



**JOÃO ANTÓNIO
ASSIS REIS**

**UM SISTEMA CONVERSACIONAL DE PESQUISA
SOBRE BASE DE DADOS MÉDICAS**

**A CONVERSATIONAL QUERY BUILDER ON
MEDICAL DATABASES**

DOCUMENTO PROVISÓRIO



JOÃO ANTÓNIO
ASSIS REIS

UM SISTEMA CONVERSACIONAL DE PESQUISA
SOBRE BASE DE DADOS MÉDICAS

A CONVERSATIONAL QUERY BUILDER ON
MEDICAL DATABASES

DOCUMENTO PROVISÓRIO

“The greatest challenge to any thinker is stating the problem in a way that will allow a solution”

— Bertrand Russell



Universidade de Aveiro
2024

**JOÃO ANTÓNIO
ASSIS REIS**

**UM SISTEMA CONVERSACIONAL DE PESQUISA
SOBRE BASE DE DADOS MÉDICAS**

**A CONVERSATIONAL QUERY BUILDER ON
MEDICAL DATABASES**

Dissertação apresentada à Universidade de Aveiro para cumprimento dos requisitos necessários à obtenção do grau de Mestre em Engenharia Informática, realizada sob a orientação científica do Doutor João Rafael Almeida, Professor auxiliar do Departamento de Eletrónica, Telecomunicações e Informática da Universidade de Aveiro, e do Doutor José Luís Oliveira, Professor catedrático do Departamento de Electrónica, Telecomunicações e Informática da Universidade de Aveiro.

o júri / the jury

presidente / president

Prof. Doutor João Antunes da Silva

professor associado da Faculdade de Engenharia da Universidade do Porto

vogais / examiners committee

Prof. Doutor João Antunes da Silva

professor associado da Faculdade de Engenharia da Universidade do Porto

Prof. Doutor João Antunes da Silva

professor associado da Faculdade de Engenharia da Universidade do Porto

Prof. Doutor João Antunes da Silva

professor associado da Faculdade de Engenharia da Universidade do Porto

Prof. Doutor João Antunes da Silva

professor associado da Faculdade de Engenharia da Universidade do Porto

Prof. Doutor João Antunes da Silva

professor associado da Faculdade de Engenharia da Universidade do Porto

**agradecimentos /
acknowledgements**

Agradeço toda a ajuda a todos os meus orientadores, amigos e família.

Palavras Chave

texto livro, arquitetura, história, construção, materiais de construção, saber tradicional.

Resumo

Um resumo é um pequeno apanhado de um trabalho mais longo (como uma tese, dissertação ou trabalho de pesquisa). O resumo relata de forma concisa os objetivos e resultados da sua pesquisa, para que os leitores saibam exatamente o que se aborda no seu documento.

Embora a estrutura possa variar um pouco dependendo da sua área de estudo, o seu resumo deve descrever o propósito do seu trabalho, os métodos que você usou e as conclusões a que chegou.

Uma maneira comum de estruturar um resumo é usar a estrutura IMRaD. Isso significa:

- Introdução
- Métodos
- Resultados
- Discussão

Veja mais pormenores aqui:

<https://www.scribbr.com/dissertation/abstract/>

Keywords

textbook, architecture, history, construction, construction materials, traditional knowledge.

Abstract

An abstract is a short summary of a longer work (such as a thesis, dissertation or research paper).

The abstract concisely reports the aims and outcomes of your research, so that readers know exactly what your paper is about.

Although the structure may vary slightly depending on your discipline, your abstract should describe the purpose of your work, the methods you've used, and the conclusions you've drawn.

One common way to structure your abstract is to use the IMRaD structure. This stands for:

- Introduction
- Methods
- Results
- Discussion

Check for more details here:

<https://www.scribbr.com/dissertation/abstract/>

Contents

Contents	i
List of Figures	v
List of Tables	vii
List of Listings	ix
Glossary	xii
1 Introduction	1
1.1 Motivation	2
1.2 Objectives	2
1.3 Dissertation outline	3
2 State of Art	5
2.1 Information Retrieval	5
2.1.1 Traditional Methods	6
2.1.2 Neural Information Retrieval Systems	8
2.1.3 Question Answering	9
2.2 Large Language Models	10
2.2.1 Definition	11
2.2.2 Architecture Overview	11
2.2.3 Optimization Techniques	13
2.2.4 Comparison between foundation Large Language Models	15
2.2.5 Limitations	16
2.3 Conversational User Assistants	17
2.3.1 Overview of Conversational User Assistants	17
2.3.2 Generative-Based Chatbot	18
2.4 Interactive Query Builder	20
2.4.1 General Query Builders	20

2.4.2	ATLAS	21
2.5	Summary	22
3	Conversational Search Assistant	25
3.1	Data	25
3.2	Information Retrieval	27
3.2.1	BioASQ Challenge 2024	28
3.2.2	BM25 Implementation	28
3.2.3	API-driven Approach	31
3.3	Large Language Model	31
3.3.1	Tasks	31
3.3.2	Implementation	32
3.4	Frameworks to streamline biomedical data discovery ¹	34
3.4.1	FlowiseAI	34
3.4.2	Langflow	35
3.5	EHDEN Chatbot Architecture ²	35
4	Query Builder	39
4.1	Implementation Strategy	39
4.1.1	Interaction between Components	39
4.1.2	Template and Questions	41
4.1.3	Template Pointer	42
4.2	Concepts Sets	42
4.2.1	Expression	42
4.2.2	Implementation	43
4.3	Cohort Definition	44
4.4	ATLAS WebAPI	45
4.4.1	Communications with ATLAS	45
4.4.2	Network Interactions	46
5	Results and Discussion	49
5.1	IR methods	49
5.2	LLM Implementation: FlowiseAI vs Langflow	50
5.3	Conversational Search Assistant	51
5.4	Query Builder	51

¹This section is mainly based in the publication *Using Flowise to Streamline Biomedical Data Discovery and Analysis, IEEE 22nd Mediterranean Electrotechnical Conference (MELECON), 2024*

²This section is mainly based in the publication *A chatbot-like platform to enhance the discovery of OMOP CDM databases, 34th Medical Informatics Europe Conference (MIE), 2024*

6	Conclusions	53
	References	55

List of Figures

2.1	Overview of an interaction-based Information Retrieval model: Retrieval and Ranker. . .	8
2.2	The Transformer Block Architecture, from Vaswani <i>et al.</i> [1].	12
2.3	Tokenization process visually explained by OpenAI [35].	13
2.4	Example of Chain-of-thought prompting. The reasoning processes are highlighted in yellow. Adapted from Wei <i>et al.</i> [2].	15
2.5	Retrieval-Augmented Generation Workflow.	20
2.6	Query builder from jQuery QueryBuilder [58].	20
3.1	The diagram of the connection between the data files.	27
3.2	Index data structure: A) All information about the databases is composed of three files, B) First approach to the structure of indexed documents, C) Current approach to the structure of indexed documents.	29
3.3	Workflow implemented in FlowiseAI tool.	34
3.4	Overview of the EHDEN chatbot architecture.	36
4.1	Interaction diagram between the user, the system components and external tools.	40
4.2	The form to create the Concept Set.	44
4.3	Network diagram illustrating the connections between the system and other institutions.	47

List of Tables

2.1	Comparison of foundation Large Language Models.	16
3.1	Overview of the LLM tasks.	32
3.2	LLM tasks parameters and output information.	32
5.1	Comparison between FlowiseAI and Langflow.	50

List of Listings

1	The configuration file of the LLM (<code>params.json</code>)	33
2	The cohort template file (<code>cohort_template.json</code>).	41
3	The cohort questions file (<code>cohort_questions.json</code>).	42
4	A Concept Set expression example.	43
5	The body to create the Concept Set in ATLAS.	46
6	The body to create the Cohort Defintion in ATLAS.	46

Glossary

EHDEN	European Health Data Evidence Network
OMOP CDM	Observational Medical Outcomes Partnership Common Data Model
OHDSI	Observational Health Data Sciences and Informatics
IR	Information Retrieval
TF	Term Frequency
IDF	Inverse Document Frequency
TF-IDF	Term Frequency - Inverse Document Frequency
DL	Document Length
BM25	Best Matching 25
QA	Question Answering
AI	Artificial Intelligence
NLP	Natural Language Processing
NLU	Natural Language Understanding
NLG	Natural Language Generation
LM	Language Models
LLM	Large Language Models
SLM	Statistical Language Models
NLM	Neural Language Models

PLM	Pre-trained Language Models
RNN	Recurrent Neural Networks
BERT	Bidirectional Encoder Representations from Transformers
GPT	Generative Pre-trained Transformer
ML	Machine Learning
RAG	Retrieval-Augmented Generation
RLHF	Reinforcement Learning from Human Feedback

Introduction

In the field of medical research, the lack of sufficient samples to conduct robust and statistically significant studies is still a challenge [3]. This problem has more impact in studies involving rare diseases, specific patient subgroups, or when analyzing rare outcomes in common diseases [4]. This lack of information can lead to underpowered studies, limiting the reliability and generalizability of the research findings [5]. Furthermore, the limitations imposed by insufficient data are significant since individual healthcare institutions often have limited patient populations, which restricts the diversity and volume of data available for research [3], [5].

Multicentre studies are a possible strategy to overcome the issue of insufficient samples used in a study [6]. By aggregating data from multiple sources, researchers may increase the number of participants and the population diversity [7]. Although this solution seems to solve the problem, it raises more challenges, namely the heterogeneity between databases [8] and privacy issues [9]. In other words, databases contain different types, formats, and/or sources, which are often not compatible with each other. The existence of these diverse data is not very effective, as they cannot be easily shared or integrated with other data.

To address the challenge of data heterogeneity in multi-institutional studies, the **Observational Health Data Sciences and Informatics (OHDSI)** initiative proposed some strategies. **OHDSI** is a collaborative network that aims to improve health by empowering a community to collaboratively generate evidence that promotes better health decisions and better care [10]. A key contribution of **OHDSI** is the development of standardized data models and tools that facilitate the harmonization of observational health data [11]. The most notable is the **Observational Medical Outcomes Partnership Common Data Model (OMOP CDM)**. The **OMOP CDM** standardizes the format and content of clinical data, enabling data from different sources to be integrated and analyzed cohesively [10], [12]. The **OHDSI** tools enable researchers to effectively overcome the barriers posed by data heterogeneity since the data would be mapped into the **OMOP CDM** format. This procedure also enhances the reproducibility and scalability of research, contributing to more reliable and impactful healthcare insights [13].

1.1 MOTIVATION

Solving the heterogeneity issues simplified the growth of multi-institutional observational databases. This raised new opportunities and challenges, considering multicentre studies [14]. One of the main challenges is the selection of the most appropriate databases for conducting a study. However, there are already some solutions to address that problem [15].

The **European Health Data Evidence Network (EHDEN)** project is one of several initiatives to simplify multicentre studies. In this project, the consortium members helped European data providers migrate their databases to the **OMOP CDM** data schema and publish their metadata in a web catalogue. This comprehensive catalogue offers researchers a centralized platform to explore the variety of available data sources [16]. With the increasing number of data partners, selecting databases only using metadata becomes complex, which leads to the creation of the **EHDEN** Network Dashboards tool. It gives statistical and aggregated information about the databases available on the network. Researchers get more context to select, manually, the databases that satisfy their study requirements.

Despite these improvements, the catalogue now hosts the information of more than 150 databases. Researchers often spend considerable time and effort to find the most suitable database for their research criteria. The current methods of database discovery often rely on manual search and evaluation, which can be time-consuming and subject to human error, impacting the quality and efficiency of their research. For these reasons, the challenge of efficiently discovering the most appropriate databases for a specific study returned.

1.2 OBJECTIVES

With the increasing adoption of chatbot-like tools, the main objective of this work is to develop a conversational query builder to help medical researchers define their study objectives. This approach offers a more personalized and efficient method of database discovery. Researchers can interact with the chatbot conversationally, specifying their needs and receiving tailored recommendations.

To achieve this objective, the present dissertation seeks to answer the following research question:

How can a conversational query builder support medical researchers when defining a study protocol?

To answer this question, the work can be addressed by focusing on different aspects, namely by dividing it into three stages:

1. Study of state-of-the-art: i) methodologies to build a conversational user interface, ii) procedures to retrieve the databases of most interest, and iii) explore the definition of a study protocol;
2. Developed a chat-like search engine to help discover the best databases for a study;
3. Enhance the engine to collect additional information to provide a query as an outcome.

1.3 DISSERTATION OUTLINE

TODO: resumir a estrutura do documento.

State of Art

This dissertation suggests the development of a conversational query builder that functions as a chat-based interface within the EHDEN project ecosystem. This interface aims to help researchers redefine their studies effectively. The conversational user assistant will return a query to help discover the best databases for a specific research. Studying state-of-the-art Information Retrieval systems to retrieve the most appropriate databases for a researcher's query is essential in this context. In addition, it is vital to explore Large Language Models, a recent advance of algorithms that are known for their language generation ability, and the actual state of conversational user assistants or chatbots. In the end, research about interactive query builder. And so these topics will be discussed below.

2.1 INFORMATION RETRIEVAL

In computing and information science, **Information Retrieval (IR)** involves retrieving information from collections of unstructured data, such as documents. According to Kumar *et al.* [17], an **IR** system requires input user queries and uses them to retrieve information from its collections, aligning with the users' needs. This process efficiently filters and narrows down the vast amount of available unstructured data, thereby preventing users from being overwhelmed by the sheer volume of data and aiding them in quickly accessing relevant and specific content.

In conformity with Hambarde and Proença [18], conventional text retrieval systems were predominant in the initial stages of the **IR** field. These systems mainly depended on matching terms between queries and documents. Nevertheless, these systems based on terms have limitations, including issues like polysemy, synonymy, and linguistic gaps, which may restrict their effectiveness.

With the advancement of technology, deep learning techniques emerged, improving conventional text retrieval systems and overcoming the constraints associated with term-based retrieval methods. For this reason, the performance of these systems increased significantly, resulting in a more accurate and streamlined retrieval of information for end-users [18].

In turn, deep learning methods have evolved. Neural Network Architectures, transfer learning, and pre-training techniques emerged. These approaches have advanced the representation of textual data and bolstered the IR system’s comprehension of natural language queries.

More recently, Transformer architectures with attention mechanisms have been implemented in IR systems to enable more effective handling of complex queries and documents. Moreover, incorporating pre-trained (masked) language models like **Bidirectional Encoder Representations from Transformers (BERT)** [19] and **Generative Pre-trained Transformer (GPT)-2** has proven to enhance the performance of IR systems, offering an advanced understanding and processing of context, capturing the relationships and nuances within the natural language query and the documents.

This field has many applications in the real world, such as search engines like Google, digital libraries, and e-commerce platforms. In compliance with Kumar *et al.* [17], IR generally functions across three main scales: searching the web, retrieving personal information, and conducting searches for enterprises, institutions, and domain-specific contexts.

This section explores the IR field, more specifically, traditional methods, neural IR systems and how IR can be joined with natural language.

2.1.1 Traditional Methods

Some successful classical methods are **Term Frequency - Inverse Document Frequency (TF-IDF)** and **Best Matching 25 (BM25)**, briefly explained next.

TF-IDF

To better understand the **TF-IDF** method, first, it is necessary to understand the concepts of **Term Frequency (TF)** and **Inverse Document Frequency (IDF)**. Manning *et al.* [20] explained these concepts in their work. It is reasonable to assume that a document containing a query term more frequently is more relevant to that query. Therefore, it should be assigned a higher relevance and/or score. So, **TF** is the number of term occurrences in a document. However, to evaluate the relevancy of a query, each term is regarded with equal importance, and this is the problem with the raw method explained above. Manning *et al.* [20] clarified this with the following example: the automotive industry is expected to include the term "auto" in nearly every document. The **IDF** calculates the rarity of a term across a set of documents. This measure is calculated as the logarithm of the inverse fraction of documents containing the term. The goal is to help prioritize some terms that are sporadic and possibly more informative.

The traditional IR method, **TF-IDF**, combines **TF** and **IDF** definitions, as the name suggests, and then produces a weight for each term in each document, as Manning *et al.* [20] and [21] mention. Equation 2.1 shows that the weight is calculated as the product of **TF** and **IDF** values, highlighting terms that are both important within a specific document and relatively uncommon in the document collection.

$$\text{tf-idf}_{t,d} = \text{tf}_{t,d} \times \text{idf}_t. \quad (2.1)$$

Manning *et al.* [20] noted the **TF-IDF** weight assigned to a term in a document is highest when the term frequently appears in a few documents, providing discriminating solid power. The weight is lower when the term occurs less frequently in a document or is widespread across many documents, indicating a less pronounced relevance signal. The weight is at its lowest when the term is present in nearly all documents. However, **TF-IDF** does not consider the semantic information of words, which limits its ability to accurately reflect the similarity between texts [21]. To address this limitation, researchers have proposed hybrid methods that combine **TF-IDF** with semantic understanding to calculate text similarity.

In summary, this **IR** method evaluates the importance of a term within a document relative to its occurrence across a collection of documents.

BM25

BM25, the short form for Best Matching 25, is a ranking algorithm for **IR** systems, especially in the context of search engines. Hambarde and Proença [18] noted that **BM25** and other initial retrievers are employed for their effectiveness in recalling pertinent documents from an extensive pool.

The core components of **BM25** include **TF**, **IDF**, **Document Length (DL)**, and tuning parameters. Recapping from the **TF-IDF** section, **TF** is the number of occurrences that a specific term is in a document, and **IDF** is a measure that indicates the importance of a term in the whole document. Equation 2.2 gives the formula to calculate the **BM25** score.

$$score(D, Q) = \sum_i^n IDF(q_i) * \frac{f(q_i, D) * (k1 + 1)}{f(q_i, D) + k1 * (1 - b + b * \frac{fieldLen}{avgFieldLen})}. \quad (2.2)$$

As Gomedé [22] explained, the **BM25** equation is composed of the i th query term (q) and the respective **IDF** and **TF** values. Also, include a division between the **DL**, represented in the formula as $fieldLen$, and the average document length, $avgFieldLen$. This ratio evaluates how much the length of the document field deviates from the average length. So, the Connolly [23] explained intuitively: a document tends to receive a lower score when it contains more terms, especially those that do not match the query. The value b is a fine-tuning parameter, and it is responsible for length normalization. When b is larger, the ratio has a more significant effect on the overall score. Finally, the $k1$ value means term frequency saturation. It is a fine-tuning parameter that prevents the term frequency component of **BM25** from having an unlimited impact on the document score.

This algorithm is simple and effective in **IR** tasks, mainly search tasks. Also, it can handle large collections. For these reasons, it is widely used and called a classic. However, **BM25** can not perform a semantic analysis of the query and the documents, so getting the context and, in turn, getting better results is challenging.

While traditional information retrieval methods like **TF-IDF** and **BM25** have been effective in matching queries to documents based on keyword frequencies and document lengths, the evolution of **Artificial Intelligence (AI)**, data complexity and user needs has led to the development of neural **IR** systems.

2.1.2 Neural Information Retrieval Systems

Neural IR systems represent a significant advancement in the field of IR. Conforming to Mitra and Craswell [24], these systems utilize neural network models and deep learning techniques to understand and interpret the semantic content of queries and documents, offering a more nuanced approach. Neural IR systems can be broadly categorized into two types: Representation-based and Interaction-based models.

Conforming to Chen *et al.* [25], the representation-based model focuses on creating representations of both queries and documents. They employ techniques to convert text into high-dimensional vector spaces. In these spaces, semantically similar terms and phrases are represented by vectors that are close to each other, capturing the underlying meaning and relationships of words beyond their surface-level appearances. This approach utilizes architectures like Recurrent Neural Networks (RNN) and allows the system to understand the content of documents and queries on a deeper level.

The interaction-based model examines how words in a query relate to words in a document, capturing complex patterns and dependencies between them [25]. They often use attention mechanisms, as seen in Transformer-based models like BERT, to dynamically weigh the importance of different parts of the text based on their relevance to the query.

However, only the interaction-based model will be explored.

There are some approaches to adopting an interaction-based model. One of them is the Retrieval and Ranker. According to Hambarde and Proença [18], this method can be separated into two stages, as the name suggests: Retrieval and Ranker. The following figure, adapted from Hambarde and Proença [18], shows us an overview of interaction-based IR systems, highlighting the two main stages.

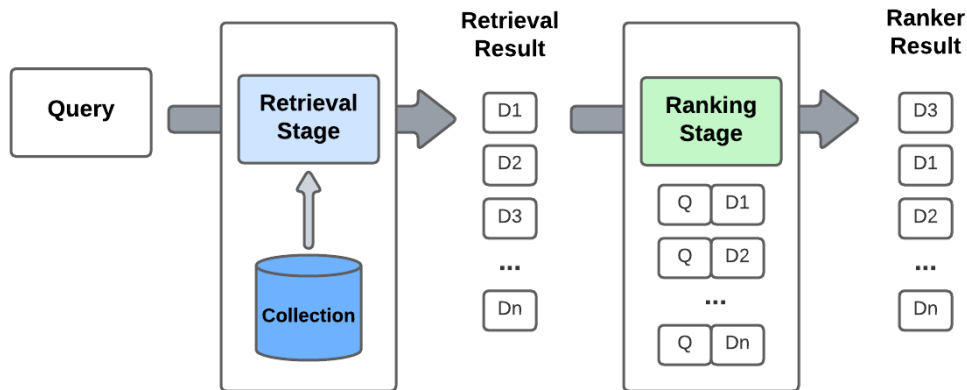


Figure 2.1: Overview of an interaction-based Information Retrieval model: Retrieval and Ranker.

After analyzing the query, the retrieval stage will select an initial set of documents that are potentially pertinent to the query, as shown in figure 2.1. Subsequently, the relevance of these documents undergoes reassessment through the similarity scores.

This is followed by the ranking stage, in which the primary objective is to adjust the order of the initially retrieved documents based on their relevance scores using a neural reranking,

like **BERT** [25]. This phase prioritizes the enhancement of result effectiveness rather than efficiency. In the end, it returns a rank of documents as close as possible to the user's query criteria.

2.1.3 Question Answering

Natural Language Processing (NLP) is the basis for building a **Question Answering (QA)** system, so it is important to give an overview of this technique. It is a field of **AI** whose primary goal is to understand, interpret, and generate human language. The **NLP** can be divided into two major components: **Natural Language Understanding (NLU)** and **Natural Language Generation (NLG)**, according to Ayanouz *et al.* [26]

The **NLU** component is focused on enabling machines to understand and interpret human language in a meaningful way. **NLU** involves processing and analyzing natural language data to comprehend its meaning, context, sentiment, and intent.

In agreement with Ngai *et al.* [27], the user queries can be processed by semantic analysis, pragmatic analysis, and syntactic analysis. Ayanouz *et al.* [26] explained these steps and added two more necessary steps to make it easier to understand: a lexical analysis and discourse integration.

- **Lexical Analysis:** This step involves analyzing and identifying the structure of words. It breaks down the text into chapters, sentences, phrases, and words. Chizhik and Zhrebtsova [28] defined lexical analysis as the pre-processing of the text following the steps: tokenization, removal of special characters, links, and punctuation, and removal of stop-words.
- **Syntactic Analysis:** The syntactic analyzer parses the grammar and arrangement of words, making the relationships between different words more explicit. Essentially, it rejects sentences with incorrect structures. This analysis can be seen as the process of normalizing tokens.
- **Semantic Analysis:** This step ensures the text is meaningful and interprets its correct meaning by mapping syntactic constructions. It ensures that only semantically valid content is retained. The recognition of entities is part of this analysis.
- **Pragmatic Analysis and Discourse Integration:** This step analyzes the overall context to derive the conclusive interpretation of the actual message in the text. It considers factors like the true meaning of a phrase or sentence based on the broader context.

The other component is **NLG**. Language generation is responsible for crafting coherent and linguistically accurate responses. Simply put, it grapples with the challenge of navigating the intricacies of natural human language [27].

Backtracking, **QA** is a subfield of **IR** and **NLP**. According to Zhong *et al.* [29], **QA** focuses on providing a single and specific answer to a question posed in natural language. Unlike **IR**, which aims to return a broad range of relevant information or documents in response to a query, **QA** seeks to pinpoint and provide one precise answer.

The traditional approach to question analysis and answering often involves mapping questions into predefined templates, such as "What-type" and "How-type". While widely utilized by existing online question-answering search engines, this template-based approach faces limitations in handling multiple questions [29].

So, with the advancement of technology, another approach emerged: deep learning-based question-answering. In contrast with the traditional approach, this approach employs deep learning techniques, like **RNN**, to offer automatic representation and analysis of questions. These neural models, trained through end-to-end approaches, excel in extracting and understanding complex characteristics in textual documents.

Recently, deep learning approaches with attention mechanisms and transfer learning have enhanced the flexibility of representation in text classification and named entity recognition. Zhong *et al.* [29] highlights **BERT** that has emerged as a powerful model, utilizing contextualized representations for transfer learning. **BERT**-based models showcase performance in question-answering tasks, even in domains like medicine.

2.2 LARGE LANGUAGE MODELS

It is crucial to trace briefly the development history to understand the concept of **Large Language Models (LLM)**. Liu *et al.* [30] explained this simply and intuitively.

Before **LLM**, there were only simple **Language Models (LM)**, a subfield of **NLP** and **AI**, that have been called foundation models. A **LM** is a statistical model used to predict the next word in a sequence of words. It calculates the probability of a given word occurring in a sequence, helping to determine which words are likely to appear next in a given context [31].

Most of these predictive models were based on probabilities and Markov assumptions, also known as **Statistical Language Models (SLM)**. This was heavily dependent on feature engineering. Afterward, as deep learning gained prominence, an architecture designed to learn data features automatically; in other words, neural networks for **NLP** emerged to enhance **LM**'s capabilities. Integrating feature learning and model training, **Neural Language Models (NLM)** established a comprehensive neural network framework applicable to diverse **NLP** tasks [30].

Most recently, the launch of the Transformer Block Architecture by Vaswani *et al.* [1] revolutionized this field. These deep-learning architectures led to the development of pre-trained models not explicitly designed for a particular task, including **BERT** and **GPT**, collectively known as **Pre-trained Language Models (PLM)**. **PLM** have shown significant performance enhancements across various **NLP** tasks.

Following this, the researchers have involved the scale of model parameters, and the paradigm of "Pre-train, Prompt, Predict" like Liu *et al.* [30] call, gained widespread acceptance. So, in terms of interaction with **LM**, the prompts became crucial. Researchers name these **PLM** with hundreds of billions of parameters as **LLM**. Prompts effectively allow **LLM** to deal with a large number of complex and diverse tasks without a lot of effort.

This section discusses mainly **LLM**, exploring briefly their architecture and comparison between foundation **LLM**. Finally, it addresses some of its limitations.

2.2.1 Definition

LLM is an advanced **LM** and belongs to generative **AI**. It is designed to comprehend and generate text that is coherent and contextually relevant, engaging in human language interactions. Essentially, these advanced **AI** systems mimic human intelligence. These models have a notable ability in natural language tasks, such as text generation and translation, **QA**, decision-making, summarization, and sentiment analysis.

These models can process and predict patterns with accuracy. Hadi *et al.* [32] combine sophisticated **SLM** and deep learning techniques to train, analyze, and understand huge volumes of data, learning the patterns and relationships among the data. For this reason, according to Naveed *et al.* [33], when provided with task descriptions and examples through prompts, **LLM** can produce textual responses to task queries.

Liu *et al.* [30] say that the release of ChatGPT 1 garnered significant social attention, and research into **LLM** triggered more interest. This has led to the development of noteworthy products like PaLM, **GPT-2**, **GPT-3**, and, most recently, **GPT-4**, and LLaMA and LLaMa-2.

LLM are revolutionizing **NLP** and **AI** research.

2.2.2 Architecture Overview

This subsection discusses the architecture of a **LLM**, which is supported by the Transformer Block architecture. Also, explains the pre-training process of this advanced **LM**.

Transformer Architecture

The development and advancement of **LLM** is thankful for the introduction of Transformers by Vaswani *et al.* [1] in 2017. Most **LLM** are built on the Transformer model, which is based on a multihead self-attention mechanism and feedforward layers. This new technology enables parallelization and efficient handling of long-range dependencies, according to Hadi *et al.* [32], and led to the development of models that have achieved enormous results, such as **GPT** by OpenAI and **BERT** by Google. The architecture of this revolutionized model is shown in figure 2.2.

The innovation of this model is due to the multihead self-attention mechanism, one of the key components [1] [32]. It allows the model to weigh the importance of different words in a sequence when processing each word. This mechanism enables the model to focus on relevant information, capturing dependencies regardless of word order.

However, the key advantage of this multihead self-attention mechanism is its highly parallelization [1]. This characteristic enables the Transformer model to be easily distributed and trained on a large scale using GPUs. The ability to parallelize computations means that Transformers can handle larger datasets and more complex tasks, unlike previous architectures like **RNN**, where sequential processing of data was required.

Since the model doesn't have recurrence and convolution to understand the order of the input sequence, another component, Position Encoding, provides some information about the position of the tokens in the sequence. This is crucial for capturing sequential information in the data.

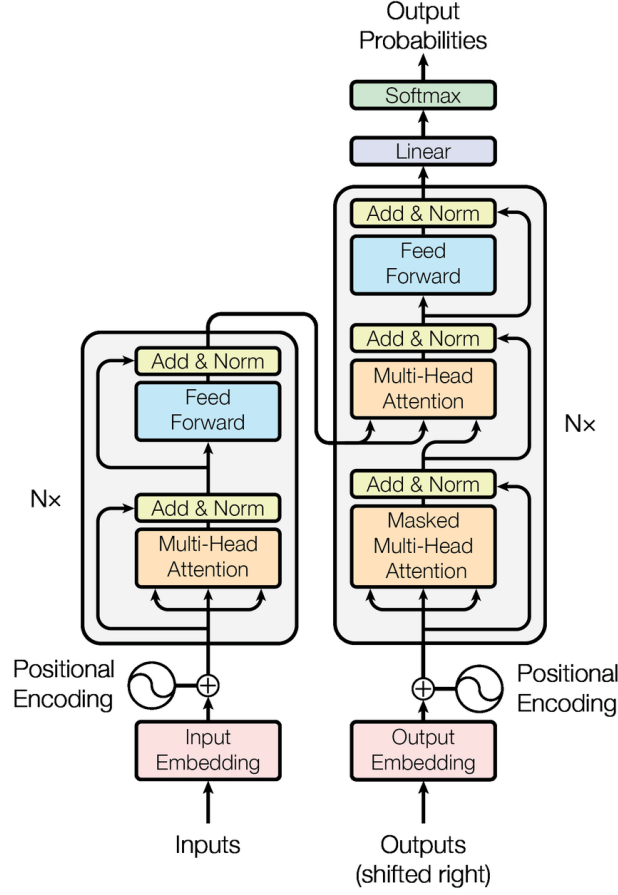


Figure 2.2: The Transformer Block Architecture, from Vaswani *et al.* [1].

Pre-training Process

Learning the patterns and relationships among the data starts with the pre-training process. In compliance with Min *et al.* [34], first, the LLM needs to access a vast volume of textual data from multiple sources. The goal of this phase is to predict the succeeding word in a sentence based on the context given by the previous words through unsupervised learning.

According to Hadi *et al.* [32], preparing and preprocessing the data before the training stage is necessary to achieve this. First, demand quality filtering from the training corpus. It is vital to remove unwanted, repetitive, duplicated, superfluous, and potentially harmful content from the massive text data. Next, it is necessary to pay attention to privacy. The data could have sensitive or personal information, so it is vital to address privacy concerns by removing this information from the pre-training corpus.

An important step, the tokenization, follows this [32]. This step aims to divide the unprocessed text into sequences of individual tokens, which are subsequently input into LLM. Moreover, it is vital in mitigating the computational load and enhancing efficiency during the pre-training phase. Figure 2.3 visually presents the tokenization process [35] carried out and explained by OpenAI.

After the pre-training process, the LLM goes through an optimization phase.

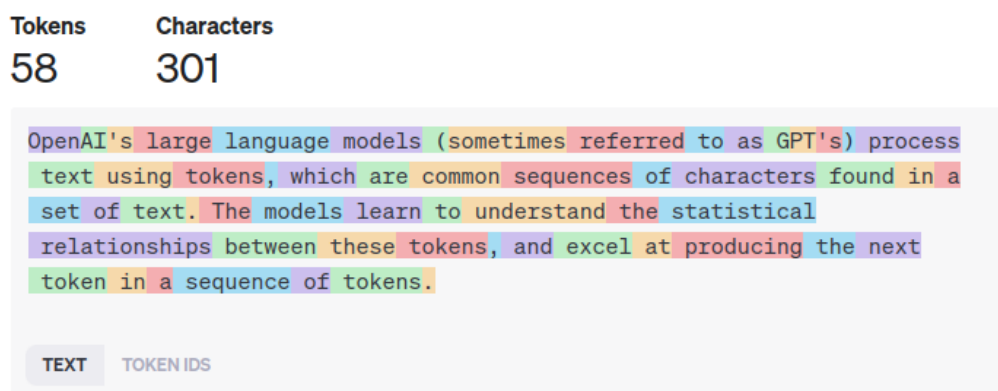


Figure 2.3: Tokenization process visually explained by OpenAI [35].

2.2.3 Optimization Techniques

There are some techniques to optimize the tasks and the accuracy of the **LLM**. Two of these techniques, Fine-tuning and Prompt Engineering, are explained next.

Fine-tuning

During pre-training, models are generally trained with the objective of next token prediction, learning the nuances of language structure and semantics. According to Kamnis [36] and Hadi *et al.* [32], the fine-tuning phase involves adapting a pre-trained model to specific tasks and aligning it with human preferences, improving the performance on particular domains.

In this stage, the model is presented with labeled data to produce more contextually accurate responses for the specific task. Fine-tuning enables the **LLM** to specialize in diverse applications, ranging from language translation and question-answering to text generation.

Some approaches could be applied to fine-tune the model. Naveed *et al.* [33] distinguishes some of them, such as Parameter-Efficient Tuning. As **LLM** typically requires a lot of computational resources, like memory and computing, the Parameter-Efficient Tuning approach is helpful because it allows the model to train by updating fewer parameters, adding new ones, or selecting existing ones. Inside this approach, there are also some different methods. The commonly used indicated by Naveed *et al.* [33] are Prompt Tuning, Prefix Tuning, and Adapter Tuning.

The Prompt Tuning method integrates trainable tokens, named soft prompts, to the beginning or within the input of a **LLM**, and only these tokens are adjusted during training to adapt the model for a specific task. This method keeps the rest of the model unchanged, ensuring the core knowledge and capabilities of the model are preserved while it learns to handle new types of requests or information.

In Prefix Tuning, a sequence of trainable tokens is introduced to transformer layers, with only the prefix parameters undergoing fine-tuning, while the remaining model parameters remain unchanged. These added prefixes function as virtual tokens, allowing input sequence tokens to attend to them during processing.

Meanwhile, in Adapter Tuning, small modules called adapters are added inside each layer

of the Transformer. These adapters can be trained to adapt the model for specific tasks. The fine-tuning process works by slightly altering the model’s internal features, allowing it to learn task-specific patterns without changing the entire model. LoRA (Low-Rank Adaptation) is one technique that implements Adapter Tuning, introduced by Hu *et al.* [37]. Instead of adding new layers like traditional adapters, LoRA learns low-rank matrices that are used to update the weights of the existing layers, maintaining the original weights of the model. This approach allows for efficient fine-tuning of the model on specific tasks while maintaining the model’s original performance and avoiding significant increases in computational costs.

Prompt Engineering

With the emergence of LLM, other research fields were born. Prompt Engineering is one of these cases and has been widely applied. In compliance with Meskó [38] and Ma *et al.* [39], this emerging field involves designing, refining, and implementing prompts or instructions to direct the generated output of LLM, aiding in diverse tasks. LLM can follow specific directions provided by users in natural language after being tuned with instructions.

There are some techniques of prompt engineering, such as Chain of Thought (CoT) and Reason-Action (ReAct).

CoT is a popular problem-solving approach for prompt engineering that aims to break complex tasks into multiple and simpler subtasks and solve them. So, Wei *et al.* [2] explained that this method involves explicitly modeling the reasoning processes that lead to a final answer, rather than directly generating an answer. This explanation of reasoning often leads to more accurate results. Figure 2.4 explains the effect of this technique.

ReAct is a prompt technique introduced by Yao *et al.* [40]. The idea behind this is to simultaneously include both reasoning and action within a single prompt. To solve a complex task, ReAct consists of three tasks for every subtask: 1) **Reason** involves analyzing the current situation and determining the necessary steps; then, 2) **Action** entails executing a task based on the reasoning. 3) **Observation** then refers to examining the outcomes following the action.

Meskó [38] raise a series of recommendations for more effective LLM prompts: it must be as precise as possible; providing the setting and the context of the question is essential; describe the goal of the prompt first; give a role to the LLM to get more context (for example, "You are a math teacher and explain the natural numbers"); continuous LLM prompt refinement; prefer open questions over close-questions. Regularly testing prompts in real-world situations is crucial, as their effectiveness is most accurately assessed through practical application.

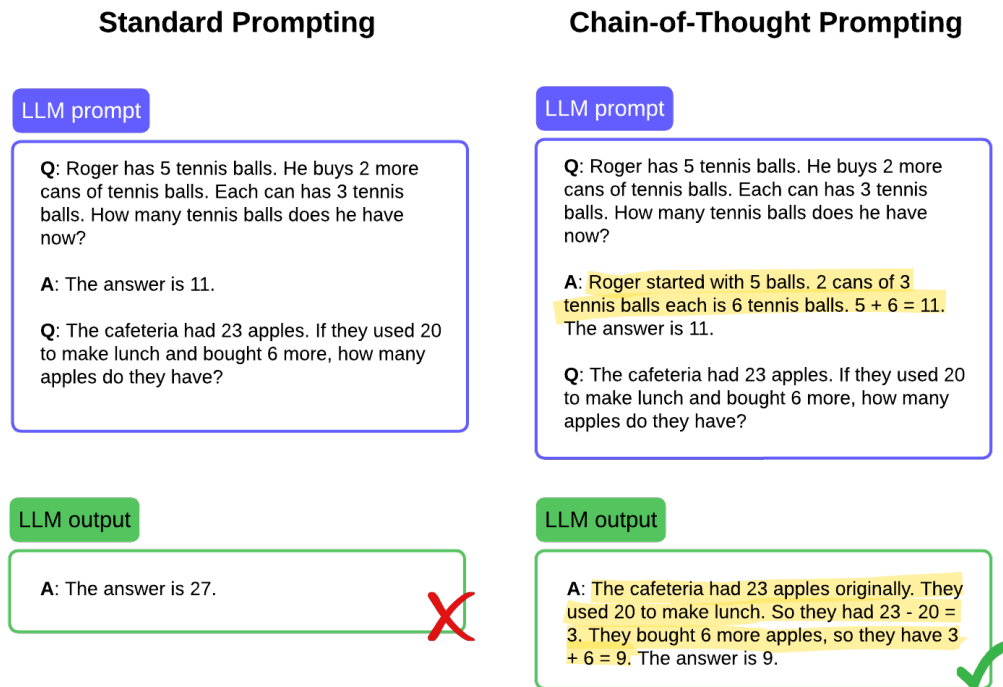


Figure 2.4: Example of Chain-of-thought prompting. The reasoning processes are highlighted in yellow. Adapted from Wei *et al.* [2].

2.2.4 Comparison between foundation Large Language Models

The best way to compare **LLM** is to evaluate the model's performance. Hadi *et al.* [32] identified five factors to make this comparison: the size of the training corpus, the quality of the training corpus, the number of parameters, the complexity of the model, and some test tasks.

The primary foundation models of **LLM** are **GPT-4** by OpenAI, **LLaMA 2** by Meta, **PaLM 2** by Google, and **Falcon** by Technology Innovation Institute (TII). These **LLM** are provided by big companies and have outstanding progress in the evolution of this area. These models gave rise to many others.

LLaMa 2 is an open source **LLM** by Meta (Touvron *et al.* [41]). **LLaMa 2** was trained on 40% more data than **LLaMa 1**, the model from which it came, and has double the context length. So, the model size of **LLaMa 2** is 7 billion, 13 billion, or 70 billion parameters. With 4096 context length and trained on 2 trillion pretraining tokens, this **LLM** is commonly fine-tuned for chat use cases. Many other models, like **Alpaca**, **Vicuna**, and **Llama-2-chat**, came from **LLaMa** and deserve further analysis. It is accessible for both research and commercial purposes

The recent **GPT** model from OpenAI, **GPT-4**, is a closed source **LLM** (OpenAI *et al.* [42]). Trained on a meticulously curated dataset from various textual sources, including books, articles, and websites, **GPT-4** exhibits remarkable performance with text and image inputs. It is the **LLM** behind **ChatGPT**. It has 32 000 context length. OpenAI has chosen to provide limited technical details about the training methodology used for this advanced

model, including specific information on parameter counts.

The Google generative chatbot, Bard, uses as **LLM** the PaLM 2 model developed by Google (Anil *et al.* [43]). It emerged from PaLM with 540 billion parameters. PaLM 2 is a closed source **LLM**, and, following the OpenAI approach, has opted to disclose limited technical specifics, including the number of parameters.

The Falcon **LLM** is an open-source model with impressive performance and scalability (Almazrouei *et al.* [44]). There are three variations of the model size: 7 billion, 40 billion, and the most recent, 180 billion of parameters. This Falcon 180B is equipped with an impressive 180 billion parameters and trained on 3.5 trillion tokens. It is accessible for both research and commercial purposes.

The table 2.1 summarizes the important aspects of comparison between this foundation **LLM**.

Model	Provider	Model size (Parameters)	Context Length	Tokens	Fine-tuneability	Open-source
GPT-4	OpenAI	-	-	-	No	No
LLaMa 2	Meta	7B, 13B, 70B	4096	2T	Yes	Yes
PaLM 2	Google	-	-	-	No	No
Falcon	TII	7B, 40B, 180B	2048	3.5T	Yes	Yes

Table 2.1: Comparison of foundation Large Language Models.

2.2.5 Limitations

It is safe to say that **LLM** are significantly impacting the world. According to Liu *et al.* [30], this is justified by their abilities, mainly in-context learning, reasoning for complex content, instruction following, and creative capacity.

However, **LLM** has some limitations. Hadi *et al.* [32] address some of them, and the most important ones are biased responses, hallucination, explainability, and cyber-attacks.

We already know that **LLM** are pre-trained with extensive training data. But suppose that data contains some biased information related to factors such as gender, socioeconomic status, and/or race. In that case, this may result in analyses and recommendations that are discriminatory or inaccurate across diverse domains. The problem of bias applies not only to training data but also to user interaction bias, algorithmic bias, and contextual bias. The user interaction bias means that, as user prompts shape responses, and if users consistently ask biased or prejudiced questions, the model may acquire and reinforce these biases in its replies.

A severe limitation that is an active area of research is hallucination. Church and Yue [45] characterized **LLM** hallucinations as when the model attempts to fill gaps in knowledge or context, relying on learned patterns during training. Such occurrences can result in inaccurate or misleading responses, detrimental to the user and the model’s reliability.

The way the **LLM** makes decisions is unknown. Comprehending the decision-making process of a complex model with billions of parameters, like **LLM**, is challenging. So, the explainability of these models is a big limitation [32]. Sometimes, it is necessary to decipher the factors that influenced an **LLM**’s decision and this limitation poses difficulties in offering a clear and concise explanation. In vital sectors like healthcare, where decisions carry substantial

consequences, ensuring transparency and the capability to elucidate the model’s predictions is essential.

Another limitation is the cyber-attacks. A **LLM** can suffer some prompt injections from a malicious user to extract sensitive information from the model, according to Kshetri [46]. This is called the Jail Break attack [32]. Another attack is Data Poisoning Attacks, which consist of data poisoning strategies to manipulate the model’s output.

Furthermore, Liu *et al.* [30] highlighted another limitation: the temporal lag of the training corpus. **LLM** cannot retrieve information in real time, and the answer generated may not be the most current.

It is important to be aware of these limitations.

2.3 CONVERSATIONAL USER ASSISTANTS

Conversational User Assistants, also known as chatbots, chatterbots, or virtual assistants, have become a vital aspect of the digital landscape. These tools are generally dialogue systems that understand, interpret, and generate human language, enabling them to communicate with users to dissolve their questions.

Chatbots are increasingly being used in various contexts due to their many benefits. These aspects that make companies bet on the use of chatbots are the continuous availability to support and assist the customer, ensuring more consistent support; the cost-efficiency by reducing the human customer support; the time-saving both for the organization and for customers due to the immediate responses to the user queries; the ease and intuitiveness of this systems; and, improve service with every interaction [47]. Because of this, the utility of the chatbots as tools is increasing as the technology advances.

The rise of conversational user assistants is underpinned by a convergence of technologies, specifically by **LLM**.

This section provides a brief overview of chatbots, followed by an in-depth focus on Generative-Based Chatbots. It encompasses important concepts from their definition to some techniques employed in their development and optimization.

2.3.1 Overview of Conversational User Assistants

To provide a comprehensive understanding of conversational user assistants, it’s crucial to first explore some of their characteristics, such as domains and method of response generation.

Nuruzzaman and Hussain [48] defined the differences between chatbots with opened or closed domains. In an open-domain environment, conversations can go in any direction without a predefined goal or intention. Conversely, in closed-domain environments, the conversation is centered on a particular topic. A closed-domain chatbot is designed with a clear objective.

It is important to differentiate chatbots in their way of giving or generating a response based on an input query. Peng and Ma [49] distinguish three main types of chatbots based on their response generation: Rule-based, Retrieval-based and Generative-based chatbots.

A **rule-based chatbot** examines fundamental features of the user’s input statement and generates a response based on a predefined set of manually crafted templates. This type is

more applicable in a closed-domain conversation. ELIZA, introduced by Weizenbaum [50], was the first chatbot that applied this primitive technique.

A **retrieval-based chatbot** picks a response from an extensive precompiled dataset. It selects the most promising reply from the top-k ranked candidates. Thus, they refrain from producing new text. It has limited flexibility regarding closed-domain and in terms of errors [51].

A **generative-based chatbot** generates a text sequence as a response rather than choosing it from a predefined set of candidates. These chatbots are very flexible and can handle open domains because they are implemented with deep learning techniques. The interactions will be more identical to those of humans, as it implements a self-learning method from a large quantity of interaction data [49] [51]. However, this could be complex and costly to implement.

The only type that will be covered will be generative chatbots, due to their capabilities.

2.3.2 Generative-Based Chatbot

Generative-based conversational user assistants are chatbots that use generative models to generate natural language responses. These chatbots utilize sophisticated deep-learning techniques, such as **LLM**. **LLM**, as described in section 2.2, has the ability to understand and generate human-like text in context. This advanced **LM** has been widely used in the modern chatbots. ChatGPT is an example of this type of chatbot.

Although an **LLM** can generate text from a query, it is not prepared to be applied to a chatbot. There are some techniques for improving and optimizing the model to behave like a chatbot, such as **Reinforcement Learning from Human Feedback (RLHF)**.

In addition, some techniques aim to combat some of the limitations of chatbots, such as hallucination, by increasing their knowledge, such as **Retrieval-Augmented Generation (RAG)**.

The **RLHF** and **RAG** are explained in more depth below.

Reinforcement Learning from Human Feedback (RLHF)

In the context of **AI**, according to Li *et al.* [52], **RLHF** is a popular approach in which an agent learns how to perform a task based on evaluative feedback provided by a human observer.

RLHF in generative-based chatbots is a topic that has been explored in the field of conversational **AI**. This technique is a transformative technique that combines reinforcement learning and supervised learning to refine **LLM** for chatbot applications. Tran *et al.* [53] stated that **RLHF** aims to align chatbot responses to human preferences, improving chatbots' performance and making them more human-like.

The process encompasses several crucial stages, in conformity with Axelsson *et al.* [54]. Initially, the **LLM** is pre-trained on a large dataset of text, which allows it to learn a wide range of language patterns and knowledge. After pre-training, the model undergoes a phase of supervised fine-tuning. In this phase, the **LLM** is trained on a dataset of conversational examples that are specifically curated to reflect the desired outputs for the chatbot.

After that, humans provide feedback on the model’s outputs. This feedback is crucial because it is used to build and train a reward model. The reward model learns to predict the quality of the model’s responses based on the human-provided feedback [54].

The **LLM** is further fine-tuned using reinforcement learning, where it learns to generate responses that maximize the predicted reward, using the reward model created in the previous phase. This stage enables the model to optimize its responses through the reward model, which is based on human feedback.

The process often involves several iterations of feedback and fine-tuning to continually improve the chatbot’s performance. This can resolve errors, improving the refining the conversational style.

Retrieval-Augmented Generation (RAG)

RAG is a subfield of **NLP** and **AI**. This approach was introduced by Lewis *et al.* [55] in 2020 and combines retrieval-based and generative models to enhance content generation and information retrieval processes. For a clearer comprehension of this method, Gao *et al.* [56] made a survey into **RAG** systems and distinguish the parametric knowledge from non-parametric knowledge.

Traditionally, **LLM** can adapt their knowledge and responses to a specific domain by fine-tuning models with parameters. This is **parametric knowledge** because the **LLM** knowledge is provided through the model’s training data. However, entirely parameterized **LLM** have limitations, including data not currently updated and hallucinations. The **non-parametric knowledge**, provided by external information sources, emerged to solve these limitations. This non-parametric knowledge approach is known as **RAG**.

So, according to Lewis *et al.* [55], **RAG** involves retrieving pertinent information from external knowledge bases, giving more context to the **LLM**. This enables the **LLM** to access and utilize up-to-date and/or domain-specific information to enhance response accuracy and relevance. This process leads to reducing hallucinations.

The figure 2.5 shows the workflow of a **RAG**. This model aims to retrieve relevant information from an extensive corpus of documents when answering questions, subsequently adding this information in the prompt as context to enhance the quality of predictions in compliance with Lewis *et al.* [55].

Gao *et al.* [56] explain simply the workflow. The first step is 1) Retrieve information from an external data source, such as a vector database, as in the example in figure 2.5. This step utilizes **IR** models, such as **BM25**, to retrieve relevant information based on the query. The second step is 2) Augment, improving the **LLM** prompt with the context retrieved in the previous stage. The last step is 3) Generation. Using the prompt with context, the **LLM** generates a response to the query, based on the external information retrieved.

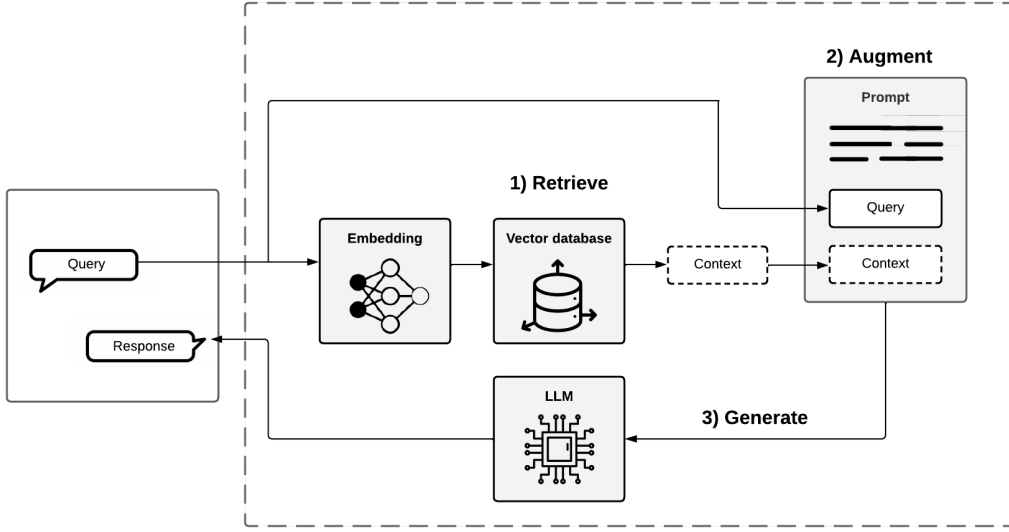


Figure 2.5: Retrieval-Augmented Generation Workflow.

2.4 INTERACTIVE QUERY BUILDER

A query builder is a user interface tool for dynamically searching and filtering database objects, constructing a query according to user preferences, as the work of Mussa *et al.* [57] shows. This query could be in different formats, such as SQL and JSON. This tool lets users construct queries visually, eliminating manual research or coding.

This section is particularly significant as there is limited documentation on conversational query builders. Therefore, it documents the general workings of query builders and delves into a detailed explanation of the functioning of the ATLAS cohort definition.

2.4.1 General Query Builders

Users interact with the query builder through a user-friendly interface. Figure 2.6 shows a query builder interface from jQuery QueryBuilder [58].

Figure 2.6: Query builder from jQuery QueryBuilder [58].

Users can add rules and conditions/groups with some clicks. Each rule typically consists of a field, an operator, and a value. The conditions/groups could be an AND or an OR. Using

figure 2.6 as an example, there is a group with two elements joined with a condition AND: a rule and another group. This other group is composed of two rules joined with a condition OR.

As users build their queries, the query builder internally represents the conditions in a structured format, often a structured JSON of rules and groups, that reflects the logical structure of the query [58].

In summary, a query builder simplifies creating complex queries by providing a visual and interactive interface, making it more accessible [59].

There are several advantages of its use [59]: offers a user-friendly interface with menus, operators, and suggestions to facilitate the creation of accurate queries; users, often without direct permissions to modify the data source, can leverage the query builder to transform datasets without making changes to the underlying database; and, the generated queries are easily modifiable, allowing for flexibility in adjustments or repetitions.

2.4.2 ATLAS

ATLAS, an open-source and web-based software application, is freely accessible and was created by the OHDSI community. It aids in helping researchers conduct scientific analyses on standardized observational data converted to the OMOP CDM.

Using healthcare claims data, researchers can define cohorts by categorizing groups of people according to their exposure to a medication or their diagnosis of a certain health condition. ATLAS offers the functionality to search medical concepts, enabling the identification of cases with particular conditions or drug exposures. Moreover, it allows for the examination of patient profiles within a given cohort, providing a way to visualize the healthcare records of specific subjects.

There are some different definitions of cohort, but, in the OHDSI research, according to OHDSI [60], a cohort is a query that defines a set of persons who meet certain inclusion criteria over a specified duration.

Cohorts serve as fundamental units for addressing research questions. A key characteristic of these cohorts is their independent definition. The distinct structure facilitates their reuse across different research contexts.

Cohort definition

There are two approaches to building a cohort: rule-based and probabilistic. The most popular one is the rule-based cohort definition, which uses explicit rules to describe when a patient is in the cohort.

In ATLAS, the process of defining a cohort is composed of 3 stages [60]: Cohort Entry Events, Inclusion Criteria, and Cohort Exit.

The creation of a cohort starts with **Cohort Entry Events**, defining the initial event criteria. This involves the primary identification of the population of interest, which might include users of a certain drug, individuals with a specific diagnosis, or a combination of factors.

The concept set needs to be specified in the Cohort Entry Events. It is a collection of standardized medical concepts used to define clinical elements like diseases, drugs, or procedures. For instance, if the study is about diabetes, the concept set will include various codes representing diabetes in different medical terminologies. These sets ensure that the cohort captures all relevant instances of the condition or exposure across different healthcare data sources.

Additional initial event criteria can also be added to refine the population further, such as the event occurring within a certain time frame.

After defining the initial event, the next step is to establish **Inclusion Criteria**. These criteria are based on a combination of domain-specific attributes to further refine and specify the cohort population, ensuring that it aligns closely with the research objectives. The inclusion criteria can be based on a range of factors such as age limits, the presence of certain symptoms, or a specified duration of medication use.

Finally, defining the **Cohort Exit** criteria is crucial for determining when individuals no longer belong to the cohort. This stage is important for studies where the duration of membership in the cohort is relevant to the research question.

In compliance with OHDSI [60], a well-defined cohort specifies how a patient enters a cohort and how a patient exits a cohort.

After defining the cohort, it is possible to generate an SQL code with the query to get the list of individuals who meet the criteria. ATLAS also facilitates the reuse of cohort definitions across different studies by allowing users to export and import cohort definitions in JSON format. This enhances the efficiency and reproducibility of research within the OHDSI network.

A cohort definition can be seen as a query builder, a little different when compared to other general query builders.

2.5 SUMMARY

To sum up, the proposed query builder is a generative-based chatbot with a closed domain. Closed-domain chatbots are specialized in specific areas, offering precise responses. Generative chatbots use advanced **LM** to create dynamic responses closer to human interactions.

The utilization of **LLM** allows improvements in the **NLG** capabilities of chatbots and guides conversations more effectively, especially in the task of defining cohorts in medical research. LLaMa-2 and Falcon appear to be good options to implement since they are open-source and have the possibility of fine-tuning the model. Fine-tune has the role of optimizing chatbot performance, alongside Prompt engineering and **RAG**. However, for this specific case, I don't think much external knowledge will be needed, so the use of fine-tuning should be enough.

In terms of **IR**, in order to retrieve the most interesting databases according to the user's needs, the **BM25** technique proves to be a good balance between effectiveness and efficiency. **BM25** is a sophisticated yet relatively straightforward algorithm that improves upon the traditional **TF-IDF** approach. Neural **IR** systems, like Interaction-based models, show good results in the **IR** tasks, but are complex and the neural networks require substantial

computational power. It is a lot simpler to implement than the Neural IR systems, as well as requires fewer computer resources.

ATLAS by OHDSI provides a cohort definition, which is a query builder to define groups of people based on the research question using healthcare claims data. However, the interface revealed not very user-friendly and intuitive because it requires the user to have a good knowledge of its use and important concepts. Therefore, a chatbot that builds a cohort definition improves the user experience by making it more intuitive and autonomous.

Conversational Search Assistant

The **EHDEN** Portal is a web-based platform that facilitates access to information, resources, and tools for data partners, researchers, and other stakeholders involved in the **EHDEN** project. The portal fosters collaboration, knowledge exchange, and innovation among diverse partners, providing a user-friendly interface to access standardized healthcare data and related resources. **EHDEN** delivers a catalogue of medical databases across Europe, offering researchers a centralized platform to explore available medical data sources. However, with the increasing number of databases in the catalogue, **EHDEN** built a tool named Network Dashboards¹ to help researchers to choose the best databases across the catalogue. **EHDEN** Network Dashboards is a complementary tool of this ecosystem that provides statistical and aggregated information about the databases available on the network.

However, the catalogue's growth is inevitable, making it difficult and time-consuming for a researcher to find his databases of interest. The proposed solution is to build a conversational search tool for the **EHDEN** catalogue.

This section describes how the **EHDEN** Catalogue search assistant and its components were implemented and the decisions and steps made over time.

3.1 DATA

This section details the structure and contents of the data used in the project, explaining their content, and illustrating how they interconnect to provide a comprehensive overview of the databases in the **EHDEN** network. The data used for this project is the real **EHDEN** data from the catalogue and the Network Dashboards. IEETA, a partner of **EHDEN**, provides this data.

There are four main files that contain the data necessary to obtain an overview of the databases available in the **EHDEN** network: the countries file, the data sources file, the medical concepts file, and the Achilles Results file. It is essential to understand the content of each file and the relationships between their data, as illustrated in Figure 3.1.

¹<https://github.com/EHDEN/NetworkDashboards>

To comprehend the purpose of the four files, each is detailed with fields, description and an example of entries. When the data is private and sensitive, not-real examples are included to better understand the content and connections between files.

Countries file (`countries.csv`).

- **Fields:** `id`, `country`, `continent`.
- **Description:** This file contains of real-world data, listing countries and their corresponding continents. It is used to provide geographical context for the analyses, and supporting studies that require demographic segmentation.
- **Example Entries:**
 - 233, Ukraine, Europe
 - 149, Montenegro, Europe

Data sources file (`data_sources.csv`).

- **Fields:** `id`, `name`, `acronym`, `hash`, `release_date`, `country_id`.
- **Description:** This file includes essential details of the databases such as the source identification, the database name and the hash code of the **EHDEN** catalogue.
- **Example Entries:**
 - 1, MEDIBASE, Medical Database for Health Information Exchange, <hash>, NA, 107
 - 2, GRAVITAS, Global Repository of Advanced Vaccine Innovations Technologies and Strategies, <hash>, NA, 228

Medical concepts file (`concepts.csv`).

- **Relevant fields:** `concept_id`, `concept_name`, `domain_id`, `vocabulary_id`.
- **Description:** The concepts file contains metadata on medical concepts, extracted from OHDSI Athena².
- **Example Entries:**
 - 2966436, Latanoprost (Apotex) 0.005% Eye Drops 2.5 Ml Bottle, Drug, AMT, Containered Pack, NA, 1009551000168106, 2016-11-01, 2099-12-31, NA
 - 2966944, Letrozole (Apotex) 2.5 Mg Tablet, Drug, AMT, Trade Product Unit, NA, 1009571000168102, 2016-11-01, 2099-12-31, NA

Database summaries (`achilles_results.csv`).

- **Relevant fields:** `data_source_id`, `analysis_id`, `stratum_n` (from 1 to 5), `count_value`, statistical information.
- **Description:** Metadata summary of metadata present in the databases providing aggregated analytics. The “`stratum_1`” field points to the concept ID.
- **Example Entries:**

²<https://athena.ohdsi.org/>

- 143, 1, NA, 8507.0, NA, NA, NA, 988296, NA, NA, NA, NA, NA, NA, NA, NA
- 138, 2, 8532.0, 8532.0, NA, NA, NA, 1354216, 526.0, 52600.0, 26300.0, 52.6, 25774.0, 52.6, 157.79, 368.2, 499.7

The combination of these four files gives the necessary information about the databases. Each row of the Achilles Results file contains statistical information about each concept available in the database, and a set of characteristics that are commonly used to filter databases in a cohort study. For instance, the number of samples segregated by range of age. The Figure 3.1 shows how these files content are joined.

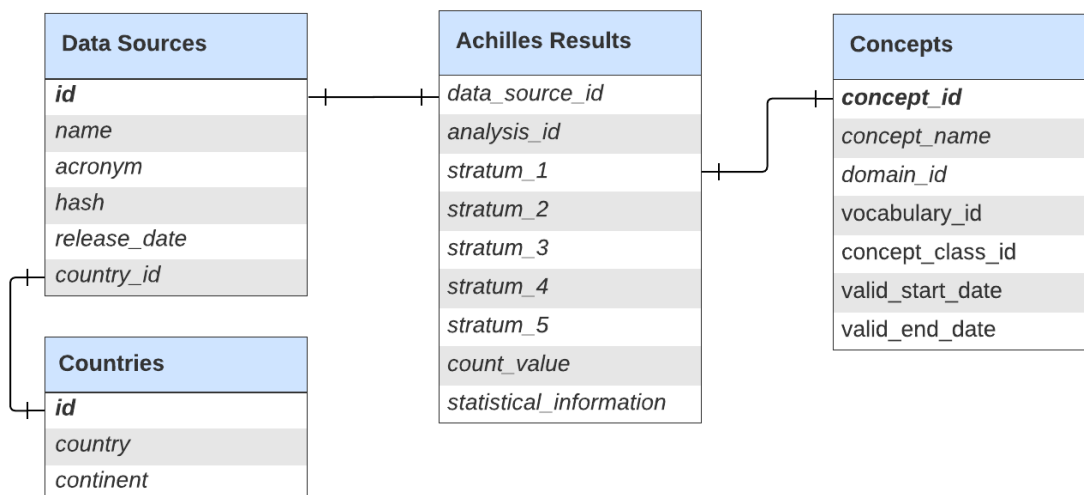


Figure 3.1: The diagram of the connection between the data files.

3.2 INFORMATION RETRIEVAL

After understanding the data and its connections, an **IR** component is crucial to this project to find the most suitable databases. This section discusses mainly the **IR** methods tested, the process of indexing and searching in the collection and an API-driven approach of the **IR** component.

BioASQ³ addresses the information access problem for biomedical experts by organizing challenges in biomedical semantic indexing and **QA**. These challenges cover tasks such as hierarchical text classification, machine learning, **IR**, **QA** from texts and structured data, and multi-document summarization.

This section outlines the **IR** component, covering the BioASQ challenge and its relevance, the implementation of BM25 for indexing and searching databases, and the API-driven approach for integrating **IR** functionalities with external systems.

³<http://bioasq.org/>

3.2.1 BioASQ Challenge 2024

The BioASQ includes several challenges. The IEETA team was involved in 'BioASQ Task 12b', focusing on biomedical semantic QA. This challenge aims to advance the development of systems capable of understanding and answering biomedical questions. Participants must respond to test questions using various types of information, including relevant concepts, articles, and snippets. The challenge involves 5,000 training questions with gold-standard answers and introduces 500 new test questions, all constructed by European biomedical experts. The challenge consists of two phases, A and B:

- **Phase A:** Participants respond to released questions with relevant articles and snippets.
- **Phase B:** Participants provide exact and ideal answers based on questions and provide relevant articles and snippets.

To choose and validate the IR method for this project, I joined the IEETA team and was involved in the 2024 edition of the BioASQ challenge. My task was implementing and testing some IR methods to determine which performs better. This task allowed me to have more concrete results. The techniques used and tested were BM25, SPLADE[61], and BGE-M3[62]. BM25 has already been explained, so here's a brief overview of the other methods tested.

SPLADE

SPLADE⁴ is a neural retrieval model that learns query/document sparse expansion through the BERT Masked Language Model head and sparse regularization, according to Formal *et al.* [61]. This technique belongs to Learned Sparse Retrieval because it combines elements of traditional sparse retrieval techniques with machine learning, particularly deep learning.

SPLADE is designed to balance the effectiveness of dense retrieval models with the interpretability and efficiency of sparse representations. So, it is advantageous because it combines the benefits of dense and sparse models [61]. This method can improve the relevance and accuracy of the search results, particularly in complex queries where understanding the context and the semantic relationships between terms is essential. The model used was `naver/splade-cocondenser-ensembledistil`, detailed by Formal *et al.* [63].

BGE-M3

In compliance with Chen *et al.* [62], BGE-M3⁵ is a multifunctional technique because it can simultaneously perform the three common retrieval functionalities of embedding model: dense retrieval, multi-vector retrieval, and sparse retrieval. Also, this technique is multilingual because it supports over 100 languages and is multi-granular because it can process inputs ranging from short sentences to long documents, accommodating up to 8192 tokens. The model used was `BAAI/bge-m3` with 1024 of dimension and 8192 tokens of sequence length.

3.2.2 BM25 Implementation

Although we have tested various methods, the implementation of BM25 has shown promising results. BM25 is the selected method for implementing the IR component in this

⁴<https://github.com/naver/splade>

⁵https://github.com/FlagOpen/FlagEmbedding/tree/master/FlagEmbedding/BGE_M3

project. The technique was implemented using the PyTerrier PISA, a python interface to PISA⁶.

This section explains the implementation of BM25, the process of creating documents based on the EHDEN data, and the indexing and searching of these documents.

Database Indexing Process

The IR component must index a collection of documents that thoroughly represent all the concepts within the databases, ensuring a clear understanding of what each database encompasses. To achieve this, the data mentioned in the section 3.1 is used to create the documents.

Figure 3.2 represents the different stages of the data. The first part of the figure, A), represents the three data files detailed in the section 3.1. These are received from the EHDEN Network Dashboard. One contains information about the data sources, another contains the metadata summary of the databases (Achilles Results), and the last contains all medical concepts used in this community. This three files data was combined and readjusted to be indexed as documents. The remaining parts of this figure (B and C) represent the strategies adopted to create the documents of each database.

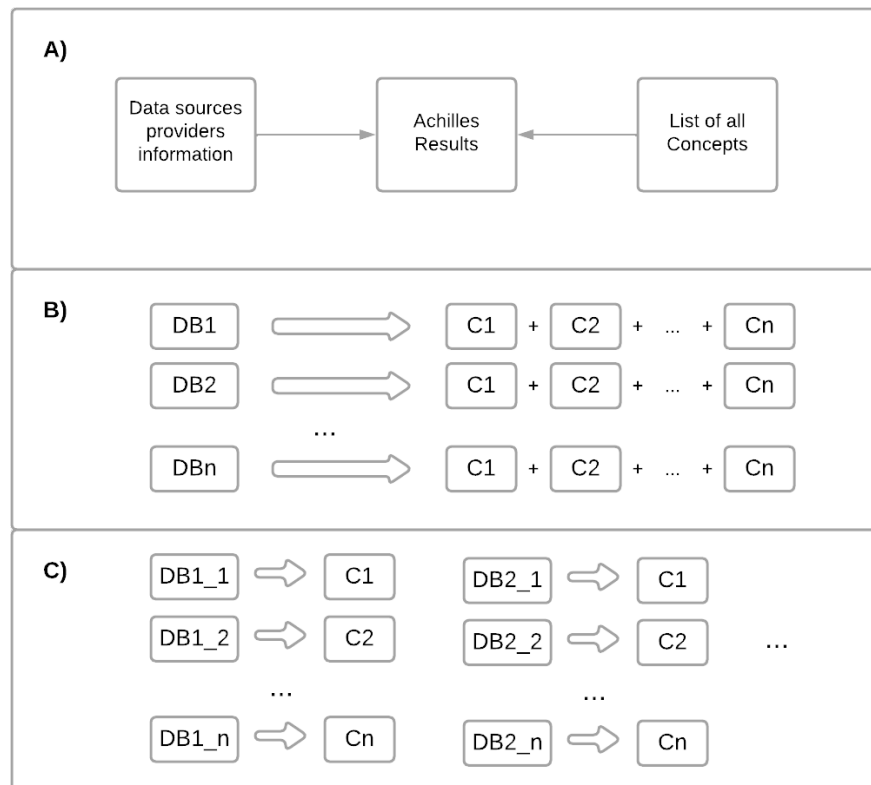


Figure 3.2: Index data structure: A) All information about the databases is composed of three files, B) First approach to the structure of indexed documents, C) Current approach to the structure of indexed documents.

⁶https://github.com/terrierteam/pyterrier_pisa

The first approach to structure of the documents to be indexed is shown in Figure 3.2, part B). The document structure consists of a pairing of database IDs with their corresponding content. The content is represented as a single string, where each concept from each database is concatenated and separated by spaces. However, this way of indexing the database contents has several drawbacks:

- **Lack of Granularity:** By concatenating all concepts into a single string, the individual relevance of each concept is lost. When a search query is made, the search engine treats the entire string as a single document, reducing search results' precision and recall.
- **Poor Search Results:** The search engine returns unsatisfactory results because the indexing does not accurately reflect the importance of each concept. This leads to low search scores, resulting in commonly retrieving irrelevant documents.

So, the approach to creating the document structure has been refined to enhance the granularity and searchability of the database content. Each concept in the database is individually indexed, as shown in Figure 3.2, part C). Each concept within a database is treated as a separate document for indexing purposes. This is achieved by assigning a document number that uniquely combines the database ID and the concept's position within its list. The document number serves as a unique identifier for each concept, ensuring that each piece of data can be individually retrieved and queried. The content of each document, represented by the concept itself, allows for reflecting the importance of concepts within the database, facilitating more precise and efficient retrieval of information.

Database Searching Process

When the IR system receives a query in free text, the system processes it using NLP techniques applied to the query to enhance the search performance. These processes can include tokenization (splitting the text into individual terms or words), stemming or lemmatization (reducing words to their base form), and stop-word removal (eliminating common words that do not contribute to the search). Standardizing the text and transforming the query into important terms improves the match between the queries and indexed concepts as documents.

After processing the query, the BM25 algorithm evaluates the relevance of each concept within the database to the query, assigning a score that reflects its relevance. The method returns the 100 most relevant results and scores each concept. Then, these concepts are grouped by databases. The total relevance score for each database is calculated by summing the scores of its concepts. This comprehensive compilation process highlights the most relevant concepts. To ensure that only the most pertinent databases are presented to the user, the engine applies a predefined threshold to filter out databases with lower cumulative scores. The remaining results are then sorted in descending order of their total scores, prioritizing those with the highest relevance to the query.

Each element in the ranked list of databases contains the following information about the databases: a unique numeric identifier, the database name, the hash that identifies the database to be added to the link to access in the EHDEN Catalogue, the total relevance score

assigned by the **BM25** algorithm, indicating how well the database matches the query, and the concepts list that has been identified.

3.2.3 API-driven Approach

This project's **IR** component is accessible through an API, allowing seamless integration with external systems. By exposing the endpoints of the **IR** component, the API facilitates the consumption of search functionalities, ensuring that various external applications can interact with the data effectively.

The API offers endpoints for submitting queries and retrieving search results. When a query is sent to the API, it undergoes the same searching steps described in the previous section 3.2.2.

This API-driven approach enhances the **IR** component's interoperability and ensures that external systems can seamlessly query and retrieve data. This enables efficient and precise identification of the most suitable databases for any given query, facilitating the integration of the **IR** component in the following project implementation steps.

3.3 LARGE LANGUAGE MODEL

To achieve the proposed goal of making a conversational search assistant, a **LLM** has a crucial role in generating coherent, contextually relevant text and engaging in human language interactions. The **LLM** used in this project is **Nous Hermes 2 Mixtral 8x7B** [64] open model, finetuned from Mixtral [65]. This model belongs to the llama family and has a 47B of parameters.

The reason for using this **LLM** model is that it is a promising model, and the IEETA has installed it for use in its projects. The **LLM** is installed a local version of the Ollama⁷ framework and deployed in a Virtual Machine to access the **LLM** using a URL.

3.3.1 Tasks

In the context of this project, between the multiple applications and capabilities, the **LLM** has two types of tasks: text generation and text analysis. So, Table 3.1 shows the different tasks given to the LLM and the respective descriptions. Also, it presents the output format and the temperature used in each task. The output format is the generated response format expected from **LLM**. The temperature of a **LLM** refers to a hyperparameter that controls the randomness of predictions made by the model. A higher temperature results in more creative and random outputs. Otherwise, a lower temperature results in more deterministic and focused outputs.

⁷<https://ollama.com/>

Type of Task	Task	Description
Generation	Response	Generate a response to a given message.
	Response with Databases	Generate a response to retrieve the best databases to the user.
	Question	Generate a question related to a given question.
Analyse	Topic	Identify if the user message is related to health. If so, identify the research topic in the field of health.
	Yes or No Answer	Identify if the answer is positive or negative for a respective question.
	Answer to Question	Identify if the answer is a response for a respective question.

Table 3.1: Overview of the LLM tasks.

3.3.2 Implementation

After defining the tasks, the following table 3.2 represents the parameters and the output format of each task.

Type of Task	Task	Output format	Temperature (0 - 1)
Generation	Response	Plain Text	0.5
	Response with Databases		0.1
	Question		0.5
Analyse	Topic	JSON object	0.1
	Yes or No Answer		0.1
	Answer to Question		0.1

Table 3.2: LLM tasks parameters and output information.

Inside the Generation main task, the LLM should generate a simple response, a response with the most suitable databases, and a question to ask the user. The temperature in this type of task is medium to provide a balance between randomness and determinism, in simple words, to be more creative in his answers and not lose the focus of the answer. The output format is Plain Text because it is intended to send messages to the user. These tasks can be viewed as NLG tasks.

In the main Analysis task, the language model should respond in a JSON object with a specific structure as specified in the prompt. This is crucial for the system to comprehend the user’s response. These tasks are related to the NLU tasks.

The model requires specific prompts to generate appropriate responses. Each task assigned to the LLM has its own specific prompt and a corresponding temperature setting, which controls the randomness of the output, as we can see in the Code 1. These parameters — prompt and temperature — are defined in a JSON file (`params.json`).

```

{
  "url": "http://<LLM_URL>:<port>/api/generate",
  "model": "nous-hermes2-mixtral:latest",
  "system": "...",
  "generate_question": {
    "prompt": "Your goal is generate a question. You already chat with the user; you don't
↪ need to welcome him. Just ask the following question (you should rewrite the question):
↪ <question>. This question aims to extract more user research requirements. You should
↪ write a message with plain text, not in JSON.",
    "temperature": 0.5
  },
  "topic": {
    "prompt": "...",
    "temperature": 0.1
  },
  "yesno": {
    "prompt": "...",
    "temperature": 0.1
  },
  "is_answer_to_question": {
    "prompt": "...",
    "temperature": 0.1
  },
  "retriever": {
    "prompt": "...",
    "temperature": 0.1
  },
  "generate_answer": {
    "prompt": "...",
    "temperature": 0.5
  }
}

```

Code 1: The configuration file of the LLM (`params.json`)

This file serves as a configuration file, ensuring that the **LLM** receives the correct instructions and settings for each task, thereby optimizing its performance and ensuring the generated responses meet the desired criteria.

The URL field contains the Ollam **LLM** endpoint to generate text. The model field indicates the **LLM** model used, and the system field refers to the system prompt. System prompt is crucial because it guides how these advanced **AI** models interpret and respond to user queries. It directs the **LLM** behavior and ensures that the generated outputs align with the intended goals.

The following fields are the tasks, with the respective prompt and temperature associated. The prompt of the question generation task, as an example, has a `<question>` string, where the cohort question replaces it later. This prompt refinement is done in every task prompt.

3.4 FRAMEWORKS TO STREAMLINE BIOMEDICAL DATA DISCOVERY⁸

There are some options to integrate the **LLM** in this system. FlowiseAI and Langflow are frameworks that enable build an **LLM** system without worrying with the orchestration flow between components. An overview of these frameworks and, in the case of FlowiseAI, an implementation are presented below.

3.4.1 FlowiseAI

FlowiseAI⁹, an open-source automation tool, plays a pivotal role by facilitating the integration of different **AI** components, combined with **IR** techniques, as Reis *et al.* [66] stated.

FlowiseAI enables the creation of customized orchestration flows for **LLM** with **AI** agents and other tools. The workflows within FlowiseAI consist of interconnected nodes or blocks that represent various actions or operations. The specific workflow implemented is illustrated in Figure 3.3.

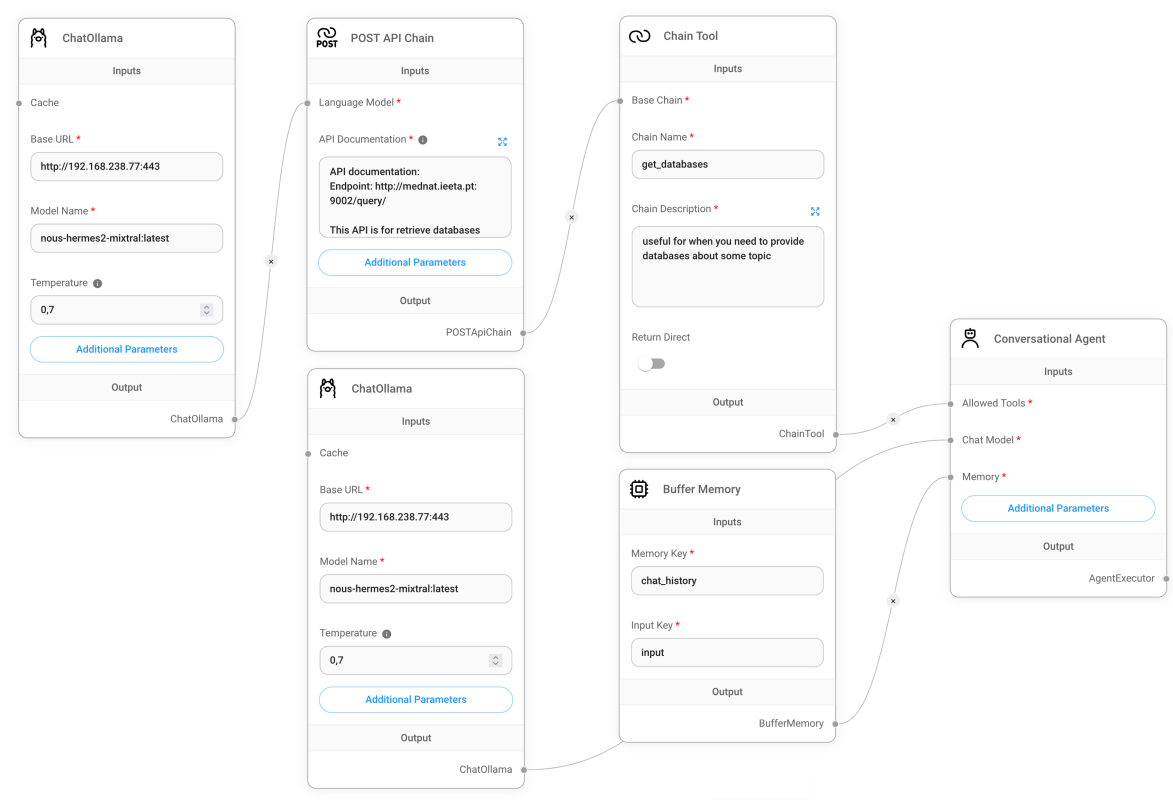


Figure 3.3: Workflow implemented in FlowiseAI tool.

The conversational agent employs a comprehensive approach to enable dynamic interactions on a healthcare **IR** platform. It orchestrates the dialogue flow through the Conversational Agent Node, utilizing an **LLM** to provide a coherent and context-aware user experience. This node is configured with parameters that control tool access, chat model specifications, and

⁸This section is mainly based in the publication *Using Flowise to Streamline Biomedical Data Discovery and Analysis*, *IEEE 22nd Mediterranean Electrotechnical Conference (MELECON)*, 2024

⁹<https://github.com/FlowiseAI/Flowise>

memory capabilities, allowing it to maintain context or state information across interactions for structured conversations.

The core of conversational dynamics is powered by the ChatOllama Node, a chatbot engine designed for processing user queries and the generation of relevant responses. The implementation rely on Ollama since other solutions like ChatGPT API possess privacy issues. The Ollama operates based on a set of parameters including a base URL, alongside model specifications and a temperature setting that modulates probabilistic distribution over the predicted tokens. Furthermore, the Conversational Agent incorporates a Buffer Memory component, essential for the retention of interaction histories and stateful data. This component ensures the persistence of conversational context, a critical feature for enhancing user engagement and response relevance.

To provide the most relevant medical databases, a **RAG** architecture was adopted, as detailed in the section 2.3.2. This approach applied to this scenario involves retrieving the best databases from the **IR** component and adding them to the **LLM** prompt. **RAG** enables the **LLM** to have up-to-date, valid, and domain-specific information to enhance response accuracy and relevance.

The integration with the Chain Tool facilitates the application of prompt engineering techniques, enabling the agent to access the list of the recommended best databases. The tool consumes the endpoint of the **IR** component, which is better described in the section 3.2.2. The conversational agent is equipped with real-time, accurate database recommendations, enhancing the quality of information provided to the user.

3.4.2 Langflow

Langflow¹⁰ is another open-source tool to build **AI** applications. It is also a low-code tool that allows the integration of **LLM** and **AI** components. This tool simplifies the process of creating flows, such as chatflows. Users can drag components from the sidebar onto the canvas and connect them to begin building their applications. The platform allows for exploration by editing prompt parameters, grouping components into high-level components, and creating custom components. This intuitive interface makes Langflow a powerful tool for developing **LLM**-based applications.

3.5 EHDEN CHATBOT ARCHITECTURE¹¹

The system was designed to be integrated as a tool in the **EHDEN** Portal. Therefore, authentication issues were solved by the current mechanisms available on this platform [67]. Therefore, the system was implemented to use the MONTRA2-SDK [15]. This also provided the Network Dashboards tool, an interface to show the metadata, when researchers want to get more details about the suggested databases.

¹⁰<https://github.com/langflow-ai/langflow>

¹¹This section is mainly based in the publication *A chatbot-like platform to enhance the discovery of OMOP CDM databases, 34th Medical Informatics Europe Conference (MIE), 2024*

The previous implementation, detailed in the section 3.4.1, tried to adopt an open-source automation tool to build Chatbot applications, FlowiseAI [66]. However, these become limited to address the new requirements. Therefore, FlowiseAI was replaced by a Python-based backend, developed for this system. In addition to the IR function, the backend also orchestrates chat flow between the various components.

Figure 3.4 represents the key components of the system and their interconnections.

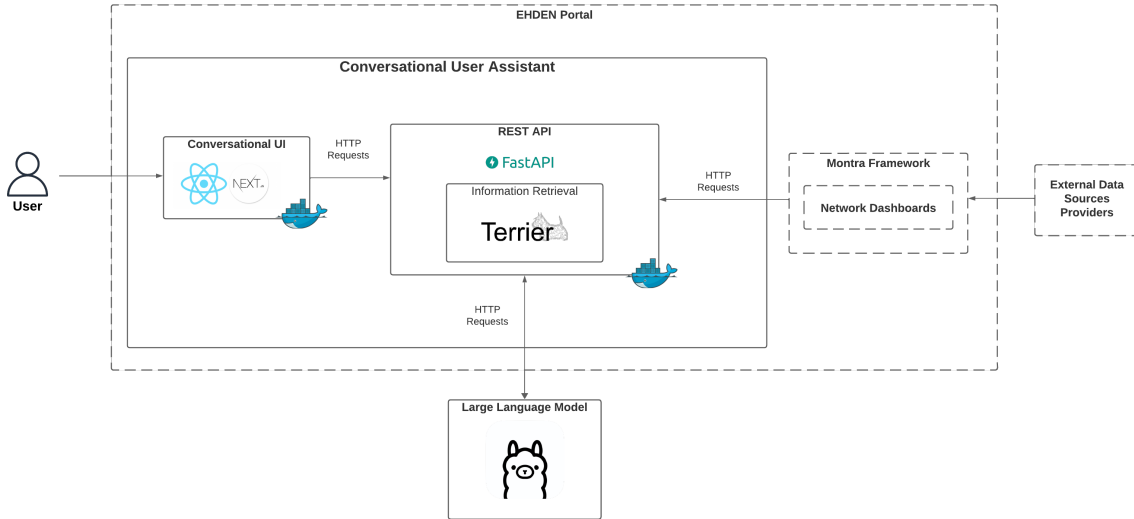


Figure 3.4: Overview of the EHDEN chatbot architecture.

The Conversational User Interface is the primary interface for user interaction, built on the React framework. The User Interface records the user questions and conveys them to the backend to be processed in the component dedicated to IR.

The backend API, built on the FastAPI framework, handles the HTTP requests from the other components. When it receives a query from the User Interface, the backend communicates with the LLM. The LLM is the same open model detailed in the section 3.3, Nous Hermes 2 Mixtral 8x7B. In this implementation, the LLM has two tasks: analyse the topic of health and generation of a response with databases or generate a simple response. If the first task returns false, *i.e.*, if the user message is not related to health, the LLM generates a simple response reminding the user of the chatbot’s purpose. Otherwise, if the first task is true, then the RAG architecture is applied. This means that the generation of a response with databases task is applied. The BM25 retrieves the best databases to the LLM, so it could generate a response with that valid information.

For research purposes, researchers may save their conversations to facilitate data analysis and enhance their research outcomes. When a user decides to save a conversation, they can do so by pressing a save button, which triggers an API endpoint designed to store the conversation. The conversations are subsequently added to a JSON file, creating a structured and accessible format for future analysis. This method ensures that valuable data is systematically archived and easily retrievable for subsequent research endeavors.

Storing these dialogues can be particularly valuable for research purposes, for example,

identify frequently asked questions by medical researchers. By analyzing the saved conversations, researchers can gain insights into common concerns, informational gaps, and the types of support often sought by medical researchers. By taking advantage of this stored information, the system can be improved, or even develop other platforms to improve the overall quality and impact of medical research.

Query Builder

This project has another goal: to enhance the conversational search assistant to collect additional information in order to provide a query as an outcome. The query is referred to as cohort definition in the ATLAS community, as detailed in section 2.4.2.

The following section discusses the implementation of the query builder, as well as the decisions and strategies involved.

4.1 IMPLEMENTATION STRATEGY

To create a system that builds a query through a conversation with the user, the first step is to get a cohort definition generated in the ATLAS platform after implementing an experimental case. Understanding how the cohort is built is crucial at this stage. The section 2.4.2 details defining a cohort. The definition of a cohort requires the definition of a concept set. So, the strategy of the cohort construction is focused first on creating the concept set and then on the rest of the cohort's parameters.

4.1.1 Interaction between Components

It is important to understand how this system interacts with the user and other external tools. Diagram 4.1 describes all the interactions between components. The interaction diagram all the interactions between components. The system components are the IR API and the Chatbot Interface, which contains the User interface and the LLM. The EHDEN Portal and ATLAS are external tools.

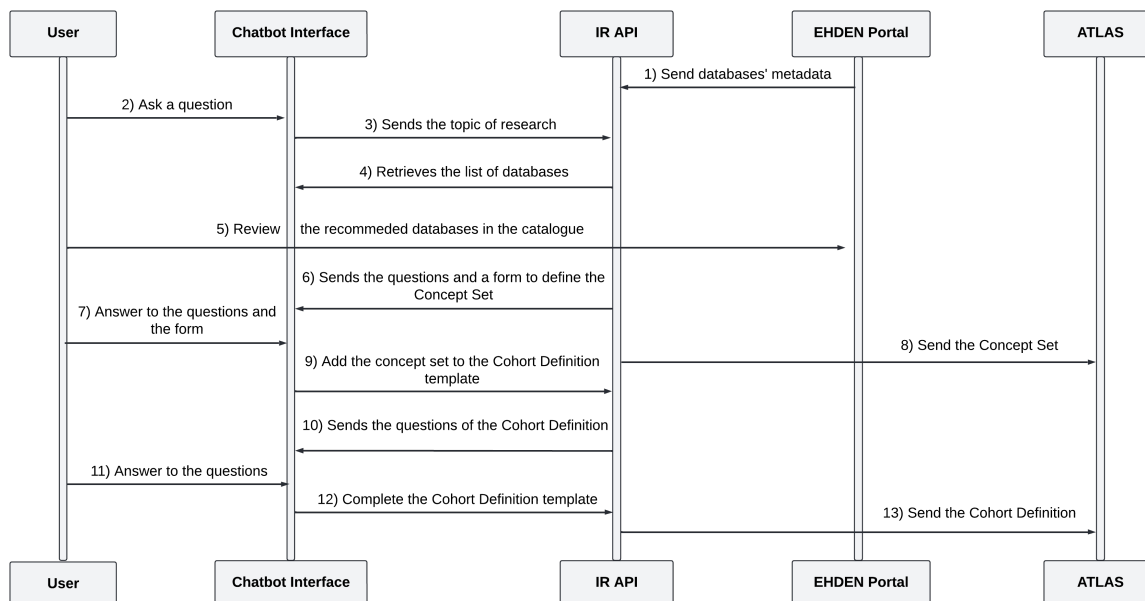


Figure 4.1: Interaction diagram between the user, the system components and external tools.

This diagram presents all the interactions of the conversational query builder, not only presenting the interactions detailed in the previous chapter 3 — steps 1) to 5) —, but also the interactions that will be performed during the detail of this chapter — steps 6) to 13).

Reminding the interactions explained in chapter 3, the **IR** component receives the databases' metadata from the **EHDEN** Portal — step 1) —, creating and indexing the documents. The user asks a question to the Chatbot Interface — step 2). Here, the **LLM** applies his task of identifying the research topic and sends it to the **IR** API — step 3). The **IR** API retrieves the most suitable databases for the research concept and sends the list of databases to the Chatbot Interface — step 4). The user can review the recommended databases in the catalogue in the **EHDEN** Portal — step 5).

Now, the query builder phase starts. As mentioned before, the cohort's construction requires a concept set. So, the **IR** API, when sending the databases to the chatbot interface, also sends questions and then a form to help define the concept set — step 6). The user answers the question and fills out the form — step 7 —, in order to create the Concept Set that could be sent to the **ATLAS** instance — step 8).

The rest of the cohort will be determined now that the concept set is defined. The Concept Set is added to the cohort definition template — step 9). The **IR** API sends a question to complete a field needed for the cohort definition — step 10). After the user answers to it — step 11) —, the **LLM**, in the Chatbot Interface, analyses if the user message is an answer to the question. If so, the answer will be saved in the cohort template — step 12). Steps 10), 11), and 12) are repeated until the cohort template is complete. Finally, the **IR** API sends the cohort definition to the **ATLAS** instance if the user requests it — step 13).

4.1.2 Template and Questions

The ATLAS platform accepts the import of a cohort definition in a JSON format. Two crucial JSON files are needed to generate this JSON object, personalized with the user's requirements information. One is the template for users to fill out during the conversation (cohort_template.json), shown in Code 2. This template is an empty cohort definition JSON structure. The other file, Code 3, is the questions associated with each key field of the template (cohort_questions.json). Each question should be simple and efficient so the medical researcher can respond to it, and its response is the value of that key. Also, each question in the template is manually inserted in the JSON file.

```
{
  "ConceptSets": null,
  "PrimaryCriteria": {
    "CriteriaList": [
      {
        "ConditionOccurrence": {
          "CodesetId": 0,
          "ConditionTypeExclude": null,
          "ConditionSourceConcept": null,
          "First": null
        }
      }
    ],
    "ObservationWindow": {
      "PriorDays": 0,
      "PostDays": 0
    },
    "PrimaryCriteriaLimit": {
      "Type": "First"
    }
  },
  "QualifiedLimit": {
    "Type": "First"
  },
  "ExpressionLimit": {
    "Type": "First"
  },
  "InclusionRules": [],
  "CensoringCriteria": [],
  "CollapseSettings": {
    "CollapseType": "ERA",
    "EraPad": 0
  },
  "CensorWindow": {}
}
```

Code 2: The cohort template file (cohort_template.json).

```

{
  "ConceptSets": "Are there any other concepts you'd like to add? If yes, please add them.
↪ If no, simply respond with 'no'.",

  "PrimaryCriteria": {
    "ObservationWindow": {
      "PriorDays": "Regarding the observation window, what is the minimum number of days
↪ needed before the continuous observation? You must choose from 0, 1, 7, 14, 21, 30, 60,
↪ 90, 120, 180, 365, 548, 730 or 1095.",
      "PostDays": "How many days after event index date? You must choose from 0, 1, 7,
↪ 14, 21, 30, 60, 90, 120, 180, 365, 548, 730 or 1095."
    }
  }
}

```

Code 3: The cohort questions file (`cohort_questions.json`).

4.1.3 Template Pointer

In a nutshell, the cohort questions file have the questions to complete the cohort template file. However, when exchanging messages with the user, tracking which questions must be asked and which question the user responded to is essential and also a problem.

The solution to this issue is creating a pointer to a template key. The pointer retrieves a template key that should be completed with the user information. The pointer points to the same key until the user responds to the question associated with the pointer. Otherwise, it moves to the following template key.

So, during the conversation, the pointer helps:

- To determine the appropriate question to ask the user.
- To identify which key of the cohort template to use to save the user's answer.
- To move on to the next question of the cohort template.

This solution keeps track of the user's responses and dynamically updates the template with the relevant information throughout the conversation.

4.2 CONCEPTS SETS

In Code 2, a cohort is defined by multiple JSON fields. One of these fields is the concepts set (the cohort template key is '*ConceptSets*'), which is a list of concepts required to meet the study requirements of the researcher. In order to properly define a cohort, the concept set must be established first.

A concept set consists of standardized medical terms that define clinical elements like diseases, drugs, and procedures. Therefore, our strategy is to define the concept set in the conversation before moving on to the remaining fields of the cohort template.

This section details the creation of the concept set needed for the cohort definition later.

4.2.1 Expression

A concept set expression is comprised of a list of concepts with the following attributes:

- **concept**: Definition of a concept, using data contained in the concepts file (`concepts.csv`) (section 3.1).
- **isExcluded**: Exclude the concept from the concept set.
- **includeDescendants**: Add all of descendants of a concept, with it also included.
- **includeMapped**: Allow the search for non-standard concepts.

The JSON code 4 below represents a concept definition within a concept set expression. This example is not real information; it is just to define a concept set visually.

```
{
  "id": 0,
  "name": "<CONCEPT SET NAME>",
  "expression": {
    "items": [
      {
        "concept": {
          "CONCEPT_ID": 231256,
          "CONCEPT_NAME": "Covid-19",
          "DOMAIN_ID": "Condition",
          "VOCABULARY_ID": "ASG67HR",
          "CONCEPT_CLASS_ID": "ASG67 code",
          "CONCEPT_CODE": "953635",
          "VALID_START_DATE": "2000-01-01",
          "VALID_END_DATE": "2099-12-31"
        },
        "isExcluded": false,
        "includeDescendants": true,
        "includeMapped": false
      },
      (...)
    ]
  }
}
```

Code 4: A Concept Set expression example.

4.2.2 Implementation

The goal is to create a JSON object identical to the one in the Code 4 with the concepts selected by the medical researcher.

When a user enters a query, the conversational search assistant retrieves the best databases related to the health topic of that query. It should also inquire whether the user has additional concepts for their study requirements. In an implementation vision, the first key to complete in the cohort template is '*ConceptsSets*', so the question regarding that key is rewritten by the LLM and then asked to the user.

The user should respond to whether there are additional concepts to add. If there are, he should indicate them. If not, he should provide a negative answer. The chatbot will continue to ask if there are more concepts until the user confirms that there are no more to add.

Now, in the IR component, the BM25 technique searches in the collection for the concepts the user indicates in his responses. After identifying each concept in the collection using the

IR technique search, additional information from the concept file is added. This includes the concept ID, concept code, and other attributes that define the concept as represented in Code 4.

The list of concepts identified is sent to the chatbot interface. The chatbot interface shows the concepts in a form. For example, Figure 4.2 shows a form to create the concept set after the user mentions his interest in COVID-19. The user gives a name to the concept set and selects the concepts and the other attributes.

Concept Set Name

Covid-19 Concept Set

ID	Concept	Select ?	Excluded ?	Descendants ?	Mapped ?
	COVID-19	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>
	Asymptomatic COVID-19	<input type="checkbox"/>			
	COVID-19 vaccine	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>
	Suspected COVID-19	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
	COVID-19 excluded	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>

Submit

Figure 4.2: The form to create the Concept Set.

Finally, the user submits the forms, creating the concept set. The chatbot will send a message informing the user that they can download the concept set JSON file, which is accepted in the ATLAS instance, to create a new concept set. Alternatively, the chatbot can send the file directly to the ATLAS instance. The pointer of the cohort template moves to the next key.

4.3 COHORT DEFINITION

Once the Concept Set is defined, the list of concepts and their attributes are saved on the template, where the template key points to. Then, the template pointer moves to the next template key. The LLM generates a specific question related to the question that the current key points to. This generated question is then presented to the user through the interface.

The user provides a response to the question presented. The user message could be a possible response to the question or, for example, can be a random message. So, the LLM checks if the user's message is indeed an answer to the question. This task is mentioned in the section 3.3. If the response is valid, the LLM extracts the value and the template is updated

with that value. The template pointer advances to the next key, where the next question will be generated by the **LLM**. Otherwise, if the response is invalid or unrelated, the pointer remains on the same key because the question is not answered yet.

This cycle of getting the question of the pointer, generating the question, presenting it, receiving a response, and verifying the response repeats until the entire cohort is complete, and all keys in the template are filled with the appropriate user responses.

At the end, users can download the cohort definition by clicking on a button sent in a message. The cohort definition is in the right structure for import into the ATLAS platform. Alternatively, the chatbot can send the Cohort Definition directly to the ATLAS instance.

4.4 ATLAS WebAPI

The **OHDSI** community provides a WebAPI¹ that contains all **OHDSI** RESTful services that can be called from **OHDSI** applications. This API has many features, such as providing a centralized API for working with databases converted to the **OMOP CDM**, searching the **OMOP CDM** standardized vocabularies for medical concepts and constructing concept sets. Also, it allows defining cohort definitions for use in identifying patient populations.

The ATLAS platform² is an open-source software provided by **OHDSI**, where can be defined not only the cohort but also the concepts sets. This software consumes the **OHDSI** WebAPI.

This section explains the communication with the external tool ATLAS, emphasizing the interaction of the IR API to assist researchers in creating their concept sets and cohort definitions in ATLAS.

4.4.1 Communications with ATLAS

When the user interacts with the tool and finalizes the definition of a concept set or cohort, the IR API facilitates seamless integration with the **OHDSI** WebAPI by consuming specific endpoints designed for these tasks.

Creating a Concept Set

Initially, for the creation of a concept set, the necessary structure for the body of the **OHDSI** WebAPI endpoint is prepared, as shown in the Code 5. This process begins by making a POST request to the **OHDSI** WebAPI (<https://api.ohdsi.org/WebAPI/conceptset/>). This POST request is responsible for creating the concept set, although at this stage it is devoid of any concept data. The response from this request includes the ID of the newly created concept set within ATLAS.

¹<https://github.com/OHDSI/WebAPI>

²<https://atlas-demo.ohdsi.org/>

```
{
  "id": 0,
  "name": "<Concept Set Name>",
  "description": "Hippocrates costum concept set"
}
```

Code 5: The body to create the Concept Set in ATLAS.

Subsequently, the concepts that were defined throughout the user's interaction are sent to the PUT endpoint (<https://api.ohdsi.org/WebAPI/conceptset/<Concept Set ID>/items>). This PUT request populates the previously created concept set with the specified concepts, completing the definition process.

Defining a Cohort

Following the concept set creation, the process of defining a cohort involves preparing the necessary structure for the body of the POST endpoint (<https://api.ohdsi.org/WebAPI/cohortdefinition/>). The body structure for this request is as shown in the Code 6.

```
{
  "id": 0,
  "name": "[Hippocrates] Cohort {current_time}",
  "expressionType": "SIMPLE_EXPRESSION",
  "description": "Hippocrates custom cohort definition",
  "expression": "<Cohort Definition>"
}
```

Code 6: The body to create the Cohort Definition in ATLAS.

Here, the "expression" field contains the cohort definition that was constructed during the user's interaction. This POST request sends the cohort definition to the OHDSI WebAPI, where it is created and stored in the ATLAS platform.

User Interaction and Submission

Upon completion of either the concept set or cohort definition, the user is presented with an option to either download the JSON object, which can be imported into ATLAS manually, or to send it directly to ATLAS through the chatbot interface. The chatbot provides a message with two buttons, each corresponding to the creation of a new Concept Set or Cohort Definition on the OHDSI platform. By clicking the appropriate button, the defined concept set or cohort is automatically submitted to ATLAS, ensuring a streamlined and efficient workflow for the researcher.

4.4.2 Network Interactions

The interaction between various components and institutions is crucial for the effective functioning of the query builder and the communication with the ATLAS platform. Figure

4.3 illustrates these interactions in detail.

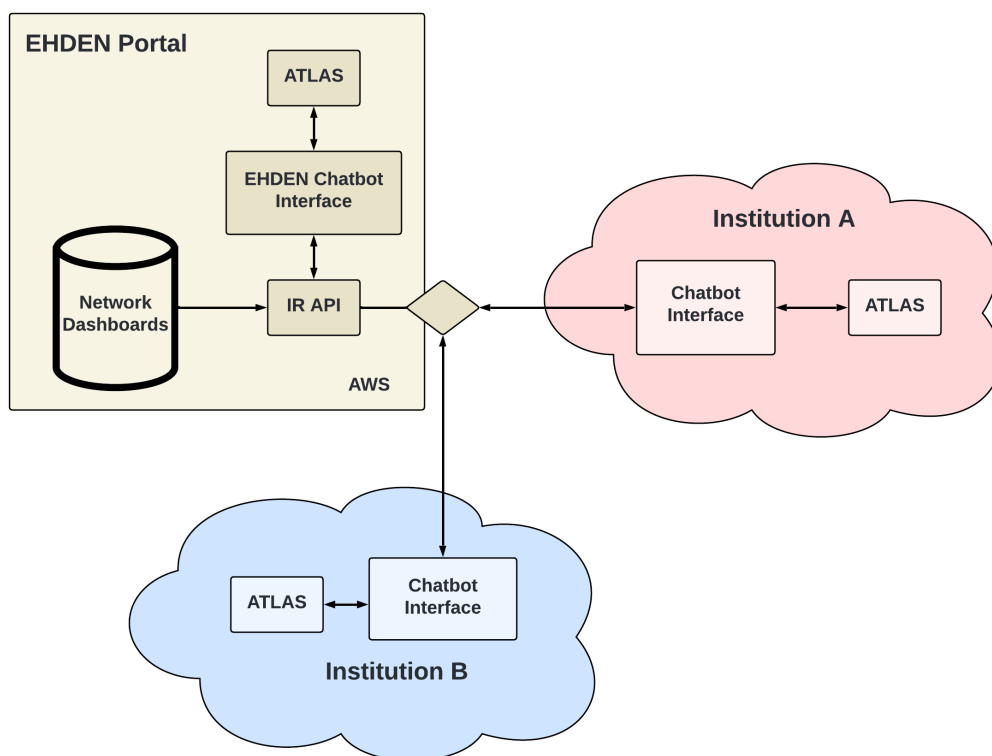


Figure 4.3: Network diagram illustrating the connections between the system and other institutions.

The **EHDEN** Portal serves as the central hub, hosting the **EHDEN** Chatbot Interface. Within this portal, the Network Dashboards provide metadata about various databases, which are crucial for the initial indexing and subsequent query building.

The **EHDEN** Chatbot Interface interacts directly with the user, guiding them through the process of defining concept sets and cohorts. It communicates with the IR API to retrieve relevant questions and forms, which are then presented to the user for completion.

The **IR** API acts as a mediator between the user inputs received via the Chatbot Interface and the ATLAS platform. It processes user queries, retrieves relevant database metadata, and prepares the necessary JSON structures for Cohort and Concept Set definitions.

Multiple institutions interact with the **EHDEN** Portal IR API, each with its own Chatbot Interface connected to its ATLAS instance. Institution A and Institution B represent different organizational units that utilize the **EHDEN** infrastructure for their research needs. The institutions' Chatbot Interfaces communicate with their respective ATLAS instances and the **EHDEN** Portal, ensuring data consistency and integration.

Institutions A and B, through their Chatbot Interfaces, interact with the **IR** API and ATLAS to perform similar tasks. Each institution can define its cohorts and concept sets, leveraging the centralized capabilities of the **EHDEN** infrastructure while maintaining individual ATLAS instances for specific research activities.

Results and Discussion

The current methods applied to help medical researchers discover databases of interest, specifically on the **EHDEN** Portal, are based on graphical and tabular views. These methods also have shown good results when dealing with a low number of databases, or concepts. However, with the increase of the databases in the network, such techniques may not be the best option.

The increasing number of databases in the **EHDEN** methods leads to a comprehensive database catalogue with information from almost 200 databases in Europe. It spans across 28 countries and offers data concerning 256 million patients, derived from 34 distinct data sources **MIE**.

The databases in the **EHDEN** Catalogue are supported by 38 unique attributes focused on exposing the core characteristics regarding data access. Additionally, the **EHDEN** Network Dashboards, it is hosted a staggering 8,504,324 unique concepts to accurately represent the scope of each database. These concepts are organized into 25 domains, including but not limited to Measurement, Condition, Procedure, Observation, and Drug. This amount of data raised unseen challenges that this project tried to address with this new approach.

The presented system aims to optimize the process of identifying and selecting relevant observational databases for medical research. The tool simplifies the discovery of medical observational databases in the **EHDEN** network, addressing the challenges of navigating through the vast and varied catalogues of medical databases. It also tackles the challenge of defining a cohort study in the **ATLAS** platform by improving the tool to allow for a more conversational approach to defining and providing a cohort query as an outcome.

This section ...

5.1 IR METHODS

mas que resultados?

5.2 LLM IMPLEMENTATION: FLOWISEAI VS LANGFLOW

Langflow and FlowiseAI are similar tools that provide multiple tools that allow the building of different workflows, architectures, and integrations with external tools. The table 5.1 shows a comparison between these two frameworks, identifying the strengths and weaknesses of each one.

FlowiseAI can run in air-gapped environments with local LLM, embeddings and vector databases. Also, it creates autonomous agents that can use tools to execute different tasks, such as Custom Tools, OpenAI Assistant, and Function Agent. FlowiseAI has proved to be a great solution to build an LLM-based chatbot. However, the FlowiseAI implementation become limited to address the requirements of the query builder. It became difficult to control the logic and the flow required by the query builder requirements.

Langflow is a dynamic graph where each node is an executable unit, so the way of development of a LLM application is very identical to FlowiseAI. But also, Langflow shown to be limited for this project because of missing basic features.

	FlowiseAI	Langflow
Strengths	falta um aqui	Flexible customization
	Free access	Bug tracking
	Visually intuitive	Visually intuitive
Weaknesses	Bug tracking	Missing basic features
	Limited documentation	Outdated documentation
	Limited features	Controlled access

Table 5.1: Comparison between FlowiseAI and Langflow.

Starting with FlowiseAI, it is a great tool to build simple LLM applications, but have some limitations:

- **Limited Features** - There are tasks that require more customization, and the custom tool that FlowiseAI provides is not enough for some of these tasks.
- **Limited Documentation** - There are some tools that not have documentation, and it became difficult sometimes to explore a tool without information about it.
- **Bug Tracking** - It is difficult to debug. When something in your system goes wrong, it is difficult to find what is causing that error.

Otherwise, Langflow offers features that either improve upon or address certain limitations of FlowiseAI, such as better features for bug tracking and more customization; for example, creating a complete custom component is possible. However, it also addresses some limitations:

- **Missing basic features** - Lacked basic functionalities such as a conversational agent. Others, requires some
- **Outdated documentation** - There was documentation of tools that didn't exist, and there wasn't documentation of some that did exist.
- **Controlled Access** - Requires an API key to access tools and features of the platform.

Although these frameworks prove to be good options, they have limitations that do not satisfy this project’s complexity. So, the **LLM** orchestration flow was implemented using a Python-based backend, developed for this project system.

5.3 CONVERSATIONAL SEARCH ASSISTANT

The integration of **IR** techniques within a conversational user interface represents a significant advancement in the field of data discovery **ritzel2019development**. Its ability to return a list of databases that are pertinent to the user’s query not only saves time but also introduces a level of precision in the selection process that was previously unattainable through manual methods. The integration of generative **AI** helps the conversation be more human-like.

Researcher can establish questions like “What is the prevalence of omeprazole?”, and the tool responds with the most relevant databases that may have the information to answer such questions. Then, the research can analyze this information in more detail, or refine the question.

Therefore, the method of simplifying the discovery of observational databases has an impact on the speed and quality of research efforts. Also, the tool has the potential to level the playing field by providing less experienced individuals with access to complex databases that they might otherwise overlook or find too challenging to navigate. However, the interaction with the proposed tool is somewhat restricted due to the data content. Users may expect to have deep insights into the databases, but that may expose sensitive data. Future enhancements could include the integration of such information available only to users with access to them, which will include the use of Rule-Based Access Control (RBAC) mechanisms over this tool.

5.4 QUERY BUILDER

?????

Conclusions

The introduction of this conversational search tool addresses the significant challenge of navigating a growing repository of medical databases. This system enables users to engage in natural language dialogue, significantly reducing the time and complexity traditionally associated with identifying suitable databases for research studies. Moreover, the system reflects a deep understanding of the needs within the medical research community for reliable, up-to-date, and accessible data.

In conclusion, the successful implementation of this system represents a significant advancement in the field of data discovery and retrieval. It not only simplifies the process of identifying the most relevant databases but also empowers researchers, regardless of their prior experience with database search tools. However, some future work remains, namely enhancing the IR methods to produce better results, considering synonyms.

References

- [1] A. Vaswani, N. Shazeer, N. Parmar, *et al.*, *Attention Is All You Need*, arXiv:1706.03762 [cs], Aug. 2023. (visited on 12/04/2023).
- [2] J. Wei, X. Wang, D. Schuurmans, *et al.*, *Chain-of-Thought Prompting Elicits Reasoning in Large Language Models*, Jan. 2023. DOI: [10.48550/arXiv.2201.11903](https://doi.org/10.48550/arXiv.2201.11903).
- [3] J. R. Rogers, G. Hripcsak, Y. K. Cheung, and C. Weng, «Clinical comparison between trial participants and potentially eligible patients using electronic health record data: A generalizability assessment method», *Journal of biomedical informatics*, vol. 119, p. 103822, 2021. DOI: [10.1016/j.jbi.2021.103822](https://doi.org/10.1016/j.jbi.2021.103822).
- [4] U. Topaloglu and M. B. Topaloglu, «Using a federated network of real-world data to optimize clinical trials operations», *JCO clinical cancer informatics*, vol. 2, pp. 1–10, 2018. DOI: [10.1200/CCI.17.00067](https://doi.org/10.1200/CCI.17.00067).
- [5] G. Lombardo, C. Couvert, M. Kose, *et al.*, «Electronic health records (EHRs) in clinical research and platform trials: Application of the innovative EHR-based methods developed by EU-PEARL», *Journal of Biomedical Informatics*, p. 104553, 2023. DOI: [10.1016/j.jbi.2023.104553](https://doi.org/10.1016/j.jbi.2023.104553).
- [6] J. R. Almeida, L. B. Silva, I. Bos, P. J. Visser, and J. L. Oliveira, «A methodology for cohort harmonisation in multicentre clinical research», *Informatics in Medicine Unlocked*, vol. Volume 27, p. 100760, 2021. DOI: [10.1016/j.imu.2021.100760](https://doi.org/10.1016/j.imu.2021.100760).
- [7] D. C. Kaelber, A. K. Jha, D. Johnston, B. Middleton, and D. W. Bates, «A research agenda for Personal Health Records (PHRs)», *Journal of the American Medical Informatics Association*, vol. 15, no. 6, pp. 729–736, 2008. DOI: [10.1197/jamia.M2547](https://doi.org/10.1197/jamia.M2547).
- [8] A. Pereira, J. R. Almeida, R. P. Lopes, and J. L. Oliveira, «Querying semantic catalogues of biomedical databases», *Journal of Biomedical Informatics*, vol. 137, p. 104272, 2023. DOI: [10.1016/j.jbi.2022.104272](https://doi.org/10.1016/j.jbi.2022.104272).
- [9] J. R. Almeida, J. P. Barraca, and J. L. Oliveira, «A secure architecture for exploring patient-level databases from distributed institutions», in *2022 IEEE 35th International Symposium on Computer-Based Medical Systems (CBMS)*, IEEE, 2022, pp. 447–452. DOI: [10.1109/CBMS55023.2022.00086](https://doi.org/10.1109/CBMS55023.2022.00086).
- [10] G. Hripcsak, P. B. Ryan, J. D. Duke, *et al.*, «Characterizing treatment pathways at scale using the OHDSI network», *Proceedings of the National Academy of Sciences*, vol. 113, no. 27, pp. 7329–7336, 2016. DOI: [10.1073/pnas.1510502113](https://doi.org/10.1073/pnas.1510502113).
- [11] K. Park, M. Cho, M. Song, *et al.*, «Exploring the potential of omop common data model for process mining in healthcare», *Plos one*, vol. 18, no. 1, e0279641, 2023. DOI: [10.1371/journal.pone.0279641](https://doi.org/10.1371/journal.pone.0279641).
- [12] J. R. Almeida, J. F. Silva, S. Matos, and J. L. Oliveira, «A two-stage workflow to extract and harmonize drug mentions from clinical notes into observational databases», *Journal of Biomedical Informatics*, vol. 120, p. 103849, 2021. DOI: [10.1016/j.jbi.2021.103849](https://doi.org/10.1016/j.jbi.2021.103849).
- [13] C. Reich, A. Ostroplets, P. Ryan, *et al.*, «OHDSI Standardized Vocabularies—a large-scale centralized reference ontology for international data harmonization», *Journal of the American Medical Informatics Association*, ocad247, 2024. DOI: [10.1093/jamia/ocad247](https://doi.org/10.1093/jamia/ocad247).
- [14] J. R. Almeida, A. Pazos, and J. L. Oliveira, «Clinical Data Integration Strategies for Multicenter Studies», in *Doctoral Conference on Computing, Electrical and Industrial Systems*, Springer, 2023, pp. 175–190. DOI: [10.1007/978-3-031-36007-7_13](https://doi.org/10.1007/978-3-031-36007-7_13).

- [15] J. R. Almeida and J. L. Oliveira, «MONTRA2: A web platform for profiling distributed databases in the health domain», *Informatics in Medicine Unlocked*, vol. 45, p. 101 447, 2024. DOI: [10.1016/j.imu.2024.101447](https://doi.org/10.1016/j.imu.2024.101447).
- [16] J. R. Almeida, J. M. Silva, and J. L. Oliveira, «A FAIR approach to real-world health data management and analysis», in *2023 IEEE 36th International Symposium on Computer-Based Medical Systems (CBMS)*, IEEE, 2023, pp. 892–897. DOI: [10.1109/CBMS58004.2023.00338](https://doi.org/10.1109/CBMS58004.2023.00338).
- [17] Kumar, T. Rao, A. R. Lakshminarayanan, and E. Pugazhendi, «An Efficient Text-Based Image Retrieval Using Natural Language Processing (NLP) Techniques», in Jan. 2021, pp. 505–519, ISBN: 9789811553998. DOI: [10.1007/978-981-15-5400-1_52](https://doi.org/10.1007/978-981-15-5400-1_52).
- [18] K. A. Hambarde and H. Proença, «Information Retrieval: Recent Advances and Beyond», en, *IEEE Access*, vol. 11, pp. 76 581–76 604, 2023, ISSN: 2169-3536. DOI: [10.1109/ACCESS.2023.3295776](https://doi.org/10.1109/ACCESS.2023.3295776). (visited on 12/06/2023).
- [19] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, *BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding*, Oct. 2018.
- [20] C. Manning, P. Raghavan, and H. Schuetze, «Introduction to Information Retrieval», en, 2009.
- [21] F. Lan, «Research on Text Similarity Measurement Hybrid Algorithm with Term Semantic Information and TF-IDF Method», *Advances in Multimedia*, vol. 2022, 2022, ISSN: 1687-5680. DOI: [10.1155/2022/7923262](https://doi.org/10.1155/2022/7923262).
- [22] E. Gomedé, *Understanding the BM25 Ranking Algorithm*, en, Sep. 2023. [Online]. Available: <https://medium.com/@evertongomedé/understanding-the-bm25-ranking-algorithm-19f6d45c6ce> (visited on 12/06/2023).
- [23] S. Connelly, *Practical BM25 - Part 2: The BM25 Algorithm and its Variables*, Apr. 2018. [Online]. Available: <https://www.elastic.co/blog/practical-bm25-part-2-the-bm25-algorithm-and-its-variables>.
- [24] B. Mitra and N. Craswell, «An Introduction to Neural Information Retrieval»,
- [25] H. Chen, Z. Dou, Q. Zhu, X. Zuo, and J.-R. Wen, «Integrating Representation and Interaction for Context-Aware Document Ranking», *ACM Transactions on Information Systems*, 2023. DOI: [10.1145/3529955](https://doi.org/10.1145/3529955).
- [26] S. Ayanouz, B. A. Abdelhakim, and M. Benhmed, «A Smart Chatbot Architecture based NLP and Machine Learning for Health Care Assistance», in *Proceedings of the 3rd International Conference on Networking, Information Systems & Security*, ser. NISS2020, New York, NY, USA: Association for Computing Machinery, 2020, pp. 1–6, ISBN: 978-1-4503-7634-1. DOI: [10.1145/3386723.3387897](https://doi.org/10.1145/3386723.3387897). (visited on 11/13/2023).
- [27] E. W. T. Ngai, M. C. M. Lee, M. Luo, P. S. L. Chan, and T. Liang, «An intelligent knowledge-based chatbot for customer service», *Electronic Commerce Research and Applications*, vol. 50, p. 101 098, Nov. 2021, ISSN: 1567-4223. DOI: [10.1016/j.elerap.2021.101098](https://doi.org/10.1016/j.elerap.2021.101098). (visited on 11/13/2023).
- [28] A. Chizhik and Y. Zhrebtsova, «Challenges of Building an Intelligent Chatbot», 2020. (visited on 11/13/2023).
- [29] B. Zhong, W. He, Z. Huang, P. E. Love, J. Tang, and H. Luo, «A building regulation question answering system: A deep learning methodology», en, *Advanced Engineering Informatics*, vol. 46, p. 101 195, Oct. 2020, ISSN: 14740346. DOI: [10.1016/j.aei.2020.101195](https://doi.org/10.1016/j.aei.2020.101195). (visited on 12/06/2023).
- [30] X. Liu, J. Wang, J. Sun, *et al.*, «Prompting Frameworks for Large Language Models: A Survey», en, *ACM Comput. Surv.*, vol. 1, no. 1,
- [31] T. A. Chang and B. K. Bergen, *Language Model Behavior: A Comprehensive Survey*, Aug. 2023.
- [32] M. U. Hadi, Q. Al-Tashi, R. Qureshi, *et al.*, «Large language models: A comprehensive survey of its applications, challenges, limitations, and future prospects», Jul. 2023. DOI: [10.36227/techrxiv.23589741](https://doi.org/10.36227/techrxiv.23589741).
- [33] H. Naveed, A. U. Khan, S. Qiu, *et al.*, *A Comprehensive Overview of Large Language Models*, arXiv:2307.06435 [cs], Nov. 2023. (visited on 12/15/2023).

- [34] B. Min, H. Ross, E. Sulem, *et al.*, *Recent Advances in Natural Language Processing via Large Pre-trained Language Models: A Survey / ACM Computing Surveys*.
- [35] *OpenAI Platform*, en. [Online]. Available: <https://platform.openai.com> (visited on 12/15/2023).
- [36] S. Kamnis, «Generative pre-trained transformers (GPT) for surface engineering», *Surface and Coatings Technology*, vol. 466, p. 129 680, Aug. 2023, ISSN: 0257-8972. DOI: [10.1016/j.surfcoat.2023.129680](https://doi.org/10.1016/j.surfcoat.2023.129680). (visited on 01/09/2024).
- [37] E. J. Hu, Y. Shen, P. Wallis, *et al.*, *LoRA: Low-Rank Adaptation of Large Language Models*, Oct. 2021. DOI: [10.48550/arXiv.2106.09685](https://doi.org/10.48550/arXiv.2106.09685).
- [38] B. Meskó, «Prompt Engineering as an Important Emerging Skill for Medical Professionals: Tutorial», en, *Journal of Medical Internet Research*, vol. 25, e50638, Oct. 2023, ISSN: 1438-8871. DOI: [10.2196/50638](https://doi.org/10.2196/50638). (visited on 12/21/2023).
- [39] X. Ma, S. Mishra, A. Liu, *et al.*, *Beyond ChatBots: ExploreLLM for Structured Thoughts and Personalized Model Responses*, arXiv:2312.00763 [cs], Dec. 2023. (visited on 12/21/2023).
- [40] S. Yao, J. Zhao, D. Yu, *et al.*, *ReAct: Synergizing Reasoning and Acting in Language Models*, Mar. 2023. DOI: [10.48550/arXiv.2210.03629](https://doi.org/10.48550/arXiv.2210.03629).
- [41] H. Touvron, L. Martin, K. Stone, *et al.*, *Llama 2: Open Foundation and Fine-Tuned Chat Models*, Jul. 2023. DOI: [10.48550/arXiv.2307.09288](https://doi.org/10.48550/arXiv.2307.09288). (visited on 01/06/2024).
- [42] OpenAI, J. Achiam, S. Adler, *et al.*, *GPT-4 Technical Report*, Dec. 2023. DOI: [10.48550/arXiv.2303.08774](https://doi.org/10.48550/arXiv.2303.08774). (visited on 01/06/2024).
- [43] R. Anil, A. M. Dai, O. Firat, *et al.*, *PaLM 2 Technical Report*, Sep. 2023. DOI: [10.48550/arXiv.2305.10403](https://doi.org/10.48550/arXiv.2305.10403). (visited on 01/06/2024).
- [44] E. Almazrouei, H. Alobeidli, A. Alshamsi, *et al.*, *The Falcon Series of Open Language Models*, Nov. 2023. DOI: [10.48550/arXiv.2311.16867](https://doi.org/10.48550/arXiv.2311.16867). (visited on 01/06/2024).
- [45] K. Church and R. Yue, «Emerging trends: Smooth-talking machines», *Natural Language Engineering*, vol. 29, 2023, ISSN: 1351-3249. DOI: [10.1017/S1351324923000463](https://doi.org/10.1017/S1351324923000463).
- [46] N. Kshetri, «Cybercrime and Privacy Threats of Large Language Models», *IT Professional*, vol. 25, 2023, ISSN: 1520-9202. DOI: [10.1109/MITP.2023.3275489](https://doi.org/10.1109/MITP.2023.3275489).
- [47] C. V. Misischia, F. Poetze, and C. Strauss, «Chatbots in customer service: Their relevance and impact on service quality», vol. 201, Jan. 2022, ISSN: 1877-0509. DOI: [10.1016/j.procs.2022.03.055](https://doi.org/10.1016/j.procs.2022.03.055). (visited on 12/15/2024).
- [48] M. Nuruzzaman and O. K. Hussain, «A Survey on Chatbot Implementation in Customer Service Industry through Deep Neural Networks», in *2018 IEEE 15th International Conference on e-Business Engineering (ICEBE)*, Oct. 2018, pp. 54–61. DOI: [10.1109/ICEBE.2018.00019](https://doi.org/10.1109/ICEBE.2018.00019). (visited on 11/13/2023).
- [49] Z. Peng and X. Ma, «A survey on construction and enhancement methods in service chatbots design», *CCF Transactions on Pervasive Computing and Interaction*, 2019, ISSN: 2524-521X. DOI: [10.1007/s42486-019-00012-3](https://doi.org/10.1007/s42486-019-00012-3).
- [50] J. Weizenbaum, «ELIZA—a computer program for the study of natural language communication between man and machine», *Communications of the ACM*, Jan. 1966. DOI: [10.1145/365153.365168](https://doi.org/10.1145/365153.365168).
- [51] R. Agarwal and M. Wadhwa, «Review of State-of-the-Art Design Techniques for Chatbots», 2020, ISSN: 2662-995X. DOI: [10.1007/s42979-020-00255-3](https://doi.org/10.1007/s42979-020-00255-3).
- [52] G. Li, R. Gomez, K. Nakamura, and B. He, «Human-Centered Reinforcement Learning: A Survey», *IEEE Transactions on Human-Machine Systems*, no. 4, Aug. 2019, ISSN: 2168-2305. DOI: [10.1109/THMS.2019.2912447](https://doi.org/10.1109/THMS.2019.2912447). (visited on 01/17/2024).
- [53] Q.-D. Tran, A.-C. Le, and V.-N. Huynh, «Enhancing Conversational Model With Deep Reinforcement Learning and Adversarial Learning», *IEEE Access*, 2023, ISSN: 2169-3536. DOI: [10.1109/ACCESS.2023.3297652](https://doi.org/10.1109/ACCESS.2023.3297652).

- [54] A. Axelsson, H. Buschmeier, and G. Skantze, «Modeling Feedback in Interaction With Conversational Agents—A Review», *Frontiers in Computer Science*, 2022, ISSN: 2624-9898. DOI: [10.3389/fcomp.2022.744574](https://doi.org/10.3389/fcomp.2022.744574).
- [55] P. Lewis, E. Perez, A. Piktus, *et al.*, «Retrieval-augmented generation for knowledge-intensive NLP tasks», in *Proceedings of the 34th International Conference on Neural Information Processing Systems*, ser. NIPS'20, Red Hook, NY, USA: Curran Associates Inc., Dec. 2020, pp. 9459–9474, ISBN: 978-1-71382-954-6. (visited on 12/24/2023).
- [56] Y. Gao, Y. Xiong, X. Gao, *et al.*, *Retrieval-Augmented Generation for Large Language Models: A Survey*, arXiv:2312.10997 [cs], Dec. 2023. DOI: [10.48550/arXiv.2312.10997](https://doi.org/10.48550/arXiv.2312.10997). (visited on 12/24/2023).
- [57] O. Mussa, O. Rana, B. Goossens, P. Orozco-TerWengel, and C. Perera, «ForestQB: An Adaptive Query Builder to Support Wildlife Research», 2022.
- [58] *jQuery QueryBuilder*. [Online]. Available: <https://querybuilder.js.org/> (visited on 01/04/2024).
- [59] *Introducing Query Builder: Write complex queries without code*, Aug. 2021. [Online]. Available: <https://www.dronahq.com/introducing-query-builder/> (visited on 01/05/2024).
- [60] OHDSI, *Chapter 10 Defining Cohorts / The Book of OHDSI*. [Online]. Available: <https://ohdsi.github.io/TheBookOfOhdsi/> (visited on 01/03/2024).
- [61] T. Formal, B. Piwowarski, and S. Clinchant, *SPLADE: Sparse Lexical and Expansion Model for First Stage Ranking*, Jul. 2021. DOI: [10.48550/arXiv.2107.05720](https://doi.org/10.48550/arXiv.2107.05720).
- [62] J. Chen, S. Xiao, P. Zhang, K. Luo, D. Lian, and Z. Liu, *BGE M3-Embedding: Multi-Lingual, Multi-Functionality, Multi-Granularity Text Embeddings Through Self-Knowledge Distillation*, Feb. 2024. DOI: [10.48550/arXiv.2402.03216](https://doi.org/10.48550/arXiv.2402.03216).
- [63] T. Formal, C. Lassance, B. Piwowarski, and S. Clinchant, *From Distillation to Hard Negative Sampling: Making Sparse Neural IR Models More Effective*, May 2022. DOI: [10.48550/arXiv.2205.04733](https://doi.org/10.48550/arXiv.2205.04733).
- [64] Teknium, theemozilla, karan4d, and huemin_art, *Nous Hermes 2 Mixtral 8x7B DPO*. [Online]. Available: <https://huggingface.co/NousResearch/Nous-Hermes-2-Mixtral-8x7B-DPO>.
- [65] A. Q. Jiang *et al.*, «Mixtral of experts», *arXiv:2401.04088*, Jan. 2024. arXiv: [2401.04088](https://arxiv.org/abs/2401.04088). [Online]. Available: <https://arxiv.org/abs/2401.04088>.
- [66] J. A. Reis, J. R. Almeida, T. M. Almeida, and J. L. Oliveira, «Using Flowise to Streamline Biomedical Data Discovery and Analysis», in *2024 IEEE 22nd Mediterranean Electrotechnical Conference (MELECON)*, IEEE, 2024.
- [67] J. R. Almeida, A. Zúquete, A. Pazos, and J. L. Oliveira, «A Federated Authentication Schema among Multiple Identity Providers», *Heliyon*, 2024. DOI: [10.1016/j.heliyon.2024.e28560](https://doi.org/10.1016/j.heliyon.2024.e28560).