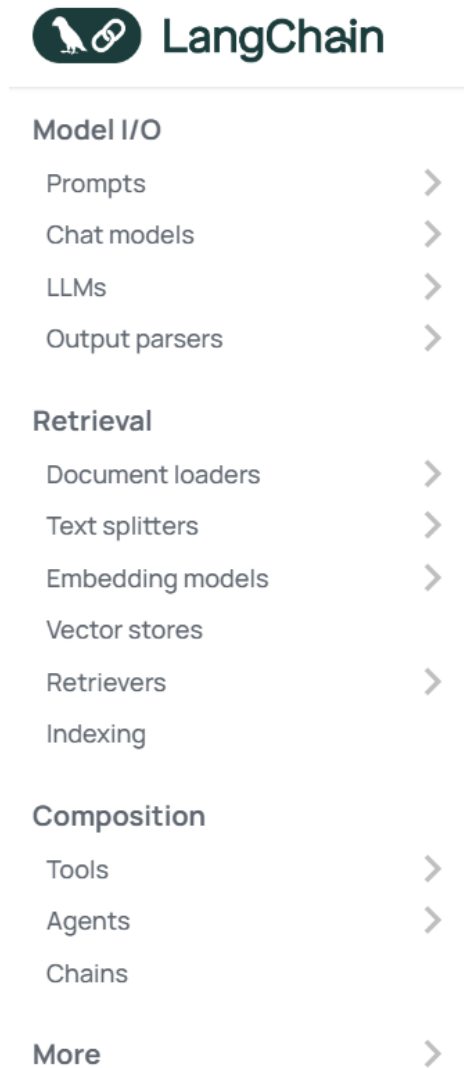


Figure 4.4: LangChain Components. ([LangChain Documentation](#))

4.2.2 LLMs

There is a plethora of models of this type, both closed and open source, with many ways to access and use them.

Locally through Hugging Face Hub

Many open-source LLM variants are available on Hugging Face Models Hub, as it is the main developer of the 'transformers' Python library.



Figure 4.5: Hugging Face logo. [[Hugging Face models Hub](#)]

”The Hugging Face Hub hosts many models for a [variety of machine learning tasks](#). Models are stored in repositories, so they benefit from [all the features](#) possessed by every repo on the Hugging Face Hub. Additionally, model repos have attributes that make exploring and using models as easy as possible.” ([[Hugging Face Models Hub documentation](#)], 2024)

We can find an abundance of pre-trained LLMs downloadable from the HF Hub.

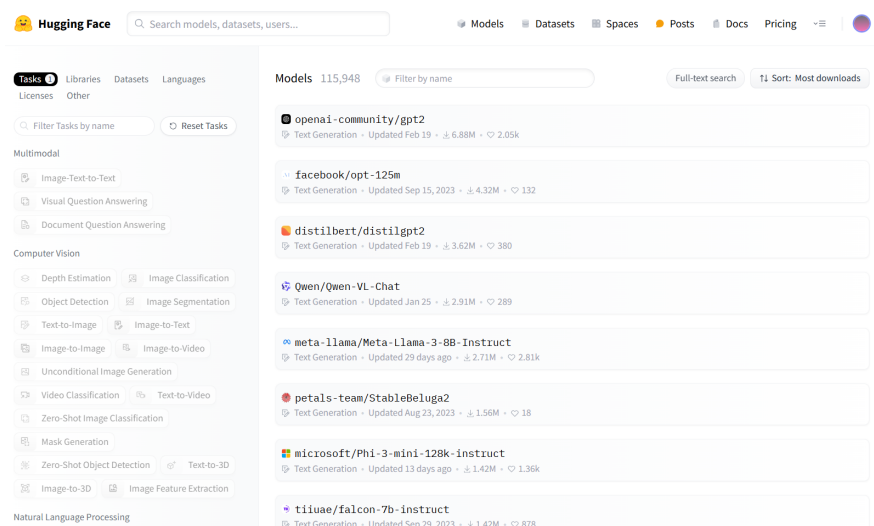


Figure 4.6: A short list of popular text-generation models available on HF Hub. [HF models filtered by ‘text-generation’ task](#)

This method allows the most customization of how an LLM work as it loads a model from locally downloaded files. These typically include a handful of customization options.

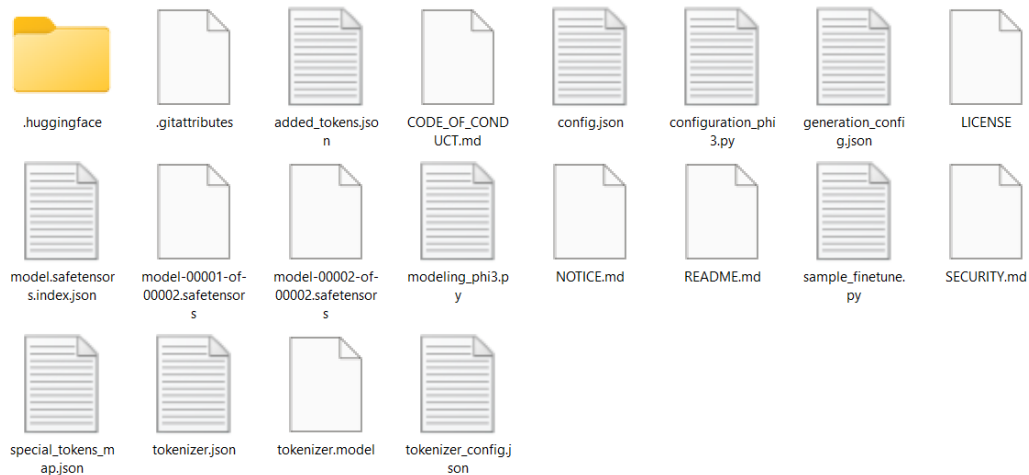


Figure 4.7: An illustration of different files of a locally downloaded model ([Phi-3-mini-4k-instruct](#)) from [Hugging Face Hub](#)

We can observe multiple files which boast many customization parameters and fine-tuning samples: `config.json`, `configuration_phi3.py`, `generation_config.json`, `tokenizer_config.json`, `sample_fineturning.json`.

While this method is suitable for fine-tuning and reinforcement learning purposes, it comes with the downsides of high hardware usage (CPU, GPU, RAM, Disk Space) and time-consuming performance of the ML models.

API-enabled services and Cloud-based solutions

Many companies behind Large Language Models development provide APIs or cloud-based environments to access their models. Some of the most popular options include OpenAI API, Google Cloud Vertex AI, Anthropic, Cohere, FireworksAI, MistralAI, TogetherAI, GroqCloud among others.

This method, in contrast to running models locally, does not come with the prerequisites of expensive hardware or delayed answer generation. Even though these solutions are paid services, most provide free trials and some of the available free plans only has a limit on daily and monthly usage, and is usually enough for personal usage. In addition to this, some of these cloud environments provide many models to use. For instance, FireworksAI allows to use Llama-3, Yi-Large, Mistral, while TogetherAI provides models such as Qwen-2, Gemma (open-source

version of Google’s model, Gemini), Phi-2, Nous Capybara, and many others, all from within a single platform. This solution also allows some customization options through their web interfaces, even though limited.

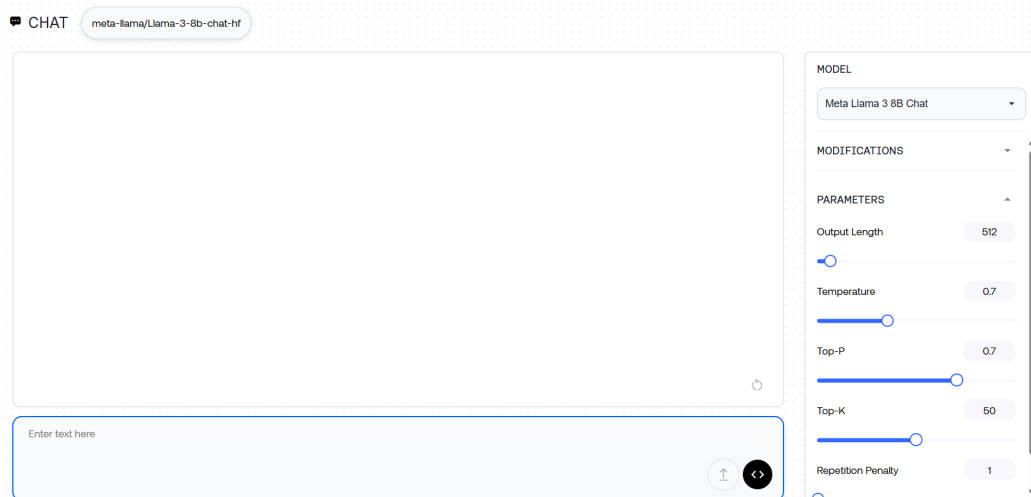


Figure 4.8: TogetherAI playground web interface, an example of a cloud-based environment allowing multiple LLM’s customization options and parameters. [TogetherAI Playground service](#)

To compare different LLMs and their providers, the ‘Artificial Analysis LLM Performance Leaderboard’ space on Hugging Face provides independent performance benchmarks and pricing across API providers of LLMs.

		CONTEXT ↕	MODEL QUALITY ↕	PRICE ↕	OUTPUT TOKENS/S ↕	LATENCY ↕		
API PROVIDER ↕	MODEL ↕	CONTEXT WINDOW ↕	INDEX <small>Normalized avg</small> ↕	BLENDED <small>USD/1M Tokens</small> ↕	MEDIAN <small>Tokens/s</small> ↕	MEDIAN <small>First Chunk (s)</small> ↕	FURTHER ANALYSIS	
OpenAI	GPT-4o	128k	100	\$7.50	86.8	0.59	Model	Providers
OpenAI	GPT-4 Turbo	128k	94	\$15.00	27.2	0.68	Model	Providers
Microsoft Azure	GPT-4 Turbo	128k	94	\$15.00	29.8	0.55	Model	Providers
OpenAI	GPT-4	8k	84	\$37.50	24.7	0.87	Model	Providers
Microsoft Azure	GPT-4	8k	84	\$37.50	20.2	0.54	Model	Providers
OpenAI	GPT-3.5 Turbo	16k	59	\$0.75	67.2	0.46	Model	Providers
Microsoft Azure	GPT-3.5 Turbo	16k	59	\$0.75	62.5	0.32	Model	Providers
OpenAI	GPT-3.5 Turbo Instruct	4k	60	\$1.63	65.3	0.37	Model	Providers
Microsoft Azure	GPT-3.5 Turbo Instruct	4k	60	\$1.63	136.5	0.61	Model	Providers
Google	Gemini 1.5 Pro	1m	95	\$5.25	62.4	1.14	Model	Providers
Google	Gemini 1.5 Flash	1m	84	\$0.53	155.2	1.17	Model	Providers
Google	Gemini 1.0 Pro	33k	62	\$0.75	88.2	2.18	Model	Providers
Fireworks AI	Gemma 7B	8k	45	\$0.20	234.1	0.26	Model	Providers
deepinfra	Gemma 7B	8k	45	\$0.07	65.9	0.29	Model	Providers
groq	Gemma 7B	8k	45	\$0.07	1,039.4	0.87	Model	Providers
together.ai	Gemma 7B	8k	45	\$0.20	141.4	0.30	Model	Providers
Replicate	Llama 3 (70B)	8k	83	\$1.18	46.8	1.66	Model	Providers
BWS	Llama 3 (70B)	8k	83	\$2.86	49.3	0.48	Model	Providers
OctoAI	Llama 3 (70B)	8k	83	\$0.90	62.0	0.29	Model	Providers
Lepton AI	Llama 3 (70B)	8k	83	\$0.80			Model	Providers
Microsoft Azure	Llama 3 (70B)	8k	83	\$5.67	17.4	2.91	Model	Providers
Fireworks AI	Llama 3 (70B)	8k	83	\$0.90	112.1	0.26	Model	Providers
deepinfra	Llama 3 (70B)	8k	83	\$0.61	20.5	0.35	Model	Providers
groq	Llama 3 (70B)	8k	83	\$0.64	350.8	0.37	Model	Providers
databricks	Llama 3 (70B)	8k	83	\$1.50	69.9	0.53	Model	Providers
perplexity	Llama 3 (70B)	8k	83	\$1.00	46.9	0.33	Model	Providers
together.ai	Llama 3 (70B)	8k	83	\$0.90	125.8	0.67	Model	Providers
Replicate	Llama 3 (8B)	8k	64	\$0.10	78.7	1.78	Model	Providers
BWS	Llama 3 (8B)	8k	64	\$0.38	79.2	0.29	Model	Providers
OctoAI	Llama 3 (8B)	8k	64	\$0.15	132.1	0.21	Model	Providers
Lepton AI	Llama 3 (8B)	8k	64	\$0.07			Model	Providers
Microsoft Azure	Llama 3 (8B)	8k	64	\$0.55	76.5	0.91	Model	Providers

Figure 4.9: Artificial Analysis LLM Performance Leaderboard [ArtificialAnalysis/LLM-Performance-Leaderboard space on Hugging Face]

The leaderboard shows some of the most common families of LL models: GPT, Gemini, Llama, Mistral, etc... Each of these models has its strengths and weaknesses, but are suitable for the RAG-context applications. For this project, some of the most popular LLMs were implemented through various cloud providers for the purpose of exploring and evaluating their performance (the selection was arbitrary in view of that multiple LLMs, rather than a single one, will generate answers in parallel).

4.2.3 Vector Stores

There are many vector stores and embedding models to choose from which we will look through their differences.

The choice of the most suitable vector store solution was based on four criteria mainly:

- **Self-hosting:** This means that the vector store will be managed locally on the same computing infrastructure as the web server. This is opposed to a cloud-based deployment, which comes with the drawbacks of higher latency, possible network errors and high-costs.
- **Latency:** This refers to the performance and speed of similarity searching algorithms which are provided by the vector store and its ability to index and handle large volumes of data with the utilization of GPU parallel computing features.
- **Accuracy:** The relevance of retrieved data to the actual searched query. Often, this has a reverse relation with latency, as more accurate results take longer to be achieved.
- **Documentation:** Online Documentation and community forums that can guide on how to use the database efficiently.

Vector Store	Self-hosting	Latency	Accuracy	Documentation
FAISS	✓	✓	✓	✗
Pinecone	✗	✓	✗	✓
Chroma	✓ / ✗	✓	✗	✓
Lance	✓	✗	✗	✗

Table 4.1: Comparison of popular vector databases

After careful consideration, FAISS vector store was selected due to its high performance and accuracy in comparison to alternatives. It leverages the GPU-enabled CUDA toolkit, and provides a state-of-the-art implementation of similarity searching algorithms based on this structure.

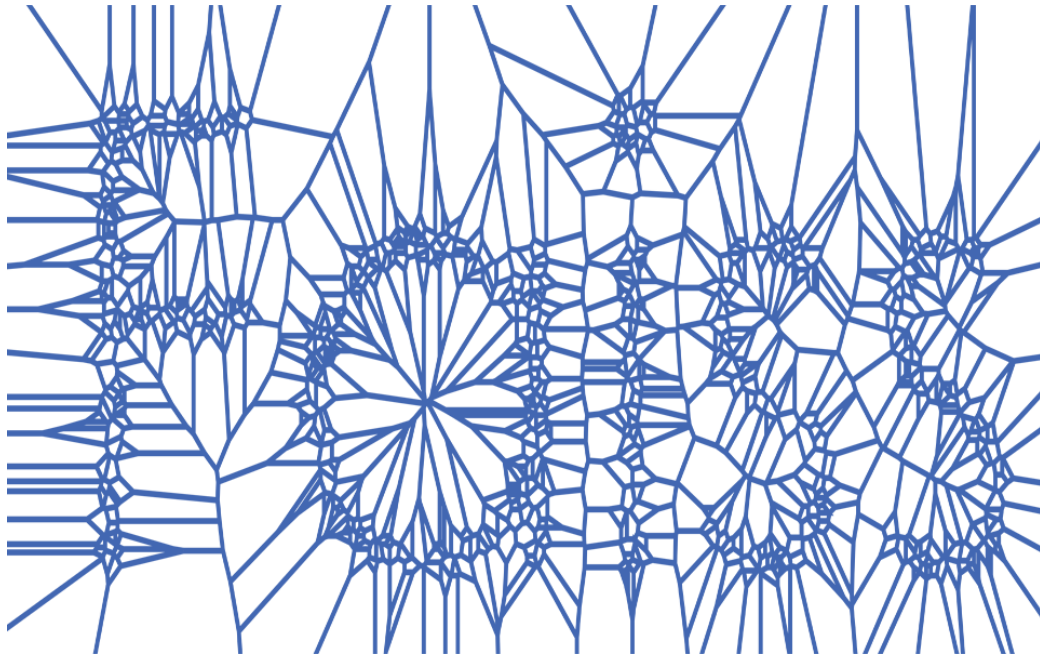


Figure 4.10: FAISS logo.

FAISS is a FOSS library, it stands for "Facebook AI Similarity Search" with an implementation of nearest-neighbor search and k-selection algorithms designed specifically to efficiently handle large data sets, 8.5x faster than previous methods.

4.2.4 Embedding Model

We have the choice to select the most suitable embeddings model.

We have highlighted in the second chapter the importance of an embedding algorithm that minimize the loss of semantics when converting textual data to embedding vectors. In addition to this, considering that a vector store once initiated with an embedding model, can no longer swap it with another (unless re-initialized from zero). For this purpose, it would be a bad choice to consider cloud solutions as these may be unavailable in some cases (network or provider failure, expiration of tokens...).

The best choice in this case, as in vector store’s choice, is to select a suitable model which is available offline (as in self-hosting). This will avoid foreseeable failures and limit problems.

This redirects us to the Hugging Face Hub where we can find many embedding models, which are available through the sentence transformers-repository.

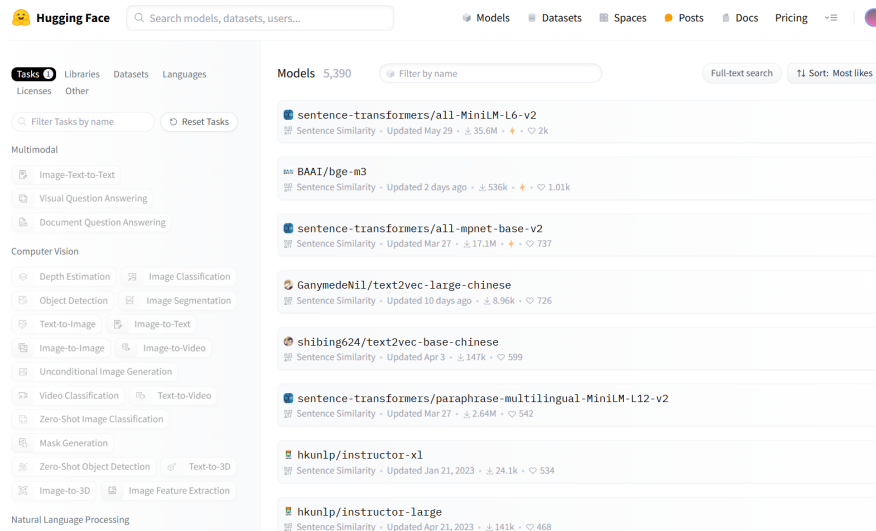


Figure 4.11: A short list of popular embedding models available on HF Hub. [HF models filtered by ‘sentence-similarity’ task](#)

Two of the most popular choices are ”all-MiniLM-L12-v2” and ”all-mpnet-base-v2”, the first of which makes a better choice when working in a limited environment while not paying much attention to retaining semantics (faster embedding generation and smaller size), while the second one, ”all-mpnet-base-v2”, is the more suitable choice for our case due to, even with its larger footprint, its capability to capture most of the semantics and meanings of sentences.


Model Name	 Performance Sentence Embeddings (14 Datasets) ⓘ	Performance Semantic Search (6 Datasets) ⓘ	Avg. Performance ⓘ	Speed ⓘ	Model Size ⓘ
all-mpnet-base-v2 ⓘ	69.57	57.02	63.30	2800	420 MB
all-distilroberta-v1 ⓘ	68.73	50.94	59.84	4000	290 MB
all-MiniLM-L12-v2 ⓘ	68.70	50.82	59.76	7500	120 MB
all-MiniLM-L6-v2 ⓘ	68.06	49.54	58.80	14200	80 MB
multi-qa-mpnet-base-dot-v1 ⓘ	66.76	57.60	62.18	2800	420 MB
multi-qa-distilbert-cos-v1 ⓘ	65.98	52.83	59.41	4000	250 MB
paraphrase-multilingual-mpnet-base-v2 ⓘ	65.83	41.68	53.75	2500	970 MB
paraphrase-albert-small-v2 ⓘ	64.46	40.04	52.25	5000	43 MB
multi-qa-MiniLM-L6-cos-v1 ⓘ	64.33	51.83	58.08	14200	80 MB
paraphrase-multilingual-MiniLM-L12-v2 ⓘ	64.25	39.19	51.72	7500	420 MB
paraphrase-MiniLM-L3-v2 ⓘ	62.29	39.19	50.74	19000	61 MB
distiluse-base-multilingual-cased-v1 ⓘ	61.30	29.87	45.59	4000	480 MB
distiluse-base-multilingual-cased-v2 ⓘ	60.18	27.35	43.77	4000	480 MB

Figure 4.12: Sentence-Transformers model performance comparison [[Pre-Trained Sentence Transformers models' performance](#)]

”The all-* models were trained on all available training data (more than 1 billion training pairs) and are designed as general purpose models. The all-mpnet-base-v2 model provides the best quality, while all-MiniLM-L6-v2 is 5 times faster and still offers good quality.” ([[Sentence Transformers documentation](#)])

4.3 Incorporation into a functional RAG system

After introducing the frameworks and tools which would allow us to build a RAG system, it is necessary to discuss the implementation details of these elements and how they were incorporated together.

This section is dedicated to showcasing how these components interact together in the implemented solution.

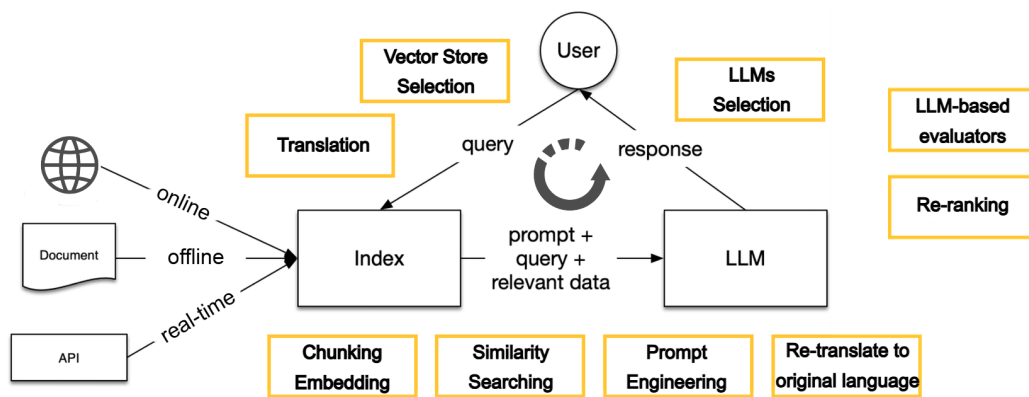


Figure 4.13: The overall pipeline of the system to be implemented, showcasing the interaction between the different components, from knowledge augmentation methods (document/data loading and chunking etc...), vector database (index), RAG pipeline (User-Index-LLM interaction), re-ranking and evaluation processes

This figure demonstrates the overall process that are executed in the system:

- Loading data from online web content, offline documents and files, or real-time data from APIs (Search engines and GPT Researcher) after which the chunking and embedding transformation occur.
- User preference selections (vector stores, LLMs, enable translation or not)
- Typical RAG processes: retrieval (similarity searching), prompt engineering and transformation, and finally generation
- Reiterate the RAG pipeline for every LLM selection (automatically retranslate the answer to the user's language if needed).
- Re-Ranking of answers based on RAG metrics and user feedback.

4.3.1 Data Ingestion Methods

The first essential part of the project is to allow the on-demand ingestion of data and new information into the vector store.

For this purpose, various file formats and online web scraping and fetching methods have been implemented; LangChain, as seen in the previous section, provides a "Document loaders" section in its "Retrieval" toolbox. These tools provide many tutorials and helpful functions to implement document loading and chunking from various sources, which has allowed to accomplish the required methods to satisfy the coordinating teams at 3S.

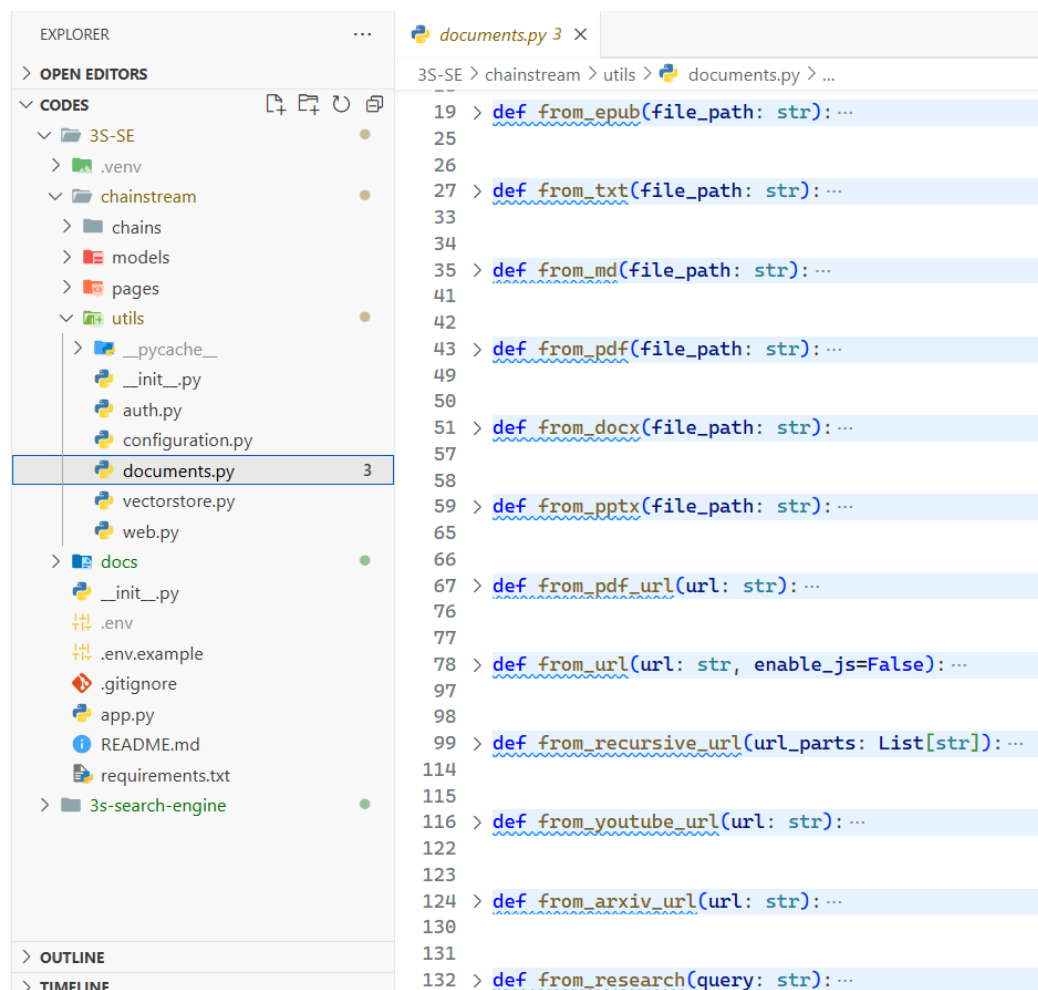


Figure 4.14: Various data ingestion methods

To make it more comprehensible, here is a short description of what each of these methods do:

- Web Pages (from_url, from_recursive_url): Allowing to read a single page from a given URL, or a page and its child pages recursively.
- Arxiv Research Papers (from_arxiv_url): Allows to read content directly from an arxiv URL, conserving metadata like author, dates, etc...
- Youtube Content (from_youtube_url): Transcribes a video uploaded on YouTube and constructs a document from a given URL.
- Generative Research (from_research): Utilizes [GPT Researcher](#), a tool that allows to, given a question, generate multiple queries to send to search engines and then generating an artificial report.
- EPUB files (from_epub): This is an ebook format suitable for storing and distributing large volumes of information.
- Text files (from_txt): A file format for storing texts and notes.
- Markdown files (from_md): An easy to use syntax to write documentation and notes.
- PDF files (from_pdf): Allowing to read PDF files.
- Same for other file formats...

4.3.2 Vector Store

As outlined in the project objectives, the implementation of a vector database should allow different teams to upload their documents into separate vector databases. To achieve this, two Python classes have been implemented: one to hold the vector store functionalities: loading and saving to local storage in addition to listing available databases, and another to manage access credentials to vector store through a name and a passphrase.

4.3. Incorporation into a functional RAG system

Chapter 4. Implementation

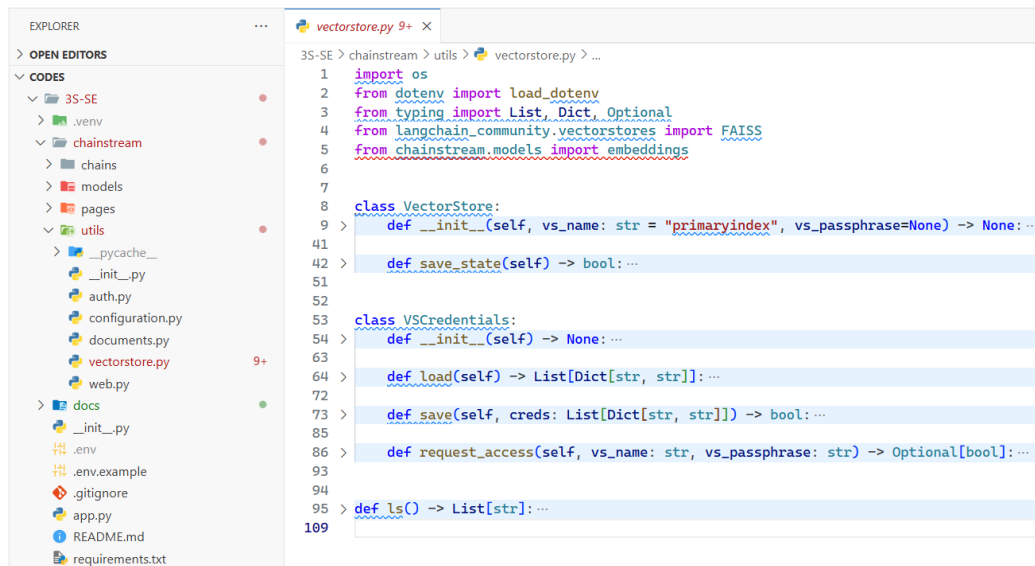


Figure 4.15: Vector Store implementation classes and methods

A vector store can be initialized with "vs.name" alone, which gives the read access to its contents. To modify it however, one should provide it with a passphrase ("vs.passphrase"), which would give users who have gained access the ability to load new content through the various data ingestion methods previously mentioned.

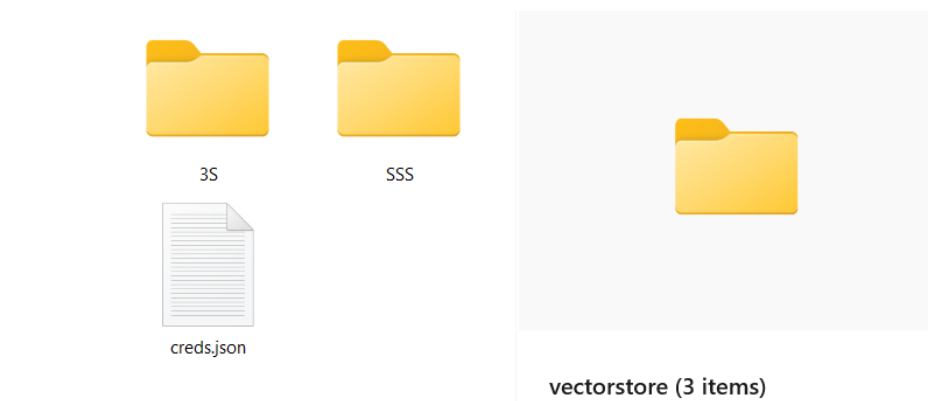


Figure 4.16: Local Directory for storing vector stores.

The "creds.json" file store access credentials to the available vector stores, while other folders ("3S" and "SSS" in this example) holds the vector store's data.

A vector store initialization typically includes an embedding model initialization. For our case, as we chose to select a model from Hugging Face Hub ("all-mpnet-base-v2"), we need to ensure that the model is downloaded and stays

up-to-date.

```
chainstream > models > utils.py > ...
import os
from huggingface_hub import snapshot_download
from dotenv import load_dotenv

load_dotenv()

def localize(model_name):
    try:
        model_path = os.environ.get(
            "PROJECT_DATA_DIR",
            os.path.expanduser("~/Downloads/3S-SE-AI/")
        ) + model_name
        snapshot_download(
            repo_id=model_name,
            local_dir=model_path
        )
        return model_path
    except BaseException as e:
        print(e)
```

```
3S-SE > chainstream > models > embeddings.py
8 def load(model_name="sentence-transformers/all-mpnet-base-v2",
9         prefer_cuda=True):
10     load_dotenv()
11     # freeze_support()
12     model_path = localize(os.environ.get("EMBEDDINGS_MODEL_NAME",
13                                         model_name))
14     try:
15         return HuggingFaceEmbeddings(
16             model_name=model_path,
17             # multi_process=True,
18             model_kwargs={
19                 "device": f"cuda:{cuda.current_device()}" if
20                     prefer_cuda and cuda.is_available() else "cpu"
21             },
22             encode_kwargs={"normalize_embeddings": True},
23         )
24     except cuda.OutOfMemoryError as e:
25         return HuggingFaceEmbeddings(
26             model_name=model_path,
27             # multi_process=True,
28             model_kwargs={
29                 "device": "cpu",
30             },
31             encode_kwargs={"normalize_embeddings": True},
32         )
```

(a) Downloading

(b) Loading

Figure 4.17: Embedding Model Implementation

The "snapshot_download" function imported from "huggingface_hub" ensures that the latest version of the model is available locally and can be loaded, while the implemented 'load' function leverages it and loads the model on the GPU with a CPU fallback.

4.3.3 LLMs and Prompts

To make the application more extensible, the implementation of LLMs and their Prompts was designed in an extendible manner.

4.3. Incorporation into a functional RAG system Chapter 4. Implementation

chainstream > utils > configuration.py > load

```
def load(args={}):  
    config = {  
        "llms": {},  
        "web_tools": {},  
    }  
  
    if "OPENAI_API_KEY" in args:  
        try:  
            config["llms"]["gpt"] = {  
                "name": "GPT 3.5 Turbo",  
                "model": ChatOpenAI(  
                    model="gpt-3.5-turbo", openai_api_key=args["OPENAI_API_KEY"]  
                ),  
            }  
        except BaseException as e:  
            print(e)  
  
    if "ANTHROPIC_API_KEY" in args:  
        try:  
            config["llms"]["claude"] = {  
                "name": "Claude 3 Sonnet",  
                "model": ChatAnthropic(  
                    model="claude-3-sonnet-20240229",  
                    anthropic_api_key=args["ANTHROPIC_API_KEY"],  
                ),  
            }  
        except BaseException as e:  
            print(e)  
  
    if "GOOGLE_API_KEY" in args:  
        try:  
            config["llms"]["gemini"] = {  
                "name": "Gemini Pro",  
                "model": ChatVertexAI(  
                    model="gemini-pro", google_api_key=args["GOOGLE_API_KEY"]  
                ),  
            }  
        except BaseException as e:  
            print(e)
```

Figure 4.18: Implementation of LLMs and Prompts, focusing on extensibility

Initialized as empty, the list of Large Language Models ("llms" key in the "config" variable) grows dynamically when the function which loads them is passed with API keys from their providers.

The list is much bigger than this illustration (and allows for integrating other tools "web_tools", which would allow the LLMs to connect to external web searching APIs), but it is easy to figure out how to add other LLMs and APIs as needed.

In addition to these APIs, this loading mechanism allowing for flexible Prompt modification when different models require different prompts.


```

config["llms"]["llama"] = {
    "name": "LlaMa 3 70B Instruct",
    "model": ChatFireworks(
        model="accounts/fireworks/models/llama-v3-70b-instruct",
        fireworks_api_key=args["FIREWORKS_API_KEY"],
    ),
    "prompt": hub.pull("rlm/rag-prompt-llama3"),
}

```

Figure 4.19: Incorporating prompts into model definition.

This approach allows to define a custom prompt template suitable for a specific model (such as in our case of Llama 3) from within a unified source code part. In theory, the "prompt" field can hold any f-string format (a concept in Python used to interpolate variables into a string), but the [LangChain Hub](#) provides many example prompts for different tasks.

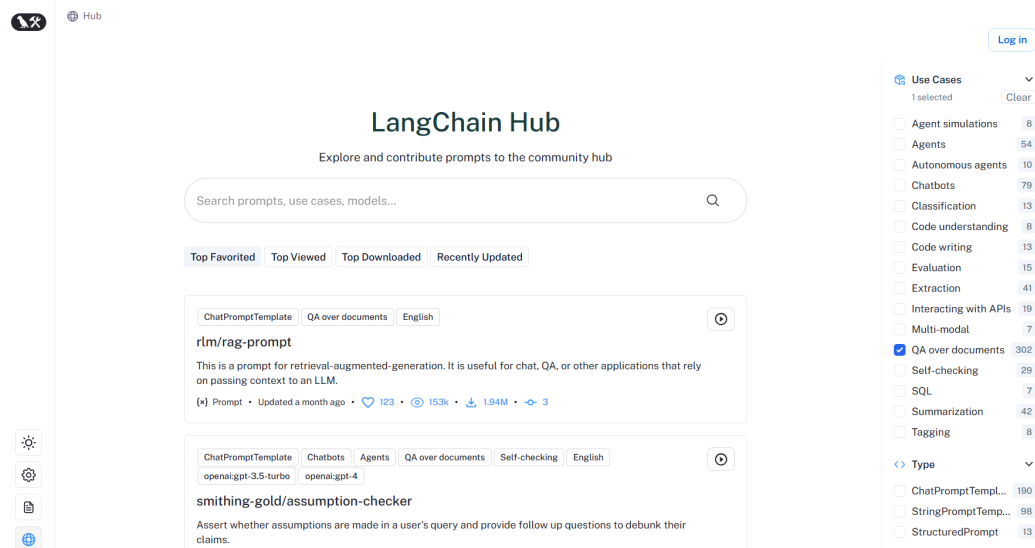


Figure 4.20: LangChain Hub web interface.

This Hub allows to find ready-to-use prompt templates suitable for various LLMs for Retrieval-augmented generation purposes.

The resulting list of implemented LLMs is as follows.



Figure 4.21: List of implemented LLMs

4.4 Testing and Validation

Chapter 5

Conclusion and Future Work

5.1 Summary of Findings

5.2 Future Research Directions

References



ESPRIT SCHOOL OF ENGINEERING

www.esprit.tn - E-mail : contact@esprit.tn

Siège Social : 18 rue de l'Usine - Charguia II - 2035 - Tél. : +216 71 941 541 - Fax. : +216 71 941 889

Annexe : Z.I. Chotrana II - B.P. 160 - 2083 - Pôle Technologique - El Ghazala - Tél. : +216 70 685 685 - Fax. : +216 70 685 454